# Arvados - Bug #3663

## [SDKs] Python CollectionReader should return at least one byte to caller per block read from Keep.

08/23/2014 04:13 PM - Tom Clegg

| | | | | |
|---|---|---|---|---|
| **Status:** | Resolved | | **Start date:** | 09/04/2014 |
| **Priority:** | Normal | | **Due date:** | |
| **Assigned To:** | Tim Pierce | | **% Done:** | 100% |
| **Category:** | SDKs | | **Estimated time:** | 0.00 hour |
| **Target version:** | 2014-09-17 sprint | | | |

**Description**

The manifest format makes it possible to describe a single file consisting of a small amount of data from each of many large data blocks.

- "foo.txt consists of the first 4 bytes of 64M blob A, followed by the first 4 bytes of 64M blob B, ..."

When this happens, clients should expect poor performance, but there are two things we should do to make this less painful.

1. (Big improvement in a future story) Retrieve partial content from Keep. Doing an HTTP request per 4-byte segment will be slow, but much faster than doing an HTTP request and 64MB of disk and network traffic per 4-byte segment!
2. (Small improvement in this story) After fetching a full (or partial) block from Keep, there is always at least one byte ready to return to the caller. Return the available data to the caller right away. Don't fetch the next block until the next read() call.

This is the offending code in sdk/python/arvados/stream.py:

```
    def read(self, size):
        """Read up to 'size' bytes from the stream, starting at the current file position"""
        if size == 0:
            return ''

        data = ''
        for locator, blocksize, segmentoffset, segmentsize in locators_and_ranges(self.segments,
self._filepos, size):
            data += self._stream.readfrom(locator+segmentoffset, segmentsize)
        self._filepos += len(data)
        return data
```

Rather than looping through locators_and_ranges, it should get as much data as it can from the *first element* of locators_and_ranges, and return that.

This mimics the [behavior of io.read()](), except that we don't [yet] support the "unspecified size" and "size=-1" shortcuts.

Related: flush outfile after each outfile.write(data) in arv-get.

**Subtasks:**

| | | |
|---|---|---|
| Task # 3816: Review 3663-collection-reader-performance | | **Resolved** |
| Task # 3745: Print output immediately after read | | **Resolved** |

**Related issues:**

| | | |
|---|---|---|
| Related to Arvados - Feature #3734: [Keep] Keepstore and keepproxy support HT... | **New** | **08/27/2014** |

## Associated revisions

### Revision 90fc7985 - 09/08/2014 02:06 PM - Tim Pierce

Merge branch '3663-collection-reader-performance'

Closes #3663.

## History

### #1 - 08/23/2014 04:33 PM - Tom Clegg

*- Description updated*

### #2 - 08/23/2014 04:33 PM - Tom Clegg

*- Category set to SDKs*

**#3 - 08/23/2014 04:38 PM - Tom Clegg**

*- Description updated*

**#4 - 08/27/2014 02:06 PM - Ward Vandewege**

*- Target version changed from Arvados Future Sprints to 2014-09-17 sprint*

**#5 - 08/27/2014 04:26 PM - Tim Pierce**

*- Assigned To set to Tim Pierce*

**#6 - 09/04/2014 02:59 PM - Tim Pierce**

*- Status changed from New to In Progress*

**#7 - 09/08/2014 02:04 PM - Brett Smith**

[d070454](#) looks good to me.  Thanks.

**#8 - 09/08/2014 02:20 PM - Tim Pierce**

*- Status changed from In Progress to Resolved*

Applied in changeset arvados|commit:90fc79852a995fd8e665cf48ae20c49a9bbc78eb.