

## Arvados - Story #4930

### [FUSE] Design: specify behavior for writable arv-mount

01/07/2015 08:52 PM - Tom Clegg

<b>Status:</b> Resolved	<b>Start date:</b> 01/28/2015
<b>Priority:</b> Normal	<b>Due date:</b>
<b>Assigned To:</b> Peter Amstutz	<b>% Done:</b> 100%
<b>Category:</b> SDKs	<b>Estimated time:</b> 0.00 hour
<b>Target version:</b> 2015-01-28 Sprint	
<b>Description</b> To address: <ul style="list-style-type: none"><li>• When do changes become visible to other Keep mounts? (what is the "commit" button?)</li><li>• When do changes become visible to other procs accessing the same mount?</li><li>• Expected good/bad performance scenarios?</li></ul>	
<b>Subtasks:</b> Task # 4939: Review and discuss <span style="float: right;"><b>Resolved</b></span>	
<b>Related issues:</b>	
Related to Arvados - Feature #4823: [SDKs] Good Collection API for Python SDK	<b>Resolved</b> 12/17/2014
Blocks Arvados - Story #3198: [FUSE] Writable streaming arv-mount	<b>Resolved</b> 04/14/2015

### History

#### #1 - 01/07/2015 08:53 PM - Tom Clegg

- Description updated

- Category set to SDKs

- Assigned To set to Peter Amstutz

#### #2 - 01/08/2015 04:34 PM - Peter Amstutz

Keep FUSE mount will work on an eventual consistency model with specific, well defined commit points where the API server collection record is updated with new manifest text.

Proposed commit points:

- Command line user uses sync(1) or application uses fsync(2). Include current state of files still open for writing. Need to implement FUSE request handlers for fsync()/fsyncdir().
- N seconds after the last file in the collection open for writing is closed (where N is relatively short, maybe 3-5 seconds)
- On unmount.

Propagating changes to other nodes (no conflicts)

1. arv-mount uses arvados.events.subscribe to listen for changes
2. If no files are open, an remotely modified collection can be safely updated (or marked "dirty" for deferred update on next access)
3. If files are open for reading but have been updated remotely, continue to read from snapshot from the point in time the file was opened.
4. Newly opened handles on updated files will see the new file contents.

Merging conflicts:

1. A file is added remotely and not present locally: keep it
2. A file is added locally and not present remotely: commit it
3. A file is deleted remotely, but modified locally. Commit the local version.
4. A file is deleted locally, but modified remotely. Keep the remote version.
5. A file is locally modified and updated remotely: rename the remote version of the file as a "conflicted copy". Commit the local version (based on the snapshot at the point in time the file was opened.)

Performance scenarios. Assume only a small number of files open simultaneously, many simultaneous readers accessing many different files is likely to result in cache thrashing. Default Keep client block cache limit is 256 MiB (holds minimum 4 blocks, may hold more blocks if they are smaller), a bigger cache should be able to support random access read/write on bigger files or larger numbers of open files for before thrashing.

1. Sequential reads/writes on large files: very efficient
2. Random access reads on small amounts of data (up to cache limit) and small number of files being accessed, or files are packed into a few blocks: very efficient
3. Random access reads on small amounts of data (up to cache limit) and large numbers of small files with one block per file: less efficient

4. Random access reads on large amounts of data (over cache limit): less efficient
5. Random access writes on small amounts of data (up to cache limit): very efficient
6. Random access writes on large amounts of data (over cache limit): mostly efficient writes, may result in inefficient reads if file becomes very fragmented

For the last one, we could compute metrics such as the ratio of file size to the total size of the blocks it is distributed over, or number of blocks that must be fetched per block of actual data, and re-pack files when they exceed certain thresholds.

### #3 - 01/19/2015 08:29 PM - Tom Clegg

rename the remote version of the file as a "conflicted copy"

Does this mean we save an extra "conflicted collection"? Or save a "conflicted file" inside this collection? Either way, this seems to produce surprising behavior, where *some* edits cause backup files/collections to accumulate, but most of the time they don't, and it's not obvious why. It might be best to skip this feature entirely: whoever would be inclined to look at these "conflicted" files could probably just look up the collection's history in the logs table -- which would be necessary anyway if you wanted to detect merges that were handled automatically but were incorrect.

I like the way the proposed read semantics align with the traditional POSIX atomic update pattern "make a temp file, close it, and rename it into place": Once you open a file for reading, you continue reading the same file. When you re-open it, you might get a different version. You'll never find yourself reading from a half-written file.

The write semantics are more troublesome, though, because there are opportunities to commit half-written data.

- Open file, write data, crash.
- Open file, write data, some unrelated process elsewhere on the system calls sync().

Some of the benefit of the "read from a snapshot" behavior -- and the symmetry between writing and reading semantics -- is lost if writers are committing new snapshots unexpectedly.

I think it would be more understandable if

- The only way to commit changes to a file is to close() the file.
- During close(), commit happens immediately, and close() returns an error if commit is unsuccessful.
- fsync(), likewise, commits the file and returns an error if commit is unsuccessful.
- If the program crashes, nothing is committed. (Or is this impossible? Can we tell the difference between an explicit close() and an implicit close() by a crashing process?)
- If the mount is detached while files are open, nothing is committed. (Again, if possible.)

This will have worse performance (it will require at least one API call per close()) but it will be predictable, and it makes it possible for the caller to learn that the data was not committed.

Ignoring sync() will cause much better performance if something like redis is calling sync() every second, which (apparently) doesn't seem that crazy in the non-arv-mount world.

Although we don't have the API support yet, putting commit inside close() gives us a chance to use etags to do atomic API updates.

Thoughts?

### #4 - 01/19/2015 10:12 PM - Peter Amstutz

Tom Clegg wrote:

rename the remote version of the file as a "conflicted copy"

Does this mean we save an extra "conflicted collection"? Or save a "conflicted file" inside this collection? Either way, this seems to produce surprising behavior, where *some* edits cause backup files/collections to accumulate, but most of the time they don't, and it's not obvious why. It might be best to skip this feature entirely: whoever would be inclined to look at these "conflicted" files could probably just look up the collection's history in the logs table -- which would be necessary anyway if you wanted to detect merges that were handled automatically but were incorrect.

My proposal is to save a "conflicted file" (the conflicted file being replaced) inside the new collection, alongside the new file. This is based on my understanding of the behavior of [Dropbox](#) (I don't actually use Dropbox, but if they have 200 million users they must be doing something right.)

While introducing new files could cause its own problems, over all it seems less surprising than "I thought I edited this file, but it looks nothing like the way I left it, so it must have deleted my work, so Keep sucks".

Unlike a git merge, we don't have an opportunity to ask the user what to do. If we let the last writer win the conflict, then the user has no way of knowing there ever was a conflict. Right now we don't have any tools for diffing and browsing a collection's log table history that would let us discover that (we should make a story for that, though). I'll add that saving backups could be a selectable/optional behavior, since we're doing the merge on the client side.

The other alternative would be to introduce a write lock on collections. I don't know how I feel about that. Locks in a distributed system have a way of going stale and jamming everything up.

I like the way the proposed read semantics align with the traditional POSIX atomic update pattern "make a temp file, close it, and rename it into place": Once you open a file for reading, you continue reading the same file. When you re-open it, you might get a different version. You'll never find yourself reading from a half-written file.

Yes, having it snapshot the file when opening for reading turns out to neatly solve several problems.

The write semantics are more troublesome, though, because there are opportunities to commit half-written data.

- Open file, write data, crash.

On the other hand, if it is a long running process that is streaming data (such as to a log file!) we probably don't want to lose that.

- Open file, write data, some unrelated process elsewhere on the system calls sync().

Yea, that could be kind of annoying.

Some of the benefit of the "read from a snapshot" behavior -- and the symmetry between writing and reading semantics -- is lost if writers are committing new snapshots unexpectedly.

Well, I'm trying to walk the line between eventual consistency and behavior that looks a little bit more like a shared file system.

I think it would be more understandable if

- The only way to commit changes to a file is to close() the file.

What happens if more than one process tries to open the file for writing? Currently my plan is that the **writers** of the same file on the same mount will make synchronized writes to the same buffer (so they would be able to see each other's writes.) Does it commit on each file handle close() or just when all file handles are closed?

- During close(), commit happens immediately, and close() returns an error if commit is unsuccessful.
- fsync(), likewise, commits the file and returns an error if commit is unsuccessful.

lfuse has fsync() and fsyncdir(). I have not yet determined if sync(1) causes those to be called or not.

- If the program crashes, nothing is committed. (Or is this impossible? Can we tell the difference between an explicit close() and an implicit close() by a crashing process?)

I don't think you can distinguish these cases. I'll have to do some tests. Also this may not be desirable, see my log file comment above.

- If the mount is detached while files are open, nothing is committed. (Again, if possible.)

Well if arv-mount crashes, obviously nothing will be committed. In general, it won't let unmount it while files are still open (without doing a force unmount as root and effectively killing arv-mount anyway.)

This will have worse performance (it will require at least one API call per close()) but it will be predictable, and it makes it possible for the caller to learn that the data was not committed.

Let's mull this over some more. I definitely considered a straightforward commit-on-close initially, but I think it brings its own set of problems.

Ignoring sync() will cause much better performance if something like redis is calling sync() every second, which (apparently) doesn't seem that crazy in the non-arv-mount world.

(Need to see what relationship sync() has with fsync() and fsyncdir())

Although we don't have the API support yet, putting commit inside close() gives us a chance to use etags to do atomic API updates.

Not sure I follow.

**#5 - 01/20/2015 04:34 PM - Peter Amstutz**

More notes/discussion items:



Storage infrastructure should support the features used by reliable programs to achieve reliability. Even if it's true that few programs are reliable, they tend to be the important ones that people... well, rely on. So I don't think that "most likely" comment is either true or especially relevant.

If a program writes 100 small files in sequence, "synchronously commit on close" means we would end up with 100 entries in the history and 100 small blocks, requiring at least 200 HTTP interactions. I think it would be better if we can combine those into a single commit and single block by deferring the commit in order to accumulate writes.

If that's the best we can do, we have to live with it. It's a storage system: reliability is a mandatory feature.

We can optimize by doing commits (and various other things) over websockets. We can write a big block to Keep even though it contains data from unclosed files. Clients can *optionally* turn on `async/unsafe` mode, and blame themselves if they forget to call `commit()` at all the necessary places in their scripts. What we *can't* do is hand users a Keep interface that looks just like POSIX except that it occasionally causes silent data loss when used that way.

## Outstanding questions

...seem to include...

How does a Python SDK consumer arrange for `close()` *not* to flush pending writes, for better performance when modifying many files at once? Setting this mode at the collection level is presumably most convenient, although it could be offered on a file level too.

How does a Python SDK consumer sync, other than `close()`? Presumably `collection.sync()`, and perhaps `file.sync()` is available too? (Would that work on a closed file? Is it a no-op to sync a file that is already closed and synced, or would it invoke a new save/merge?)

How does a FUSE consumer ask for fast-and-loose `close()` behavior? (`arv-mount` could accept a flag. Any other ways to switch on the fly? Per-collection?) When running in that mode, how does a FUSE consumer obtain a guarantee that all writes have been committed? (Presumably, crunch tasks will do this before setting `state=Complete`.)

How does one get the new `portable_data_hash` of a collection after writing to it via FUSE mount?

What do conflicted files look like, exactly? Perhaps a leading "." will make them less intrusive? Or should we go out of our way to make them ugly, like emacs autosave files? (Hope not.) Ending them with "~" would help many tools understand that they're just backup files.

Under which conditions does a conflicted file appear? E.g., while deleting a file I discover that someone else has modified it. In POSIX my delete would win. If I had written garbage to it just before deleting it, I would win (and get a conflicted file). So, should I win and get a conflicted file even if I didn't write garbage?

What is the recommended procedure for cleaning out conflicted files? Is there a way for a client to turn them off, or does a client just have to do an explicit clean-up step after each commit? (Some use cases will want to strip out the conflicted files as a matter of course. `portable_data_hash` will be different for (otherwise) identical manifest content, depending on whether a conflicted file is present.)

### #8 - 01/22/2015 07:38 PM - Peter Amstutz

Tom Clegg wrote:

#### Snapshots / who sees buffered writes?

After looking at `write(2)`, I don't think we can use the "snapshot read" and "snapshot write" behavior within a single mount point, because they will break POSIX-expectant programs in ways that are painful to diagnose (and sometimes hard to detect). Consider a program that opens the same file three times, two `fds` for `r+` and one for `r` (at different points in the file). POSIX requires that the written data is visible to the reader immediately after `write()` returns. Also, if two distinct sections of the file are written by the two writing `fds`, both writes must be visible to the reader while the file is open, and to subsequent readers after the file is closed.

This was my original design (before we went off on the "snapshot read/snapshot write" tangent), for basically the same reasons you state here. The biggest drawback with this approach is the need to take a lock on every read and write. This can impose significant overhead, in <https://arvados.org/issues/4823#note-15> I noted a 30% slowdown between taking and releasing a lock on every read and no locking. So I would like to find a way to avoid unnecessary locking in single-thread crunch scripts.

If we mark the SDK Collection object explicitly as writable or read-only, we can skip taking a lock when the collection is read only.

Fortunately "snapshot read" can be obtained easily enough when desired: open the collection through its `portable_data_hash` path instead of its name or UUID path. (The Python SDK could also offer a way to enable snapshot mode when opening via UUID, `snapshot=True` or something.)

At the SDK level I think we might want to default creating collections as read-only / snapshotted, with a flag to indicate when you want it to be writable / updatable.

#### Commit on close()

It seems to me flushing on `close` *must* be the default behavior in both Python SDK and FUSE because otherwise, standard good programming practices -- which work reliably elsewhere -- will occasionally result in silent data loss when using Keep. This would be contrary to the goal of providing a familiar API, i.e., let programmers do things the way they're used to, and expect them to work.

- Python's IOBase says close does "Flush and close".
- Pipes flush on close().
- close(2) says, in the paragraph preceding the one quoted above, "Not checking the return value of close() is a common but nevertheless serious programming error" because close() can report errors about previous writes.
- Even NFS flushes on close() rather than expecting everyone to add "if might\_be\_on\_nfs { sync() }" to their programs.

Also from close(2):

```
A successful close does not guarantee that the data has been successfully saved to disk, as the kernel defers writes. It is not common for a filesystem to flush the buffers when the stream is closed. If you need to be sure that the data is physically stored, use fsync(2). (It will depend on the disk hardware at this point.)
```

We should at least have an option for close(2) to schedule a background commit of the manifest (individual blocks may have already been uploaded) to the API server, as opposed to synchronously doing the commit and blocking until complete. As noted above, this is already standard practice on many file systems. The drawback is that if it *does* fail, close() won't be able to return an error, but clearly close(2) won't return an error if there is a disk error on a deferred write to a conventional file system.

(On the other hand, network failures are a lot more common than disk failures.)

Here is an example of a program that is perfectly reliable in Linux, but unreliable if close() is not a guaranteed commit:

- [...]

That's a fair point.

If a program writes 100 small files in sequence, "synchronously commit on close" means we would end up with 100 entries in the history and 100 small blocks, requiring at least 200 HTTP interactions.

If that's the best we can do, we have to live with it. It's a storage system: reliability is a mandatory feature.

We can optimize by doing commits (and various other things) over websockets. We can write a big block to Keep even though it contains data from unclosed files. Clients can *optionally* turn on async/unsafe mode, and blame themselves if they forget to call commit() at all the necessary places in their scripts. What we *can't* do is hand users a Keep interface that looks just like POSIX except that it occasionally causes silent data loss when used that way.

Ok, I concede that defaulting to safe behavior is the right thing to do.

## Outstanding questions

...seem to include...

How does a Python SDK consumer arrange for close() *not* to flush pending writes, for better performance when modifying many files at once? Setting this mode at the collection level is presumably most convenient, although it could be offered on a file level too.

This setting should be at the Collection level, because it syncs at the collection level.

How does a Python SDK consumer sync, other than close()? Presumably collection.sync(), and perhaps file.sync() is available too? (Would that work on a closed file? Is it a no-op to sync a file that is already closed and synced, or would it invoke a new save/merge?)

Right now, there is Collection.save() and Collection.save\_as() (save\_as() always creates a new collection record, you can only use save() if you already have a collection record).

Currently it does not call save() on close(). The intended pattern is:

```
with Collection("abc") as c:
    with c.open("foo", "w") as f:
        f.write("hello world")
```

f.\_\_exit\_\_() calls close(), but doesn't save the whole collection.

c.\_\_exit\_\_() calls c.save() which does save the whole collection.

This is no worse than the current CollectionWriter API, which requires that you call finish() to indicate when you want to commit your collection to the API server.

Note that `save()` currently commits all outstanding changes, including still open, partially written files.

How does a FUSE consumer ask for fast-and-loose `close()` behavior? (`arv-mount` could accept a flag. Any other ways to switch on the fly? Per-collection?) When running in that mode, how does a FUSE consumer obtain a guarantee that all writes have been committed? (Presumably, crunch tasks will do this before setting `state=Complete`.)

I was thinking we could add some additional pseudo-files to the root of the Collection directory:

- `.arv-mount.commit_policy` : One of "on\_close", "async", "none". Read from the file to get the current setting, write to the file to change the setting.
- `.arv-mount.commit` : Reading from this will immediately return "1" or "0" indicating whether it is currently synchronized. Writing a "1" to this will trigger a synchronization (if modified) and block until complete (returning an error code if there was an error.)

How does one get the new `portable_data_hash` of a collection after writing to it via FUSE mount?

There is the `.arvados#collection` pseudo-file that contains the most recent API response, including the `portable_data_hash`.

What do conflicted files look like, exactly? Perhaps a leading "." will make them less intrusive? Or should we go out of our way to make them ugly, like emacs autosave files? (Hope not.) Ending them with "~" would help many tools understand that they're just backup files.

This is definitely on the ugly side, but how about:

`XYZ.txt` becomes one of

- `XYZ.txt~conflict-2015-01-22-13:53:22~` to end with ~
- `XYZ~conflict-2015-01-22-13:53:22~.txt` to preserve the file extension

Under which conditions does a conflicted file appear? E.g., while deleting a file I discover that someone else has modified it. In POSIX my delete would win. If I had written garbage to it just before deleting it, I would win (and get a conflicted file). So, should I win and get a conflicted file even if I didn't write garbage?

Yes, it would be more consistent for the local delete to result in the delete "winning" and the remote file becoming a conflict.

What is the recommended procedure for cleaning out conflicted files? Is there a way for a client to turn them off, or does a client just have to do an explicit clean-up step after each commit? (Some use cases will want to strip out the conflicted files as a matter of course. `portable_data_hash` will be different for (otherwise) identical manifest content, depending on whether a conflicted file is present.)

We could add a `conflict_policy` flag (and a `.arv-mount.conflict_policy` pseudo-file?) of either "keep" or "discard".

I see the default policy as being intended as a safety net for people who are updating a collection manually and might be using multiple computers or collaborating with someone else. More sophisticated applications which write collections should either be avoiding conflicts entirely (only produce new collections, don't update existing ones) or be able to apply a custom merge policy (such as diff3 for text files).

## #9 - 01/22/2015 09:10 PM - Tom Clegg

Thoughts

Put all special magic files in a directory, `.arvados` or `.arv-mount` or whatever

- This way you can "ls `.arv-mount`" to get the list of special magic files

Prefer not to end the conflict filename with the original file's extension. Like emacs backup files, `foo.txt~`, and `rsync` temp files, this avoids breaking programs that glob `*.txt` and expect to get only the "real" files, etc.

- `rsync: .ubuntu-14.04-desktop-amd64.iso.xw9PMu`
- emacs autosave via tramp: `#_ascp_balcan_b_ahome_atom_asrc_aarvados_asdk_apython_aarvados_acollection.py#` (such escape...)
- normal backup file: `foo.txt~`

Select mode by using a different class (`BufferedCollection`, `Collection`, `CollectionSnapshot`)? By passing a mode argument to the constructor?

- "mode" is probably not the greatest argument name, since it already means various other things (`r/r+/w/w+/a/a+` or file permissions) in this context. Perhaps `sync=Collection.TRANSACTION`, `Collection.SNAPSHOT`, `Collection.LIVE...`?

It might be OK not to support `Collection.LIVE` fully in the SDK, initially. We could start by requiring the caller to push websocket updates into the collection object in order to stay up-to-date with remote updates. This would be convenient for `arv-mount`, and would just be a documented shortcoming (extra work for the caller) until we add automatic websocket updates to the SDK.

- Perhaps this is `Collection.LIVE_PUSH_ONLY` vs. `Collection.LIVE`, and the latter throws "not implemented" until we implement it?

**#10 - 01/29/2015 05:53 PM - Tom Clegg**

- Status changed from New to In Progress

**#11 - 01/29/2015 05:55 PM - Tom Clegg**

- Target version changed from 2015-01-28 Sprint to 2015-02-18 sprint

**#12 - 01/29/2015 05:56 PM - Tom Clegg**

- Subject changed from [FUSE] Specify behavior for writable arv-mount to [FUSE] Design: pecify behavior for writable arv-mount

- Story points changed from 1.0 to 0.5

**#13 - 01/29/2015 05:56 PM - Tom Clegg**

- Subject changed from [FUSE] Design: pecify behavior for writable arv-mount to [FUSE] Design: specify behavior for writable arv-mount

**#14 - 01/29/2015 08:11 PM - Tom Clegg**

- Status changed from In Progress to Resolved

- Target version changed from 2015-02-18 sprint to 2015-01-28 Sprint