

```

1.  返回 -EAGAIN struct seq_file {
2.      char *buf;
3.      size_t size;
4.      size_t from;
5.      size_t count;
6.      size_t pad_until;
7.      loff_t index;
8.      loff_t read_pos;
9.      u64 version;
10.     struct mutex lock;
11.     const struct seq_operations *op;
12.     int poll_event;
13. #ifdef CONFIG_USER_NS
14.     struct user_namespace *user_ns;
15. #endif
16.     void *private;
17. };

```

[buf, buf + size)是seq_file 分配的memory。

初始只分配PAGE_SIZE size,但如果要输出的内容装不下one page , 则seq_file会在原来size的基础上加一倍大小申请。

```

1.      if (!m->buf) {
2.          m->buf = seq_buf_alloc(m->size = PAGE_SIZE);
3.          if (!m->buf)
4.              return -ENOMEM;
5.      }

```

```

1.         error = m->op->show(m, p);
2.         if (error < 0)
3.             break;
4.         if (unlikely(error)) {
5.             error = 0;
6.             m->count = 0;
7.         }
8.         if (seq_overflow(m))
9.             goto Eoverflow;
10.
11.     .....
12.
13. Eoverflow:
14.     m->op->stop(m, p);
15.     kvfree(m->buf);
16.     m->count = 0;
17.     m->buf = seq_buf_alloc(m->size <= 1);
18.     return !m->buf ? -ENOMEM : -EAGAIN;

```

如果在.show()的过程中，通过seq_printf()写入的内容造成buf overflow，则会释放原来的buf，然后把buf size加倍。

这里有个注意的地方，释放掉原来buf中内容前并不会保留原来通过.show() callback循环调用生成的内容。当扩大buf后，seq_file会重新读取这些内容。

比如,原来有item 0 to item 9共10个iteratable item要输出.

第一次读取是这样的

```

item = seq_file.start()

if(item)

{

    seq_file.show() // 往buf中写入

    item = seq_file.next()

}

seq_file.stop()

```

如果在item 6时，seq_file.show()把buf撑满了，这时seq_file会调用seq_file.stop()，虽然实际上并没有枚举完所有item，

但seq_file.stop()必须与seq_file.start()匹配调用。然后释放原buf。

在扩大一倍buf后，seq_file会再次启动，从item 0重新开始枚举,就象上面的枚举没有发生过一样。

由此可见，为了输出整个virtual file的内容，buf实际上要象输出的内容一样大！

如果输出内容非常多，是否会令kernel失败呢？

目前seq_file这么处理的。

```
1. static void *seq_buf_alloc(unsigned long size)
2. {
3.     void *buf;
4.
5.     buf = kmalloc(size, GFP_KERNEL | __GFP_NOWARN);
6.     if (!buf && size > PAGE_SIZE)
7.         buf = vmalloc(size);
8.     return buf;
9. }
```

如果size大于one page，则从vmalloc区域分配。因为kmalloc分配的是物理上连续的kernel memory，

而vmalloc分配的只是virtual address连续的memory,这样如果申请的size比较大，以不会浪费宝贵的

物理上连续的kernel memory，而是用物理上完全不连着的kernel memory来拼出virtual address连续

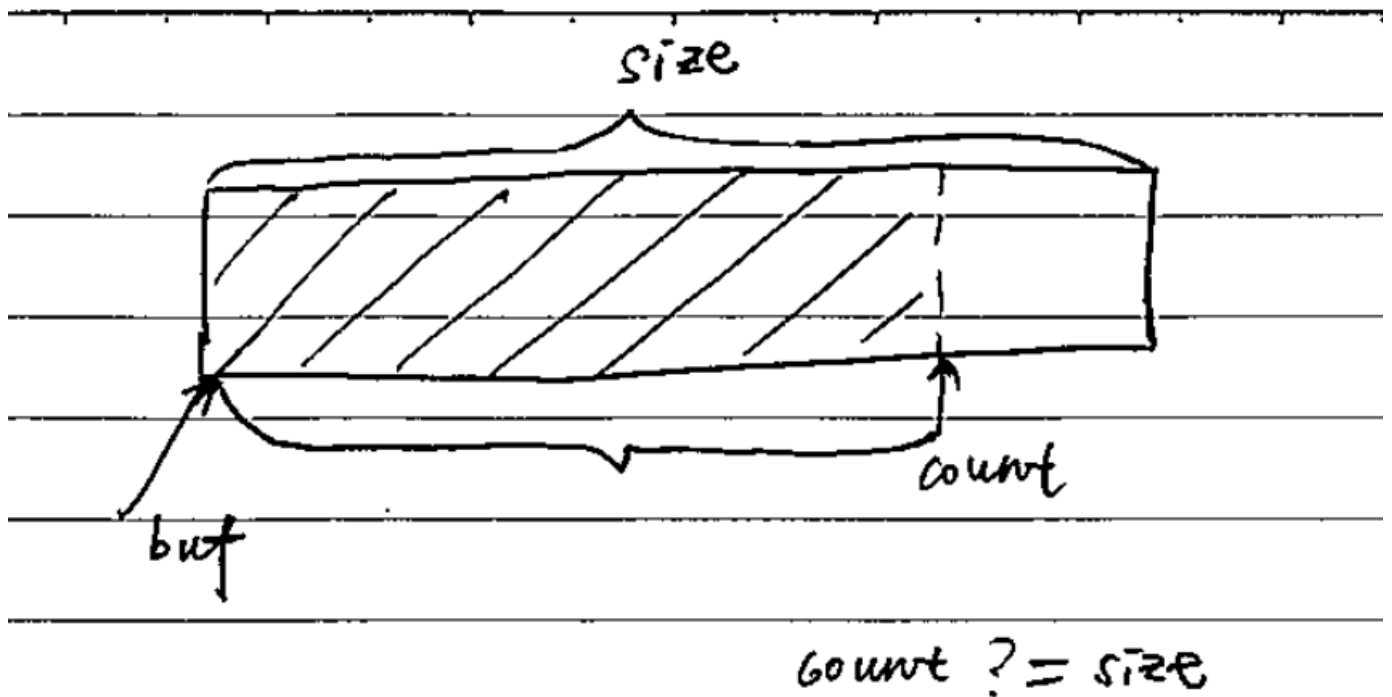
的kernel memory。这对seq_file功能没有丝毫影响。

另外如果连vmalloc也分配失败，则会返回ENOMEM，而不是崩溃。即只是对virtual file的access失败，

读取不出内容而已。

```
return -ENOMEM;
```

count应该是buf中通过.show() callback已经fill的内容。



```
1.  /*
2.   * seq_files have a buffer which can may overflow. When this happens a larger
3.   * buffer is reallocated and all the data will be printed again.
4.   * The overflow state is true when m->count == m->size.
5.   */
6.  static bool seq_overflow(struct seq_file *m)
7.  {
8.      return m->count == m->size;
9.  }
```

通过判断m->count是否等于m->size来判断show的内容是否填满了buf

```

1.  /**
2.   * seq_get_buf - get buffer to write arbitrary data to
3.   * @m: the seq_file handle
4.   * @bufp: the beginning of the buffer is stored here
5.   *
6.   * Return the number of bytes available in the buffer, or zero if
7.   * there's no space.
8.   */
9.  static inline size_t seq_get_buf(struct seq_file *m, char **bufp)
10. {
11.     BUG_ON(m->count > m->size);
12.     if (m->count < m->size)
13.         *bufp = m->buf + m->count;
14.     else
15.         *bufp = NULL;
16.
17.     return m->size - m->count;
18. }

```

返回buf中可以写入的首地址和还可以写入多少字节。

```

1.  *bufp = m->buf + m->count;

```

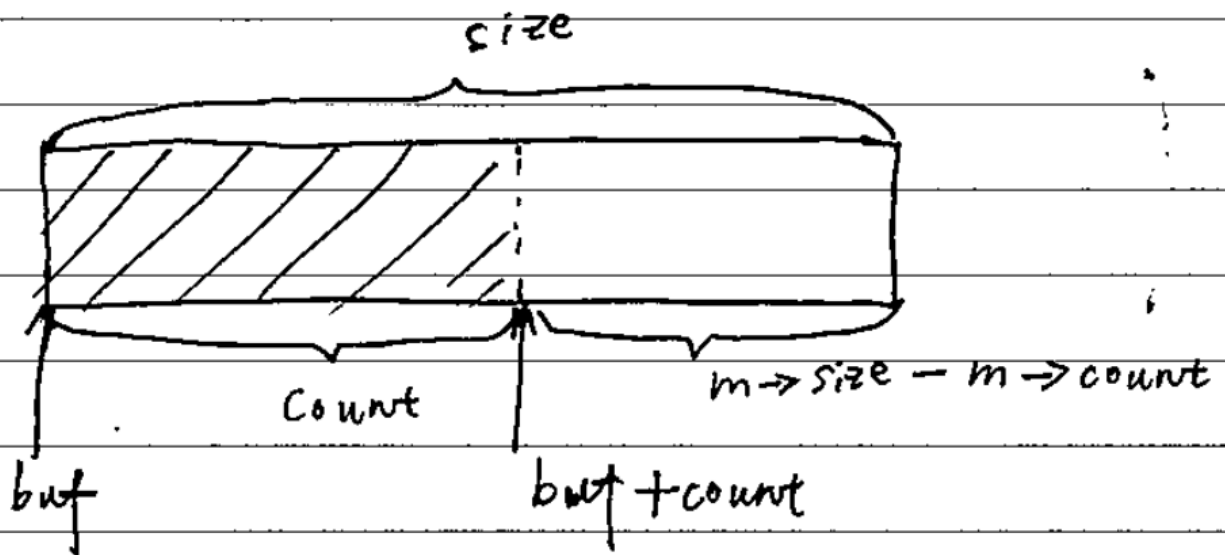
如果满了，则返回NULL

```

1.  return m->size - m->count;

```

还可以写入多少bytes。



read_pos为当前已经读过的内容在整个virtual file中的offset

```

1.  ssize_t seq_read(struct file *file, char __user *buf, size_t size, loff_t *ppos)
2.  {
3.      struct seq_file *m = file->private_data;
4.      size_t copied = 0;
5.      loff_t pos;
6.      size_t n;
7.      void *p;
8.      int err = 0;
9.
10.     mutex_lock(&m->lock);
11.
12.     /*
13.      * seq_file->op->..m_start/m_stop/m_next may do special actions
14.      * or optimisations based on the file->f_version, so we want to
15.      * pass the file->f_version to those methods.
16.      *
17.      * seq_file->version is just copy of f_version, and seq_file
18.      * methods can treat it simply as file version.
19.      * It is copied in first and copied out after all operations.
20.      * It is convenient to have it as part of structure to avoid the
21.      * need of passing another argument to all the seq_file methods.
22.      */
23.     m->version = file->f_version;
24.
25.     /* Don't assume *ppos is where we left it */
26.     if (unlikely(*ppos != m->read_pos)) {           ①
27.         while ((err = traverse(m, *ppos)) == -EAGAIN) ②
28.             ;
29.         if (err) {                                     ③
30.             /* With prejudice... */
31.             m->read_pos = 0;
32.             m->version = 0;
33.             m->index = 0;
34.             m->count = 0;
35.             goto Done;
36.         } else {
37.             m->read_pos = *ppos;                       ④
38.         }
39.     }
40.     .....

```

①

if (unlikely(*ppos != m->read_pos))

这里*ppos就是virtual file的file pointer，如果不等于seq_file当前读取到的在virtual file

中的position，则实际上要作废原来seq_file所有的读取，要重头开始读取，也就是traverse到指定的file pointer,即这里的*ppos。

从此可看出seq_file对随机读是极其极其极其低效的，它只适合sequence read，所以它的name是seq_file

的原因。

②

如果traverse()返回-EAGAIN，则要重新开始traverse()。

```
1.      m->buf = seq_buf_alloc(m->size <= 1);  
2.      return !m->buf ? -ENOMEM : -EAGAIN;
```

返回-EAGAIN的原因是在移动file pointer到*ppos的过程中，发觉buf太小，要加大buf,并再次重头开始读取。

比如一开始buf为one page, 4K (4096)，而*ppos为 17K = 16K + 1K

则第一次traverse()，读了4K，发觉overflow了，释放buf，申请8K

对使用seq_file的client端而言

```
seq_file.start()
```

```
while(seq_file.next())
```

```
{
```

```
    seq_file.show()  到fill 满4K的时候退出loop
```

```
}
```

```
seq_file.stop()
```


第二次traverse(), 读了8K, 发觉overflow了, 释放buf, 申请16K

```
seq_file.start()
while(seq_file.next())
{
    seq_file.show() 到fill 满16K的时候退出loop
}
seq_file.stop()
```

第三次traverse(), 读了16K, 发觉overflow了, 释放buf, 申请32K

```
seq_file.start()
while(seq_file.next())
{
    seq_file.show() 到fill 满17K的时候退出loop
}
seq_file.stop()
```

这次总算成功了, seq_file中size = 32K。

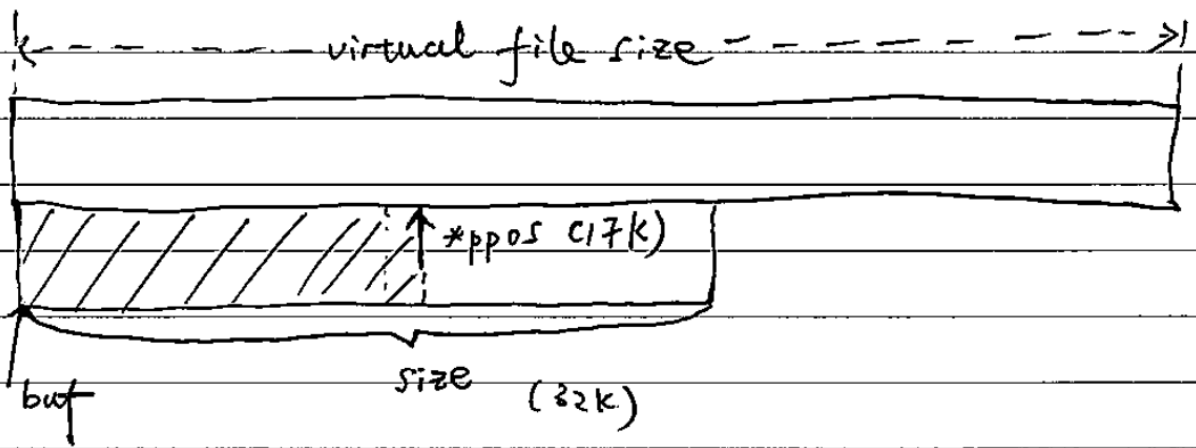
③

如果出错了, 成功的话, traverse()返回0

④

m->read_pos = *ppos

反应了在virtual file中的position。



traverse() 3 遍, 才读到 17k.

关于from field

```

1.  ssize_t seq_read(struct file *file, char __user *buf, size_t size, loff_t *ppos)
2.  {
3.      struct seq_file *m = file->private_data;
4.      size_t copied = 0;
5.      loff_t pos;
6.      size_t n;
7.      void *p;
8.      int err = 0;
9.
10.     mutex_lock(&m->lock);
11.
12.     /*
13.      * seq_file->op->..m_start/m_stop/m_next may do special actions
14.      * or optimisations based on the file->f_version, so we want to
15.      * pass the file->f_version to those methods.
16.      *
17.      * seq_file->version is just copy of f_version, and seq_file
18.      * methods can treat it simply as file version.
19.      * It is copied in first and copied out after all operations.
20.      * It is convenient to have it as part of structure to avoid the
21.      * need of passing another argument to all the seq_file methods.
22.      */
23.     m->version = file->f_version;
24.
25.     /* Don't assume *ppos is where we left it */
26.     if (unlikely(*ppos != m->read_pos)) {
27.         while ((err = traverse(m, *ppos)) == -EAGAIN)
28.             ;
29.         if (err) {
30.             /* With prejudice... */
31.             m->read_pos = 0;
32.             m->version = 0;
33.             m->index = 0;
34.             m->count = 0;
35.             goto Done;
36.         } else {
37.             m->read_pos = *ppos;
38.         }
39.     }
40.
41.     /* grab buffer if we didn't have one */
42.     if (!m->buf) {
43.         m->buf = seq_buf_alloc(m->size = PAGE_SIZE);
44.         if (!m->buf)
45.             goto Enomem;
46.     }
47.     /* if not empty - flush it first */
48.     if (m->count) {
49.         n = min(m->count, size);
50.         err = copy_to_user(buf, m->buf + m->from, n);
51.         if (err)
52.             goto Efault;
53.         m->count -= n;

```

```

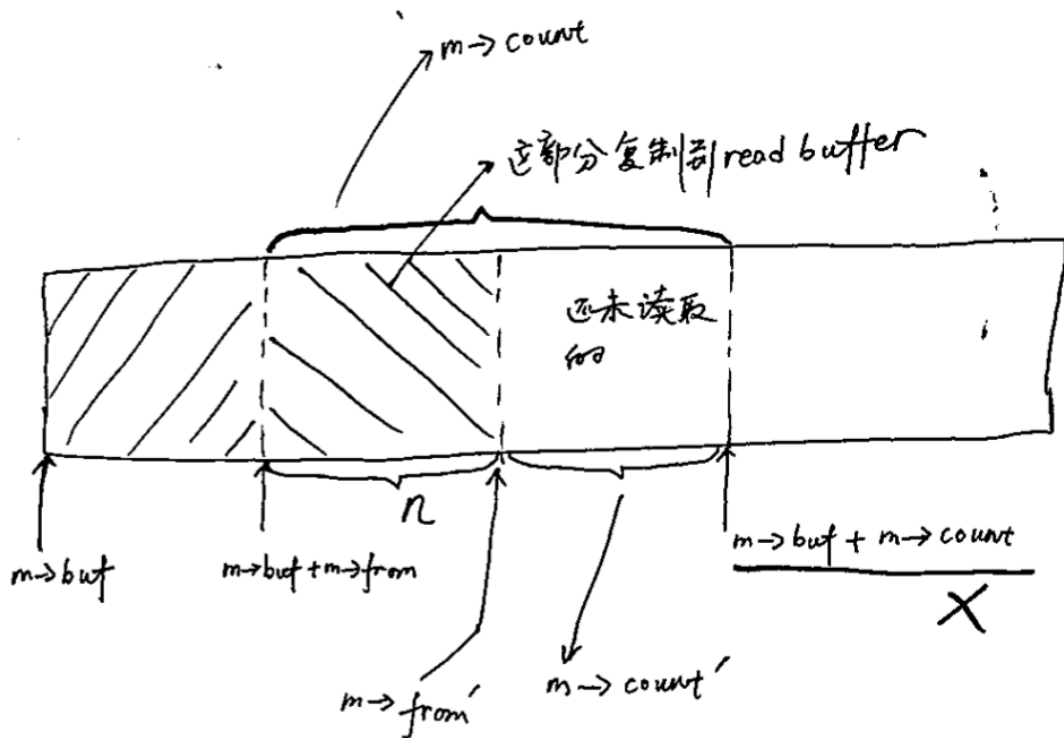
54.         m->from += n;
55.         size -= n;
56.         buf += n;
57.         copied += n;
58.         if (!m->count)
59.             m->index++;
60.         if (!size)
61.             goto Done;
62.     }
63.
64.     .....

```

①

这里 $m \rightarrow \text{count} > 0$ ，表示 $m \rightarrow \text{buf}$ 中有内容。

则需要从 $m \rightarrow \text{buf}$ 中复制内容到 read buffer。稍微有点复杂。



$$\left. \begin{array}{l} m \rightarrow \text{count} -= n \\ m \rightarrow \text{from} += n \end{array} \right\} \Rightarrow \begin{array}{l} m \rightarrow \text{count}' = m \rightarrow \text{count} - n \\ m \rightarrow \text{from}' = m \rightarrow \text{from} + n \end{array}$$

$m->from + n = m->from'$ (new read position in $m->buf[]$)

$m->count' = m->count - n$

这样整个 $[m->buf, m->buf + m->buf + m->size)$ 的空间被分成了3部分

1. $[m->buf, m->buf + m->from)$

这是已经被读取的内容所占空间

2. $[m->buf + m->from, m->buf + m->from + m->count)$

这是还未被读取的内容所占空间

3. $[m->buf + m->from + m->count, m->buf + m->size)$

这是没有valid内容的空间

$m->count$ 的含义这张图比较准确。

index的含义

index用来记录当前seq_file读取到的item index，主要给seq_file->op->start()和seq_file->op->stop()使用。

seq_file假设读取的virtual file中的内容是enumeratable，即不是stream(流式)的。可以item_1, item_2, ...,

item_n这样一个个数出来。有点象自然数和实数一样的感觉。

index就是用于计数这些item的！

```
void * (*start) (struct seq_file *m, loff_t *pos);
```

```
void * (*next) (struct seq_file *m, void *v, loff_t *pos);
```

这两个callback function中的loff_t *pos就是这里index的作用。传递指针是为了callback可以修改这个值。

比如start() callback function , *pos可以告诉它从哪个item开始enumerate,而不一定必定从0开始。

next() callback function可以修改*pos的值，这样下次在调用next()时，就可以接受到上次next()中的修改。

由于enumeration可能被中断，比如m->buf[]空间不够等，所以在seq_file structure中需要记录这个enumerate count。

