

# ALMA science data management

Felix Stoehr  
ALMA Science Archive

credit goes to the entire ALMA team

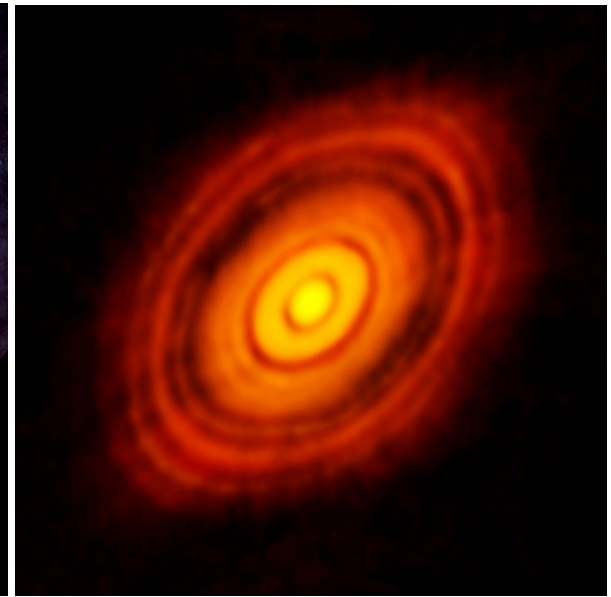
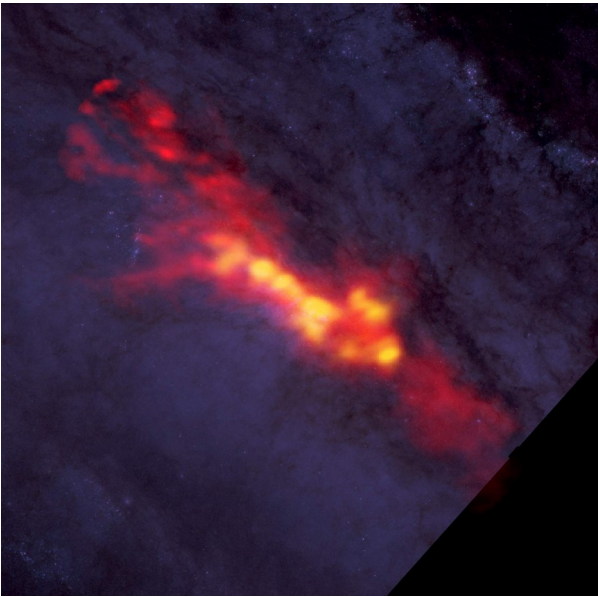
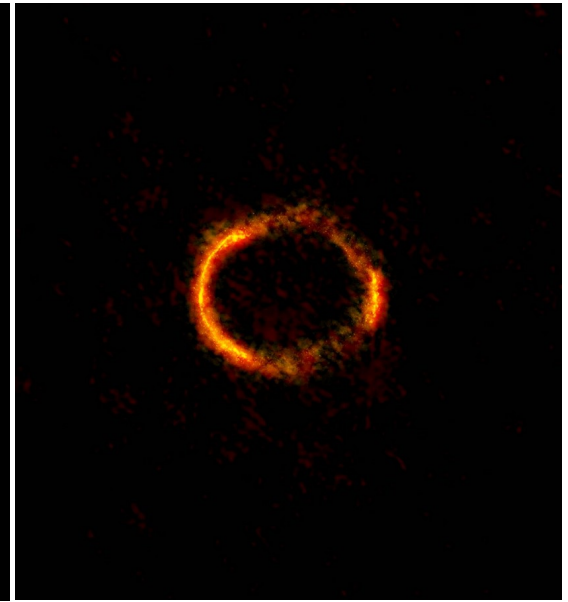
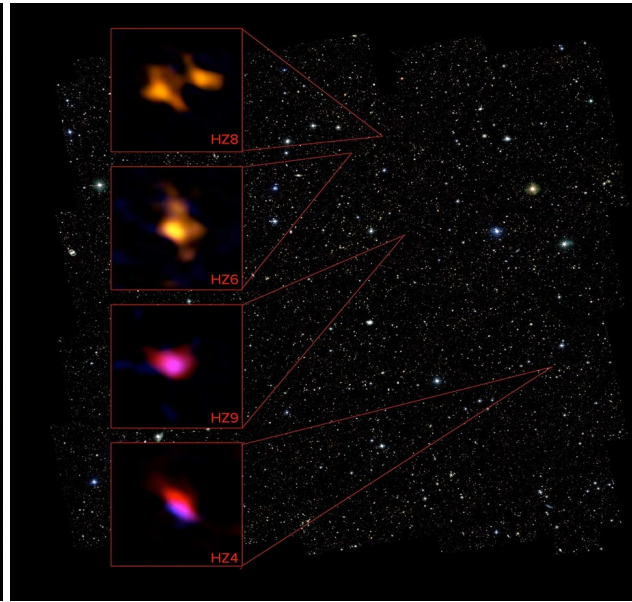
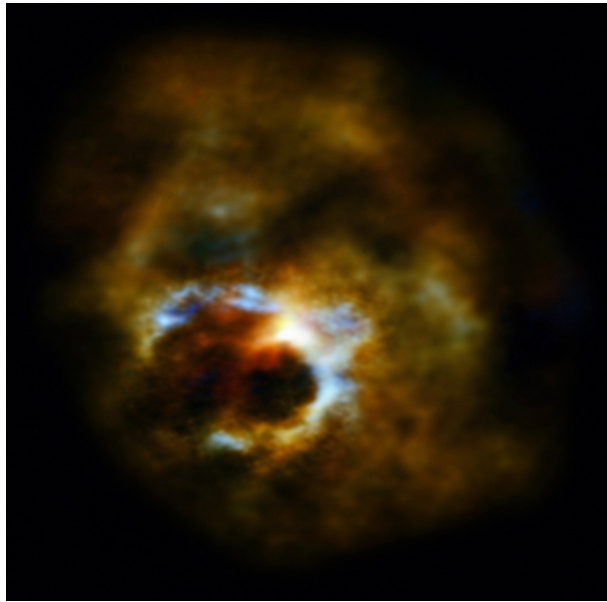




C. Malin

- built by ESO, NRAO, NAOJ in cooperation with Chile
- 66 antennas at 5000m elevation in the Atacama Desert
- interferometry at 84-702 GHz with 16km baselines
- full operations: 200TB/yr=6.6Mbytes/s
- Cycle 3: observing, call for Cycle 4: spring 2016

# ALMA



- demands of astronomers **scale faster than funding**
- collaboration over **continents** is required for the largest projects
- **challenges**
  - John Hibbard: Who does train the next generation of experts who have hands-on end-to-end expertise? Maybe there is not only room for smaller telescopes but astronomy needs those to train the next generation?
  - communication/organization

# workflow



# workflow



# workflow

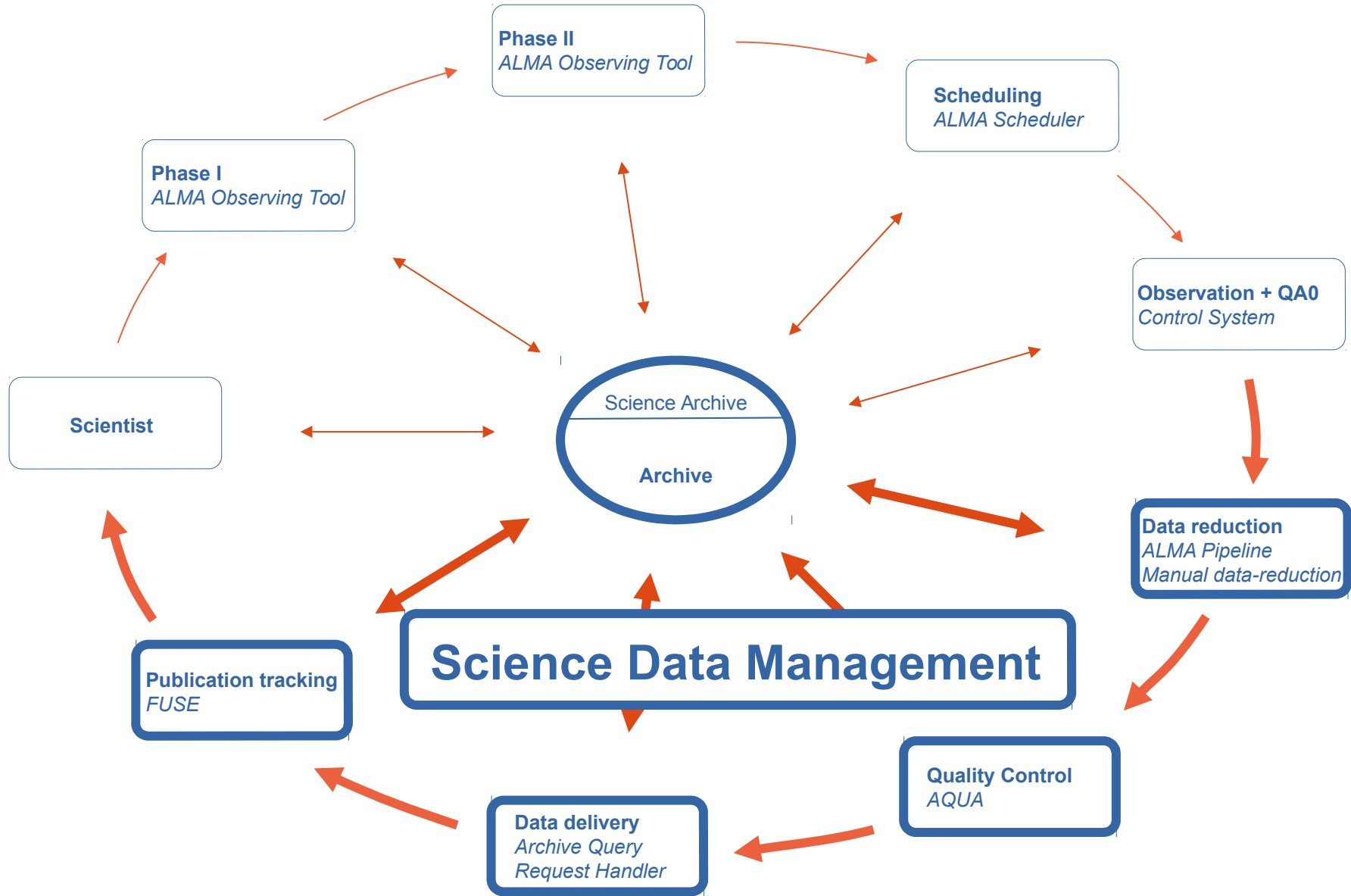


# workflow





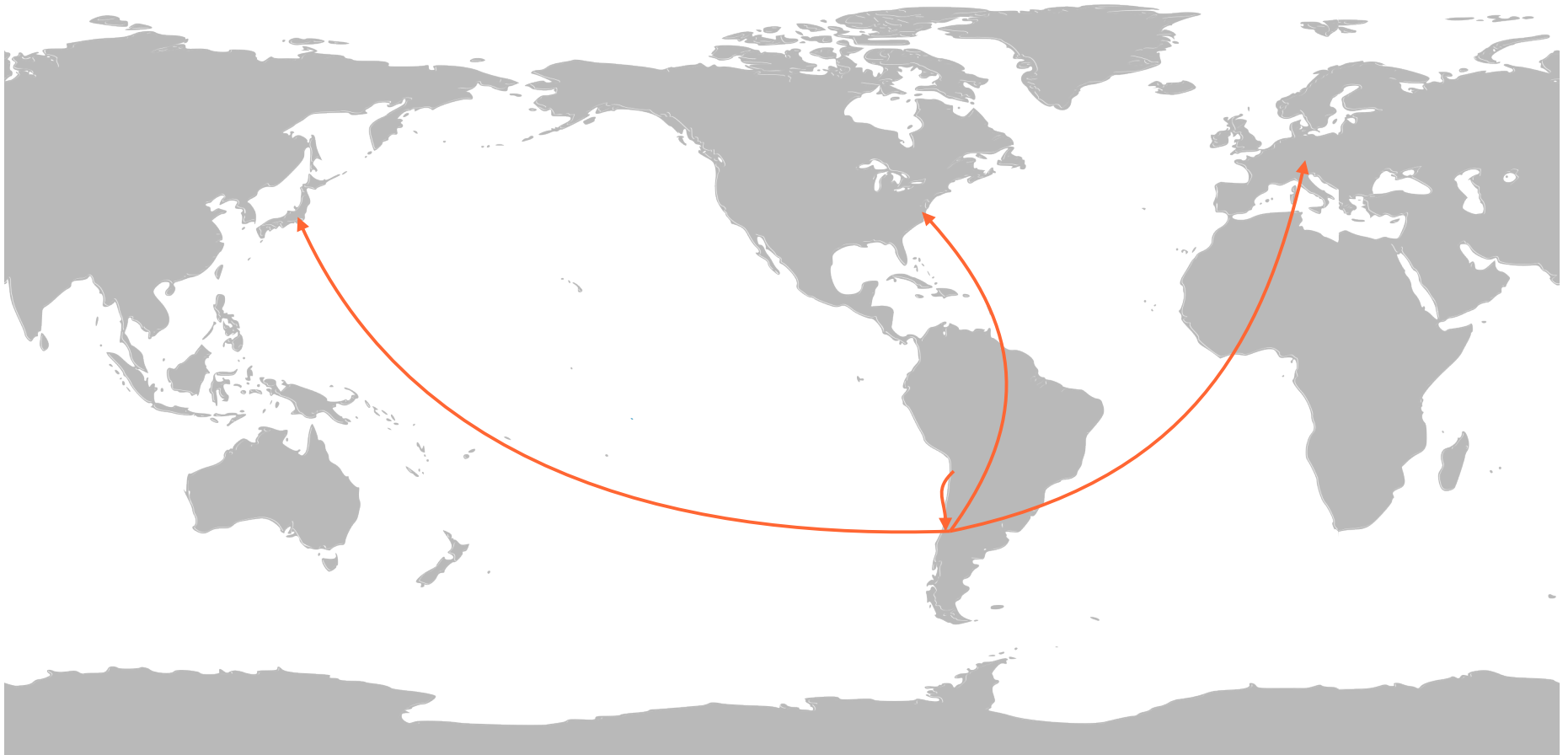
# science data management



- **success = science output of the community**
- **someone else is responsible** for our success
- we have to create the best **end-to-end user-experience possible**

# dataflow





- ALMA will produce about the same amount of data in one year as ESO has produced in its first 50 years. (And ESO will, too!)
- LOFAR, MWA, Gaia, PanStars, LSST, SKA, Euclid, ELTs
- T. Tyson: “Astronomy is **transformed** from being a data-starved science to one where data is overabundant”
- multi-wavelength science: **less time** per wavelength regime
- **astronomers do not scale**: bytes/astronomer grow exponentially
- my prediction:
  - now: astronomers **compete for observing time**
  - future: observatories will **compete for astronomers** to work with their data

- technical evolution will be ahead of astronomy
  - networks, storage, CPU
- however:
  - transfers, storage, processing will need to be **parallel**
- **code-to-data** will become more and more prominent

# data reduction



# data reduction

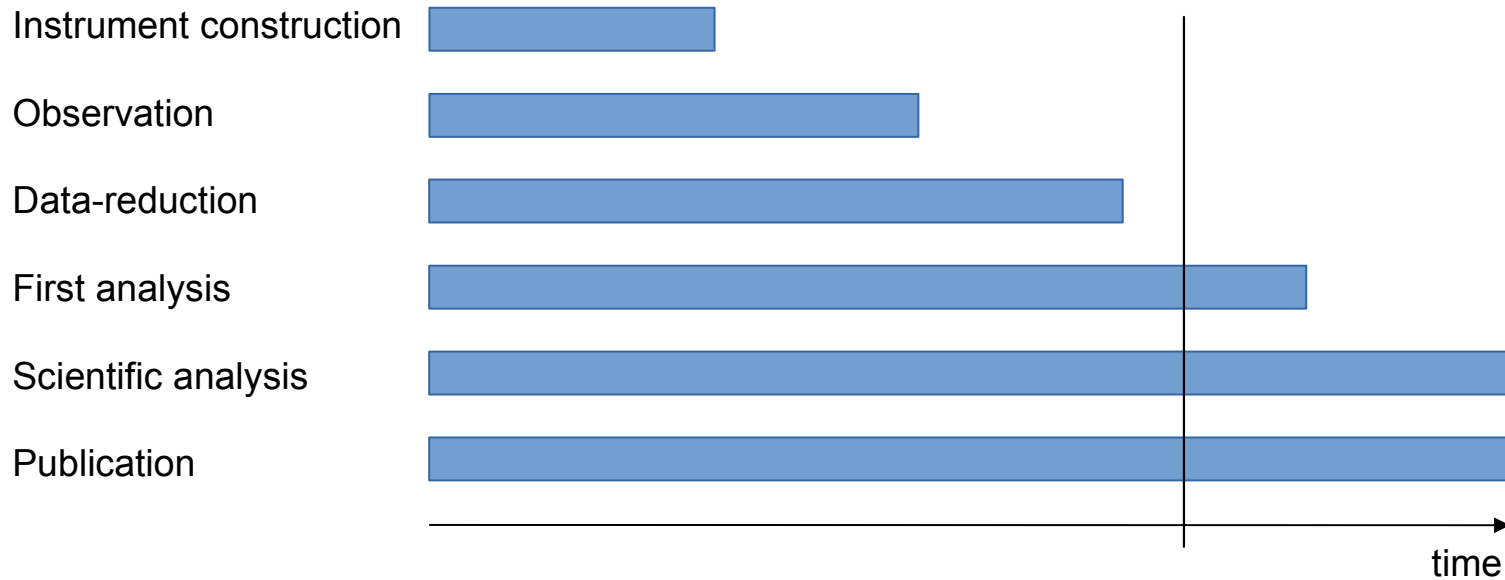
- science-grade data
- pipeline
- manual data-reduction
- first science analysis



# science-grade data

- ALMA policy: will provide **science grade products** for all raw data
- we develop a full **data-reduction package** CASA (python and C++)
- **extreme usersupporting**: face-to-face support from proposal preparation to science analysis in three continents

- what do PIs do?



- workload is **shifting from PIs to observatories**
- more complex telescopes and more data: **pipelines and science-grade data-products** will become an integral part of the telescope/instrument design

- deep impact on how science is done
- astronomers will become consumers of data instead of being co-producers
  - good: astronomy makes best use of astronomers, its future rare resource. → The science output will increase.
  - risk: astronomers may not understand data limitations
- observatories will have a much **larger responsibility**
- in some areas: **data-mining** will gain importance
- it will be (even) easier to have an **archival career**

- ALMA pipeline is **deployed with CASA**
- is **data-driven** and organized into **tasks**
- status: **calibration has been accepted**, imaging is in commissioning
  
- **reprocessing?** under discussion
- challenge: **backwards compatibility** of python code

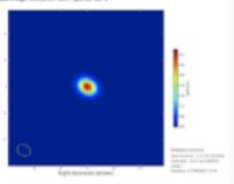
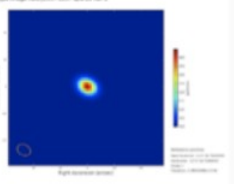
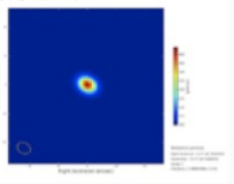
## WebLog: Calibrator Images

slides courtesy of Liz Humphreys

Home By Topic By Task

TASKS IN EXECUTION ORDER

1. hifa\_importdata
2. hifa\_flagdata
3. hifa\_fluxcallflag
4. hif\_rawflagchans
5. hif\_refant
6. hifa\_tsyscal
7. hifa\_tsysflag
8. hifa\_wvgcallflag
9. hif\_lowgainflag
10. hif\_setjy
11. hifa\_bandpass
12. hif\_bpflagchans
13. hifa\_spwphaseup
14. hifa\_gfluxscale
15. hifa\_timegaincal
16. hif\_applycal
17. hif\_makecleanlist
18. hif\_cleanlist
19. hif\_exportdata

field	spw	pol	image details	image result												
J1337-1257 (BANDPASS)	16	I	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>frequency</td><td>137.9893GHz</td></tr> <tr><td>beam</td><td>10.70x8.26arcsec</td></tr> <tr><td>beam p.a.</td><td>60.4deg</td></tr> <tr><td>image maximum</td><td>4.50e-00 Jy/beam</td></tr> <tr><td>residual rms</td><td>5.26e-04 Jy/beam</td></tr> <tr><td>channels</td><td>1 x 1937.44MHz</td></tr> </table>	frequency	137.9893GHz	beam	10.70x8.26arcsec	beam p.a.	60.4deg	image maximum	4.50e-00 Jy/beam	residual rms	5.26e-04 Jy/beam	channels	1 x 1937.44MHz	
frequency	137.9893GHz															
beam	10.70x8.26arcsec															
beam p.a.	60.4deg															
image maximum	4.50e-00 Jy/beam															
residual rms	5.26e-04 Jy/beam															
channels	1 x 1937.44MHz															
J1337-1257 (BANDPASS)	18	I	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>frequency</td><td>139.9267GHz</td></tr> <tr><td>beam</td><td>10.62x8.00arcsec</td></tr> <tr><td>beam p.a.</td><td>61.2deg</td></tr> <tr><td>image maximum</td><td>4.49e-00 Jy/beam</td></tr> <tr><td>residual rms</td><td>5.43e-04 Jy/beam</td></tr> <tr><td>channels</td><td>1 x 1937.44MHz</td></tr> </table>	frequency	139.9267GHz	beam	10.62x8.00arcsec	beam p.a.	61.2deg	image maximum	4.49e-00 Jy/beam	residual rms	5.43e-04 Jy/beam	channels	1 x 1937.44MHz	
frequency	139.9267GHz															
beam	10.62x8.00arcsec															
beam p.a.	61.2deg															
image maximum	4.49e-00 Jy/beam															
residual rms	5.43e-04 Jy/beam															
channels	1 x 1937.44MHz															
J1337-1257 (BANDPASS)	20	I	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>frequency</td><td>149.9888GHz</td></tr> <tr><td>beam</td><td>9.87x7.48arcsec</td></tr> <tr><td>beam p.a.</td><td>61.3deg</td></tr> <tr><td>image maximum</td><td>4.41e-00 Jy/beam</td></tr> <tr><td>residual rms</td><td>6.18e-04 Jy/beam</td></tr> <tr><td>channels</td><td>1 x 1937.44MHz</td></tr> </table>	frequency	149.9888GHz	beam	9.87x7.48arcsec	beam p.a.	61.3deg	image maximum	4.41e-00 Jy/beam	residual rms	6.18e-04 Jy/beam	channels	1 x 1937.44MHz	
frequency	149.9888GHz															
beam	9.87x7.48arcsec															
beam p.a.	61.3deg															
image maximum	4.41e-00 Jy/beam															
residual rms	6.18e-04 Jy/beam															
channels	1 x 1937.44MHz															

- importance of pipelines for the success of a facility can not be overestimated
- pipeline as **integral part** of the design of an instrument
- challenges:
  - maintenance and **evolution** of the pipelines
  - reprocessing
  - Most of all:  
cost **very easily underestimated**

# manual data-reduction

- ALMA's **deliver science-grade data products** principle is kept at all cost
- manual data-reduction JAO and the three ARCs until pipeline is ready
- some modes might need manual reduction also in the future
- benefits
  - helps train expert ALMA staff
  - gained knowledge can be used to improve the pipeline

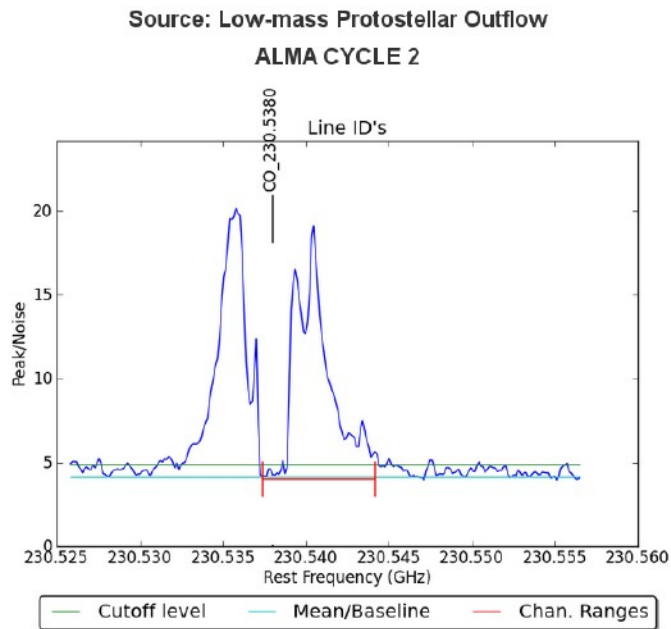
- very high cost (but also high quality)
- products can be tailored to the needs of the PI
- not complete/homogeneous products
- reprocessing is essentially infeasible



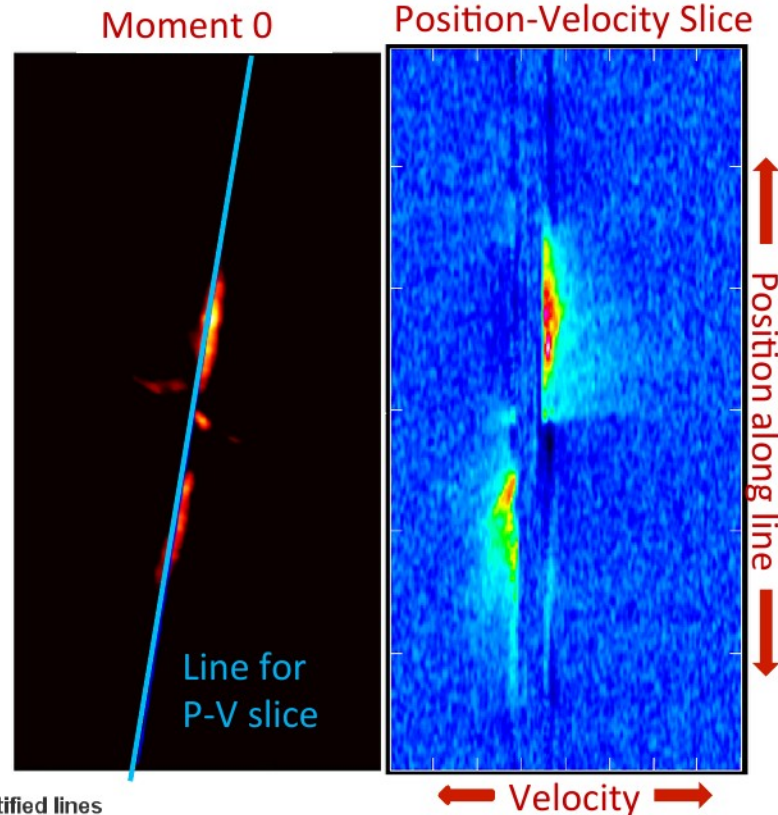
# first data analysis

- ADMIT (ALMA data mining toolkit) development program
- CASA add-on, will run at JAO on all products

plots courtesy of ADMIT team



Spectrum: peak emission over noise per channel



Identified lines

frequency	formula	transition	velocity	fwhm	startchan	endchan
230.53800	CO	2-1	1.200E+01	3.503E+00	100	155

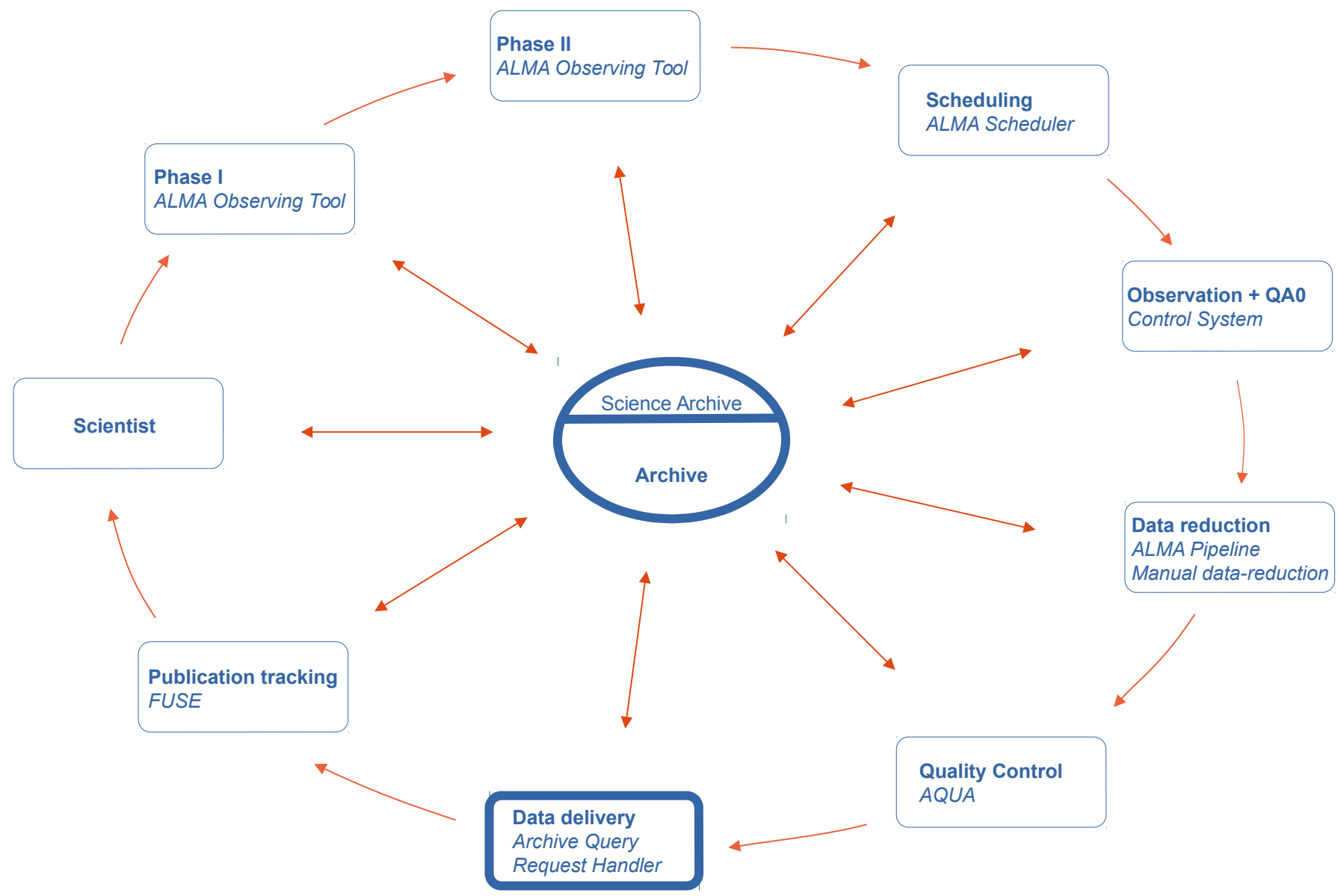
- so far most facilities describe their **observations**
- next frontier: first data-analysis allows to describe the **science content** of those observations
  
- PI ground-based facilities are “following” space facilities and surveys
- again: more workload and responsibility shift from scientists to observatories

# quality control



- ALMA quality assessment
  - QA0: raw data quality assessment after one observation
  - QA1: trend analysis of the facility
  - QA2: data product quality assessment for one dataset using pipeline weblog
  - QA3: analysis of problems reported by users after data delivery
  
- QA0 and QA1 are expected to be automatized
  
- QA2 and QA3 will remain tasks of a real person

# science archive



- complete, correct, consistent, (homogenized)
- speak the **language of the astronomers**. Query by **physical concepts**

**DANIEL DURAND TEST** of science archives:

You know have one, when you can find your own data  
by only giving physical parameter constraints

- reduce **interaction cost**
- unscope search: give access to the full **n-dimensional parameter space**

## ALMA Science Archive Query

Query Form

Results Table

Search

Reset

[Query Help](#)

### Position

Source name (Resolver)  
Source name (ALMA)  
RA Dec  
Spatial resolution  
Largest angular scale

### Energy

Frequency  
Bandwidth  
Spectral resolution  
Band

### Time

Observation date  
Integration time

### Polarisation

Polarisation type

### Observation

Water vapour

### Project

Project code  
Project title

- ... lensed z=2.3 sub-mm gala...
- ... , Luminous Infrared Galaxy
- ... medium in the inner galaxy
- ... ns of a rare triple galaxy ...
- A dusty dwarf galaxy at z=7.37
- ... f five SZE-selected galaxy...
- ... minous Star-Forming Gala...
- ... of the z=8.23 Host Galaxy...
- ... in the Luminous IR Galax...

**Project Title**  
ALMA Project Title

**Description**  
Case-insensitive search over the project title

**Example**  
\*GALAXY\*  
\*?ebula\*

### Options

View:

- raw data
- project
- publication
- public data only
- science observations only

- future: products, CARTA, PPI, HiPS+AladinLite, VO (ObsCore, SIAPv2, TAP)

## ALMA Request Handler

Anonymous User: Request #413184767 ✓  
 Request Title: [Click to edit](#)

Download Selected

Include Raw

Project / OUSet / Executionblock

- Request 413184767
  - Project 2012.1.00661.S
    - Science Goal OUS uid://A002/X6444ba/X1b2
      - Group OUS uid://A002/X6444ba/X1b3
        - Member OUS uid://A002/X6444ba/X1b4
          - product
          - raw
          - raw
    - Project 2012.1.00762.S
      - Science Goal OUS uid://A002/X5a9a13/X687
        - Group OUS uid://A002/X5a9a13/X688
          - Member OUS uid://A002/X5a9a13/X689
            - product
            - raw
            - raw
            - raw
            - raw
            - raw

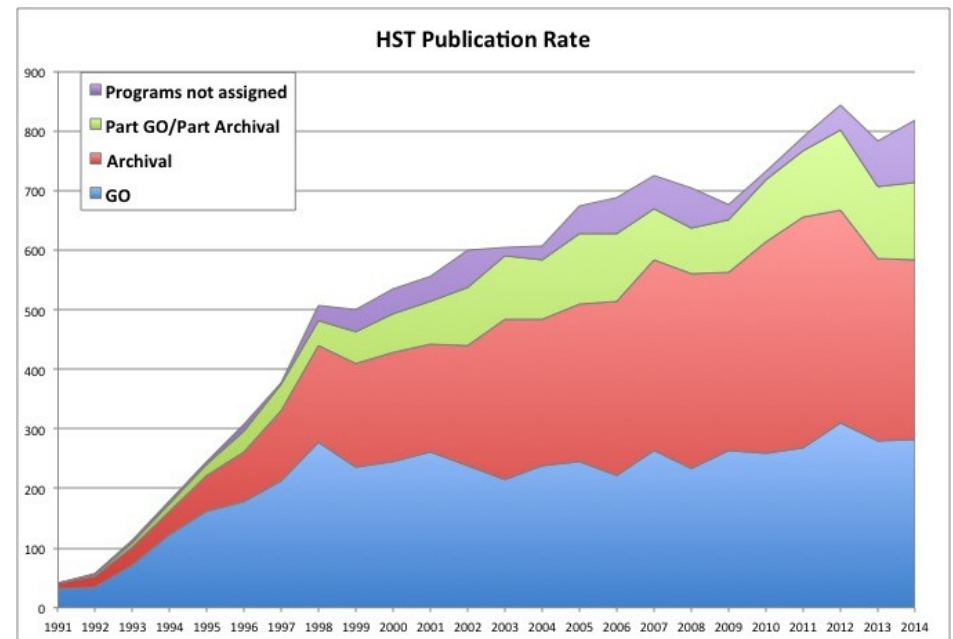
Choose one of the following download methods:

<b>Download Script</b>	The downloads are scripted for you. You just need to execute the script from the command line.
<b>Download Manager</b>	ALMA's download manager is launched as a browser applet. This is a simpler, more user-friendly way to download files in parallel, allowing you to pause and resume.
<b>Web Start Download Manager</b>	ALMA's download manager is launched as a desktop application via Java Web Start. It will not stop if you close your browser.
<b>File List</b>	View a text file containing a list of URLs. This is useful for using third-party download manager's such as <i>DownThemAll</i> .

[2012.1.00762.S\\_uid\\_A002\\_X74821d\\_X16bf.asdm.sdm.tar](#)



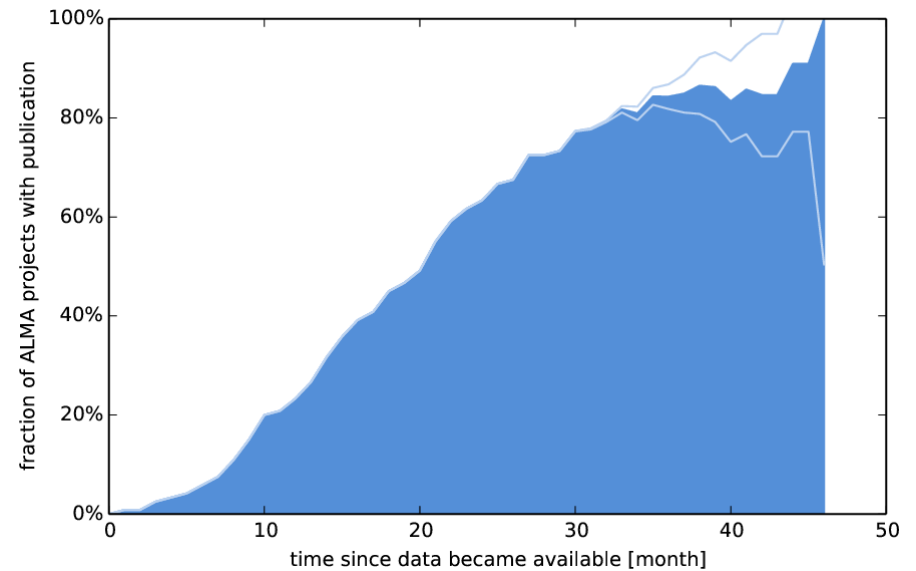
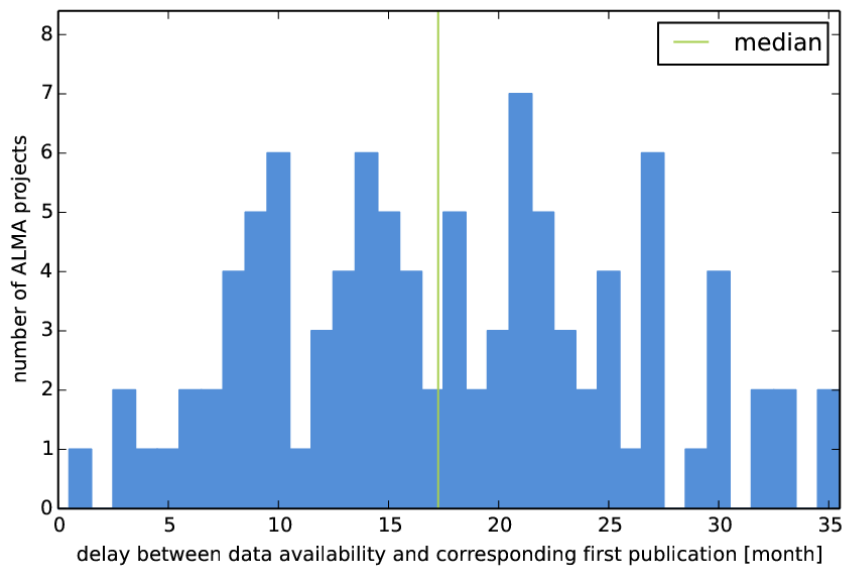
- success? → best possible **end-to-end user-experience**
- creating a good Science Archive **helps maximise the scientific return**
- great **return-for-investment** ratio
- interoperable: VO



# publications



- ALMA policy: users **must** cite the data they use
- ESO/NRAO/NAOJ libraries **track** publications
- we check arXiv and contact PIs if needed before articles go into press



# publications

## ALMA project status questionnaire

ALMA has great interest in providing data and support in a way that enables PIs and the community at large to advance science and to publish their findings in astronomical journals.

While we attempt to provide data that is useful for expert and non-expert PIs, we may not be reaching this goal effectively in some cases.

According to our records, it has been two years since the end of the proprietary period of your project, and we have been unable to identify a related publication.

We would like to use this occasion to ask you to help us provide better data and support for yourself and future users of ALMA, by giving us your feedback on this very short, anonymous questionnaire.

**\* Please describe the status of your project.  
Choose one of the following answers**

- There is a publication. (Please tell us below which journal we need to add monitoring)
- A publication is in press.
- A publication is in preparation.
- Only part of the requested project was actually observed. The data were not complete enough for the expected science.
- The project was complete, but the quality of the data was not good enough.
- The project was complete, and the quality of the data was good enough but we had problems producing the data products.
- Although the data were complete and good enough, the expected science was not contained in the data. For example, the experiment was a detection experiment.
- Although the data were complete, good enough and the expected science was contained in the data, the scientific field had moved on in the meantime. For example a competitor has already published similar results.
- Although the data were complete, good enough and still relevant, no effort was available (any more) to analyze the data.
- Although the data were complete and good enough, still relevant and effort was available we have to wait for data from other facilities. (Please tell us which ones)
- Personal reasons.
- Other:

**We would very highly appreciate any other comments or suggestions you might have for us with respect to the survey above but also with the general view to improving the publishable quality of ALMA data as well as our support to PIs and the community at large.**

Submit

Exit and clear survey

- publication tracking is
  - essential to **measure the success** and
  - a crucial tool to **improve operations**

to make sure that the investment really gets converted into science

- ALMA's has a modern science data management concept
- goal: **science-grade products** for all raw data
- extreme **user-support**
- so far ALMA is **very successful**
- “**transformational facility**”
- still **a lot of work** is required to develop ALMA to its full potential

- **success** = maximizing the **end-to-end user experience**
- **astronomers** not data will be the rare resource
- more **work and responsibility** will shift **to observatories**
- the way we do science is changing
  
- This will be **golden times for astronomers**