

```
1 /*
2  * malloc.c --- a general purpose kernel memory allocator for Linux.
3  *
4  * Written by Theodore Ts'o (tytso@mit.edu), 11/29/91
5  *
6  * This routine is written to be as fast as possible, so that it
7  * can be called from the interrupt level.
8  *
9  * Limitations: maximum size of memory we can allocate using this routine
10 *      is 4k, the size of a page in Linux.
11 *
12 * The general game plan is that each page (called a bucket) will only hold
13 * objects of a given size.  When all of the object on a page are released,
14 * the page can be returned to the general free pool.  When malloc() is
15 * called, it looks for the smallest bucket size which will fulfill its
16 * request, and allocate a piece of memory from that bucket pool.
17 *
18 * Each bucket has as its control block a bucket descriptor which keeps
19 * track of how many objects are in use on that page, and the free list
20 * for that page.  Like the buckets themselves, bucket descriptors are
21 * stored on pages requested from get_free_page().  However, unlike buckets,
22 * pages devoted to bucket descriptor pages are never released back to the
23 * system.  Fortunately, a system should probably only need 1 or 2 bucket
24 * descriptor pages, since a page can hold 256 bucket descriptors (which
25 * corresponds to 1 megabyte worth of bucket pages.)  If the kernel is using
26 * that much allocated memory, it's probably doing something wrong.  :-)
27 *
28 * Note: malloc() and free() both call get_free_page() and free_page()
29 *      in sections of code where interrupts are turned off, to allow
30 *      malloc() and free() to be safely called from an interrupt routine.
31 *      (We will probably need this functionality when networking code,
32 *      particularly things like NFS, is added to Linux.)  However, this
33 *      presumes that get_free_page() and free_page() are interrupt-level
34 *      safe, which they may not be once paging is added.  If this is the
35 *      case, we will need to modify malloc() to keep a few unused pages
36 *      "pre-allocated" so that it can safely draw upon those pages if
37 *      it is called from an interrupt routine.
38 *
39 *      Another concern is that get_free_page() should not sleep; if it
40 *      does, the code is carefully ordered so as to avoid any race
41 *      conditions.  The catch is that if malloc() is called re-entrantly,
42 *      there is a chance that unnecessary pages will be grabbed from the
43 *      system.  Except for the pages for the bucket descriptor page, the
44 *      extra pages will eventually get released back to the system, though,
45 *      so it isn't all that bad.
46 */
47
/*
 *      malloc.c - Linux 的通用内核内存分配函数。
 *
 *      由 Theodore Ts'o 编制 (tytso@mit.edu), 11/29/91
 *

```

* 该函数被编写成尽可能地快，从而可以从中断层调用此函数。

*

* 限制：使用该函数一次所能分配的最大内存是 4k，也即 Linux 中内存页面的大小。

*

* 编写该函数所遵循的一般规则是每页（被称为一个存储桶）仅分配所要容纳对象的大小。

* 当一页上的所有对象都释放后，该页就可以返回通用空闲内存池。当 malloc() 被调用

* 时，它会寻找满足要求的最小的存储桶，并从该存储桶中分配一块内存。

*

* 每个存储桶都有一个作为其控制用的存储桶描述符，其中记录了页面上有多少对象正被

* 使用以及该页上空闲内存的列表。就象存储桶自身一样，存储桶描述符也是存储在使用

* get_free_page() 申请到的页面上的，但是与存储桶不同的是，桶描述符所占用的页面

* 将不再会释放给系统。幸运的是一个系统大约只需要 1 到 2 页的桶描述符页面，因为一

* 个页面可以存放 256 个桶描述符（对应 1MB 内存的存储桶页面）。如果系统为桶描述符分

* 配了许多内存，那么肯定系统什么地方出了问题©。

*

* 注意！malloc() 和 free() 两者关闭了中断的代码部分都调用了 get_free_page() 和

* free_page() 函数，以使 malloc() 和 free() 可以安全地被从中断程序中调用

* （当网络代码，尤其是 NFS 等被加入到 Linux 中时就可能需要这种功能）。但前

* 提是假设 get_free_page() 和 free_page() 是可以安全地在中断级程序中使用的，

* 这在一旦加入了分页处理之后就很可能不是安全的。如果真是这种情况，那么我们就

* 需要修改 malloc() 来“预先分配”几页不用的内存，如果 malloc() 和 free() 被

* 从中断程序中调用时就可以安全地使用这些页面。

*

* 另外需要考虑到的是 get_free_page() 不应该睡眠；如果会睡眠的话，则为了防止

* 任何竞争条件，代码需要仔细地安排顺序。关键在于如果 malloc() 是可以重入地

* 被调用的话，那么就会存在不必要的页面被从系统中取走的机会。除了用于桶描述

* 符的页面，这些额外的页面最终会释放给系统，所以并不是象想象的那样不好。

*/

```

48 #include <linux/kernel.h> // 内核头文件。含有一些内核常用函数的原形定义。
49 #include <linux/mm.h> // 内存管理头文件。含有页面大小定义和一些页面释放函数原型。
50 #include <asm/system.h> // 系统头文件。定义了设置或修改描述符/中断门等的嵌入式汇编宏。
51
// 存储桶描述符结构。
52 struct bucket_desc { /* 16 bytes */
53     void *page; // 该桶描述符对应的内存页面指针。
54     struct bucket_desc *next; // 下一个描述符指针。
55     void *freeptr; // 指向本桶中空闲内存位置的指针。
56     unsigned short refcnt; // 引用计数。
57     unsigned short bucket_size; // 本描述符对应存储桶的大小。
58 };
59
// 存储桶描述符目录结构。
60 struct bucket_dir { /* 8 bytes */
61     int size; // 该存储桶的大小(字节数)。
62     struct bucket_desc *chain; // 该存储桶目录项的桶描述符链表指针。
63 };
64
65 /*
66  * The following is the where we store a pointer to the first bucket
67  * descriptor for a given size.
68  *
69  * If it turns out that the Linux kernel allocates a lot of objects of a

```

```

70 * specific size, then we may want to add that specific size to this list,
71 * since that will allow the memory to be allocated more efficiently.
72 * However, since an entire page must be dedicated to each specific size
73 * on this list, some amount of temperance must be exercised here.
74 *
75 * Note that this list must be kept in order.
76 */
/*
* 下面是我们存放第一个给定大小存储桶描述符指针的地方。
*
* 如果 Linux 内核分配了许多指定大小的对象，那么我们就希望将该指定的大小加到
* 该列表(链表)中，因为这样可以使内存的分配更有效。但是，因为一页完整内存页面
* 必须用于列表中指定大小的所有对象，所以需要总做总数方面的测试操作。
*/
// 存储桶目录列表(数组)。
77 struct bucket\_dir bucket\_dir[] = {
78     { 16, (struct bucket\_desc *) 0}, // 16 字节长度的内存块。
79     { 32, (struct bucket\_desc *) 0}, // 32 字节长度的内存块。
80     { 64, (struct bucket\_desc *) 0}, // 64 字节长度的内存块。
81     { 128, (struct bucket\_desc *) 0}, // 128 字节长度的内存块。
82     { 256, (struct bucket\_desc *) 0}, // 256 字节长度的内存块。
83     { 512, (struct bucket\_desc *) 0}, // 512 字节长度的内存块。
84     { 1024, (struct bucket\_desc *) 0}, // 1024 字节长度的内存块。
85     { 2048, (struct bucket\_desc *) 0}, // 2048 字节长度的内存块。
86     { 4096, (struct bucket\_desc *) 0}, // 4096 字节(1 页)内存。
87     { 0, (struct bucket\_desc *) 0}}; /* End of list marker */
88
89 /*
90 * This contains a linked list of free bucket descriptor blocks
91 */
/*
* 下面是含有空闲桶描述符内存块的链表。
*/
92 struct bucket\_desc *free\_bucket\_desc = (struct bucket\_desc *) 0;
93
94 /*
95 * This routine initializes a bucket description page.
96 */
/*
* 下面的子程序用于初始化一页桶描述符页面。
*/
///// 初始化桶描述符。
// 建立空闲桶描述符链表，并让 free_bucket_desc 指向第一个空闲桶描述符。
97 static inline void init\_bucket\_desc()
98 {
99     struct bucket\_desc *bdesc, *first;
100     int i;
101
102     // 申请一页内存，用于存放桶描述符。如果失败，则显示初始化桶描述符时内存不够出错信息，死机。
103     first = bdesc = (struct bucket\_desc *) get\_free\_page();
104     if (!bdesc)
105         panic("Out of memory in init\_bucket\_desc()");
106     // 首先计算一页内存中可存放的桶描述符数量，然后对其建立单向连接指针。

```

```

105     for (i = PAGE\_SIZE/sizeof(struct bucket\_desc); i > 1; i--) {
106         bdesc->next = bdesc+1;
107         bdesc++;
108     }
109     /*
110      * This is done last, to avoid race conditions in case
111      * get_free_page() sleeps and this routine gets called again....
112     */
113     /*
114      * 这是在最后处理的，目的是为了避免在 get_free_page() 睡眠时该子程序又被
115      * 调用而引起的竞争条件。
116     */
117     // 将空闲桶描述符指针 free_bucket_desc 加入链表中。
118     bdesc->next = free\_bucket\_desc;
119     free\_bucket\_desc = first;
120 }
121
122 ///// 分配动态内存函数。
123 // 参数: len - 请求的内存块长度。
124 // 返回: 指向被分配内存的指针。如果失败则返回 NULL。
125 void *malloc(unsigned int len)
126 {
127     struct bucket\_dir      *bdir;
128     struct bucket\_desc    *bdesc;
129     void                    *retval;
130
131     /*
132      * First we search the bucket_dir to find the right bucket change
133      * for this request.
134     */
135     /*
136      * 首先我们搜索存储桶目录 bucket_dir 来寻找适合请求的桶大小。
137     */
138     // 搜索存储桶目录，寻找适合申请内存块大小的桶描述符链表。如果目录项的桶字节数大于请求的字节
139     // 数，就找到了对应的桶目录项。
140     for (bdir = bucket\_dir; bdir->size; bdir++)
141         if (bdir->size >= len)
142             break;
143     // 如果搜索完整整个目录都没有找到合适大小的目录项，则表明所请求的内存块大小太大，超出了该
144     // 程序的分配限制(最长为 1 个页面)。于是显示出错信息，死机。
145     if (!bdir->size) {
146         printk("malloc called with impossibly large argument (%d)\n",
147             len);
148         panic("malloc: bad arg");
149     }
150     /*
151      * Now we search for a bucket descriptor which has free space
152     */
153     /*
154      * 现在我们来搜索具有空闲空间的桶描述符。
155     */
156     cli(); /* Avoid race conditions */ /* 为了避免出现竞争条件，首先关中断 */
157     // 搜索对应桶目录项中描述符链表，查找具有空闲空间的桶描述符。如果桶描述符的空闲内存指针

```

```

// freeptr 不为空，则表示找到了相应的桶描述符。
139     for (bdesc = bdir->chain; bdesc; bdesc = bdesc->next)
140         if (bdesc->freeptr)
141             break;
142     /*
143     * If we didn't find a bucket with free space, then we'll
144     * allocate a new one.
145     */
146     /*
147     * 如果没有找到具有空闲空间的桶描述符，那么我们就需要新建一个该目录项的描述符。
148     */
149     if (!bdesc) {
150         char        *cp;
151         int         i;
152         // 若 free_bucket_desc 还为空时，表示第一次调用该程序，或者链表中所有空桶描述符都已用完。
153         // 此时就需要申请一个页面并在其上建立并初始化空闲描述符链表。free_bucket_desc 会指向第一
154         // 个空闲桶描述符。
155         if (!free_bucket_desc)
156             init_bucket_desc();
157         // 取 free_bucket_desc 指向的空闲桶描述符，并让 free_bucket_desc 指向下一个空闲桶描述符。
158         bdesc = free_bucket_desc;
159         free_bucket_desc = bdesc->next;
160         // 初始化该新的桶描述符。令其引用数量等于 0；桶的大小等于对应桶目录的大小；申请一内存页面，
161         // 让描述符的页面指针 page 指向该页面；空闲内存指针也指向该页开头，因为此时全为空闲。
162         bdesc->refcnt = 0;
163         bdesc->bucket_size = bdir->size;
164         bdesc->page = bdesc->freeptr = (void *) cp = get_free_page();
165         // 如果申请内存页面操作失败，则显示出错信息，死机。
166         if (!cp)
167             panic("Out of memory in kernel malloc()");
168         /* Set up the chain of free objects */
169         /* 在该页空闲内存中建立空闲对象链表 */
170         // 以该桶目录项指定的桶大小为对象长度，对该页内存进行划分，并使每个对象的开始 4 字节设置
171         // 成指向下一对象的指针。
172         for (i=PAGE_SIZE/bdir->size; i > 1; i--) {
173             *((char **) cp) = cp + bdir->size;
174             cp += bdir->size;
175         }
176         // 最后一个对象开始处的指针设置为 0(NULL)。
177         // 然后让该桶描述符的下一描述符指针字段指向对应桶目录项指针 chain 所指的描述符，而桶目录的
178         // chain 指向该桶描述符，也即将该描述符插入到描述符链链头处。
179         *((char **) cp) = 0;
180         bdesc->next = bdir->chain; /* OK, link it in! */ /* OK, 将其链入! */
181         bdir->chain = bdesc;
182     }
183     // 返回指针即等于该描述符对应页面的当前空闲指针。然后调整该空闲空间指针指向下一个空闲对象，
184     // 并使描述符中对应页面中对象引用计数增 1。
185     retval = (void *) bdesc->freeptr;
186     bdesc->freeptr = *((void **) retval);
187     bdesc->refcnt++;
188     // 最后开放中断，并返回指向空闲内存对象的指针。
189     sti(); /* OK, we're safe again */ /* OK, 现在我们又安全了*/

```

```

172     return(retval);
173 }
174
175 /*
176  * Here is the free routine. If you know the size of the object that you
177  * are freeing, then free_s() will use that information to speed up the
178  * search for the bucket descriptor.
179  *
180  * We will #define a macro so that "free(x)" is becomes "free_s(x, 0)"
181  */
182
183 // 下面是释放子程序。如果你知道释放对象的大小，则 free_s() 将使用该信息加速
184 // 搜寻对应桶描述符的速度。
185 //
186 // 我们将定义一个宏，使得 "free(x)" 成为 "free_s(x, 0)"。
187 //
188 ///// 释放存储桶对象。
189 // 参数: obj - 对应对象指针; size - 大小。
190
191 void free_s(void *obj, int size)
192 {
193     void          *page;
194     struct bucket_dir *bdir;
195     struct bucket_desc *bdesc, *prev;
196
197     /* Calculate what page this object lives in */
198     /* 计算该对象所在的页面 */
199     page = (void *) ((unsigned long) obj & 0xffff000);
200     /* Now search the buckets looking for that page */
201     /* 现在搜索存储桶目录项所链接的桶描述符，寻找该页面 */
202
203     //
204     for (bdir = bucket_dir; bdir->size; bdir++) {
205         prev = 0;
206         /* If size is zero then this conditional is always false */
207         /* 如果参数 size 是 0，则下面条件肯定是 false */
208         if (bdir->size < size)
209             continue;
210         // 搜索对应目录项中链接的所有描述符，查找对应页面。如果某描述符页面指针等于 page 则表示找到
211         // 了相应的描述符，跳转到 found。如果描述符不含有对应 page，则让描述符指针 prev 指向该描述符。
212         for (bdesc = bdir->chain; bdesc; bdesc = bdesc->next) {
213             if (bdesc->page == page)
214                 goto found;
215             prev = bdesc;
216         }
217     }
218     // 若搜索了对应目录项的所有描述符都没有找到指定的页面，则显示出错信息，死机。
219     panic("Bad address passed to kernel free_s()");
220 found:
221     // 找到对应的桶描述符后，首先关中断。然后将该对象内存块链入空闲块对象链表中，并使该描述符
222     // 的对象引用计数减 1。
223     cli(); /* To avoid race conditions */ /* 为了避免竞争条件 */
224     *((void **)obj) = bdesc->freeptr;
225     bdesc->freeptr = obj;
226     bdesc->refcnt--;

```

```

// 如果引用计数已等于 0，则我们就可以释放对应的内存页面和该桶描述符。
208     if (bdesc->refcnt == 0) {
209         /*
210          * We need to make sure that prev is still accurate. It
211          * may not be, if someone rudely interrupted us...
212          */
213         /*
214          * 我们需要确信 prev 仍然是正确的，若某程序粗鲁地中断了我们
215          * 就有可能不是了。
216          */
217         // 如果 prev 已经不是搜索到的描述符的前一个描述符，则重新搜索当前描述符的前一个描述符。
218         if ((prev && (prev->next != bdesc)) ||
219             (!prev && (bdir->chain != bdesc)))
220             for (prev = bdir->chain; prev; prev = prev->next)
221                 if (prev->next == bdesc)
222                     break;
223         // 如果找到该前一个描述符，则从描述符链中删除当前描述符。
224         if (prev)
225             prev->next = bdesc->next;
226         // 如果 prev==NULL，则说明当前一个描述符是该目录项首个描述符，也即目录项中 chain 应该直接
227         // 指向当前描述符 bdesc，否则表示链表有问题，则显示出错信息，死机。因此，为了将当前描述符
228         // 从链表中删除，应该让 chain 指向下一个描述符。
229         else {
230             if (bdir->chain != bdesc)
231                 panic("malloc bucket chains corrupted");
232             bdir->chain = bdesc->next;
233         }
234         // 释放当前描述符所操作的内存页面，并将该描述符插入空闲描述符链表开始处。
235         free\_page((unsigned long) bdesc->page);
236         bdesc->next = free\_bucket\_desc;
237         free\_bucket\_desc = bdesc;
238     }
239     // 开中断，返回。
240     sti();
241     return;
242 }
243
244
245

```
