

Аллокации в linux: от new до физической памяти

Created by: Пресняков Арсений и Бережной Аким

2023-12-28

Содержание

Аллокации в linux: от new до физической памяти	1
1. Глава 1. Введение	2
2. Глава 2. malloc()	3
2.1. Основные структуры	3
2.2. int_malloc()	6
2.3. sysmalloc()	6
3. Глава 3. Системные вызовы brk() и mmap().	10
3.1. sbrk()	11
3.2. mmap()	12
4. Глава 4. Аллокатор физических страниц	16
4.1. Аллокация	16
4.2. Освобождение	16
5. Глава 5. Итоги	17
6. Список источников	18

Глава 1. Введение

Целью данной работы является разбор дерева вызовов, которые получит пользователь при запросе на выделении памяти. Задача дойти от вызова `new` до основ ядра операционной системы. Мы рассмотрим как это происходит в C/C++. **Аллокатор** - функция, динамически распределяющая память для пользователя, которая выделяет или освобождает память при необходимости. При вызове `new` происходит вызов функции `malloc(size_t)` из `glibc.h`, которая является одним из представителей «семейства аллокаторов» и относится к `general-purpose` аллокатором (аллокатор общего предназначения, который подходит для решения любых задач). Существует огромное количество алгоритмов распределения памяти, преследующих различные цели. Они выбирают компромиссы между основными требованиями к аллокатору такими как:

Максимальная мобильность

Использование как можно меньшего количества системно-зависимых функций (таких как системные вызовы), при этом обеспечивая дополнительную поддержку других полезных функций, доступных только в некоторых системах; соответствие всем известным системным ограничениям на правила выравнивания и адресации.

Минимизация пространства

Аллокатор не должен тратить пространство впустую: он должен получать как можно меньше памяти от системы и поддерживать память таким образом, чтобы минимизировать фрагментацию — «дыры» в смежных участках памяти, которые не используются программой.

Минимизация времени

Подпрограммы `malloc()`, `free()` и `realloc()` в среднем должны работать как можно быстрее.

Максимизация возможностей настройки

Дополнительные функции и поведение должны контролироваться пользователями либо статически (с помощью макросов и т.п.), либо динамически (с помощью команд управления, таких как `mallopt`).

Максимизация локальности

Выделение фрагментов памяти, которые обычно используются вместе, рядом друг с другом. Это помогает свести к минимуму промахи страниц и кэша во время выполнения программы.

Максимизация обнаружения ошибок

Аллокторы должны предоставлять средства для обнаружения повреждений из-за перезаписи памяти, многократного освобождения и т. д.

Минимизация аномалий

Аллокатор, использующий настройки по умолчанию, должен хорошо работать в широком диапазоне реальных нагрузок, которые сильно зависят от динамического распределения — приложения с графическим интерфейсом, компиляторы, интерпретаторы, инструменты разработки, программы с интенсивным использованием сети (пакетов), пакеты с интенсивным использованием графики, веб-браузеры, приложения для обработки строк и т. д.

Глава 2. MALLOC()

2.1. Основные структуры

Спускаемся ниже. Теперь рассмотрим реализацию **malloc()** в glibc.h. Рассмотрим используемый здесь алгоритм. При небольшом количестве потоков каждый поток получает свой **struct malloc_state**. Сделано это для скорости работы, поскольку тогда нет необходимости в примитивах синхронизации. В случае очень большого количества потоков к одному **malloc_state** будут иметь доступ несколько потоков, поэтому внутри ему необходим mutex. Далее для простоты будем рассматривать только однопоточные программы. В них будет только один **malloc_state**, который называется **main_arena**. Рассмотрим **malloc_state**:

```

1  /* malloc.c */
2  struct malloc_state
3  {
4      /* Serialize access. */
5      __libc_lock_define(, mutex);
6
7      /* Flags (formerly in max_fast). */
8      int flags;
9
10     /* Set if the fastbin chunks contain recently inserted free blocks. */
11     /* Note this is a bool but not all targets support atomics on booleans. */
12     int have_fastchunks;
13
14     /* Fastbins */
15     mfastbinptr fastbinsY[NFASTBINS];
16
17     /* Base of the topmost chunk -- not otherwise kept in a bin */
18     mchunkptr top;
19
20     /* The remainder from the most recent split of a small request */
21     mchunkptr last_remainder;
22
23     /* Normal bins packed as described above */
24     mchunkptr bins[NBINS * 2 - 2];
25
26     /* Bitmap of bins */
27     unsigned int binmap[BINMAPSIZE];
28
29     /* Linked list */
30     struct malloc_state *next;
31
32     /* Linked list for free arenas. Access to this field is serialized
33        by free_list_lock in arena.c. */
34     struct malloc_state *next_free;
35
36     /* Number of threads attached to this arena. 0 if the arena is on
37        the free list. Access to this field is serialized by
38        free_list_lock in arena.c. */
39     INTERNAL_SIZE_T attached_threads;
40
41     /* Memory allocated from the system in this arena. */
42     INTERNAL_SIZE_T system_mem;
43     INTERNAL_SIZE_T max_system_mem;
44 };

```

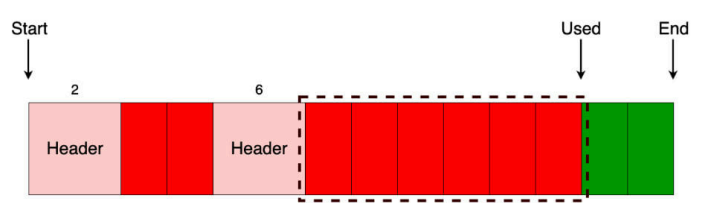
В нем много непонятных нам полей. Рассмотрим основные. **mfastbinptr** и **mchunkptr** - разные typedef одного struct'a **malloc_chunk**:

```
1 /* malloc.c */
2 typedef struct malloc_chunk* mchunkptr;
```

```
1 /* malloc.c */
2 typedef struct malloc_chunk *mfastbinptr;
```

```
1 /* malloc.c */
2 struct malloc_chunk {
3
4     INTERNAL_SIZE_T    mchunk_prev_size; /* Size of previous chunk (if free). */
5     INTERNAL_SIZE_T    mchunk_size;      /* Size in bytes, including overhead. */
6
7     struct malloc_chunk* fd;              /* double links -- used only if free. */
8     struct malloc_chunk* bk;
9
10    /* Only used for large blocks: pointer to next larger size. */
11    struct malloc_chunk* fd_nextsize; /* double links -- used only if free. */
12    struct malloc_chunk* bk_nextsize;
13 };
```

Это основообразующая структура всего аллокатора. Так называемый «Header» перед contiguous блоком памяти. Это выглядит примерно вот так:



Чтобы разобраться в его полях давайте введем два понятия: свободный блок и занятый блок.

Свободный блок - contiguous блок памяти, которым владеет **malloc** (он уже есть в виртуальном адресном пространстве программы), но который еще не выделен для пользователя, или блок, к которому был применен **free**).

Занятый блок - contiguous блок памяти, который был передан пользователю с помощью **malloc**.

Структура **malloc_chunk** хранит размер предыдущего блока (который лежит contiguous перед ним) - **mchunk_prev_size**, размер текущего блока - **mchunk_size**. Размер каждого чанка кратен 8 байтам, следовательно последние 3 бита поля **mchunk_size** не важны для определения размера и их можно использовать под флажки:

- A (0x04)

Allocated Arena - **main_arena** использует heap программы. Другие **malloc_state** аллоцируют память под их собственную кучу с помощью системного вызова **mmap()**. Если значения бита 0, этот чанк относится к **main_arena**, иначе этот блок из куска памяти, выделенного **mmap()** для другого **malloc_state**.

- M (0x02)

MMap'd chunk - этот блок был выделен для с помощью системного вызова **mmap()** и не является частью кучи. При его освобождении вызывается парный вызов - **munmap()**.

- P (0x01)

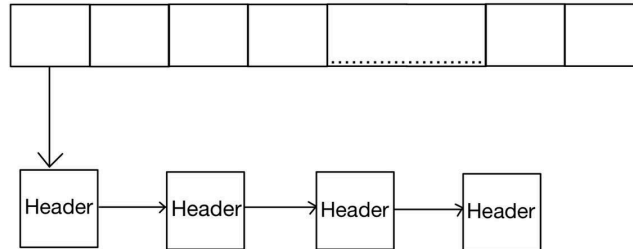
Previous chunk is in use - предыдущий блок занят. Этот флаг нам нужен для объединения двух соседних свободных чанков в один.

Далее поделим блоки на быстрые, маленькие и большие. Их границы зависят от **sizeof(size_t)** (будем считать его равным 8). Быстрые (≤ 160 байт) – кэшируются.

Вернемся к **malloc_state** и рассмотрим **fastbins**. Это статический массив поинтеров на **malloc_chunk**, которые указывают на головы односвязных списков, хранящих равные по размеру быстрые **malloc_chunk**. Эти спис-

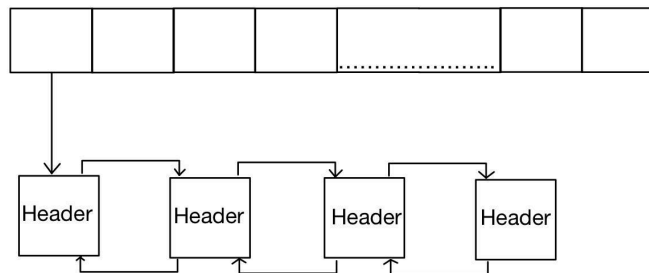
ки работают по принципу LIFO(то есть последний попавший в него блок будет лежать в голове списка), что повышает вероятность нахождения адреса в TLB-кэше при запросе на выделение из этого списка. Так же есть специальный макрос, позволяющий моментально найти список с подходящим размером блоков, что экономит значительное время, которое было бы потрачено на поиск.

• fastbinsY



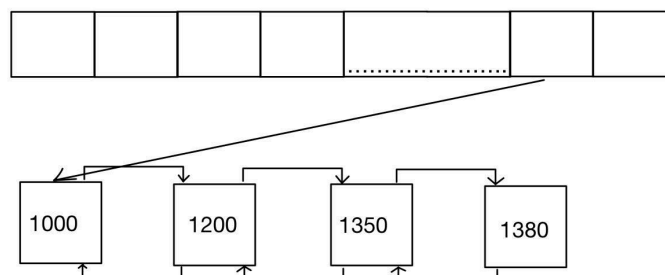
Теперь рассмотрим как хранятся блоки маленького размера(≤ 512 байт). Свободные чанки одинакового размера хранятся в двусвязном списке. В **malloc_chunk** поле **fd** - указатель на следующий элемент, а поле **bk** на предыдущий. Массив указателей на головы этих списков также хранится в **malloc_state**, а именно в поле **bins**. Под хранение маленьких блоков выделены первые 64 ячейки данного массива. Блоки выделяются по FIFO. Здесь мы все еще имеем быстрый доступ к спискам с помощью **#define**, однако чанки требуют уже куда более долгой и тяжелой обработки.

• bins



В больших блоках(все остальные) мы начинаем использовать последние два поля **malloc_chunk** - **fd_nextsize** и **bk_nextsize**. Оставшиеся 64 ячейки массива **bins** выделяются под указатели со списками на большие чанки. Мы отходим от концепции хранения списков с одинаковым размером блоков. Теперь мы храним двусвязные списки с заранее известным диапазоном значений, а блоки в них отсортированы - **fd_nextsize** ссылается на минимальный блок имеющий больший или равный нам размер, а **bk_nextsize**, соответственно, на максимальный блок с меньшим или равным размером.

• bins



Помимо этого, есть еще один двусвязный список, хранящий блоки маленького и большого размера в не отсортированном порядке. Из него достаются блоки по принципу FIFO и добавляются в **bins**.

2.2. `int_malloc()`

Теперь, когда мы знакомы с тем, как хранятся **malloc_chunk**, мы можем перейти к разбору алгоритма **int_malloc()**, внутренней функции, которая вызывается в **malloc()**. Мы не будем вдаваться во все подробности, опустив некоторые детали (например, оптимизацию TCACHE в glibc). Наша задача – понять основные принципы работы алгоритма:

- При запросе на аллокацию менее 160 байт мы находим индекс в массиве **fastbinY** и выделяем голову этого односвязного списка, новая голова - следующий за ним блок.
- При запросе на аллокацию среднего блока (160 - 512 байт) мы получаем в **bins** индекс двусвязного списка с блоком нужного нам размера. Если такой список пуст, то обращаемся к списку, хранящему не отсортированные блоки, и пробуем среди них найти подходящий, параллельно добавляя блоки в соответствующие списки **bins**. Если мы нашли подходящий блок то его мы и вернем, иначе возьмем больший блок и разделим его на два: один, который мы займем и остаток. Остаток добавим в список свободных блоков.
- При аллокации более 512 байт мы очищаем **fastbinY**, объединяя свободные соседние блоки в один. Это необходимо для более эффективного использования памяти и борьбы с фрагментацией. Затем мы добавляем один чанк из списка неотсортированных блоков в **bins** и пробуем достать подходящий блок из соответствующего размеру запроса списка из **bins**. Если он пуст, то пробуем найти подходящий чанк в неотсортированном списке.
- Запрос не получилось удовлетворить. В таком случае вызывается вспомогательная функция **sysmalloc()**.

Снова вернемся к структуре **malloc_state** и посмотрим на поле **top**. Оно хранит последний доступный блок памяти (если представить память contiguous, то самый правый блок, указатель на конец нашего heap), он не лежит в **bins**. Это некий указатель на конец доступной нам памяти, и когда мы не можем удовлетворить запрос пользователя необходимо обратиться к операционной системе за новой порцией. В этом нам поможет функция **sysmalloc()** Рассмотрим алгоритм ее работы более подробно в следующем пункте.

2.3. `sysmalloc()`

Как было оговорено выше, **sysmalloc()** обрабатывает те случаи, когда от системы требуется выделить несколько больше памяти, чем предусмотрено в обработке случаев с небольшими аллокациями. Рассмотрим алгоритм работы **sysmalloc()** на примере его исходного кода:

```

1  /* malloc.c */
2  if (av == NULL ||
3      ((unsigned long) (nb) >= (unsigned long) (mp_.mmap_threshold)
4       && (mp_.n_mmaps < mp_.n_mmaps_max)))
5  {
6      char *mm;
7      if (mp_.hp_pagesize > 0 && nb >= mp_.hp_pagesize)
8      {
9          /* There is no need to issue the THP madvise call if Huge Pages are
10             used directly. */
11          mm = sysmalloc_mmap (nb, mp_.hp_pagesize, mp_.hp_flags, av);
12          if (mm != MAP_FAILED)
13              return mm;
14      }
15      mm = sysmalloc_mmap (nb, pagesize, 0, av);
16      if (mm != MAP_FAILED)
17          return mm;
18      tried_mmap = true;
19  }
```

Сначала смотрим на трешхолд (константу, которая ограничивает то минимальное количество байт, когда мы делаем **sysmalloc_mmap()** – обёрточную функцию, вызывающую **mmap()**) и, если он достигнут, пользуемся **sysmalloc_mmap()**. Здесь стоит сделать оговорку, что код написан с учётом так называемых **Transparent Huge Pages (THP)** – альтернативой **HugeTLB** (в современных версиях OS, они включены по умолчанию). Здесь используются два подхода: **system-wide** и **per-process**. При **system-wide** подходе, система автоматически пытается назначить большие страницы любому процессу, когда есть на то возможность, и процесс использует большие непрерывные участки виртуальной памяти. В случае с **per-process**, ядро отдаёт большие страницы отдель-

ным процессам в те участки памяти, которые специфицирует системный вызов **madvise()**. Пользователь может управлять этим поведением используя `/sys/kernel/mm/transparent_hugepage/enabled`. В этот файл могут быть записаны такие опции, как:

- **always** - для **system-wide** распределения;
- **madvise** - для **per-process** распределения;
- **never** - для отключения **THP**.

Далее, в случае, когда доступных арен нет и **sysmalloc_mmap()** провалился, возвращаем пустой указатель:

```
1 /* malloc.c */
2 /* There are no usable arenas and mmap also failed. */
3 if (av == NULL)
4     return 0;
```

Далее рассчитываем количество памяти, которое нам может потребоваться, учитывая выравнивание и минимальный порог, записывая его в **size**. Мы также заранее надеемся на то, что память, запрашиваемая нами сейчас попадет в contiguous кусок сразу после прошлого запроса на **sysmalloc()**, если таковой был. Поэтому, мы вычитаем тот размер, который мы могли уже выделить на куче до текущего вызова:

```
1 /* malloc.c */
2 size = nb + mp_.top_pad + MINSIZE;
3
4 if (contiguous (av))
5     size -= old_size;
```

Здесь **old_size** – это округленный **top** из аренды, которая подалась аргументом в данный системный вызов.

После того, как мы посчитали размер, учитываем **THP** оптимизацию и выравниваем память по размеру большой страницы. Иначе - по размеру обычной:

```
1 /* malloc.c */
2 #ifdef MADV_HUGEPAGE
3     /* Defined in brk.c. */
4     extern void *__curbrk;
5     if (__glibc_unlikely (mp_.thp_pagesize != 0))
6     {
7         uintptr_t top = ALIGN_UP ((uintptr_t) __curbrk + size,
8                                   mp_.thp_pagesize);
9         size = top - (uintptr_t) __curbrk;
10    }
11    else
12 #endif
13        size = ALIGN_UP (size, GLRO(dl_pagesize));
```

В случае, когда размер посчитался корректно, вызываем **sbrk()**. На самом деле, в дальнейшем коде вызывается функция **MORECORE()**, но вы можете убедиться, что это тот же **sbrk()**, просто обёрнутый в другую функцию. Название изменено семантически, так как **MORECORE()** используется при добавление микро-кусочков памяти на куче при устранении проблем с выравниванием и сторонними вызовами **sbrk()**:

```
1 /* morecore.c */
2 void * __default_morecore (ptrdiff_t increment)
3 {
4     void *result = (void *) __sbrk (increment);
5     if (result == (void *) -1)
6         return NULL;
7
8     return result;
9 }
```

Если **sbrk()** провалился, нам не удалось выделить contiguous кусок и мы пытаемся использовать **sysmalloc_mmap_fallback()** (функцию, похожую на **sysmalloc_mmap()**, но предназначенную для случаев, когда **sbrk()** не смог выделить память), чтобы выделить размер большой страницы, если это возможно или, если и такой вызов провалится, хотя бы кусок требуемого размера:

```

1  /* malloc.c */
2  if (size > 0)
3  {
4      brk = (char *) (MORECORE (size));
5      if (brk != (char *) (MORECORE_FAILURE))
6          madvise_thp (brk, size);
7      LIBC_PROBE (memory_sbrk_more, 2, brk, size);
8  }
9
10 if (brk == (char *) (MORECORE_FAILURE))
11 {
12     char *mbrk = MAP_FAILED;
13     if (mp_.hp_pagesize > 0)
14         mbrk = sysmalloc_mmap_fallback (&size, nb, old_size, mp_.hp_pagesize,
15                                         mp_.hp_pagesize, mp_.hp_flags, av);
16     if (mbrk == MAP_FAILED)
17         mbrk = sysmalloc_mmap_fallback (&size, nb, old_size,
18                                         MMAP_AS_MORECORE_SIZE, pagesize, 0, av);
19     if (mbrk != MAP_FAILED)
20     {
21         /* We do not need, and cannot use, another sbrk call to find end */
22         brk = mbrk;
23         snd_brk = brk + size;
24     }
25 }

```

Если же **sbrk()** не провалился, то всё гораздо лучше: мы можем расширить кучу и получить кусок памяти, который лежит непрерывно с тем куском, что мы могли аллоцировать ранее (кэшу теперь будет приятно). При этом, не забываем передвинуть указатель на голову кучи подальше:

```

1  /* malloc.c */
2  if (mp_.sbrk_base == 0)
3      mp_.sbrk_base = brk;
4  av->system_mem += size;
5  if (brk == old_end && snd_brk == (char *) (MORECORE_FAILURE))
6      set_head (old_top, (size + old_size) | PREV_INUSE);

```

Если же **brk** стал меньше чем **old_end**, то произошло что-то странное: где-то произошла ошибка и мы ходим выйти как можно скорее. Выводим это в поток ошибок:

```

1  /* malloc.c */
2  else if (contiguous (av) && old_size && brk < old_end)
3      /* Oops! Someone else killed our space.. Can't touch anything. */
4      malloc_printerr ("break adjusted to free malloc space");

```

Остался третий случай: когда все указатели в корректном состоянии, но у нас остались проблемы с выравниванием. Далее в исходном коде идет большой пласт иффов, которые отвечают только за то, чтобы поправить выравнивание в куче различных случаев. Для того, чтобы не захламлять общее понимание работы **sysmalloc()**, эту деталь мы опустим.

Итак, осталось только проверить, что соблюлись все инварианты и отдать указатель на память:

```

1  /* malloc.c */
2  if ((unsigned long) av->system_mem > (unsigned long) (av->max_system_mem))

```



```

3     av->max_system_mem = av->system_mem;
4     check_malloc_state (av);
5
6     /* finally, do the allocation */
7     p = av->top;
8     size = chunksize (p);
9
10    /* check that one of the above allocation paths succeeded */
11    if ((unsigned long) (size) >= (unsigned long) (nb + MINSIZE))
12    {
13        remainder_size = size - nb;
14        remainder = chunk_at_offset (p, nb);
15        av->top = remainder;
16        set_head (p, nb | PREV_INUSE | (av != &main_arena ? NON_MAIN_ARENA : 0));
17        set_head (remainder, remainder_size | PREV_INUSE);
18        check_malloted_chunk (av, p, nb);
19        return chunk2mem (p);
20    }
21
22    /* catch all failure paths */
23    __set_errno (ENOMEM);
24    return 0;

```

Подводя итог, разреем работу **sysmalloc()** поэтапно:

1. Смотрим на трэшхолд:
 - если перевалили его, вызываем **sysmalloc_mmap()** и при успешном завершении → см. пункт 3
 - если нет, пробуем вызвать **sbrk()**, чтобы выделить contiguous кусок памяти
2. Как отработал **sbrk()**?
 - провалился → вызываем **sysmalloc_mmap_fallback()** и при успешном завершении → см. пункт 3
 - отработал без ошибок → учитываем выравнивания и идем дальше
3. Удостоверяемся, что память выделяется успешно: проверяем выравнивания, количество получишейся памяти, инварианты указателей и **malloc_state**:
 - если всё ок → отдаем указатель на выделенную память
 - если где-то ошибка, возвращаем пустой указатель и ставим флажок ENOMEM, если память кончилась.

Глава 3. Системные вызовы `brk()` и `mmap()`.

Теперь рассмотрим как работают `sbrk()` и `mmap()` в linux. Для начала дадим определение:

Регион памяти - contiguous кусок памяти лежащий не на куче и выделяемый с помощью системного вызова `mmap()`.

Познакомимся с двумя структурами:

vm_area_struct(vma) - структура, хранящая информацию о регионе памяти. Ее основные поля:

- `vm_start`, `vm_end` - начало и конец зоны памяти, соответственно
- `vm_next`, `vm_prev` - `vm_area_struct` каждого процесса завязаны в двусвязный список

mm_struct - структура, включающая в себя все `vm_area_struct`, относящиеся к данному процессу. Ее основные поля:

- `mmap` - указатель на голову двусвязного списка `vm_area_struct`
- `mmap_cache` - lru `vm_area_struct`
- `mmlist` - список всех `mm_struct`
- `start_brk`, `brk` - начало и конец кучи, соответственно
- `start_stack` - начало стека

Ниже приведен код со всеми полями данной структуры(это не настоящий код ядра linux, здесь добавлены комментарии и вырезаны некоторые детали, усложняющие читаемость кода):

```

1  /* linux/mm_types.h */
2  struct mm_struct {
3      struct vm_area_struct *mmap;           /* list of memory areas */
4      struct rb_root mm_rb;                  /* red-black tree of VMAs */
5      struct vm_area_struct *mmap_cache;     /* last used memory area */
6      unsigned long free_area_cache;         /* 1st address space hole */
7      pgd_t *pgd;                           /* page global directory */
8      atomic_t mm_users;                     /* address space users */
9      atomic_t mm_count;                     /* primary usage counter */
10     int map_count;                          /* number of memory areas */
11     struct rw_semaphore mmap_sem;          /* memory area semaphore */
12     spinlock_t page_table_lock;            /* page table lock */
13     struct list_head mmlist;               /* list of all mm_structs */
14     unsigned long start_code;               /* start address of code */
15     unsigned long end_code;                 /* final address of code */
16     unsigned long start_data;               /* start address of data */
17     unsigned long end_data;                 /* final address of data */
18     unsigned long start_brk;                /* start address of heap */
19     unsigned long brk;                      /* final address of heap */
20     unsigned long start_stack;              /* start address of stack */
21     unsigned long arg_start;                /* start of arguments */
22     unsigned long arg_end;                  /* end of arguments */
23     unsigned long env_start;                /* start of environment */
24     unsigned long env_end;                  /* end of environment */
25     unsigned long rss;                      /* pages allocated */
26     unsigned long total_vm;                 /* total number of pages */
27     unsigned long locked_vm;                /* number of locked pages */
28     unsigned long def_flags;                /* default access flags */
29     unsigned long cpu_vm_mask;              /* lazy TLB switch mask */
30     unsigned long swap_address;             /* last scanned address */
31     unsigned dumpable:1;                   /* can this mm core dump? */
32     int used_hugetlb;                       /* used hugetlb pages? */
33     mm_context_t context;                   /* arch-specific data */
34     int core_waiters;                       /* thread core dump waiters */
35     struct completion *core_startup_done;   /* core start completion */
36     struct completion core_done;            /* core end completion */
37     rwlock_t ioctx_list_lock;              /* AIO I/O list lock */
38     struct kiocx *ioctx_list;              /* AIO I/O list */

```

```

39     struct kioctx          default_kioctx;    /* AIO default I/O context */
40 };

```

3.1. sbrk()

Сейчас, когда мы познакомились с основными структурами управления виртуальной памятью в ядре, мы можем перейти к разбору кода функции **sbrk()** в linux.

```

1  /* linux/tools/include/nolibc/sys.h */
2  void *sys_brk(void *addr)
3  {
4      return (void *)my_syscall1(__NR_brk, addr);
5  }
6
7  void *sbrk(intptr_t inc)
8  {
9      /* first call to find current end */
10     void *ret = sys_brk(0);
11
12     if (ret && sys_brk(ret + inc) == ret + inc)
13         return ret + inc;
14
15     SET_ERRNO(ENOMEM);
16     return (void *)-1;
17 }

```

sbrk() дважды вызывает **sys_brk()** - обертку над системным вызовом **brk()**, который пробует заменить конец кучи процесса(поле **brk** структуры **mm_struct**) на переданный ему указатель. Первым вызовом **sys_brk(0)** мы получаем текущий **brk**, вторым вызовом **sys_brk(ret + inc)** мы пробуем увеличить **brk** на **inc** и, в случае неудачи, возвращаем ошибку. Теперь рассмотрим реализацию системного вызова **brk()**:

```

1  /* linux/mm/nommu.c */
2  SYSCALL_DEFINE1(brk, unsigned long, brk)
3  {
4      struct mm_struct *mm = current->mm;
5
6      if (brk < mm->start_brk || brk > mm->context.end_brk)
7          return mm->brk;
8
9      if (mm->brk == brk)
10         return mm->brk;
11
12     /*
13      * Always allow shrinking brk
14      */
15     if (brk <= mm->brk) {
16         mm->brk = brk;
17         return brk;
18     }
19
20     /*
21      * Ok, looks good - let it rip.
22      */
23     flush_icache_user_range(mm->brk, brk);
24     return mm->brk = brk;
25 }

```

Мы получаем указатель на **mm_struct** текущего процесса. Далее, если принятый указатель меньше начала кучи(**start_brk**), больше максимально возможного конца кучи(**context.end_brk**) или равен текущему **brk** мы не меняем **brk** и возвращаем его. Таким образом, **sys_brk(0)** возвращает текущий **brk**. Данный системный вызов

допускает уменьшение размера кучи, поэтому, если принятый указатель больше **start_brk** и меньше **brk** мы уменьшим текущий **brk**. В единственном оставшемся случае наш указатель увлечит **brk**, расширив объем кучи, и вернет его. Это **brk()** делает, вызывая функцию **flush_icache_user_range()**, которая вызывает аллокатор физических страниц памяти, о котором мы расскажем в следующей главе. Так, **sbrk()** понимает, что системный вызов **brk()** сработал ожидаемым образом, увеличив кучу.

3.2. mmap()

```
1 /* mm/mmap.c */
2 void *mmap(void *addr, size_t length, int prot, int flags, int fd, off_t offset);
3 int munmap(void *addr, size_t length);
```

Обратите внимание, что здесь и далее слово "регион" и словосочетание "участок памяти" будут пониматься как полные синонимы.

Системный вызов **mmap()** создаёт новый регион памяти внутри адресного пространства процесса. Конкретно на архитектуре x86, происходит вызов цепочки функций **sys_mmap2()** → **do_mmap2()** → **do_mmap_pageoff()**, последняя из которых и делает всю работу, при этом являясь общей для всех архитектур.

Сначала выполняются проверки входных параметров: мы проверяем корректность используемой файловой системы, убеждаемся, что размер запрашиваемой страницы выровнен и мы не пытаемся выделить страницы в адресном пространстве ядра. Помимо этого, мы проверяем, не слишком ли много на данный момент у процесса уже выделенных страниц.

После вышеперечисленных этапов подготовки происходят следующие шаги:

1. Поиск свободного куска адресного пространства, достаточно большого, чтобы быть использованным в рамках данной аллокации
2. Проверки флагов доступа выбранного региона
3. Аллокация **vm_area_struct** (структура описана выше) с помощью slab-аллокатора.
4. Провязывание страницы в **VMA**

Далее будут рассмотрены несколько процессов, необходимые для полного понимания работы **mmap**.

Первый – нахождение региона памяти, уже выделенного **mmap()**:

Для того, чтобы найти **VMA**, к которому принадлежит определенный адрес, например, во время таких операций, как **page_fault**, обычно используется функция **find_vma()**.

```
1 /* mm/mmap.c */
2 struct vm_area_struct *
3 find_vma(struct mm_struct * mm, unsigned long addr)
4 {
5     struct vm_area_struct *vma = NULL;
6
7     if (mm) {
8         /* Check the cache first. */
9         /* (Cache hit rate is typically around 35%. ) */
10        vma = mm->mmap_cache;
11        if (!(vma && vma->vm_end > addr && vma->vm_start <= addr)) {
12            rb_node_t * rb_node;
13
14            rb_node = mm->mm_rb.rb_node;
15            vma = NULL;
16
17            while (rb_node) {
18                struct vm_area_struct * vma_tmp;
19
20                vma_tmp = rb_entry(rb_node,
21                                struct vm_area_struct, vm_rb);
22
```

```

23     if (vma_tmp->vm_end > addr) {
24         vma = vma_tmp;
25         if (vma_tmp->vm_start <= addr)
26             break;
27         rb_node = rb_node->rb_left;
28     } else
29         rb_node = rb_node->rb_right;
30 }
31 if (vma)
32     mm->mmap_cache = vma;
33 }
34 }
35 return vma;
36 }

```

Сначала она проверяет поле кэша **mmap**, которое запоминает результат последнего вызова **find_vma()**, поскольку вполне вероятно, что один и тот же регион может потребоваться несколько раз подряд. Если это не желаемый регион, то совершается поиск по красно-черному дереву (в поле **mm_rb**). Если желаемый адрес не содержится ни в одном **VMA**, функция аллоцирует с помощью физического аллокатора (см. Главу 4) и возвращает вернет **VMA**, ближайший к запрошенному адресу, поэтому вызывающей стороне важно перепроверить, что возвращаемый **VMA** содержит переданный адрес. Также используется другая функция – **find_vma_prev()**,

```

1  /* mm/mmap.c */
2  struct vm_area_struct *
3  find_vma_prev(struct mm_struct * mm, unsigned long addr, struct vm_area_struct **pprev)
4  {
5      if (mm) {
6          /* Обход красно-черного дерева */
7          struct vm_area_struct * vma;
8          rb_node_t *rb_node, *rb_last_right, *rb_prev;
9
10         rb_node = mm->mm_rb.rb_node;
11         rb_last_right = rb_prev = NULL;
12         vma = NULL;
13
14         while (rb_node) {
15             struct vm_area_struct * vma_tmp;
16
17             vma_tmp = rb_entry(rb_node, struct vm_area_struct, vm_rb);
18
19             if (vma_tmp->vm_end > addr) {
20                 vma = vma_tmp;
21                 rb_prev = rb_last_right;
22                 if (vma_tmp->vm_start <= addr)
23                     break;
24                 rb_node = rb_node->rb_left;
25             } else {
26                 rb_last_right = rb_node;
27                 rb_node = rb_node->rb_right;
28             }
29         }
30         if (vma) {
31             if (vma->vm_rb.rb_left) {
32                 rb_prev = vma->vm_rb.rb_left;
33                 while (rb_prev->rb_right)
34                     rb_prev = rb_prev->rb_right;
35             }
36             *pprev = NULL;
37             if (rb_prev)
38                 *pprev = rb_entry(rb_prev, struct vm_area_struct, vm_rb);
39             if ((rb_prev ? (*pprev)->vm_next : mm->mmap) != vma)

```

```

40     BUG();
41     return vma;
42 }
43 }
44 *pprev = NULL;
45 return NULL;
46 }

```

которая в общем то аналогична `find_vma()`, за исключением того, что она также возвращает указатель на VMA, предшествующий желаемому VMA. Это нужно сделать, потому что мы хотим провязывать VMA в односвязный список, чтобы при возможности, смёрджить несколько страниц в одну. Последней функцией для поиска VMA является `find_vma_intersection()`,

```

1  /* include/linux/mm.h */
2
3  static inline struct vm_area_struct *
4  find_vma_intersection(struct mm_struct * mm, unsigned long start_addr, unsigned long end_addr)
5  {
6      struct vm_area_struct * vma = find_vma(mm, start_addr);
7
8      if (vma && end_addr <= vma->vm_start)
9          vma = NULL;
10     return vma;
11 }

```

которая используется для поиска VMA, перекрывающей заданный диапазон адресов. Наиболее заметное использование этой функции - во время вызова функции `do_brk()`, когда регион увеличивается. Важно убедиться, что растущий регион не будет перекрывать старый регион.

Второй – нахождение еще не выделенного региона памяти:

Когда требуется выделить новую область в памяти, необходимо найти свободный регион, который достаточно велик, чтобы вместить запрашиваемое количество страниц. Функция, ответственная за поиск такой свободной области, называется `get_unmapped_area()`.

```

1  /* mm/mmap.c */
2  unsigned long get_unmapped_area(struct file *file, unsigned long addr, unsigned long len,
3                                unsigned long pgoff, unsigned long flags)
4  {
5      if (flags & MAP_FIXED) {
6          if (addr > TASK_SIZE - len)
7              return -ENOMEM;
8          if (addr & ~PAGE_MASK)
9              return -EINVAL;
10         return addr;
11     }
12
13     if (file && file->f_op && file->f_op->get_unmapped_area)
14         return file->f_op->get_unmapped_area(file, addr, len, pgoff, flags);
15     return arch_get_unmapped_area(file, , len, pgoff, flags);
16 }

```

Функции передается ряд параметров, такие как:

- структура с метаданными по файлу, которому мы сейчас выделяем страницы
- адрес, с которого выделять страницы и длина
- смещение внутри файла
- restriction-флаги области памяти

Если выделение страниц происходит для внешнего устройства, например видеокарты, используется соответствующий ей `f_op->get_unmapped_area()`. Это связано с тем, что внешние устройства могут иметь дополни-

тельные требования к выделению страниц, о которых общий код ядра не может знать, например, выравнивание адреса по определенному виртуальному адресу. Если особых требований нет, вызывается специфичная для архитектуры функция **arch_get_unmapped_area()**. Не все архитектуры предоставляют свои собственные функции. Для тех, кто так не делает, в `mm/mmap.c` предоставляется универсальная версия.

Третий – вставка региона:

Основной функцией для вставки новой области памяти является **insert_vm_struct()**. На самом деле, это очень простая функция. Если разбирать поэтапно, сначала она вызывает **find_vma_prepare()**, чтобы найти подходящие **VMA**, между которыми должна быть вставлена новая область. Помимо этого, мы также спускаемся по вышеупомянутому красно-черному дереву и находим узел для вставки. Затем вызывается **vma_link()**, которая просто провязывает нас в найденные точки `linked_list`'а и красно-черного дерева.

Многие функции ядра не используют **insert_vm_struct()**, предпочитая вместо этого самим вызывать **find_vma_prepare()**, за которой сразу вызывается **vma_link()**, чтобы избежать многократного обхода дерева.

Связывание в **vma_link()** состоит из двух основных этапов, которые содержатся в трех отдельных функциях. Сначала **__vma_link_list()** вставляет **VMA** в вышеупомянутый односвязный список. Если это первое наименьший адрес в адресном пространстве (то есть мы произвели вставки в голову связного списка), оно помимо всего прочего станет корневым узлом красно-черного дерева. Второй этап - привязка узла к красно-черному дереву с помощью **__vma_link_rb()**.

Четвертый – слияние участков памяти:

В Linux раньше была функция под названием **merge_segments()**, которая отвечала за объединение соседних областей памяти вместе, если файл и разрешения совпадали. Цель состояла в том, чтобы уменьшить требуемое количество **VMA**, особенно потому, что многие операции приводили к выделению непомерно большого количества страниц. Это была очень дорогостоящая операция, поскольку она приводила к тому, мы много раз проходили почти по всем **VMA**. Именно поэтому, позже она была удалена, так как программы, которым результат её работы был критично важен, часто пытались её вызывать и проводили ну уж очень много времени в поисках страниц, которые можно сгладить. Взамен ей, сейчас мы имеем эквивалентную функцию **vma_merge()**, но она используется только в двух местах:

1. В **sys_mmap()**, если сопоставляется анонимный участок памяти (потому что анонимные регионы часто поддаются объединению).
2. Во время выполнения функции **brk()**, которая расширяет один участок памяти во вновь выделенный, где два участка должны быть объединены.

Вместо линейного прохода по всем регионам и поиска тем, что могут быть объединены, **vma_merge()** проверяет, может ли существующий регион быть расширен в соответствии с новым распределением, что сводит на нет необходимость создания нового региона. Следует также отметить, что расширение участка памяти может произойти, только если разрешения для двух областей совпадают.

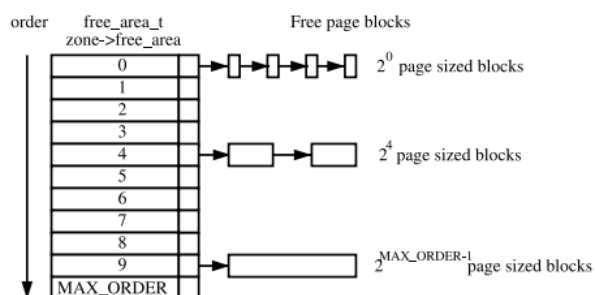
Глава 4. Аллокатор физических страниц

В этой главе мы рассмотрим как работает аллокация физических страниц памяти. В linux за основу взят Binary Buddy Allocator, который выделяет 2^n физических страниц памяти. В нем довольно большая внутренняя фрагментация, однако он идеально подходит под эту задачу, поскольку запрашиваются в основном большие блоки памяти, которые редко освобождаются.

Рассмотрим реализацию buddy аллокатора. Все выделяемые свободные блоки имеют размер в степень двойки: 4, 8, 16, 32, ... страниц физической памяти. Для каждой степени двойки есть специальная структура, хранящая его - **free_area_t**.

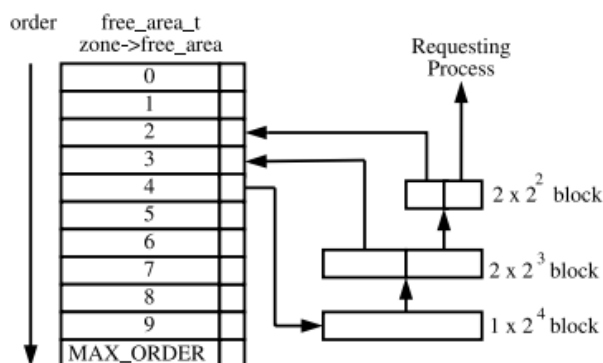
```
1 typedef struct free_area_struct {
2     struct list_head free_list;
3     unsigned long *map;
4 } free_area_t;
```

Хранится статический массив **free_area_t**, где i -ый элемент отвечает за свободные блоки размера 2^i . **free_list** - двусвязный список свободных блоков одного размера. Выглядит это так:



4.1. Аллокация

При запросе на аллокацию N страниц физической памяти мы находим минимальную степень 2 больше N . Далее проверяем наличие свободных страниц такого размера и отдаём если есть. В случае, если двусвязный список оказался пустым, то находим минимальный блок большего размера и рекурсивно разбиваем его пополам до нужного размера. Получившиеся половины пополняют списки свободных блоков нужного размера. Ниже приведен пример разбиения блока при запросе на аллокацию двух физических страниц:



Далее полученный физический адрес транслируется в виртуальный и возвращается пользователю.

4.2. Освобождение

При освобождении блока мы добавляем его в список свободных, а далее проверяем: если соседний блок так же свободен, то мы склеиваем эти два блока и добавляем в список свободных блоков следующего уровня. Эту процедуру мы продолжаем до тех пор, пока соседний блок пуст.

Глава 5. Итоги

Исходя из проведённого исследования, можно сказать, что, вопреки всеобщему мнению, прерывания при вызове **malloc()** вызываются не всегда, а довольно редко и в основном в случаях аллокаций довольно больших кусков памяти. В принципе этапы аллокации на других операционных системах схожи с тем, что описали мы. Отличаются отдельные детали реализации, но идеи и выбранные структуры данных идентичны.

В ходе исследования было также выяснено, что выгоднее не чередовать выделение маленьких кусков памяти с большими, поскольку большие аллокации затирают кэши.

Авторы считают, что данная работа имеет большую ценность и высокую образовательную составляющую. Она может помочь углубиться в понимание того, как работают аллокаторы в целом. Было разобрано более 4000 строк исходного кода и выделены ключевые моменты, позволяющие расширить понимание линуксовых аллокаций.

Список источников

1. Исходный код glibc [Электронный ресурс] : репозиторий. URL : <https://sourceware.org/git/?p=glibc.git;a=tree;hb=refs/heads/master;hb=refs/heads/master>. (дата обращения: 20.12.23)
2. Исходный код malloc.c [Электронный ресурс] : отдельный файл репозитория в частности. URL : <https://sourceware.org/git/?p=glibc.git;a=blob;f=malloc/malloc.c;h=78a531bc7a04bde7032886d08321d913a5541eb3;hb=refs/heads/master>. (дата обращения: 20.12.23)
3. Исходный код linux [Электронный ресурс] : персональный репозиторий Линуса Торвальдса. URL : <https://github.com/torvalds/linux>. (дата обращения: 24.12.23)
4. «Understanding the Linux® Virtual Memory Manager» by Mel Gorman [Электронный ресурс] : техническая документация. URL : https://pdos.csail.mit.edu/~sbw/links/gorman_book.pdf. (дата обращения: 25.12.23)