# Efficient data reduction and analysis of DECam images using multicore architecture Poor man's approach to Big data

## Roberto Muñoz

### Instituto de Astrofísica
### Pontificia Universidad Católica de Chile

Thomas Puzia, Maren Hempel, Paul Eigenthaler, Matt Taylor, Yasna Ordenes, Simón Angel

Ariane Lançon, Steffen Mieske, Michael Hilker, Patrick Côté, Laura Ferrarese & NGFS team.

**Tools for Astronomical Big Data**
**Tucson, Arizona 2015**

# Challenges of Big Data

- Computer scientists have been aware of big data for many years, maybe decades. They paved and keep paving the road for dealing with big amounts of data.

- Three concepts are repeated: Volume, Variety and Velocity

**Volume**: Quantity of data. GB, TB, PT.
　　　　　How large will be the data release?
**Variety**: Different types of data.
　　　　　Infrared images, High-res spectra
**Velocity**: Rate at which generate data.
　　　　　GB, TB or PT per night?
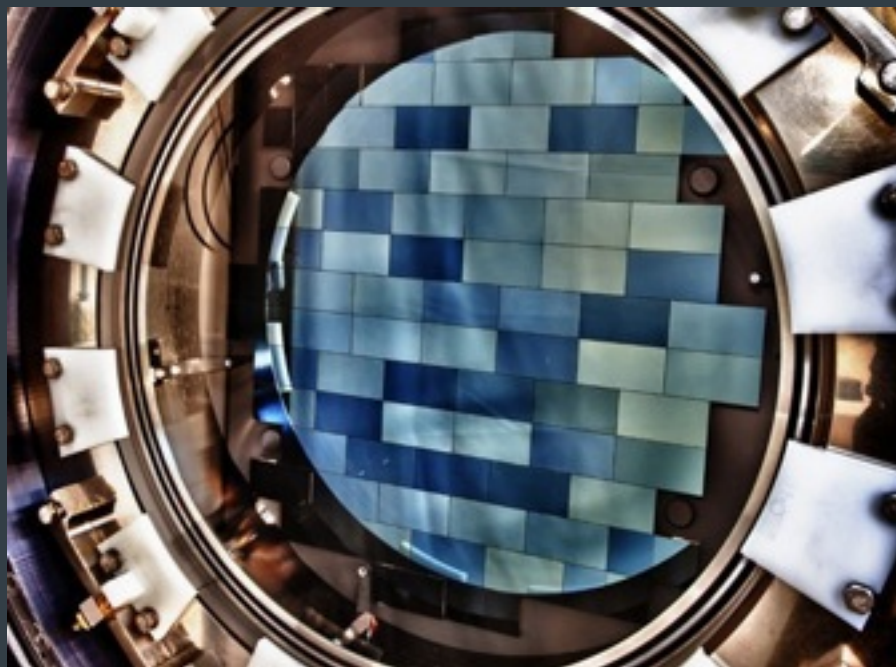
Roberto Muñoz
Tools for Astronomical Big Data

# Astronomy and big data

- Science is the main driver of any scientific project where astronomers are involved.

- We were not looking for big amounts of data in the Universe, it just happened that cameras became bigger and bigger and the data processing became more complex.

- Maybe the first project where astronomers were forced to face big data challenges was the SDSS survey. Final data release (DR12) consists of 100 TB.

# What about DECam?

- Each image is about 1 GB and raw data per night could be up to 0,5 TB. At this stage 16 bit signed integer.

- All the raw data is transferred to NOAO servers, processed by the Community pipeline and then the PI is able to download different types of products. At this stage 32 bit floating point.

- Just downloading the images and weight maps, the 0,5 TB/night becomes 3 TB/night of just Community pipeline processed data.
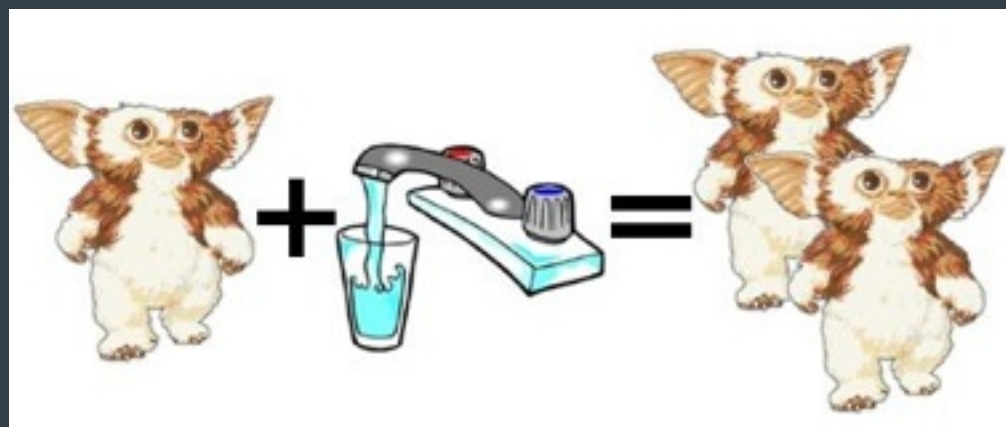
Roberto Muñoz
Tools for Astronomical Big Data

# What about post processing?

- Depending on the science you want to do with DECam, maybe you need to download all the individual images processed by the Community pipeline and do some post processing.

- Let's say you have a very particular observational strategy and want to model the sky from each individual image. But for doing that, you need to know what pixels belong to the sky.

- The easiest way is building a first stack, detecting the sources (SExtractor) and masking them. You end up with a mask per image.

- Since you can't load the entire 3 TB of data into memory, you need to create the individual sky and sky subtracted images, and then you can stack again.



0,5 TB ➡ 5 TB per night

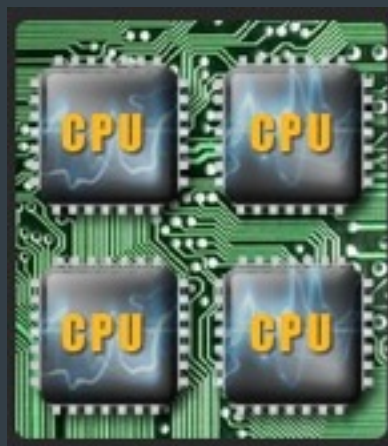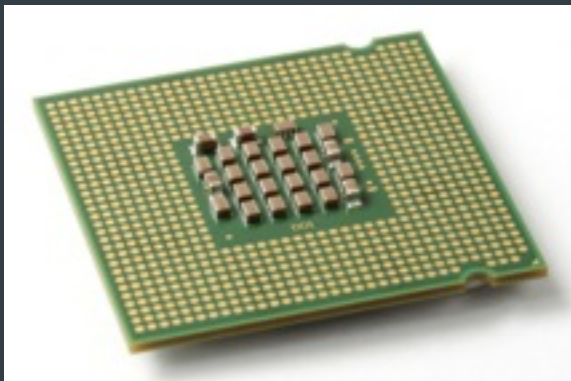Roberto Muñoz
Tools for Astronomical Big Data

# RAID systems

- Our programs consists of about 5 nights per year. We need about 25 TB of space for post processing.

- We found that video professionals (movies) have been using small and medium sizes RAID systems for a while. Nowadays, the technology is cheaper and small groups can setup a system.
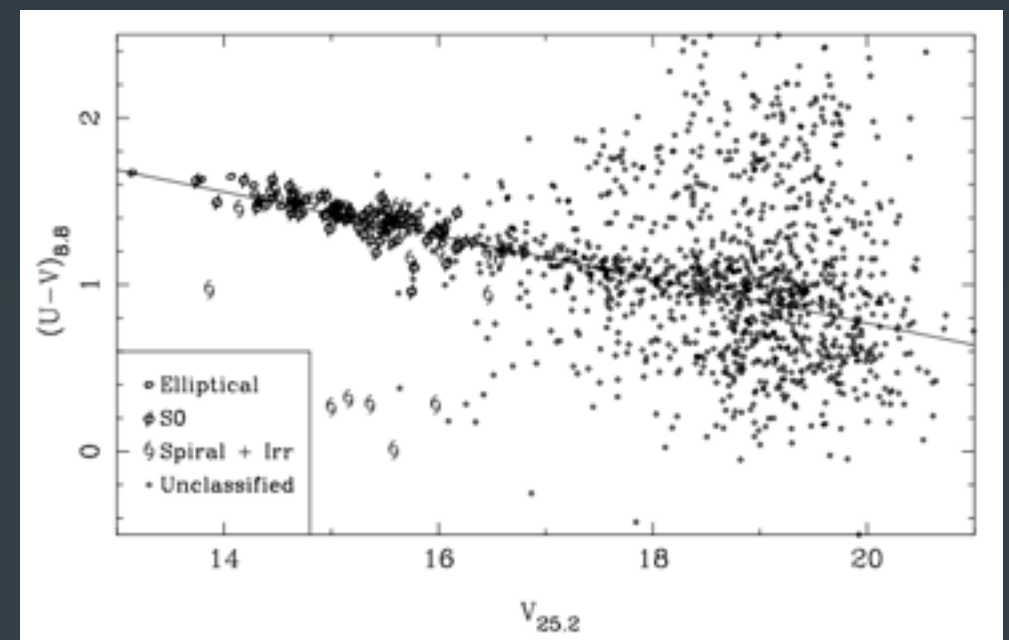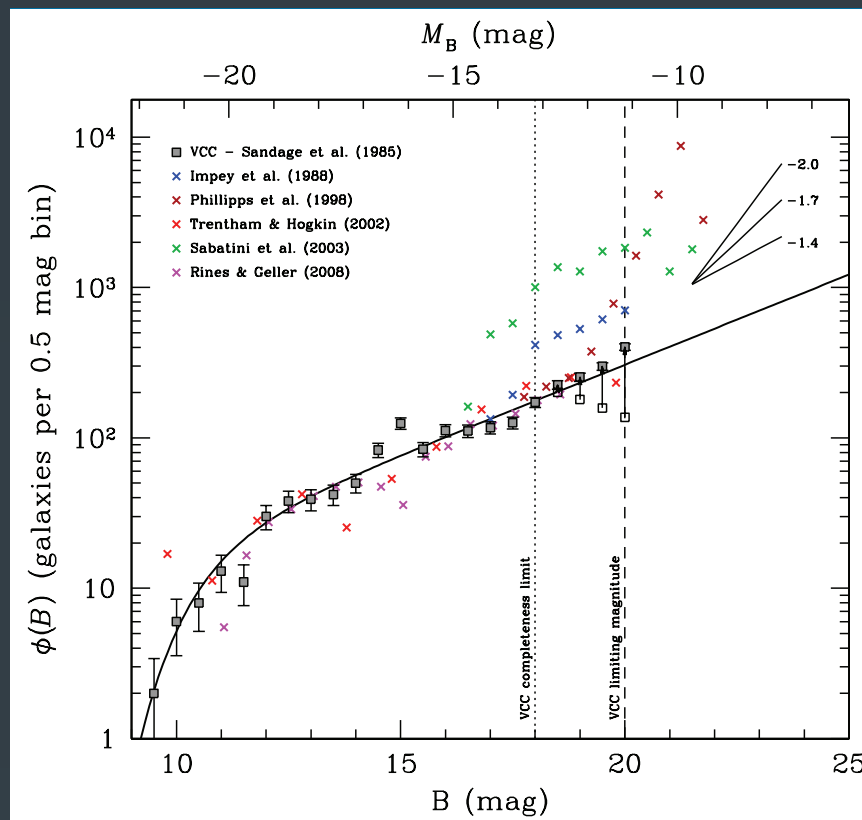
# Multicore systems

- Multicore systems started becoming popular around 2005. Most of the computers and smartphones we use nowadays are multicore.

- Large storage capacity + High I/O data transfer + Fast CPU processing + Low cost = Short execution times and accessibility for testing new algorithms for image processing
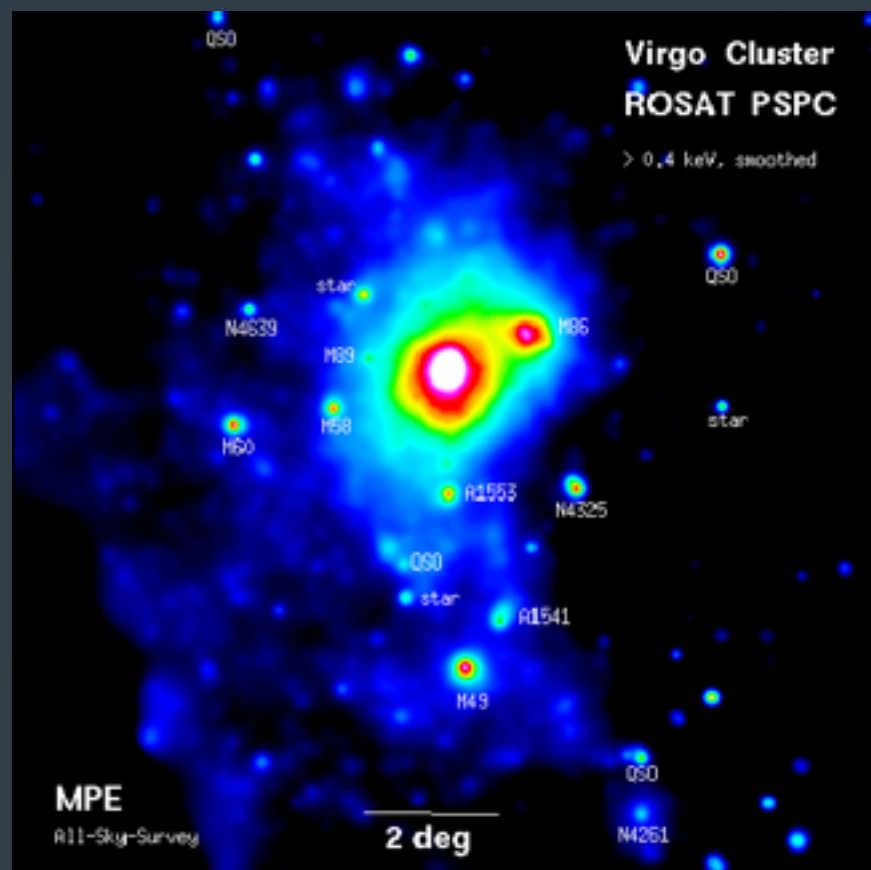
Roberto Muñoz

Tools for Astronomical Big Data

# Science: groups and clusters

- Our team is using DECam to observe several galaxy groups and clusters closer than 20 Mpc.

- We are interested in studying the stellar mass function down to $10^6$ $M_O$, how environment affects galaxy evolution and how the GCs are distributed around galaxies and ICM, among several other goals.

- Similar questions were asked by Binggeli et al. (1985) about the Virgo galaxy cluster and by Ferguson et al. (1989) about the Fornax galaxy cluster.
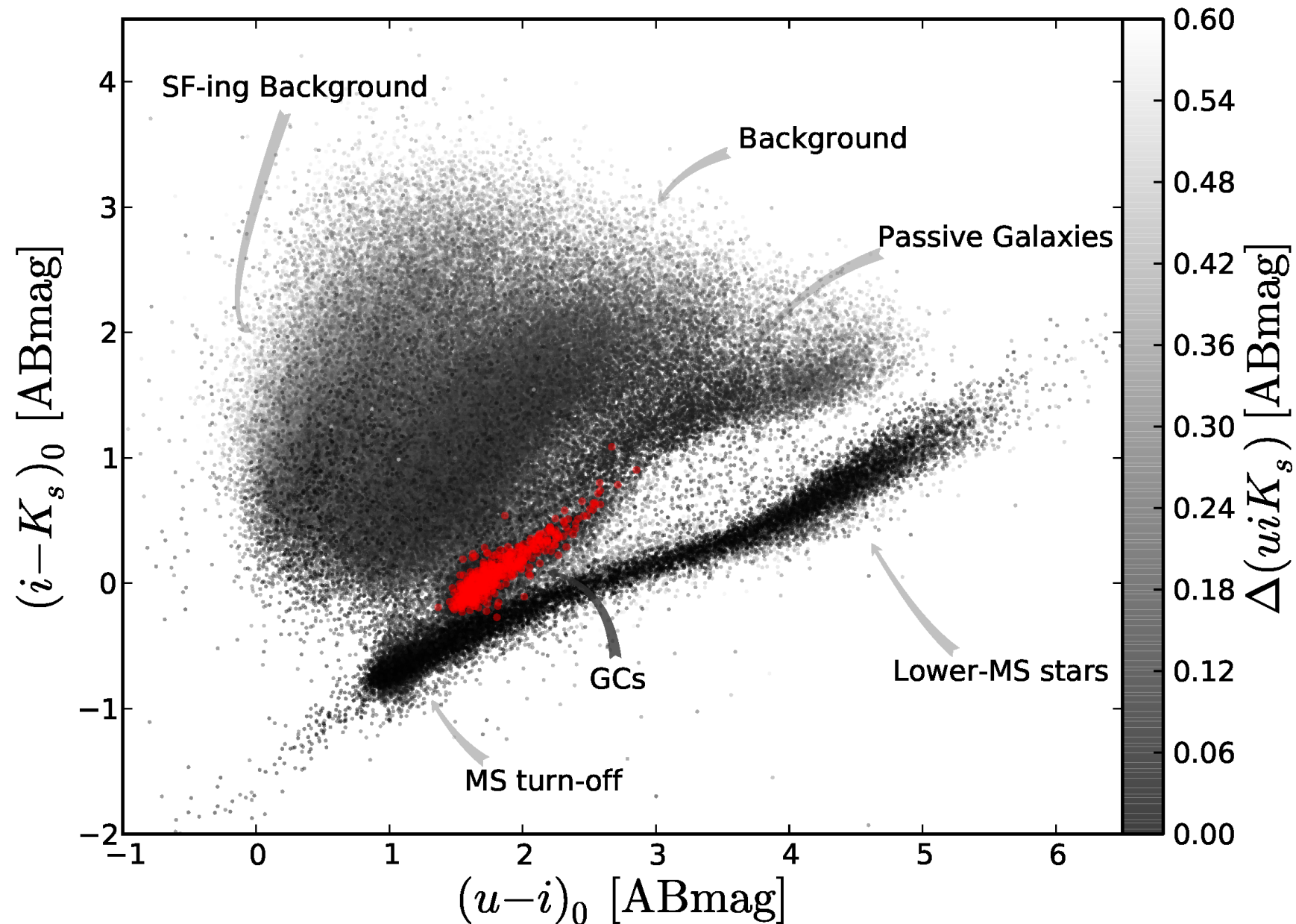




Roberto Muñoz
Tools for Astronomical Big Data

# NGVS: Virgo cluster

- The Next Generation Virgo Survey (NGVS; PI: L. Ferrarese) is designed to provide deep, high spatial resolution, and contiguous coverage of the Virgo cluster from its core to virial radius.



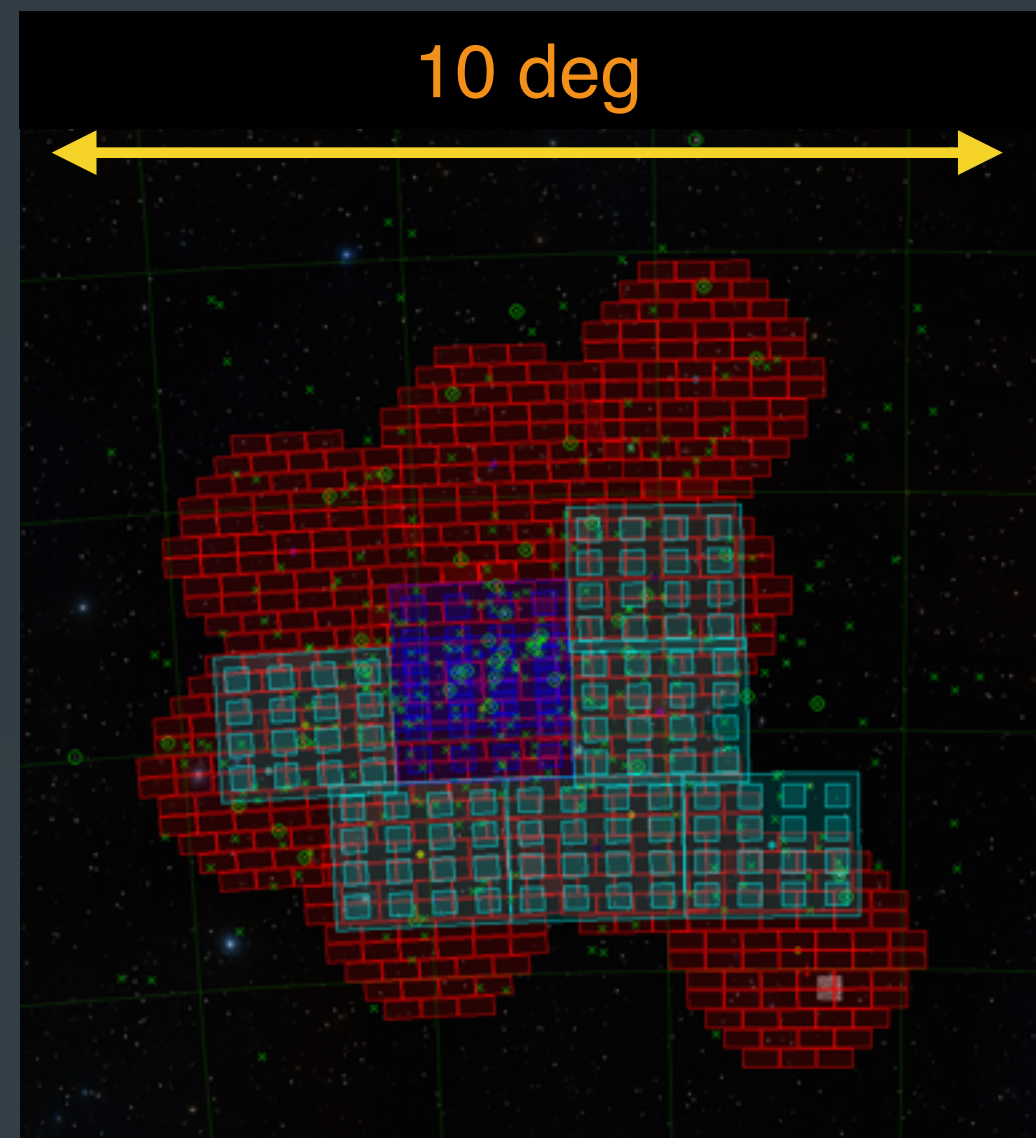Virgo cluster in X-ray from ROSAT (Böhringer et al. 1994)

# The uiK diagram



Muñoz et al. 2014

Roberto Muñoz
Tools for Astronomical Big Data

# NGFS: Fornax

- The Next Generation Fornax Survey (NGFS; PI: R. Muñoz) is an ongoing multipassband optical and NIR survey of the Fornax galaxy cluster. It will cover the central 30 deg$^2$ out to the virial radius and will allow to study the galaxy and GC populations.

1.8 deg



Astrophotography by Marzo Lorenzini

10 deg



Roberto Muñoz
Tools for Astronomical Big Data

# SCABS: Centaurus A

- The Survey of Centaurus A's Baryonic Structures (SCABS; PI: M Taylor) is mapping 72 deg$^2$ around Cen A.



7 arcmin



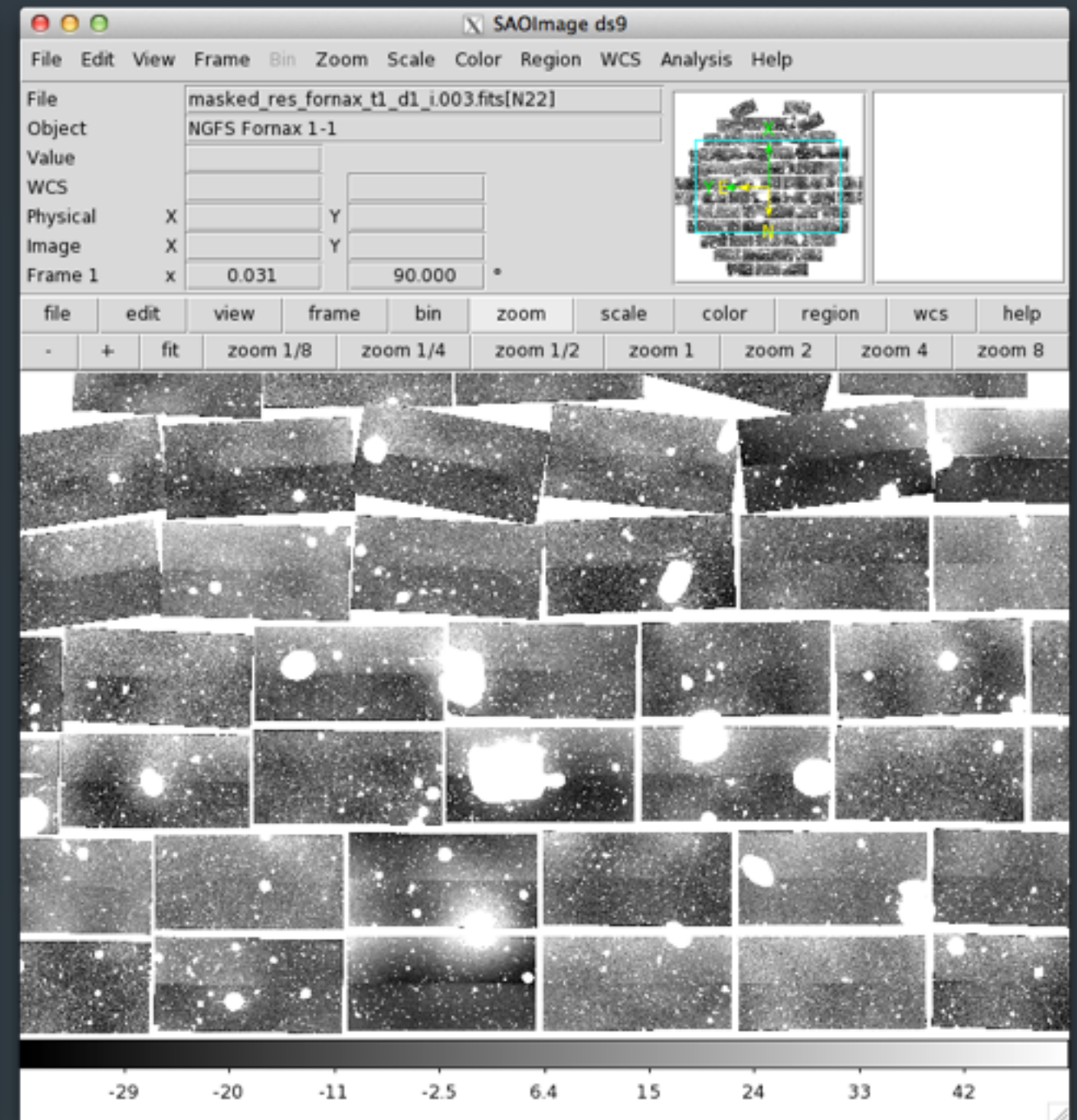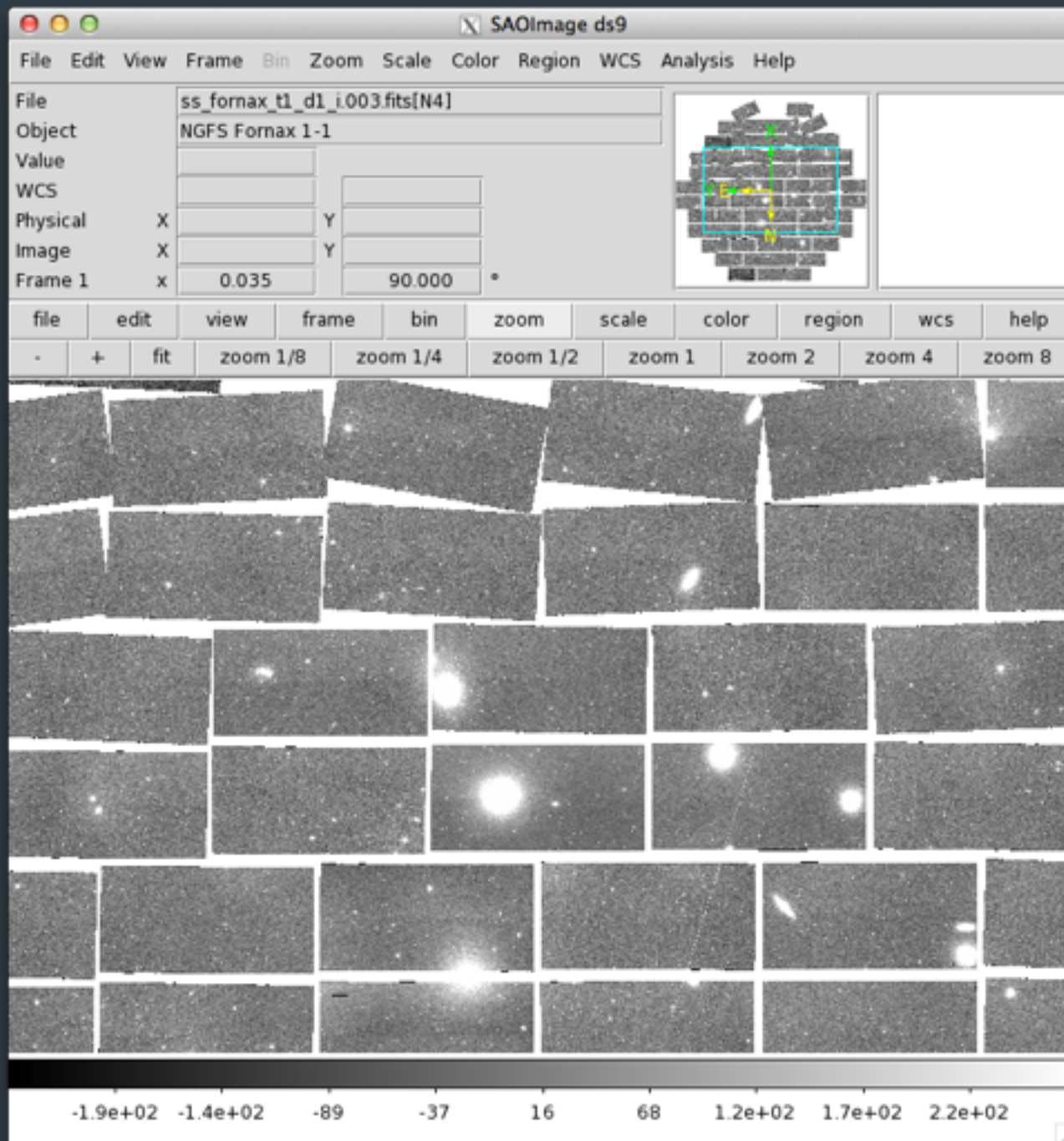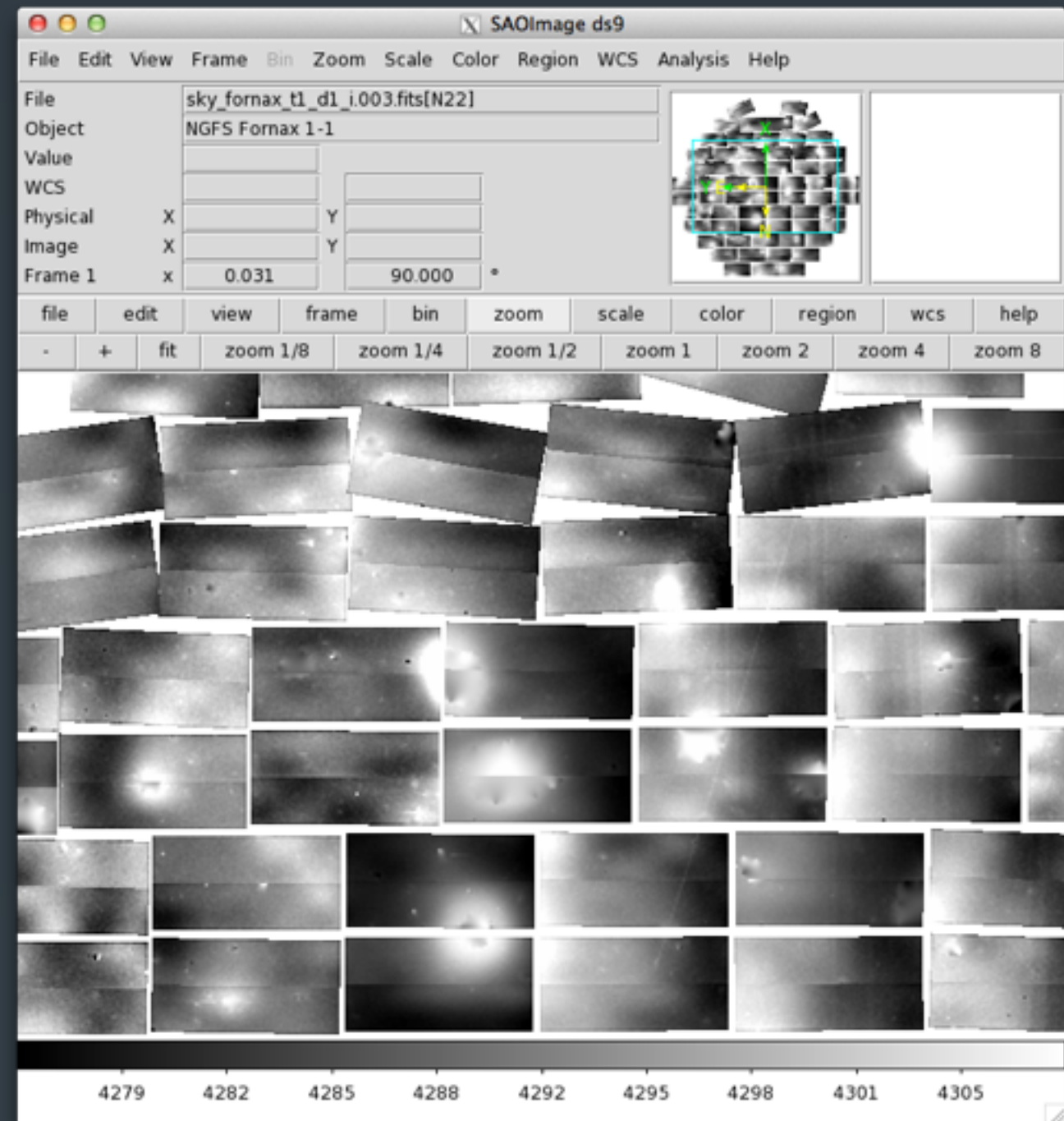20 deg

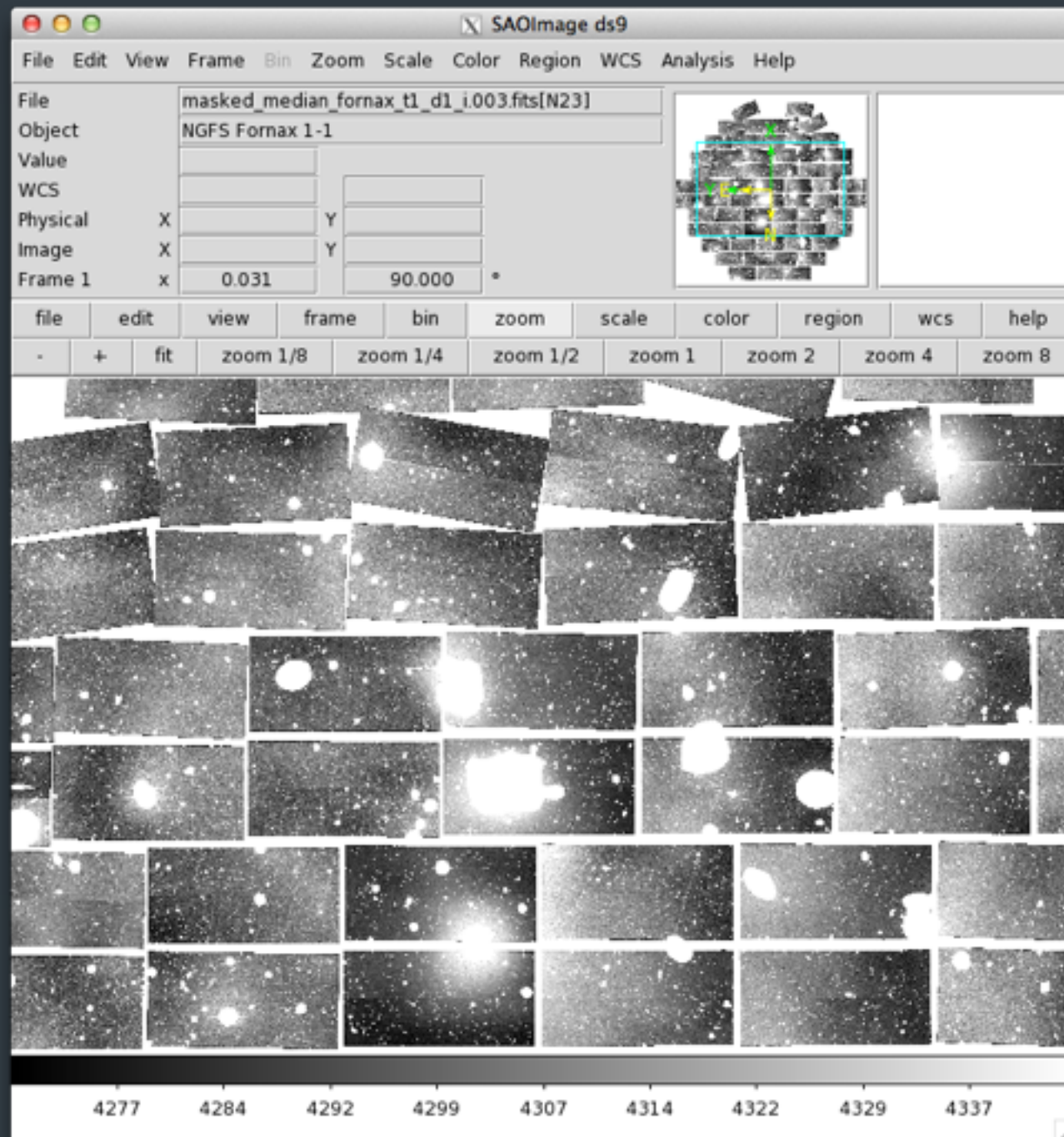Roberto Muñoz
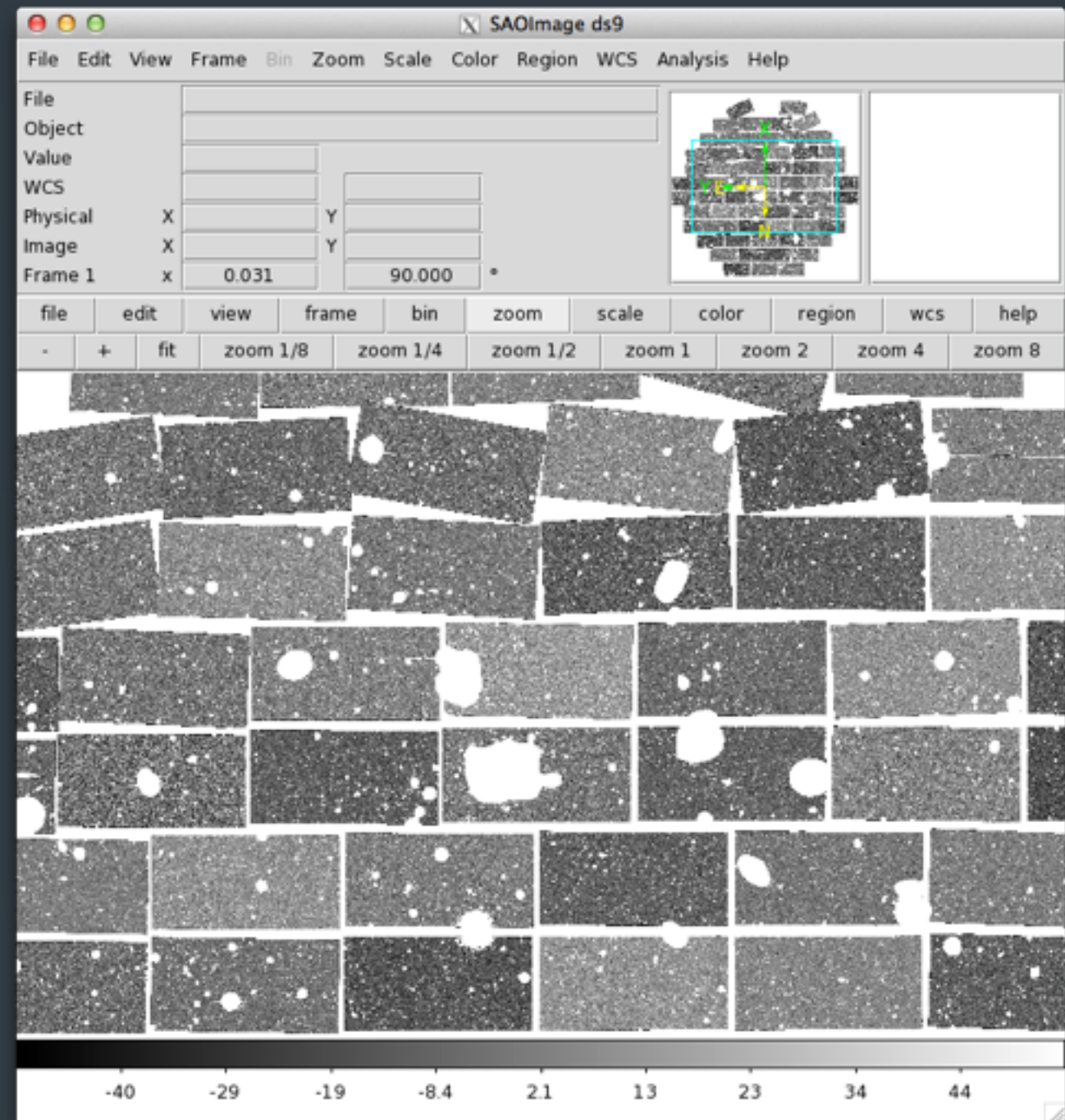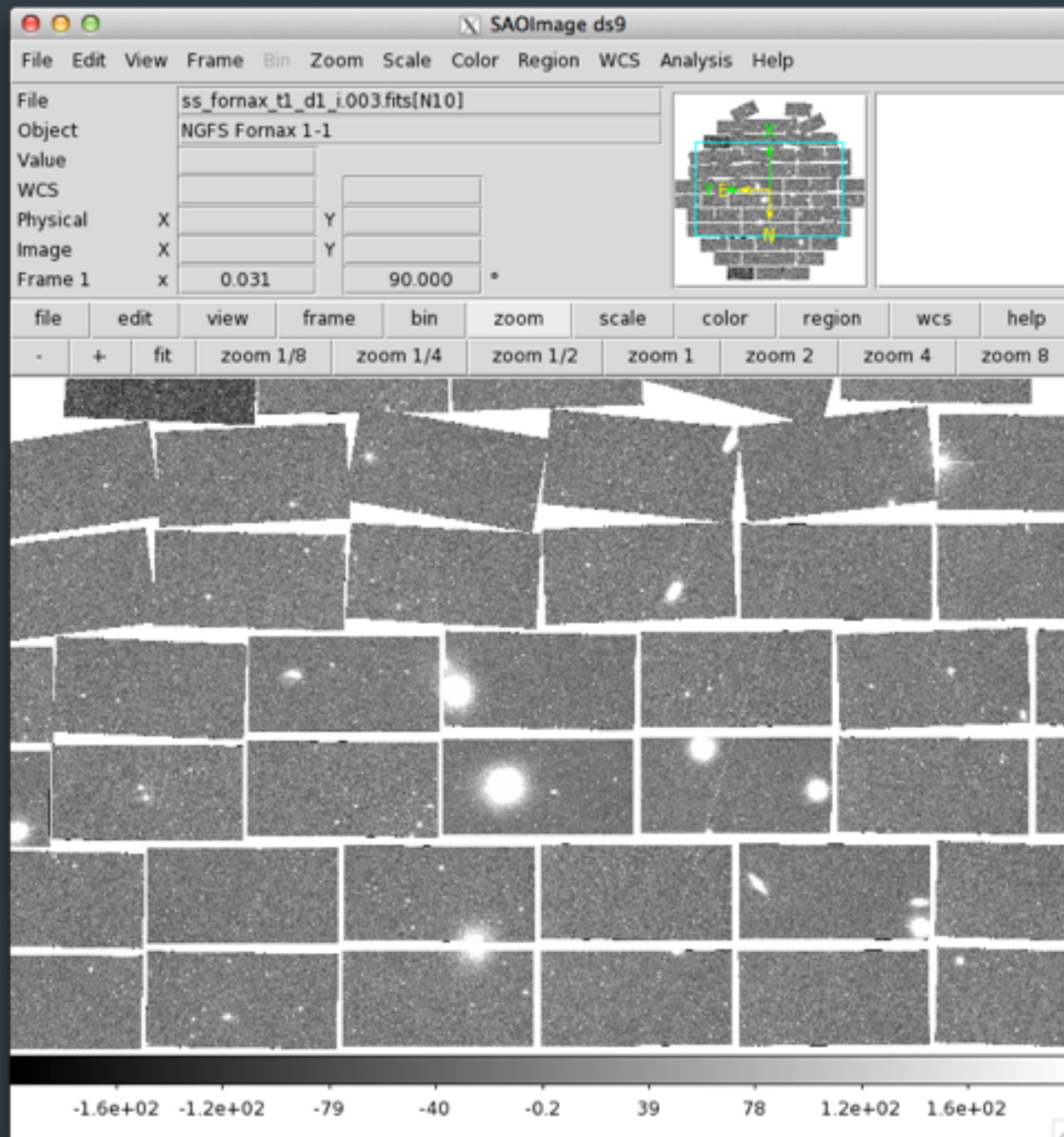Tools for Astronomical Big Data

# Sky subtraction: Polynomia

# Sky subtraction: Polynomia

# Sky subtraction: Splines
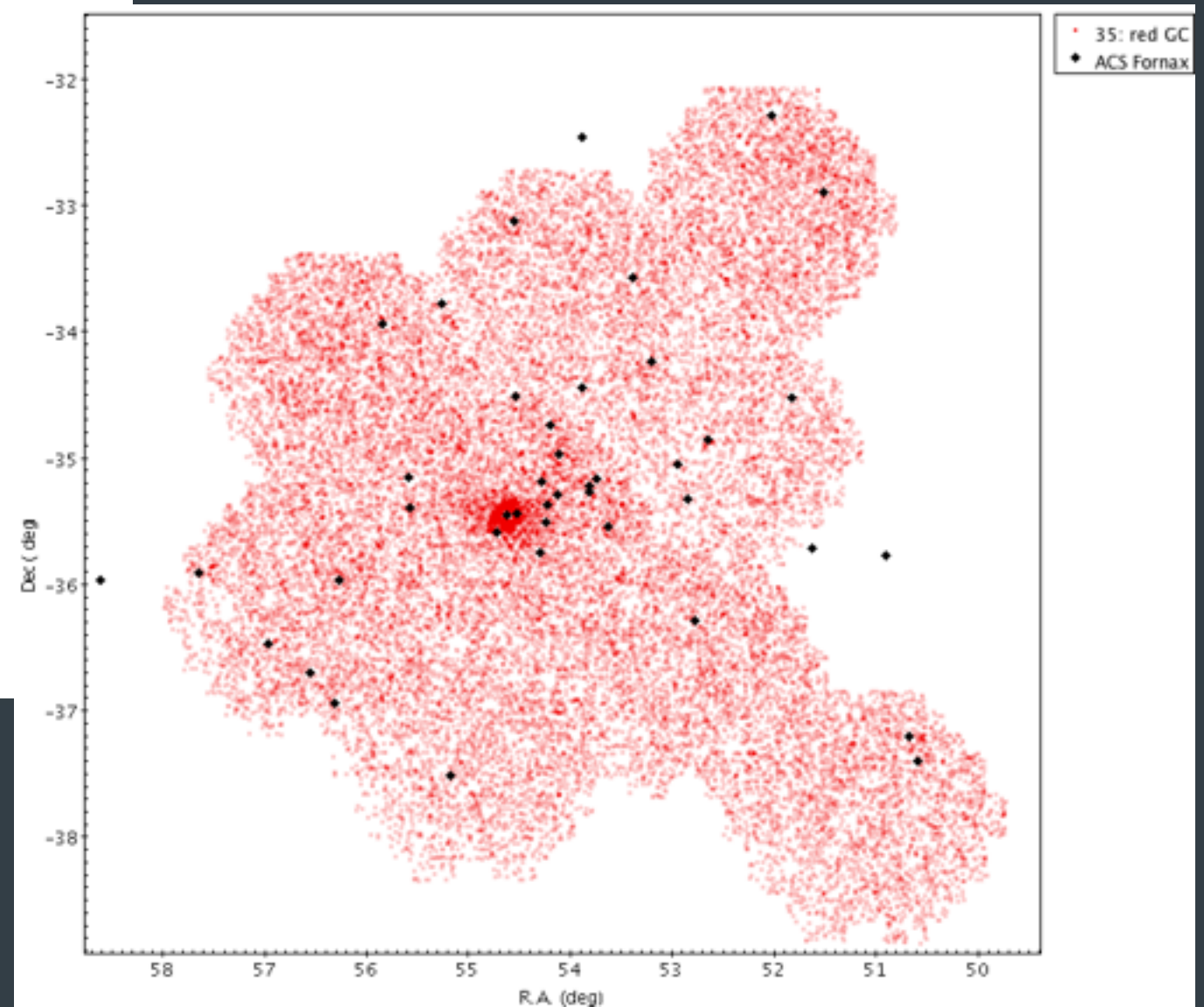


Roberto Muñoz
Tools for Astronomical Big Data

# Sky subtraction: Splines

# Stacks



Roberto Muñoz
Tools for Astronomical Big Data

# Fornax GCs



Roberto Muñoz
Tools for Astronomical Big Data

# Dwarf galaxies



Roberto Muñoz
Tools for Astronomical Big Data
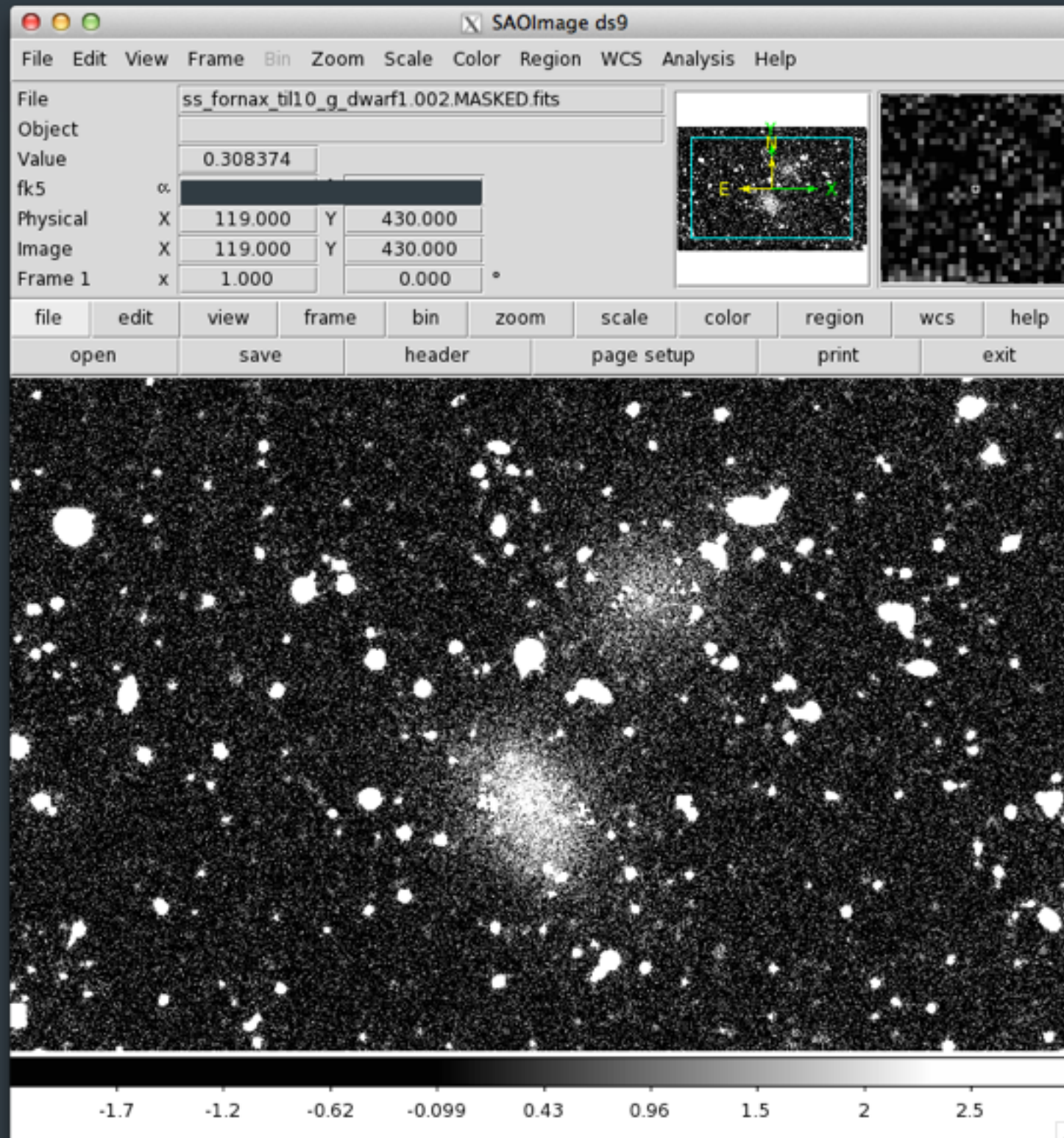
# Summary

- Setting up a large storage system in your office or institute to serve all your team 24 hours a day / 7days a week is possible.

- Plenty of disk space, high I/O and fast cpu processing for testing new algorithms and designing new methods. Different type of users in your team, some of them have lot of experience and others ones are still learning.

- Doing image processing and developing new methods take time, but the payback is good because you learn about your dataset and push it to the limits.

- Different projects and users require different solutions.
  We should support different approaches and look for the best solution depending on the needs and resources.