# 2 Data Storage, Backup and Archiving Architecture

HCS systems scan a multiwell plate with cells or cellular components in each well, acquire multiple images of cells, and extract multiple features (or measurements) relevant to the biological application, resulting in a large quantity of data and images. The amount of data and images generated from a single microtiter plate can range from hundreds of megabytes (MB) to multiple gigabytes (GB). One large-scale HCS experiment, often resulting in billions of features and millions of images that needs multiple terabytes (TB) of storage space. High content informatics tools and infrastructure is needed to manage the large volume of HCS data and images.

## 2.1 Organization of HCS Informatics Infrastructure

There are many rules that are common for the image based HCS informatics infrastructure in academic or non academic organization. Answering the following questions analyzed by entire organization tells one exactly which strategy and organization setup has to be taken and what type of work has to assign to experts and researchers. In choosing the strategy and organization setup one needs to answer the following questions:

- Is the required analysis software available off-the-shelf or must it be written in-house? This decision has to be taken in collaboration between IT and scientists, based on the defined requirements.
- What kind of data will be acquired (how many screens in year)?
- How is the data stored, managed, and protected for short-, medium-, and long-term use?
- What type of desktop clusters and servers are required for HCS computing? (brand, type, speed, and memory)
- How do the computer systems interface with the necessary data collection instrumentation and connect to the network and servers at the same time?
- Can allowances and accommodations be made for external collaborations and programs shared among scientists?
- Are we interested in setup a safety buffered zone outside of our firewalls to allow this external data exchange?

After analysis of those questions one would think to have dedicated IT person from IT department working together with the scientists to allow IT professionals to take over responsibility for informatics tasks. The side-by-side person would allow the informatics organization to understand needs of HCS unit. For example the servers processes could be placed inside of HCS pipeline or infrastructure and not be placed as usual and forced to add extra steps to the workflow. It is also important to decide what will be operated by informatics department and what by HCS unit within organization. It makes better sense for informatics department to own, operate, and manage a data center because they have overview on this and they can provide the service for the researchers also. Some advantages of placing the data center in a central IT department:

- Physical data center
- Facility management (electricity, air conditioning, fire protection, physical security).
- Networking Infrastructure (rules for sockets, switches)
- Security infrastructure
- Existing backup and recovery mechanism.
- Standards for PCs and peripherals, servers, desktop applications, middleware applications, web standards.
- nvestment in informatics infrastructure elements.

All those items allows HCS unit to take advantage also of economies of scale.

## 2.2     Hardware and Network Infrastructure

It is obvious that the High Content Screening units have very specialized computing needs, which are very different from other facilities, research labs of the academic institute or companies. Research laboratories which are using an HCS unit are not able to use off-the-shelf items in most cases because of their specific requirements. They have various instruments so there are requirements for special hardware to connect to instruments and special software for collecting, manipulating, synthesizing and managing the data coming from those instruments.

There are a wide variety of ever evolving options for server hardware, storage hardware, and networking capabilities. The number of HCS instruments, number of users, the number of sites, and the network bandwidth within a site (i.e., Local Area Network), are a few of the key factors impacting the hardware requirements for an informatics solution. Sizing and scoping the optimal hardware for an informatics solution is an area where professional IT support is critical. Nevertheless, each organization is unique in their HCS usage scenarios, which directly impacts the requirements put on an informatics solution. Academic and industrial units have completely different setups. In general, it is best to identify an informatics solution with a system architecture that can scale as the organization's HCS needs development over time. The best way for a new unit is to start with one reader (microscope or scanner) and one robotic platform which will be used in ongoing experiment. The setup for one microscope (reader) has to be organized and connected together in one data flow pipeline (Fig 1). Follow items are demand for pipeline:

- intermediate server (buffer) for data acquisition
- image storage server
- image processing application
- data mining and visualization
- archiving procedure
- databases or laboratory information management systems (LIMS)

Download free eBooks at bookboon.com

Later on, the number of readers may be increased. In non-academic institutions the entire architecture is mostly based on external vendor's informatics applications, so it should be flexible and scalable to fit a variety of hardware configurations and usage scenarios. In academic units the setup is usually less complicated and very often the connection between components is based on a customized development. For both, academic and not academic institution one of the key factors is the bandwidth of the network that connects multiple computers with various operating systems. HCS instruments typically generate data and images at a rate of 1–20 GB (or more in the future) per hour and there are limits to current network and server technology that can support writing this amount of data across networks at multiple sites. The scale of experiments determines the network infrastructure including buffer servers, storage systems and a high speed network. In case of one experiment with one single plate only the microscope local storage space and an external hard drive are required. That is why, balance between networks bandwidth, server, and storage system configurations, and each organization's unique determines how information will be accessed and shared, all need to be taken into account in order to optimize overall system performance.
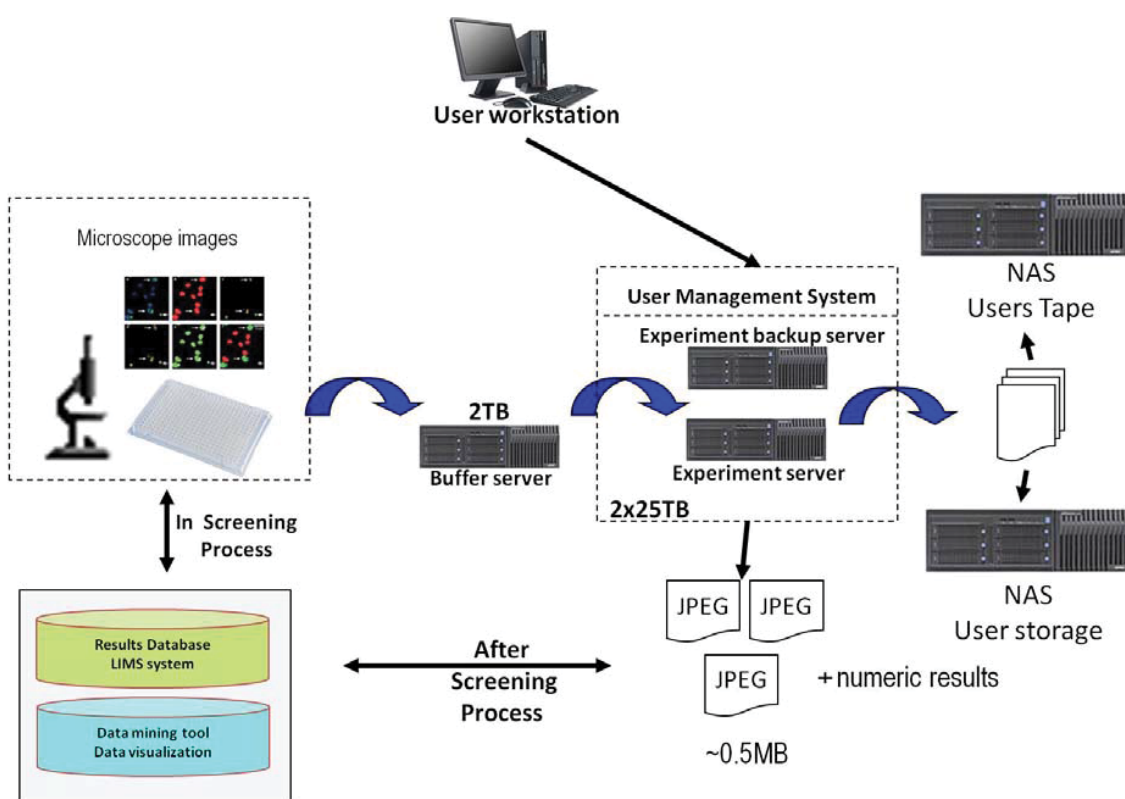


**Fig. 1**. Example of High Content Screening Informatics system architecture.

To manage a tremendous amount of HCS data collected over a period of experiment, an effective and automated data-flow must be developed. During the data-flow setup, the following questions may arise: who is allowed to view and manage the data, how the data will be backed up, and when are data archived or deleted.

Rules or procedures for storing HCS data need to be determined by each organization. Policies and communication between IT experts must be formulated by organizations. In many cases they just decide to play it safe and store everything. Regarding who can view or manage the data, some forethought must occur just to get the system up and running. This is the most difficult part of the entire setup.

HCS unit operation teams and users usually ask the following questions: Are my data secured with backup? Do I have full and easy access to my data from my workstation? Can we protect our data against loss? Therefore having access to professional IT personnel with experience in assigning and managing permissions is extremely important. Managing permissions is also a key point that reveals why user management (UM) applications are so important. The most powerful and commonly used is a UNIX user management system including groups, users where read/write permission set on folder level. Thanks to this type of UM architecture users can be assigned permissions to the file storage. UM also has to be organized to access relational databases which are used as storage for LIMS applications. UM greatly simplifies deploying and managing the system, especially when trying to share data across multiple sites or different domains. Backing up the data is another area where professional IT support is very valuable. In large organizations a dedicated IT department usually helps with archiving of data. The key feature of a successful HCS backup strategy has is preventing the volume of data that needs to be backed up from growing beyond the manageable range of the backup solution. One of the best approaches to achieve this, is to store the HCS data in different locations based on time. A location's time may then be used to determine whether the data has already been backed up. Once a particular location is no longer having data added, a final backup of this location may be completed. A backup policy depends on internal agreement. The data may be archived based on simple criteria such as creation date, storage location, or creating user based. Other option, if biological metadata like projects, compounds, or hits could drive the archive process, then scientists will need the ability to archive data. Regardless of who actually performs the archiving, coordination among users, knowledge of IT department rules, understanding communication between different IT experts and IT staff is vitally important to effectively manage HCS data.

## 2.3 Do we need robust Data Movers (DM) in High Content Screening for data-flow automation?

Manual data movement between hardware elements of HCS informatics is challenging task. Can a data-flow to be automated? Programs running in the background, that take care of moving file-based raw data produced by readers or any other measurement device to a (remote) central storage are called "Data Movers". Here is a list of data movers program freely available and used in HCS:

- DataMover, ETH Zurich, http://www.cisd.ethz.ch/software/Data_Mover (Windows, Linux)
- Rsync, available on every Linux/Unix operating system (Linux)
- RoboCopy, SH Soft, http://www.sh-soft.com (Windows)
- AllwaysSync, Allway Sync, http://allwaysync.com/ (Windows)

In very simple scenario DM would just copy a single plate/run/screen/experiment directory or can recursively copy a directory and its subdirectories. Such tools classify files by whether they exist in the source directory, in the destination directory, or both. Such tools should classify files by comparing time stamps and file sizes between the source file and the corresponding destination file. Users must be able to specify those copies that are restarted in the event of a failure which saves even more time when your network links are unreliable. In general DM should fulfill the following functionality:

- Selectively copy data files.
- Copy data files with full accuracy
- Deletion of plate files and directories from microscope computer after copying is possible
- Allowing the user to control the number of times the program retries an operation after encountering a recoverable network error.
- Allowing the user to control recovery program handle the state when program retries an operation after encountering network error.
- Usage of plate file names, wildcard characters, paths, or file attributes to include or exclude source files as candidates for copying.
- Allow plate file names, wild card characters, paths or file attribute to be included or excluded in source file as candidates for copying.
- Able to exclude directories by name or by path.
- Allowing the user to schedule DM jobs to run automatically.
- Allowing the user to specify when copying is to be performed.
- Monitor directory tree for changes to detect new plate data
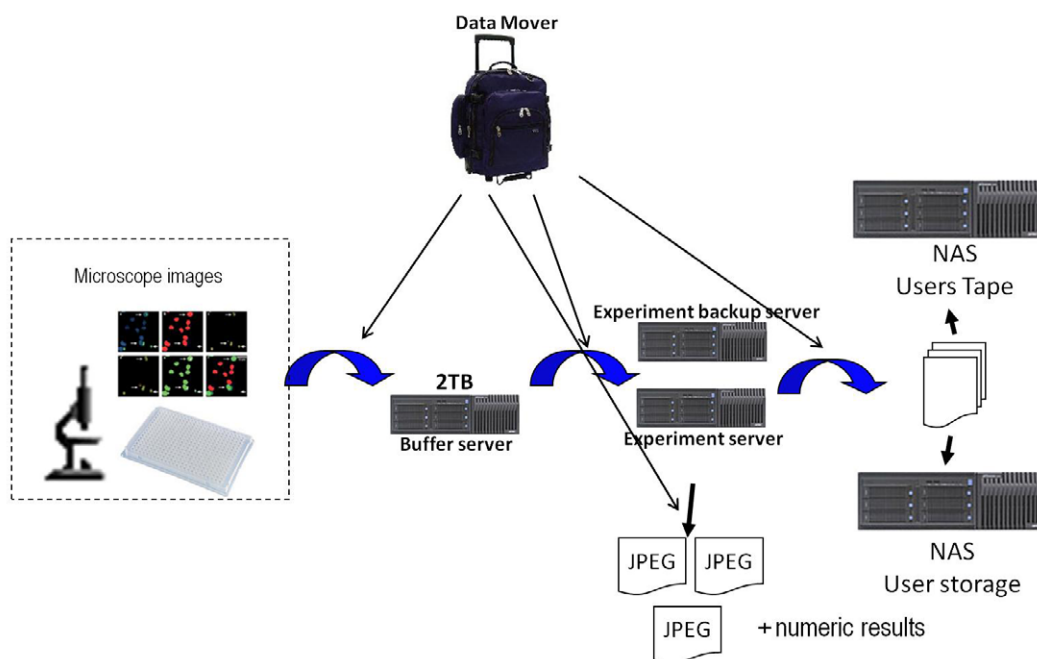- Report in form of logging mechanism errors appeared during data movement



**Fig. 2.** Illustration of Data Mover function in High Content Screening informatics infrastructure.

A role of DM is illustrated on Figure 2. Usually data movement procedures will be repeated in a configurable period of time. Since the buffer server or central storage is connected via a network, the copy process can create network trouble, i.e.: get terminated or stuck. The DM take cares of these problems to the certain extent and let the user know when the problem persists. What other problems have to be addressed in data movement process and what other options DM provides?

**Protection against disk capacity being exceeded.**

For the storage or buffer system, it should be possible to specify a "high-water mark". A high-water mark is the lowest level of free disk space for a given directory. Once the high-water mark is reached (the available free disk space lies below the specified high-water mark), the administrators should be automatically notified (for example via email) and the DM should stop moving files, waiting until sufficient disk space is available again.

**Data completion check**

Data Mover should contain an option which determines if an item (multiwell plate folder) on the microscope local storage has been completed and ready to be moved. This mechanism should start after an incoming file or folder that has not changed during the specified quiet period (see quiet-period option). With the option "data completed timeout" one can specify a time-out (in seconds) for the data movement process. If the script does not finish before the time out it will be killed. Using a "handshake" robustness of the system can increase as it eliminates the need for Data Mover to "guess" when an incoming item is ready to be moved.

**File Cleansing on Microscope Side**

Data cleansing the feature that the removes certain files before moving experiment images to the storage system or buffer. The rationale behind this feature is that sometimes user cannot prevent the microscope from creating certain files that are not needed but would eat up time and network bandwidth when moving to the central storage. It is quite important to remove these files before moving the directory that contains them to the remote side.

**Prefixing Incoming Data Sets**

The prefix option allows setting a prefix for all imaged plates that have to be moved to storage system for analysis and later for archiving. Two examples when this would be needed are given:

1. Assuming that a microscope could produce files or directories with the same name (e.g. derived from the same barcode), prefixing for example with the time point will make it unique.
2. If a screening unit has more than one microscope of same type running in parallel, where data moved to the same outgoing directory and one would still know from which microscope the data are derived

**Extra Local Copy of Produced Images**

If there is a need to access the raw data from an intermediate server (buffer server) or additional storage and the user does not want to transfer these data from the central storage, it is useful to automatically generate an extra temporary copy in a specified a target directory. However it is very useful in extra copy option to use hard links to save disk space, if the file system supports it.

**Dealing with Failures**

When a plate can not be successfully copied to the storage or buffer system even after a specified number of retries. It is important to have a notification system in which, depending on the logger configuration, will send an alarm to an administrator per email or sms. This failure notification enables timely intervention by an administrator.

**Robustness with Respect to Clock Mismatch**

When the microscope computer is located on a different host than the final storage system, there the clocks of the two hosts might be not synchronized. In order to avoid this problem, Data Mover should use an algorithm that ensures that this condition does not lead to a premature transfer and deletion processes (which might even lead to data loss). To this end, the Data Mover should never compare time from the microscope computer and the storage system directly. Instead, the last modification time of a plate (which may be an image file or a directory) is compared to the last modification time of the same item at an earlier time, where the time difference that decides on when to compare last modification times is determined from the storage system clock.

**Robustness with Respect to Program Restarts**

What happened if the program will be restarted and at the same during movement of data?? The directory-based communication has to be preferred over a memory-based one, because it is more robust with respect to restarting the program. This is because with the directory-based approach information is kept on the file system instead of the memory. Thus, restarting the program, or even the server, will restart an operation where it was terminated. It has to be ensured, that after restarting the program recovers properly, finishing all operations that were stopped in the middle. Such mechanism is called a *recovery cycle* which is run automatically after each program start.

Download free eBooks at bookboon.com

**Conclusion**

The organization of on HCS informatics infrastructure and management of highly IT skilled personal is the most difficult part of the entire screening operation. For example, informatics people never thought it would be possible to gain the acceptance, respect and ultimate responsibility to take computing away from the scientists. On the other hand, the scientists did not believe that anyone could do as good a job as they had been doing, no matter what level of expertise people had. There are a number of people in the organization that were PhD biologists in HCS units and are now system administrators or Matlab, R, java programmers, IT project managers. Below is a summary of general recommendation for running HCS IT infrastructure:

- Scientific hardware/software suppliers should be partnering with other computer hardware/software vendors
- All storage systems should be located in the corporate data center and monitored by dedicated expert.
- Data storage system should be part of the enterprise storage program used by the entire institution.
- Backups should be done on a scheduled "off hours" basis to minimize work disruption.
- Side-by-side partnership between research informatics and corporate informatics department with a new understanding of requirements and demands in both parts of the organization.
- Placing proper equipment to proper role of HCS unit: super users are equipped with very high-end pc or workstations, normal users, users traveling/working from home are equipped with laptops or notebooks.
- Standardized desktop systems (hardware and software) should be in place including the same application suite for all (i.e., Microsoft Windows XP, Microsoft Office, antivirus software, and Web access and portals).
- Evaluation of availability of computer equipment located on the actual HCS laboratory (flat panel monitors, 100 megabyte or gigabyte network connections).

It should be mentioned that completion of all rules and setting up an IT infrastructure can take approximately 2–3 years to put in place. Setup period can depend on the numbers of involved IT experts and input from HCS researchers. It is very clear that there is a need for different experts in different areas and although scientists are highly skilled and learned, they do not have the specific professional knowledge that people in the computer industry have about making the correct business decisions related to IT.