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

```
53.                    if (!m->buf)
54.                            goto Enomem;
55.            }
56.            /* if not empty - flush it first */
57.            if (m->count) {                                    ⑤
58.                    n = min(m->count, size);
59.                    err = copy_to_user(buf, m->buf + m->from, n);
60.                    if (err)
61.                            goto Efault;
62.                    m->count -= n;
63.                    m->from += n;
64.                    size -= n;
65.                    buf += n;
66.                    copied += n;
67.                    if (!m->count)
68.                            m->index++;
69.                    if (!size)
70.                            goto Done;
71.            }
72.            /* we need at least one record in buffer */
73.            pos = m->index;                           ⑥
74.            p = m->op->start(m, &pos);
75.            while (1) {
76.                    err = PTR_ERR(p);
77.                    if (!p || IS_ERR(p))
78.                            break;
79.                    err = m->op->show(m, p);
80.                    if (err < 0)                       ⑦
81.                            break;
82.                    if (unlikely(err))                 ⑧
83.                            m->count = 0;
84.                    if (unlikely(!m->count)) {      ⑨
85.                            p = m->op->next(m, p, &pos);
86.                            m->index = pos;
87.                            continue;
88.                    }
89.                    if (m->count < m->size)         ⑩
90.                            goto Fill;
91.                    m->op->stop(m, p);                  (A)
92.                    kvfree(m->buf);
93.                    m->count = 0;                       (B)
94.                    m->buf = seq_buf_alloc(m->size <<= 1);        (C)
95.                    if (!m->buf)
96.                            goto Enomem;
97.                    m->version = 0;
98.                    pos = m->index;             (D)
99.                    p = m->op->start(m, &pos);  (E)
100.           }
101.           m->op->stop(m, p);
102.           m->count = 0;                    (F)
103.           goto Done;
104.   Fill:
105.           /* they want more? let's try to get some more */
106.           while (m->count < size) {        (G)
```

```
107.                    size_t offs = m->count;        (H)
108.                    loff_t next = pos;             (I)
109.                    p = m->op->next(m, p, &next);
110.                    if (!p || IS_ERR(p)) {         (J)
111.                            err = PTR_ERR(p);
112.                            break;
113.                    }
114.                    err = m->op->show(m, p);
115.                    if (seq_overflow(m) || err) {  (K)
116.                            m->count = offs;              (L)
117.                            if (likely(err <= 0))
118.                                    break;
119.                    }
120.                    pos = next;                         (M)
121.            }
122.            m->op->stop(m, p);
123.            n = min(m->count, size);          (N)
124.            err = copy_to_user(buf, m->buf, n);
125.            if (err)
126.                    goto Efault;
127.            copied += n;
128.            m->count -= n;
129.            if (m->count)                            (O)
130.                    m->from = n;
131.            else
132.                    pos++;                           (P)
133.            m->index = pos;
134.    Done:
135.            if (!copied)                         (Q)
136.                    copied = err;
137.            else {
138.                    *ppos += copied;                 (R)
139.                    m->read_pos += copied;           (S)
140.            }
141.            file->f_version = m->version;
142.            mutex_unlock(&m->lock);
143.            return copied;
144.    Enomem:
145.            err = -ENOMEM;
146.            goto Done;
147.    Efault:
148.            err = -EFAULT;
149.            goto Done;
150.    }
```

①

\*ppos != m->read_pos

表示read system call的*ppos与seq_file中记录当前的pos不一致。由于是sequence file，不支持random access，所以需要traverse()来从头开始读取文件，以便把file pointer移动到read()指定的*ppos。

②

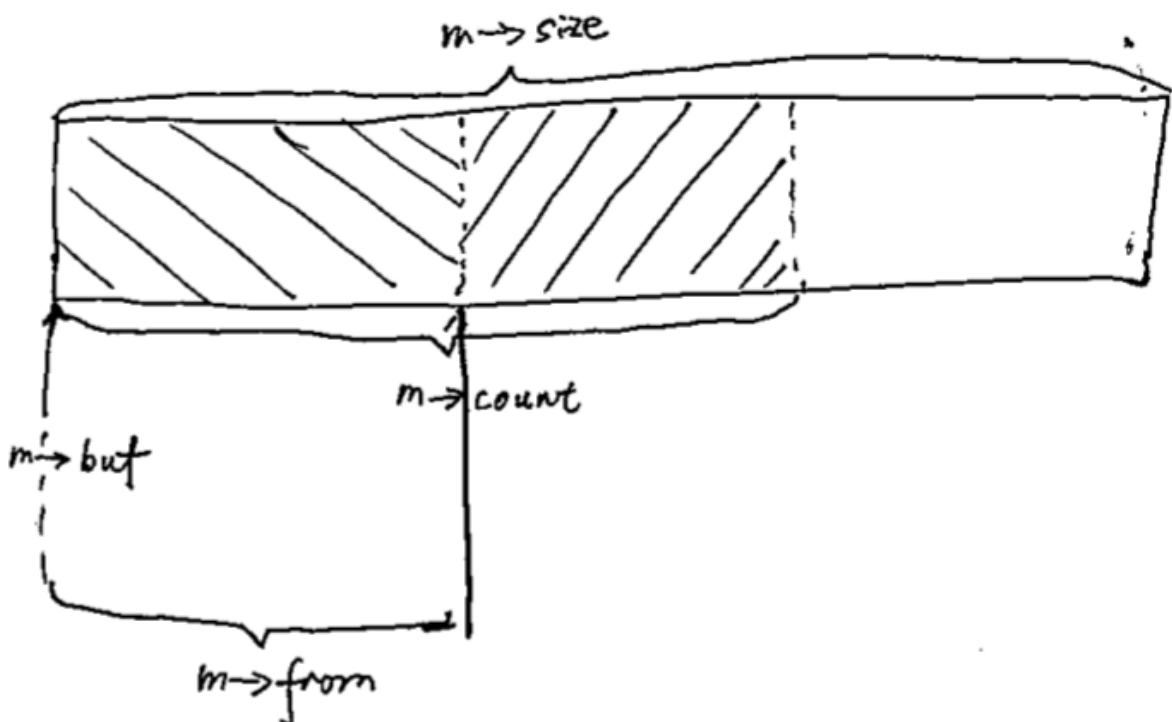在sequence virtual file中移动file pointer是非常低效和费力的，见traverse()的comments。

③

如果移动出错了，则reset seq_file变量m中的fields

④

travserse()移动file pointer正确，则令m->read_pos记录当前virtual file的position。

⑤



m->count > 0

表示m->buf[]中有内容。这里需要复制m->buf[]中内容到read buffer。

```
n = min(m->count, size);

err = copy_to_user(buf, m->buf + m->from, n);

if (err)

        goto Efault;

m->count -= n;

m->from += n;

size -= n;

buf += n;

copied += n;

if (!m->count)

        m->index++;

if (!size)

        goto Done;
```

如果read buffer size < m->count，则m->buf[]中还会留有部分内容。上图中[m->buf, m->buf + m->from)已经被复制到read buffer，

而[m->buf + m->from, m->buf + m->count)则还有效。


⑥

这里的m->index有可能是在traverse()后的值，即比如traverse()跳过了6个item，则m->start()

显然应该接受的参数是7，而不是依然是0。

⑦

err < 0表示m->op->show() fail，所以要break出while loop，并把err值做为错误值返回。

⑧

err非零，也表示m->op->show()出错了，所以令m->count = 0;这实际上就是令

[m->buf, m->buf + m->count)中的内容都无效了。

⑨

从这里一下，m->op->show()返回0,表示正常返回。但如果m->count为0，表示m->buf中是空的。

即上面的m->op->show()放回成功，但其实什么都没有产生。那就通过m->op->next()来enumerate

下一个item。

```
        m->index = pos;

        continue;
```

并且把下一个enumerate的index记录在m->index，然后就再次去.show()该新的item。

⑩

m->count < m->size

如果在m->op->show()以后，m->buf中的还没有填满整个buf,则跳转到去填满整个buf的逻辑。

(A)

运行到这里，表示[m->buf, m->buf + m->size)已经被充满。所以发起.stop(),暂停enumerate

item。

(B)

无效m->buf中的内容

(C)

m->buf[]空间扩大一倍

(D)

这里的m->index还保留着⑥处的值，即从⑥到(A)的整个.start(), .show(), .next(). stop

都白费了。

pos = m->index;

从原来老地方(old index)处再开始enumerate，但m->buf已经被扩大了。

(E)

再次开始新一轮enumeration。

紧接着⑥的while(1) loop的退出，只有两种情况。

1. 出错了

2. 在m->buf[]中存放了所有的enumeratable items,即m->op->start()或m->op->next()返回

NULL,表示没有进一步的输出了，才会退出。如果在此期间，由于m->buf[]空间耗尽，但还有输出，

则把m->buf[]扩大一倍，并从该while lop开始前的index开始enumerate。由此可见比较低效，

但也没办法，因为是sequence virtual file。

(F)

invalid m->buf[]中的内容

(G)

当m->count < m->size,即m->buf[]还由空闲，就会运行到此。

m->count < size,表示在m->buf[]中的内容少于read buffer.

(H)

offs记录下m->buf[]当前的内容的边界，以便在需要时恢复，比如(L)处

(I)

loff_t next = pos;

这里的pos记录的是上次m->op->next()对index的返回值，即下次希望读取的item的index

(J)

m->op->next()出错了

(K)

seq_overflow(m) || err

最新的m->op->show()填满了m->buf或出错了

(L)

m->count = offs;

恢复(H)处的m->buf[]中内容边界，因为???


(M)

与(I)处呼应，保持enumeration index


(N)

要复制到read buffer中的内容长度


(O)

m->count非零，表示m->buf[]中东西比read buffer要求的多


m->from = n;


这样从[m->buf, m->buf + n)已经被复制到read buffer，而[m->buf + n, m->buf + m->count)

还未被读取。m->from就记录了下次该从哪儿读取。


(P)

m->count为零的情况下，即在m->buf[]中的item都被读取了，下次枚举的item index应该从下

一个开始。


(Q)

copied为复制到read buffer的字节数，如果为零，表示出错了

(R)

没有出错

*ppos += copied;

更新virtual file的file pointer

(S)

m->read_pos += copied;

更新seq_file struct中记录的读取到的file pointer。