Image credit: Alex Cherney

# Ingest pipeline for ASKAP

**Max Voronkov**

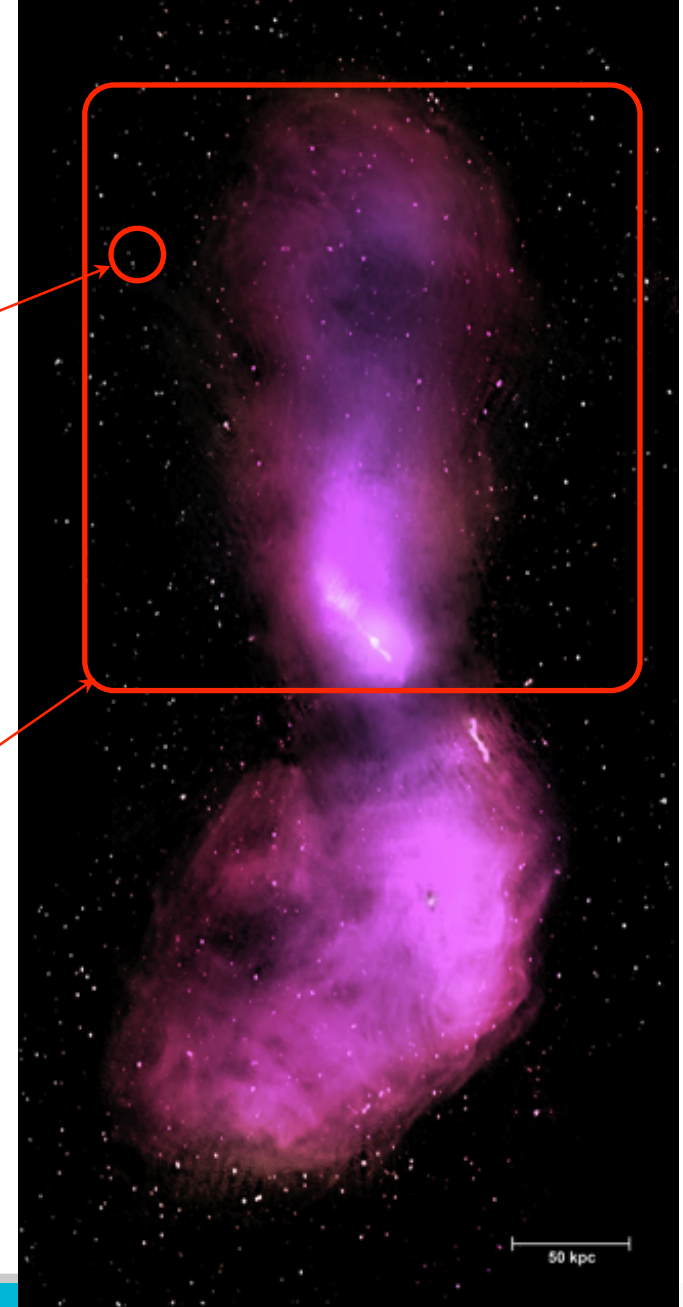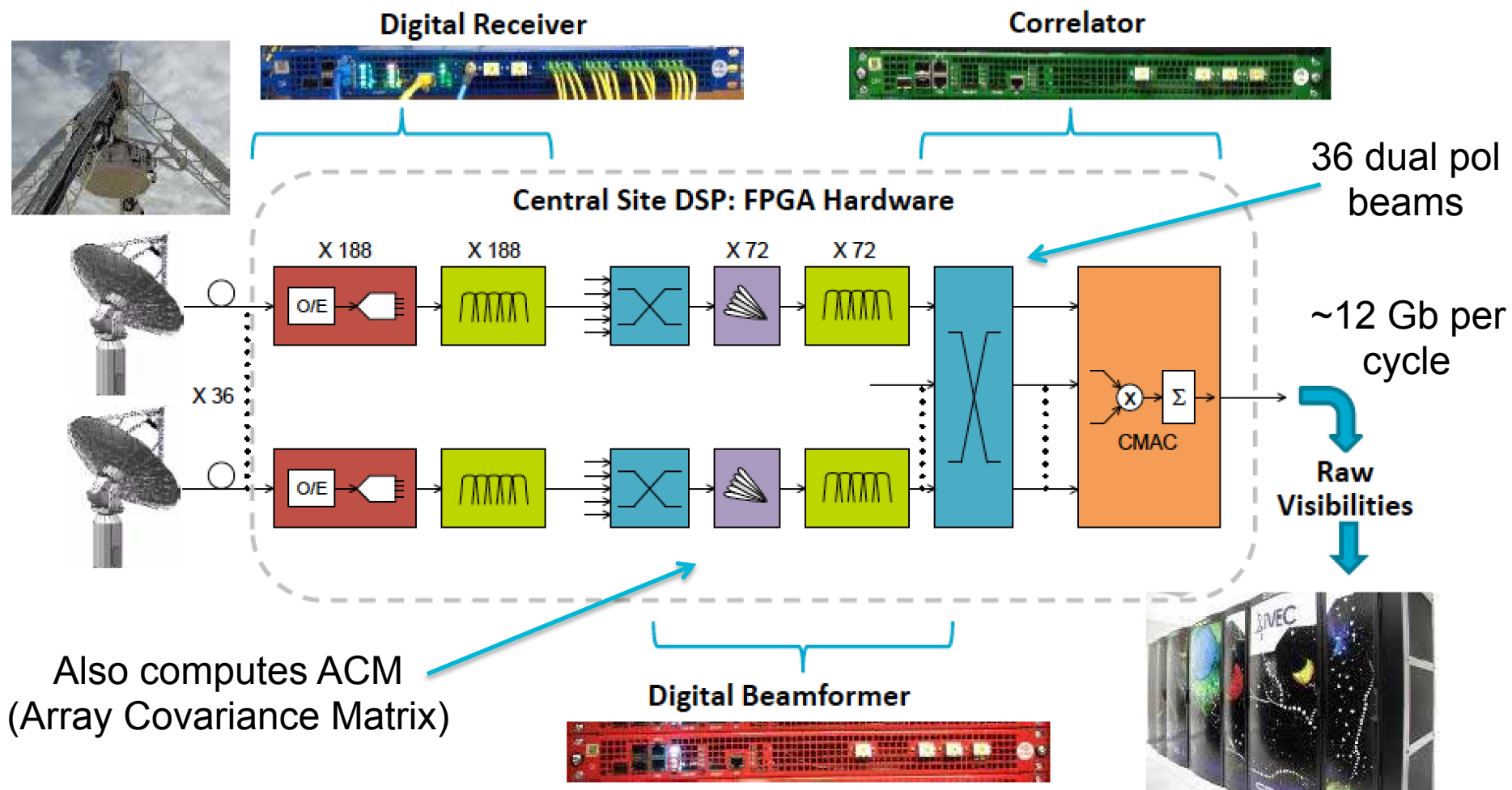Computing for SKA, Auckland, NZ – 14 Feb 2019

CSIRO

# ASKAP: Wide Field of View

- Required 1200 hours observing with the Australia Telescope Compact Array



- ASKAP will take about 10 minutes





50 kpc

Cen A image credit: Tim Cornwell / Ilana Feain

# ASKAP – system architecture (I)



Digital Receiver

Correlator

36 dual pol beams

Central Site DSP: FPGA Hardware

X 188    X 188    X 72    X 72

O/E

X 36

O/E

CMAC

~12 Gb per cycle

Raw Visibilities

Also computes ACM
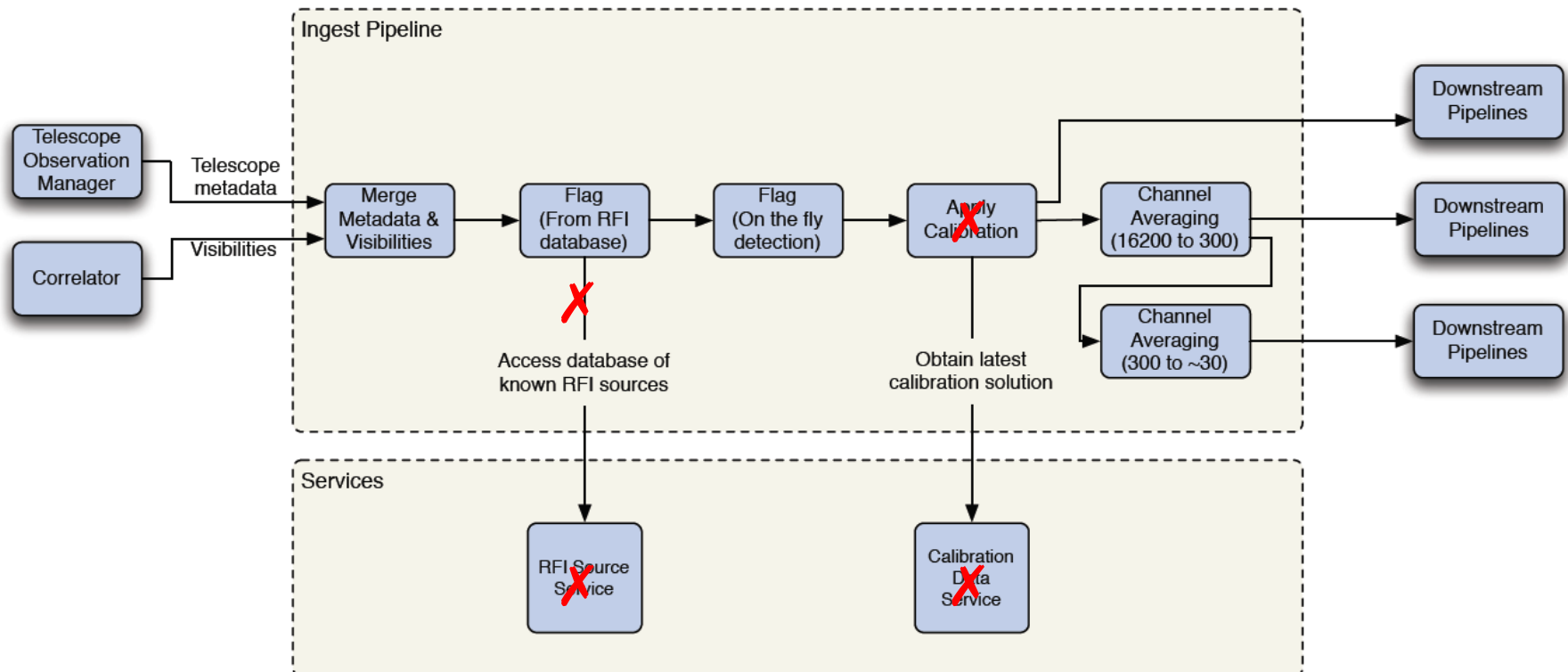(Array Covariance Matrix)

Digital Beamformer

CSIRO

# Plans vs. reality

- Initial design aimed at automatic processing of data
  - Largely due to inability to store data (run out of current storage in 1 week)
  - One size fits all approach with on-the-fly calibration – just 1 read of dataset
- Traditional processing model for now
  - User preference, more storage/better algorithms, lesser push for commensality
- Commissioning, staged deployment and support of BETA
  - Additional requirements never envisaged at the design stage
  - Lots of data inconsistencies in various ways, sometimes transient
  - Intermediate s/w solutions allowed to get science faster

# The role of ingest

- Prepare/re-format data for calibration & imaging pipelines
  - Synchronisation is the main part (especially if something doesn't work right)
  - Standard formats allow us to have more generic imaging & calibration tools
- Some processing which has to be done on-the-fly to reduce I/O
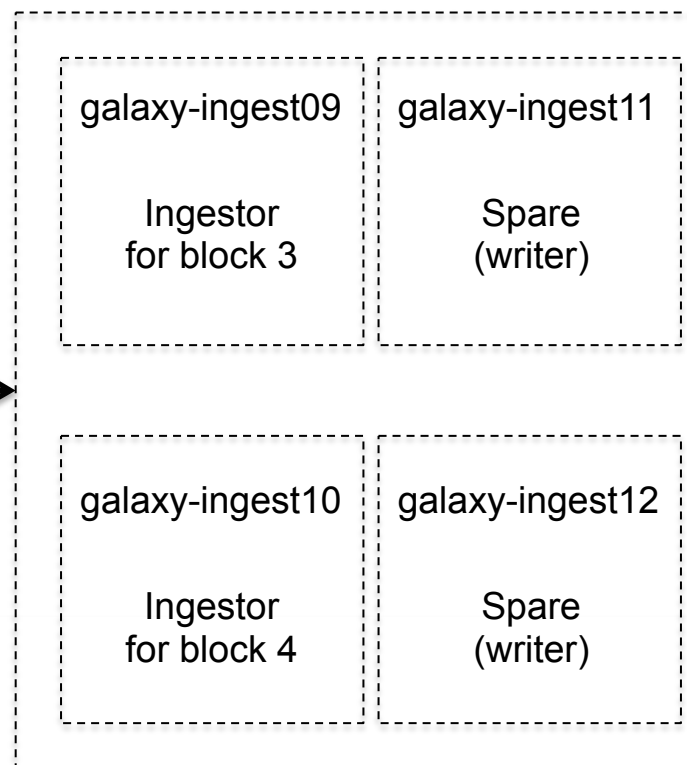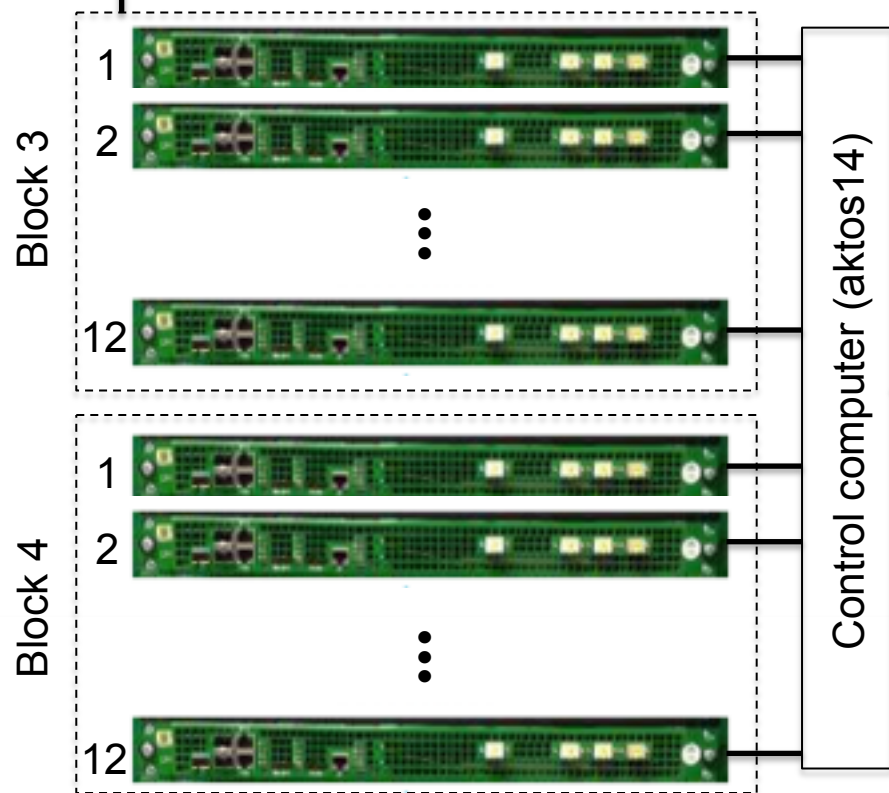- Crucial part of the online calibration loop (when/if we do it)
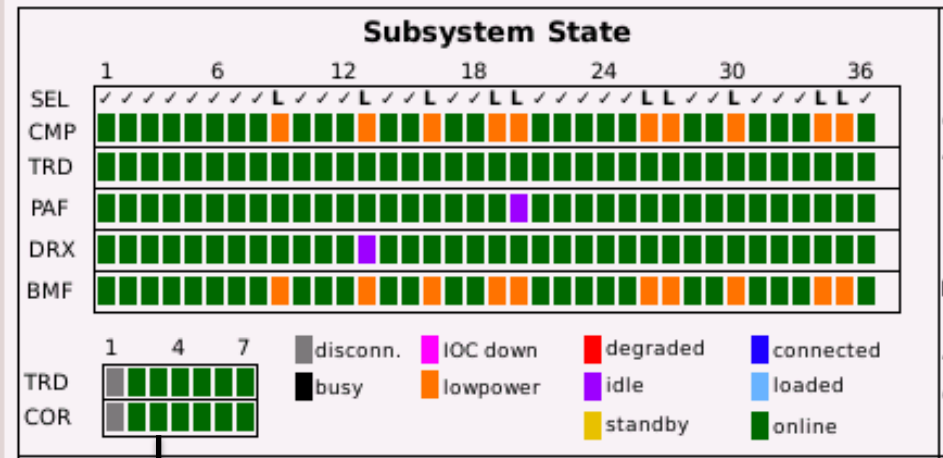


Diagram from the SW-0020 document (http://www.atnf.csiro.au/projects/askap/ASKAP-SW-0020.pdf)

# ASKAP - system architecture (II)



**Ingest cluster @ Pawsey**

- 16 nodes, 2 sockets per node
- 8 cores CPUs, 64 Gb of RAM per node
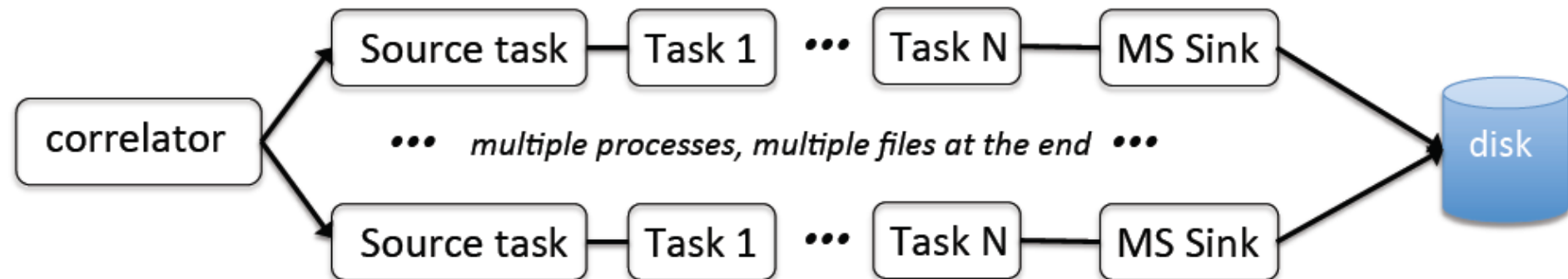- Shared Lustre storage with 30 Gb/s peak I/O performance

VLAN covering this fibre & 4 nodes

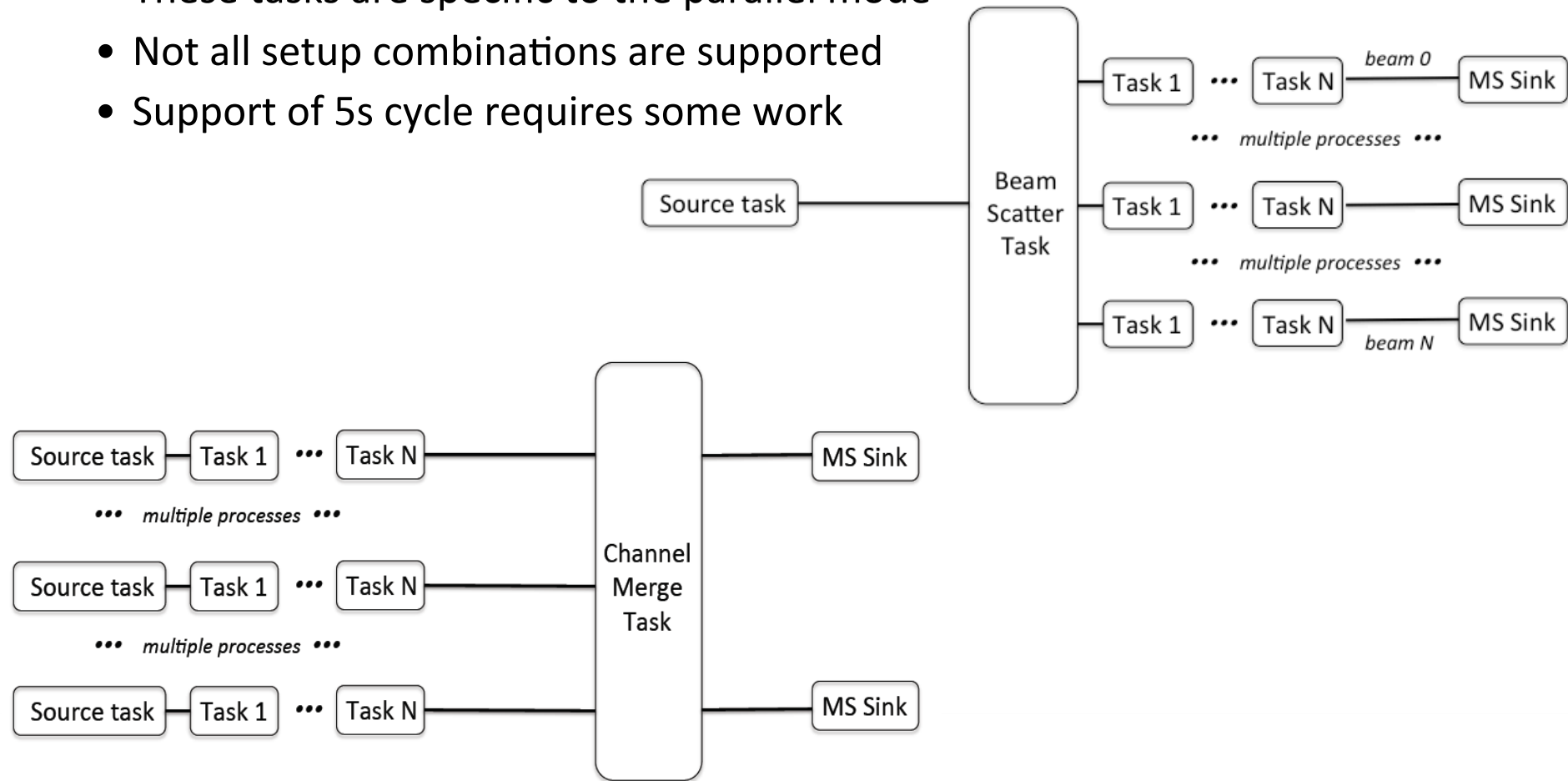| galaxy-ingest09 | galaxy-ingest11 |
| --- | --- |
| Ingestor for block 3 | Spare (writer) |

| galaxy-ingest10 | galaxy-ingest12 |
| --- | --- |
| Ingestor for block 4 | Spare (writer) |

x4

Block 3

Block 4

Control computer (aktos14)

# Tasks and data streams

- Ingest can be viewed as a chain of "processing" tasks
  - Configurable (via **tasklist** Facility Control Manager parameter)
  - Processing chain should always start with a source task (two options available)
  - Any task except source can occur in the chain any number of times (with the same or different parameters)
  - Sink task doesn't have to be the last (or even doesn't have to be present)
- Both parallel and serial modes are supported
  - Source tasks are rank-aware and would listen different UDP ports
  - Some tasks (merging, splitting) are specific to the distributed mode
  - Service ranks (i.e. those which do not run source task) are now also supported
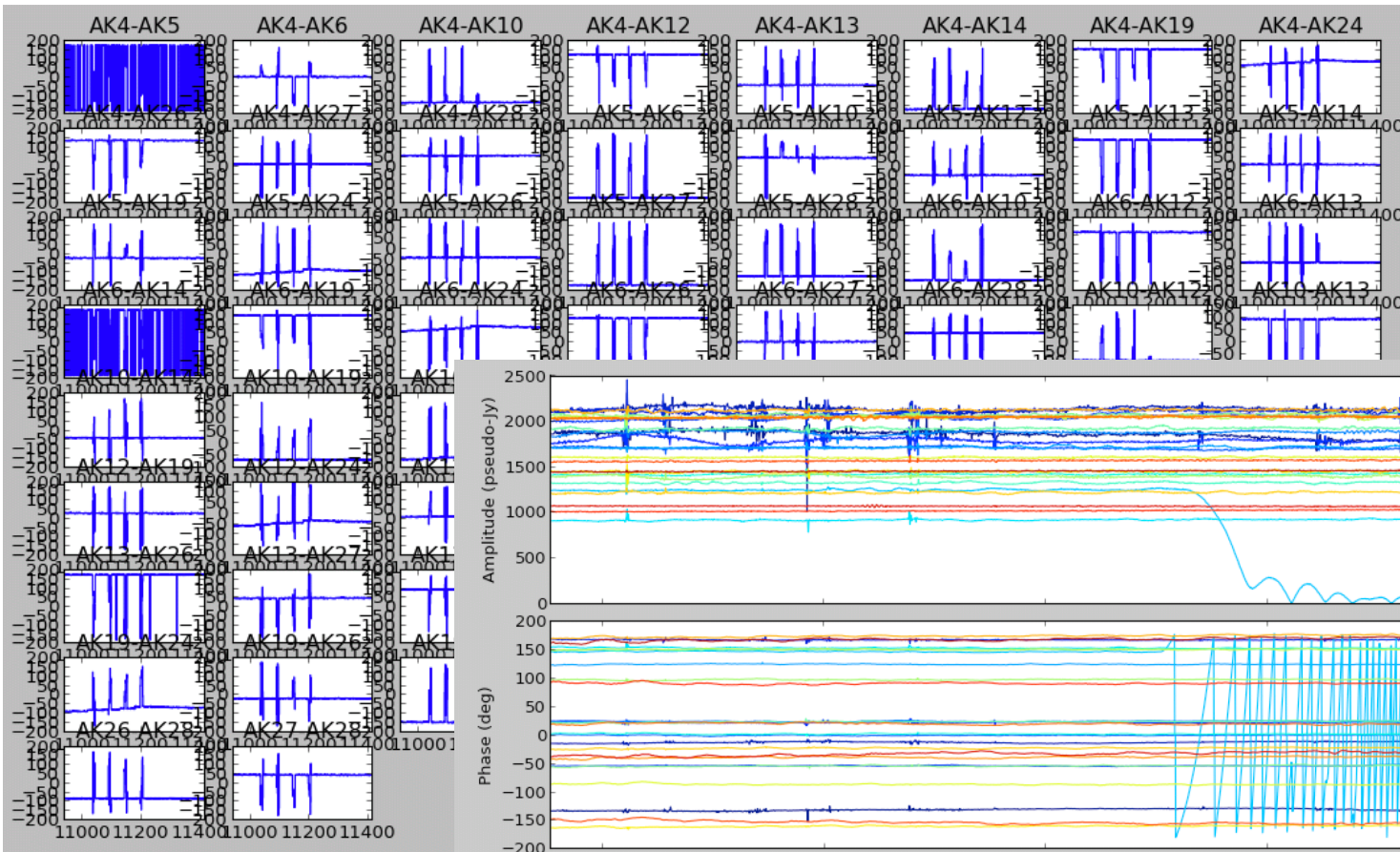
# Splitting and merging data streams

- Each stream can be active or deactivated (behind the scene)
  - Splitting/merging tasks can change the state + can use service ranks
  - These tasks are specific to the parallel mode
  - Not all setup combinations are supported
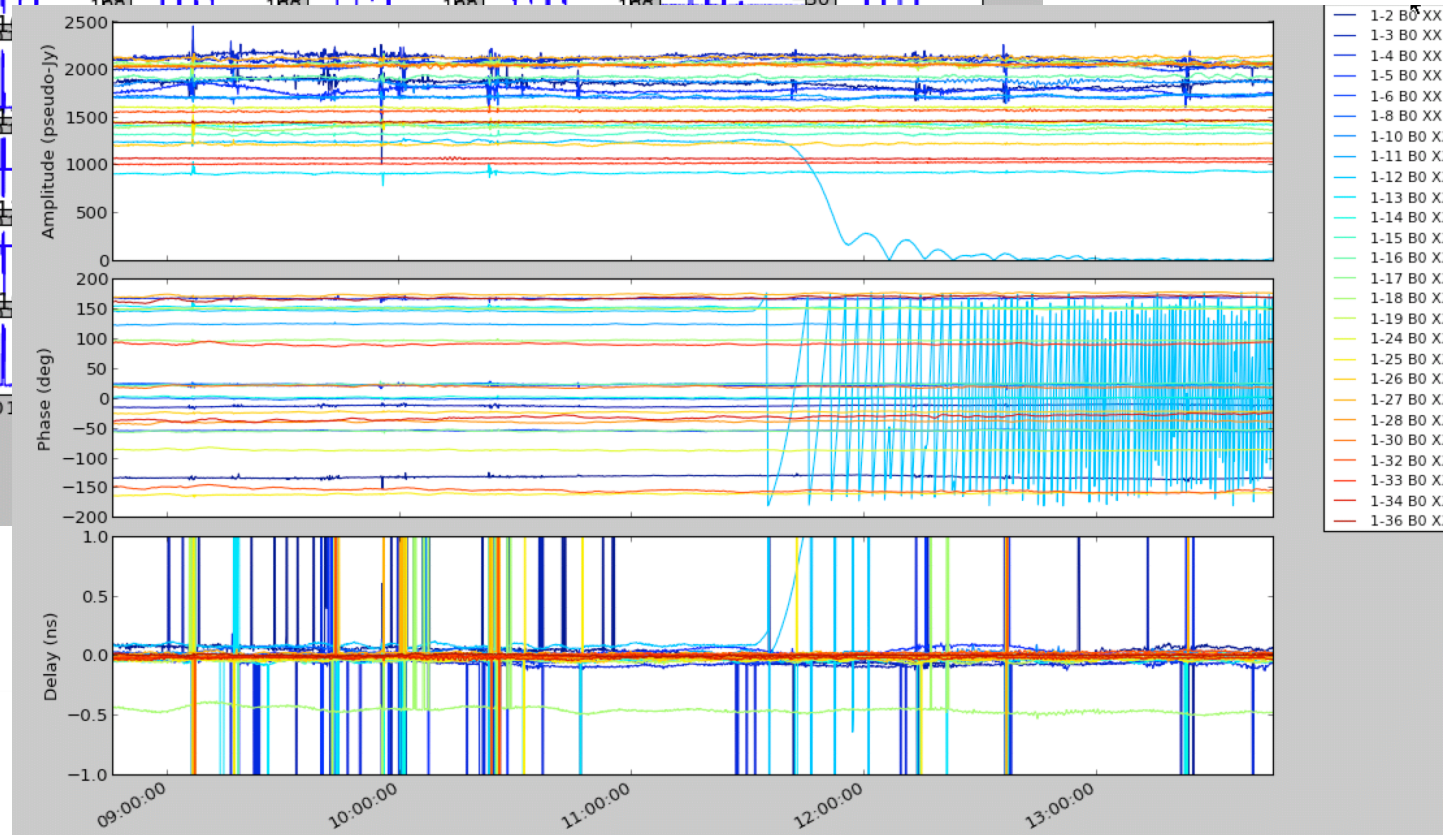  - Support of 5s cycle requires some work
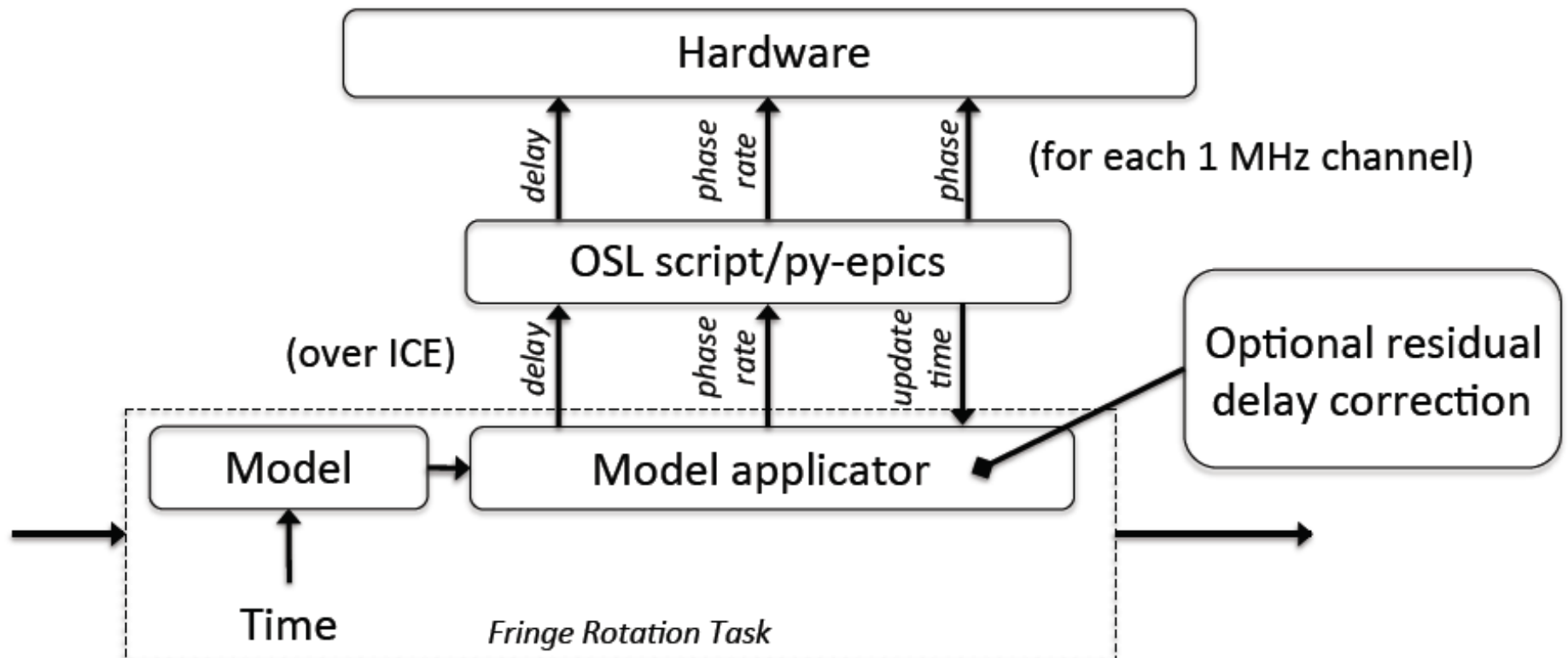
# Data monitoring via TCPSink/vispublisher



← Spectra
(spd)

Frequency
average ➔
(vis)
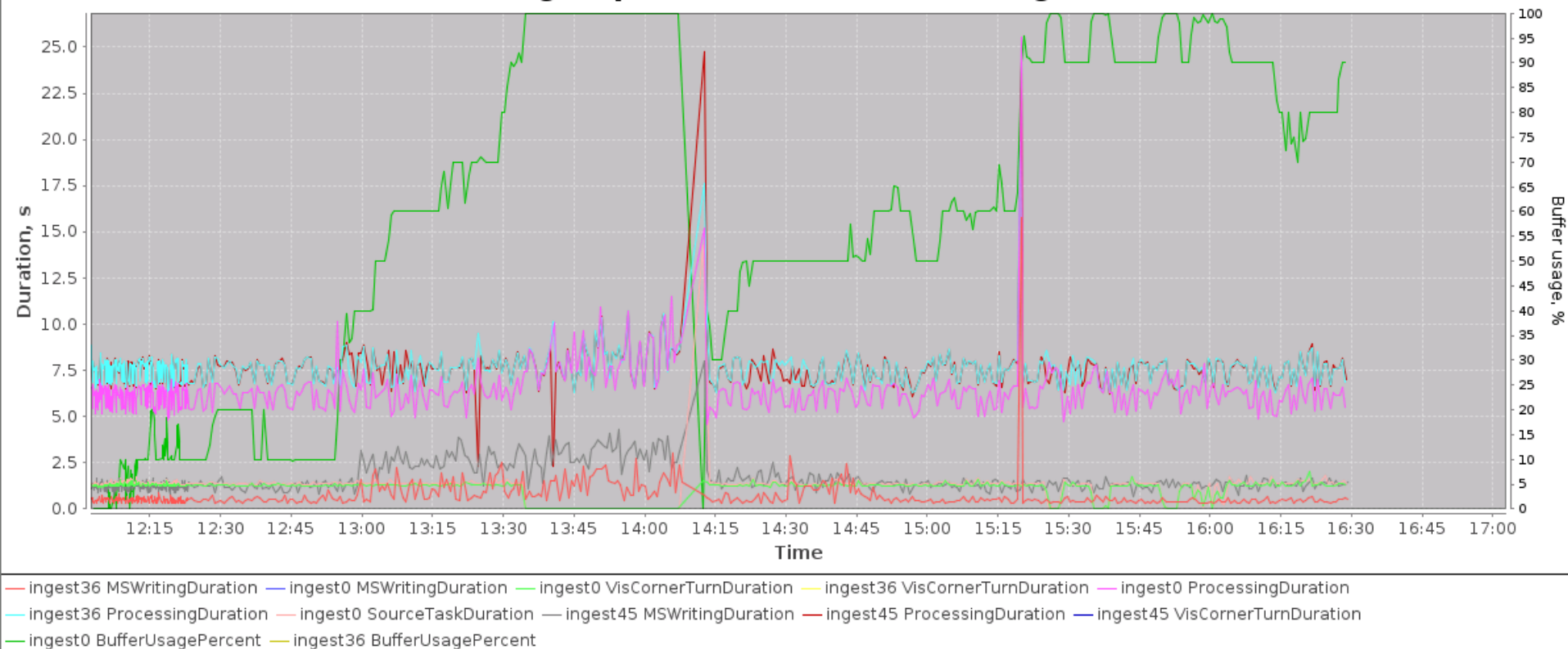
# Commissioning experiments

- Intermediate mode to control fringe rotator from ingest
  - Grew up from the need to support early BETA array, extended to early ASKAP
  - Investigate relative timing + well-controlled baseline system
  - Early science done a few years earlier

# Performance monitoring

- Many performance metrics are published every cycle via Ice
  - Figures for representative ranks can be monitored (monica or Grafana)
  - As many tasks have implicit synchronisation barriers, interpretation requires knowledge of the system architecture and configuration of ingest pipeline
  - Buffer usage and writing times are the most straightforward metrics



Ingest performance and buffer usage

# Various performance lessons

- Logical vs. physical isolation of ingest
  - Writing data to shared lustre filesystem (with dedicated OST/metadata nodes)
  - It took surprisingly long time to get to an acceptable level
    - May be not over yet – we are not pushing the system the same way now
    - Our own processes/processing can also affect performance
    - If it locks up, it happens inside 3[rd] party libraries, so can't really time out
- Real-world astronomy data formats vs. idealistic I/O benchmarks
- Implicit barriers
  - Data sent in staggered fashion (if you need more than one chunk, you wait)
  - Metadata are not available until the end of the cycle
  - Logging at scale may be non-trivial, especially for a  synchronous system
  - Consistency cross-checks may require additional communication
- Not in the regime how HPC is typically used – watch out for bugs
  - It matters where each rank goes

# What wasn't in design, but in use now

- On-the-fly averaging as a separate mode
  - Instead of full-resolution mode, not in addition to
  - This is largely to save disk space / for projects which don't need full resolution
  - Prior flagging is essential
- Per-beam partitioning of the data (one beam in one MS)
  - Single beam mode is a particular case
- Flexible partitioning in frequency (merge/split)
- Real-time monitoring of data after ingest (i.e. vis/spd)
- Various data consistency cross-checks
- Flexible configuration options
  - Zoom modes (user-controlled)
  - Adding correlator hardware (FCM controlled, but requires expert knowledge)
  - Changing antennas included in the array
    - We used to have the main and commissioning array

CSIRO

# Summary

- Ingest is a flexible adapter s/w between correlator and processing
  - Allows us to debug/test synthesis jobs from a standard MS
  - Synchronise parallel data and metadata streams
  - Aggregate/split data as required
  - On-the-fly flagging, if necessary
  - On-the-fly calibration application in the future
  - Optional on-the-fly averaging in frequency
  - Interface to on-the-fly data monitoring (vispublisher -> vis and spd)
- Invaluable tool for commissioning
  - Detect oddities in data stream
  - Non-standard experiments
  - Intermediate solutions (e.g. fringe rotation) to get science results faster

CSIRO

*We acknowledge the Wajarri Yamatji people as the traditional owners of the Observatory site*

# Thank you

**Astronomy and Space Science**
Max Voronkov
Senior Research Scientist

**t** +61 2 9372 4427
**e** maxim.voronkov@csiro.au
**w** www.narrabri.atnf.csiro.au/people/Maxim.Voronkov

CSIRO