ASTRON

Netherlands Institute for Radio Astronomy

LOFAR data

Yan Grange Vishambhar Nath Pandey



ESCAPE WP2 telecon 2020-04-06



LOFAR telescope and data chain



Netherlands Institute for Radio Astronomy

Long-Term Archive

• Somewhat averaged data products go to the Long-Term Archive (LTA)

80

- 1 year embargo, then 'public'
- All data on tape, access through grdftp, maybe webdav
- Post-processing by users on local clusters
 - Survey using grid infra to create images



Description of the parameter space

- Parameters we need to consider in this discussion
- Properties of the files
 - File sizes, number of files in a single data set, number of data sets in an observation
- Properties of the application
 - Random access vs sequential
 - I/O bound versus compute bound (number of operations per read size)
 - Amount of parallelism and inter-process communication



LOFAR data wrt parameter space

- Typical LOFAR observation consists of hundreds of data sets (by frequency)
- Data in Measurement set format:
 - Self-contained with metadata and data in separate files
 - Several ~kB MB sized metadata files, few (often 1) GB-TB sized data files
- Processing:
 - Processing is mostly I/O bound
 - Processing of each frequency is embarrassingly parallel in the first steps
 - At one stage, data and solutions combined per 10 frequencies



Prefactor use case (using the LOFAR sw stack)

- Prefactor pipeline
 - Use a calibrator to calibrate the observed target
 - De facto standard of an imaging pipeline
 - This is a first step in imaging but for now it is the most standardised we have
- Currently using custom pipeline framework
 - We can't experiment with frameworks that could make it
- Everything speaks POSIX
 - But all access to MS through casacore library
- Processing properties
 - Processing is mostly I/O bound
 - Processing of each frequency is embarrassingly parallel in the first steps
 - At one stage, data and solutions combined per 10 frequencies



LOFAR to SKA

- SKA aims to deliver much more science-ready data
 - Reprocessing of visibilities can still be part of the SRC (e.g. EOR project)
- Postprocessing use cases could cover a broader range of the parameter space
- SKA data will be larger so scale-up effects need to be taken into account



Mapping to WP2

- Caching
 - Access is fully linear. Files are written to disk during pipeline. Also link between POSIX and data lake?
- QOS:
 - Large data sets, in the LOFAR case all data is stored on bulk storage (tape). However for processing it needs to be on faster storage
 - Meta data on faster storage, bulk data on slower?
- Network
 - One observation is typically tens of TB to be moved at once over the network.

