

Tiered Storage for CAMS

Proposal:

As a general rule, we'd like to establish a standard, or some guidelines, for minimum recommended storage capacity for "new" CAMS systems or CAMS systems that are "upgrading" to more cameras. The main goal is to be able to keep 30 days of ALL data before culling and migrating to the archive storage. This paper is not intended to get current CAMS stations to expend additional money to upgrade only their storage... that is, unless storage management is becoming a problem.

One of the issues we've been dealing with for the last several months on all the cams stations around the world is **storage management**. In this paper, I would kind of like to open the discussion and brainstorm the recommendations we should be making to new or expanding cams stations for storage. I want to present to you some of the real numbers so we are all under the same base understanding to start with.

Tier-1 Storage planning:

- A 1 camera station requires over 187 GB for 30 days of storage or 6.2 GB per camera per day.
- A 16 camera station requires over 3.0 TB for 30 days of storage.
- A 20 camera system requires over 3.8 TB for 30 days of storage.

Tier-2 Storage planning:

- A single camera station will require, on average, over 10 MB – 40 MB per camera per night for archival storage, depending on the conditions. On average, we are experiencing about 20 MB on some systems and 41 MB on others.
- A single camera station will require, on average, around 15 GB to store 365 days of archives, providing that the cameras don't have excessive noise.
 - Example 1: Jim Wray's 16 camera system is using 213 GB / yr of Tier-2 storage.
 - Example 2: Lick's 20 camera system consumed 1.33 TB for 04/01/2018 through 09/24/2018 (6 months) for "archived_SubmissionFiles" data. Lick has had some noisy cables, increasing the storage requirements for archiving. These noisy cameras were consuming 365 MB of storage per camera per night.

Overview

For starters, I've looked at various CAMS systems out there that have been running more than a year. It seems to me that several of them have undersized storage. But this depends on our retention goals. In part, it seems that we are not recommending people to purchase the correct storage for the number of cameras they have or are eventually are expanding to. My goal is to **recommend the correct hard drive size for new and expanding stations**.

There are 5 main types of configurations:

- Jim Wray's 16 x 1/3" camera system is a good example of a single-drive system that is working.
- Lick, with 20 Watec cameras is another example.
- Another configuration is Pete Gural's system, which uses SSD and HD combination. *(This would be an ideal configuration if it also had a removable USB3 or NAS storage for Tier-2 storage. However, we have not yet optimized the scripts to take advantage of this ideal configuration)*
- LOCAMS has four 16 camera stations, each with a dedicated Tier-3 archive drive.
- Some systems are using one or more EZCap dongles.

How long to retain data and why?

One thing that determines how long we need to keep the Tier-1 data is our time (or delay) in responding to events. **Very bright fireballs and bolides escape immediate detection.** They can be so bright that they saturate these sensitive cameras, and the FTP_DetectMultipleFF.exe can't determine the centroid. Sometimes we don't get to tracking them down for days or weeks after the event – hopefully before we archive them (the significance here is that part of the archiving step is to **cull** all FF files that don't have detections. If a bolide was not automatically detected, then we cannot do any research or manual reduction on it after the FF file has been culled. Therefore, the plan is to keep all the CapturedFiles data for at least 30 days and to size our Tier-1 CAMS storage accordingly). To handle these undetected events, we need to use the FF files from the CapturedFiles dir. Therefore, it is important that we preserve the FF files that are in the SubmissionFiles\...\CapturedFiles directory for those sessions long enough to be able to go back and do that analysis. My thinking is that **30 days** is ideal.

New CAMS stations should be made aware of the information in this paper to help them in making a decision on the size of their storage. Much of the thinking on this topic revolves around remote systems or those systems that aren't maintained by "extremely" knowledgeable and reliable station operators. This storage management should be automated as much as possible because of that. An extra hundred dollars or so of disk storage in the beginning can make the difference between having to babysit the system all the time or letting it hum away on its own.

CAMS is very storage hungry. On a winter night of capture, each camera will capture for over 14 hours, or over 3 million VGA frames of data. CAMS combines groups of 256 frames and stores them into a single compressed FF file, which is just under 1.2 MB NTSC (or 1.6 MB PAL). For those 14 hours, it totals to about 5,600 FF files per camera per night (just under 5,000 for PAL), or 6.2 GB of FF data files per camera per night. A 16 camera station will consume an average of about 100 GB per night of storage. Therefore, **we have to be smart about storage and storage management.**

Modern storage management for such a system as a CAMS station should include a combination of Tier-0 and/or Tier-1, Tier-2, and possibly Tier-3 storage.

- **Tier-1** storage is where the working directories are kept and number crunching is done from. A fast computer system with a HD as the Tier-1 CAMS storage can usually handle the requirements for 16 cameras without dropping "many" capture frames.
- **Tier-0** storage is essentially used like Tier-1 storage, but it is much faster memory-based (SSD-like) storage. A fast computer system with an SSD as the Tier-0 storage can often handle the requirements for 16 cameras (or more) with a likelihood of not dropping "any" frames.

Note: In CAMS, the OS does not need to be on the SSD drive for boot performance. This is not a gaming computer, so the user's experience is not what we're trying to optimize. A perceivably fast boot of windows doesn't mean anything for CAMS. Putting the OS on the expensive Tier-0 storage makes it so that CAMS processing has less room to do that work. Also, for SSD systems, it makes sense to configure the pagefile as a fixed pagefile and on the larger Tier-1 storage. More about pagefiles later..

- **Tier-2** storage is where items are stored that are no longer necessary to keep live (CAMS Archive). Tier-2 storage is typically a hard drive (HD or Hybrid HD abbreviated as SSHD) and is often slower (cheaper) but it often needs to be larger than Tier-1 storage. They store those items that need to be kept close at hand, but can be brought back into working Tier-1 or Tier-0 storage for analysis. In CAMS, a Tier-2 storage device *could be* an internal drive, external drive (DAS), or network drive (NAS).
- **Tier-3** storage means a more permanent archive. Typically, Tier-3 storage is housed off-site. In CAMS, the zip files in the "F:\Cams_Archives\archive_SubmissionFiles" dir qualify as Tier-3, except that they are not kept off-site by the default configuration. We compress each session into a zip file and place it onto the Tier-2 hard drive. Since the zip requires some special effort to restore back to either of the other tiers, we can call it Tier-3 for the purpose of this discussion. Furthermore, the zip files in the cams_Archive drive should be copied to off-site storage at least once per year to a safe location (if you desire to preserve them). See the paper/discussion on the 01/2020 policy on automatic archive purging. As of 01/2020, the archive scripts now purge archive files older than the number of years specified in the cams2global.ini file.

Pagefiles should not typically be kept on the SSD drive for a few reasons:

- Pro – Accessing the pagefiles is quicker
- Con – Pagefiles take up a lot of space
- Con – If you have 8 GB or more of memory, the system won't use the pagefile during CAMS anyway. In Windows, the pagefile is more of a "last resort" virtual memory paging system for overcommitted read/write memory. Even on a 20 camera system, CAMS would never overcommit 8 GB of RAM. Windows' paging system is not like Unix in that it has what's called "discardable memory", which causes windows systems to use the pagefile less than Unix does. There are more details that I can explain if you care to understand this more thoroughly.
- Con – A pagefile on the SSD drive would consume storage that pushes you closer to a drive that is "nearly full", which performs much worse than one that is less full.

The ideal system would capture to a Tier-0 SSD drive, perform post-capture processing, queue the results to a Tier-1 HD, upload the queued results to the network coordinator, and then "move" the data to either Tier 1 or Tier 2 storage. Thus, clearing the path for more capture to Tier-0 SSD the next evening. The archives would be on a removable external hard drive that could easily be swapped out every 1-3 years.

The actual numbers from real systems might be a little surprising. I've included numbers based on Lick (which had noisy cameras at this time and produces large ArchivedFiles sets) and Forest Hill, which has normal file sizes.

CAMS CapturedFiles data consumes about 6.2 GB/camera/night. Averages from stations that include the CapturedFiles and the ArchivedFiles over the past 30 days consumed between 5.4 GB – 6.3 GB/camera/night (5400 FF files because it's not Dec 21 yet). So, based on that, look what 30 days of storage for various systems would require:

- A 1 camera station requires over 187 GB for 30 days of storage or 6.2 GB per camera per day.
- A 16 camera station requires over 3.0 TB for 30 days of storage.
- A 20 camera system requires over 3.8 TB for 30 days of storage.

Archiving data should be placed on Tier-2 storage devices. When we archive a capture session, we first delete (or "cull") the FF*.bin files from its CapturedFiles dir and keep the ones in its ArchivedFiles and Confirmed files dir. The Archiving data is a zip file that is stored out-of-line from the CAMS working directories. The zip file consumes almost exactly 50% of the disk space as the unzipped files in the set. The zip file includes all of the ArchivedFiles (including its FF*.bin files), log files from the CapturedFiles dir, the FTP dir (includes 2 FF*.bin files for each camera on that board plus the detect and cal files and the weather forecast file), and the Log dir, which is under 50 kb.

The Archiving data consumes between 41 MB/camera/night to 3.75 GB/camera/night, depending on camera noise. Archiving data does not include the FF files that did not have detections (no CapturedFiles FF files). Zip compression reduces the occupied disk space to about to ½ of the size of the original files. The archiving data should be kept around much longer than the CapturedFiles data. It would be nice if **archiving data could be kept around for about a year** in case they are needed in ways they have been needed in the past.

- A single camera station will require, on average, over 10 MB – 40 MB per camera per night for archival storage, depending on the conditions. On average, we are experiencing about 20 MB on some systems and 41 MB on others.
- A single camera station will require, on average, around 15 GB to store 365 days of archives, providing that the cameras don't have excessive noise.
 - Example 1: Jim Wray's 16 camera system is using 213 GB / yr of Tier-2 storage.
 - Example 2: Lick's 20 camera system consumed 1.33 TB for 04/01/2018 through 09/24/2018 (6 months) for "archived_SubmissionFiles" data. Lick has had some noisy cables, increasing the storage requirements for archiving. These noisy cameras were consuming 365 MB of storage per camera per night.

Note: Since most of our systems don't have properly sized Tier-1 storage for 30 days, we have reduced the number of days to keep to about 7 in many cases. This prevents us from being able to go back and examine events that happened more than 7 days prior. One could easily argue that the proper storage management should be to move the files to Tier-2 storage as soon as we can verify that the session has been processed and successfully queued to the server. However, we'd need to come up with a way to preserve the CapturedFiles FF*.bin files for the 30 days. I don't know how to accomplish that right now because of some issues that haven't been solved yet.

One thing that comes to mind is to temporarily compress those files into a zip file. So, we'd have two zip files on the Tier-2 device, one in the "archived_SubmissionFiles" dir and one in the "archived_CapturedFiles" dir. That would only save 50%. So, a 20 camera system might save

1.9 TB by doing that. As those CapturedFiles zip files age past 30 days, we can delete them while preserving the normal “archived_SubmissionFiles” zips. Restoring to Tier-1 working directories would be a two-step process. (1) unzip an “archived_SubmissionFiles” zip to the working drive; (2) unzip the corresponding “archived_CapturedFiles” zip to the working drive. One real problem with this is that it could take several hours to perform the zip operation for a single CapturedFiles session. For those stations with 3 Sensoray capture boards, this could take the better part of the day. Another option for this would be to use the NTFS compressed directory option, which would perform zip-like compression on the files as they are moved onto the “archived_CapturedFiles” dir. Benchmarks I’ve done on this prove that this could be at least as slow as zipping and then deleting all the individual files also takes a very long time. Having “archived_CapturedFiles” dir that contains zip files for individual sessions could treat those sessions as special “containers”, which could be backed up, copied, etc. much easier. One of the problems with this is that you’re going to have zip files that are 30 GB in size. It might be necessary to enable split zipping where the zip file can be split into segments. Another option would be to create archived_CapturedFiles zips for each camera in the capture session. This option might be the easiest way to split the CaptureFiles zips and preserve the original FF files.

To handle the size objection, we could use the “archived_CapturedFiles” zip files option like that mentioned above, but to zip individual cameras. Zipping the CapturedFiles FF files for individual cameras that have the fireball could be de-archived.

Obviously, there is more thought and experimentation that needs to be done with all this.

I/O Performance: One thing that I haven’t mentioned is the time it takes to “move” large sets of files from one tier to the next. It varies from system to system. One thing I have noticed is that a copy operation is handed off by the program to the system and it is the system that actually does the work. This can interfere with, or block, other operations. So, we would not want to be moving data across tiers during capture at the risk of dropping frames.

There are various configurations and drive types.

- SSD drives are essentially storage arrays of memory. They are not recoverable, so when they go bad, everything on them is lost. They are very fast as compared to a standard 5400 RPM HD hard drive. There are several different types of SSD, including SATA, PCIe, M.2, U.2, mSATA, SATA Express, and more. The cost of storage is much higher than for HD and the available storage is less than HD. The cost is much greater. Good for Tier-0 storage. Prices for these have come down. I’m looking at a Samsung 1 TB 2.5” SATA III drive on amazon for as low as \$130.
- SSHD drives are called “Hybrid” drives because they have a hard drive, which is buffered with a small SSD component and internal processing logic to “learn” usage for optimal and predictive usage of the SSD. The typical SSD component is only 8GB. Some have 16 GB. These are more ideal than HD for single drive systems, like laptops, as they offer better performance plus larger storage. Good for Tier 1 storage. The size options of SSHD are not as flexible as HD, typically in the 2 TB or under range.
- HD drives come in various and flexible storage sizes and configurations to suit many needs. This is good for Tier 1 or Tier 2 storage. The price is lower per GB. The rotation speed and buffer

generally govern the speed of the drive. There are 5400 rpm, 7200 rpm, and even 10,000 rpm drives. Each level is more expensive.

Ideal Configuration for 16 cameras:

Drive 1:	1 TB SSD:	C:	(Tier-0) Windows OS, CAMS capturing/working
Drive 2:	4 TB	E:	(Tier-1) CAMS (SubmissionFiles)
Drive 3:	1-4 TB external USB 3.0	F:	(Tier-3) CAMS Archive and Configuration Backup Under \$100

Good Configuration for 16 cameras (but not yet automated):

Drive 1:	2 TB SSHD Partitioned:	C:	Windows OS, etc.
		D:	(Tier 1) CAMS working
Drive 2:	2 TB	E:	(Tier 2) CAMS
Drive 3:	1-4 TB external USB 3.0	F:	(Tier 3) CAMS Archive and Configuration Backup Under \$100

With this configuration, the SSHD is an HD platter with a small SSD component for speed. The SSD part is temporary SSD storage that contains copies of parts of the HD. The HD is a physical drive that is recoverable.

LOCAMS has this configuration, which works well:

C:	512 GB	SSD	Windows OS, etc.
E:	4 TB	HD	CAMS
F:	4 TB	HD	Cams Archive

Lick has this configuration, which has issues from time to time because the Tier-2 and Tier-1 storage are on the same drive.

C:	512 GB	SSD	Windows OS, etc.
D:	6 TB	HD	CAMS and Archive

About once a month, these systems should be checked to ensure there are no issues causing them to not be moving their storage between tiers. About once per year, someone should determine whether the archive storage should be swapped out or purged.

Conclusion:

1. The amount of time that we want to give ourselves for going back and doing analysis, for example, looking for bolides that escaped normal detection is 30 days. Let's call this retention period "MaxDays_Captured". Retention for this is configurable in the Cams2Global.ini file.
2. We should retain the FF*.bin files on the Tier-1 storage for "MaxDays_Captured" days before deleting/culling the FF*.bin files that don't have detections. I will attempt to eventually modify

the code so that the CapturedFiles are zipped and that the capture sessions are queued to the queue dir as soon as possible. Then MaxDays_Captured will trigger purging of the archived_CapturedFiles zip sessions. Retention for this is configurable in the Cams2Global.ini file.

3. We should retain the zip files that we've transmitted up to the SETI server on the Tier-1 storage for at least 60 days, then archive to the Tier-2 storage. We kind of need to keep them at least as long as "MaxDays_Captured" days so that we don't re-upload. I want to refer to this as "MaxDays_Transmitted". Retention for this is configurable in the Cams2Global.ini file.
4. How long should we keep the CAL files on Tier-1 storage? I want to refer to this as "MaxDays_Cal". Retention for this is configurable in the Cams2Global.ini file.

In the meantime, we should be recommending at a minimum that new or expanding CAMS stations should have adequate storage for at least a minimum of "MaxDays_Captured=30" based on:

- 6.3 GB/camera/night storage requirements for Tier 0 or Tier 1 storage with enough storage to accommodate 30 days of this storage.
- 15 MB/camera/night requirement for external USB Tier-2 archive storage with enough storage to accommodate about 1 year of this storage.

This is an open discussion, so please contact us with details before making a final purchasing decision.

Dave Samuels