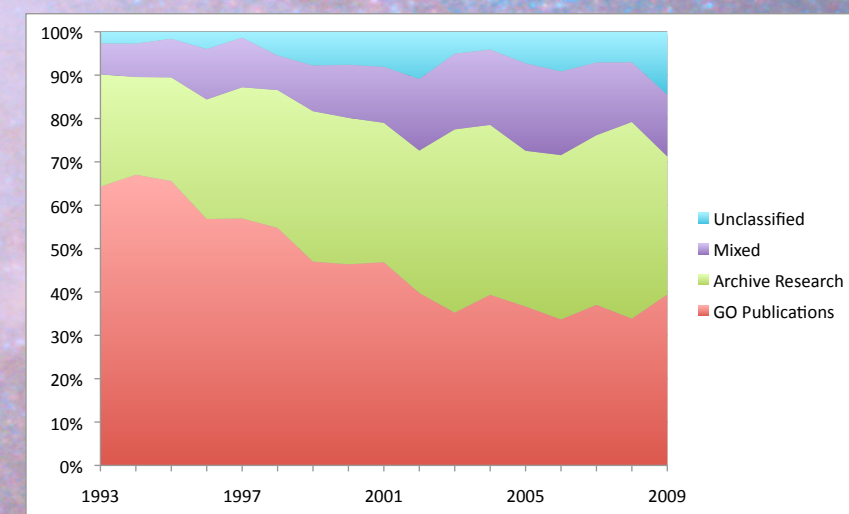# The History and Future of STScI DADS

## Niall Gaffney, Richard Kidwell, Mark Kyprianou, and Faith Abney (STScI)

## An Introduction to DADS

DADS (Data Archive and Distribution System) is used for archive and retrieval of the Hubble, FUSE, Kepler, and the JWST data from STScI. DADS was initially developed by Loral to support the HST archive in 1993. In its current incarnation it provides

- Near-line and deep storage of all data on archival class storage media
- A catalog for all data in the system and key associated metadata
- A flexible XML interface from which different applications can request data
- Options for different methods of data delivery (CD, DVD, ftp, sftp, or staged for ftp)
- A proprietary rights system to allow only approved users to retrieve data during its proprietary period.

As technologies evolved over the sixteen years of this system, so has both the system and users expectations. Computations and storage have largely followed Moore's Law while network performances have also improved. We have also witnessed the growth of archival data publications, growing from a small publication rate to that which currently rivals that of data published by PIs from GO programs.
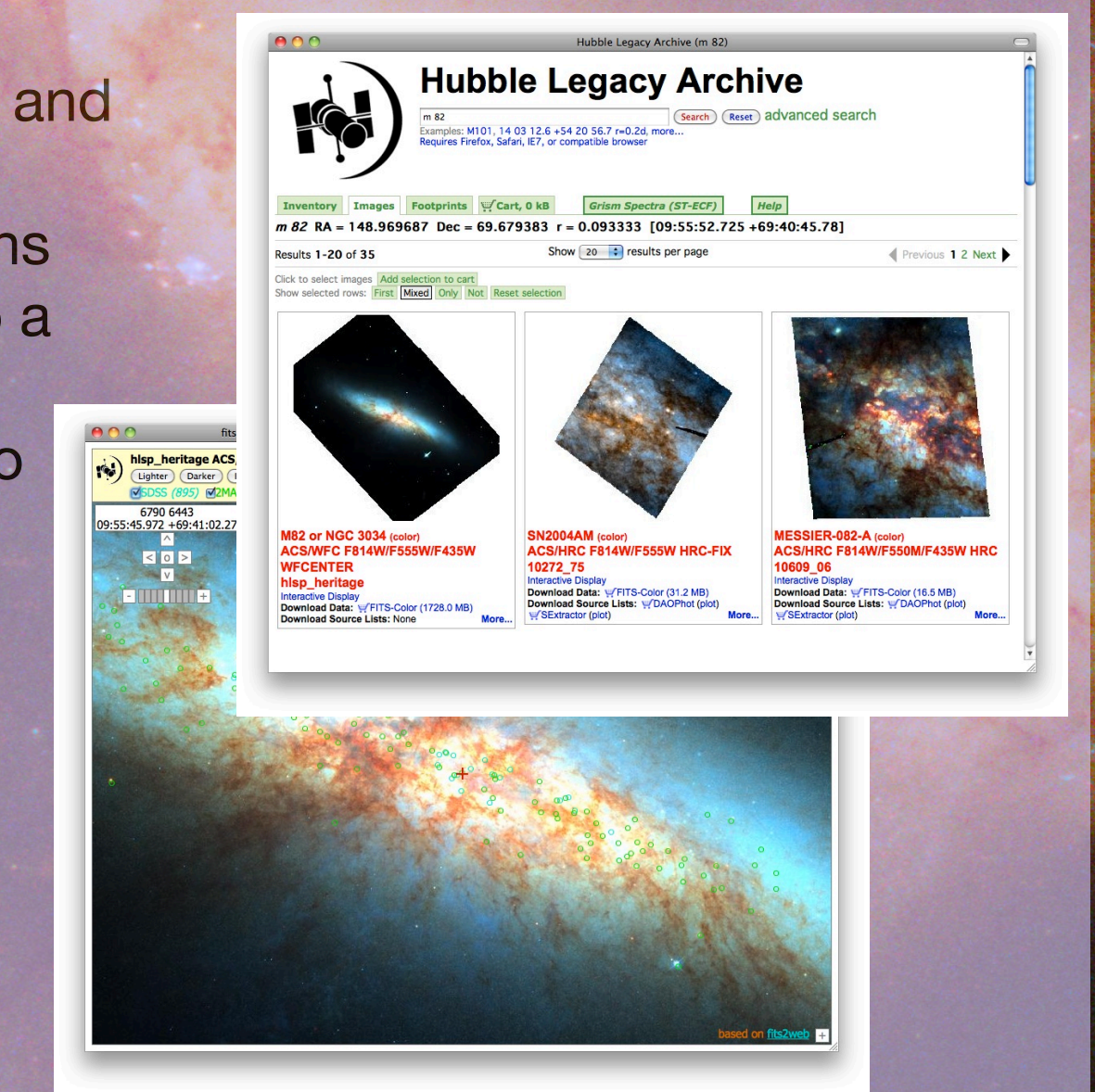
We discuss that evolution here with a focus on what lessons have been learned to craft a flexible, extensible, and reusable astronomical data archiving system from our experiences with DADS and its current (HST, FUSE, KEPLER) and future missions (JWST).

## Changing with Changing Times

In 1999, the on the fly reprocessing system (OTFR) was added to DADS. In this system, data for active instruments are generated at the time of request using the current calibration software and reference files. This allowed STScI to deliver the most up to date data for all instruments that had not been finally calibrated without storing each version on expensive archive media. Initially this increased our average data delivery time by up to a day, but with improvements in hardware this has decreased to less than an hour. Given the reliability of the current system, even the most skeptical of users eventually approved of the change.

With computational hardware following Moore's law and archive media becoming significantly cheaper while hardware on orbit remains limited by communications bandwidth, it has become possible to move back to a recalibrate when needed archive. We are investigating a system inspired by CADC and ECF to return to a static archive and reprocess the data as changes in reference files or calibration files dictate.

The advent of web 2.0 and the VO has created a need for more immediate data responses. These technologies allow archives to create more interactive interfaces with applications such as the Hubble Legacy Archive (HLA) and Aladin.

## DADS round one

The DADS system was originally a custom software package designed and developed by the Loral corporation in 1993. While initially meeting the requirements set forth for the system, the system did not perform up to expectations mainly due to its optimization for large files, which significantly hindered the ingest of the smaller files. This shortcoming was not revealed until the vendor delivered the completed software system for validation. This optimization also pinned the archives storage hardware on the Sony platter system, neglecting any potential future growth. From these mistakes, we learned that

- Numerous small files can lead to suboptimal file storage issues
- Designing a system without an eye towards future technology improvements leads to an inflexible and stagnant system
- Designing purely to requirements without use-cases will often lead to significant shortcomings in the system delivered
- Transparency during the design and development process is important as it can reveal potential problems before they become immutable and expensive to fix

This system ran from 1993 to 2004 when it was superseded by the next generation of DADS, designed and developed completely in house at STScI using these lessons to guide us.

After operating for some time, it became clear that the HST archive, with its one year proprietary data periods for data, was being used more like a library than a classic archive; with most data being retrieved multiple times during the course of the mission. Hence the design of the distribution system is at least as important as data ingest.

## Life with New Missions and Systems

One of the other major drivers for the redesign of DADS was to liberate it from being a VMS-only system to one that can adapt to changes in hardware, OSes, and capabilities driven by the needs of new missions. Over this time, we have

- run under 4 OSes including operationally on multiple versions of Solaris and Linux
- run on 5 different hardware platforms with different core architectures
- upgraded to 3 different long term archive storage media (Sony Platters, MO, UDO)
- ported data between 5 different file systems (UFS, EXT3, QFS, XFS, and ZFS)
- incorporated 3 operational missions, 2 engineering missions, and 1 future mission
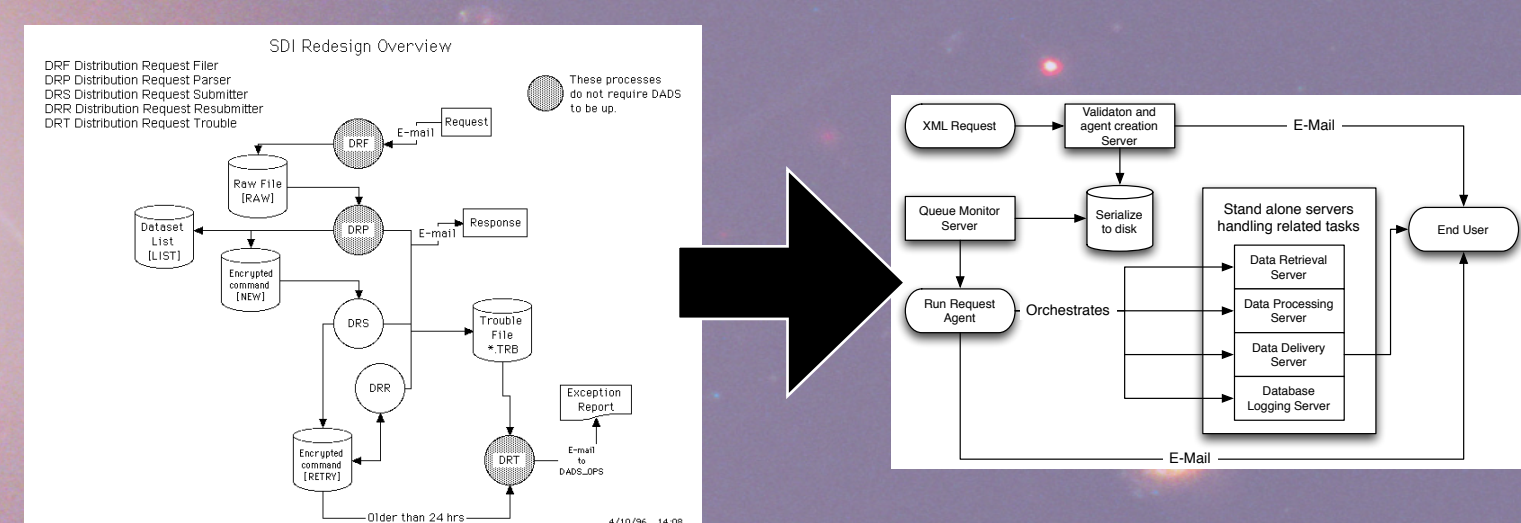- accomplished this with a current development staff of less than 3 FTEs.

This would not have been possible without the modular design of the system that isolates system and mission specifics within the code. Nor could it be done without the virtual and object oriented nature of the Java system; removing us from the bonds of many OS specific problems.

Currently we are also migrating our database servers from Sybase to SQLServer. This migration has been as simple as bcping the current tables from one server to the other and simply changing the JDBC driver within DADS. The code neutral nature of the JDBC has saved us significant porting efforts in migration, as illustrated by the efforts needed to port our other archive subsystems.

## Migration from FORTRAN to Java

From 2001 to 2004 the DADS system was completely redesigned and migrated from FORTRAN to Java. The intent was to create a more generic astronomical data archiving system based on software agents interacting with stand alone servers coordinating individual tasks. While needed at the time, much of the effort went into object serialization and interprocess communications.

Were this developed today, the code could be significantly simplified using technologies such as J2EE to handle object serialization, IPCs, and queue management. That said, the agent based system has lead to more flexibility and extensibility of the system than previously possible due to the isolation of mission specifics within the code. Such mission specific code need only be in the agent creation step. During this redesign, as many of the HST specifics were isolated as possible. However the initial DB design could not be fully restructured while the mission was active.

In addition to mission abstraction, storage media was also abstracted in the archive. By encapsulating the storage requirements in a factory interface, we are now able to easily incorporate new storage media as the media technologies improve over time without significant impact to any other part of the archive.

## DADS 2010 and Beyond

As the needs of the astronomical community grows, so will DADS. With the future missions at STScI planned to be archived using DADS, improvements still need to be made. During this time we learned that

- Design without attention to use cases and potential future growth can lead to systems that do not meet current and future expectations

- While more expensive to develop, designs that encapsulates mission and hardware specifics in smaller components that lead to systems that are cheaper to expand and modify once operational.

- It is more expensive and complicated to refractor a system once it is operational

- Science archives that make their data public behave more like libraries than classic archives and should focus heavily on data discovery and distribution.

With redesigns inspired by these lessons learned, DADS has been able to support missions with very different data requirements (e.g. HST and Kepler) with little to no changes needed to its core architecture. Many of these lessons were learned while the system was actively archiving and distributing HST mission data which lead to more complexity in their implementation. However, these lessons have been key to the creation of the current flexible and reusable archive system making it viable for the needs of the JWST project and beyond.