## Data challenges of modern HI Surveys Attila Popping



The HI/Story of the Nearby Universe 10-12 September 2018 Groningen, The Netherlands









#### **CHILES**

#### COSMOS HI Large Extragalactic Survey

2.7

#### single pointing in COSMOS 1000 hours integration VLA B-configuration







|                            | OLD   | PILOT  | NEW                       |
|----------------------------|---|--|---------------------------|
| Bandwidth (MHz)            | 6.25  | 240  | 480                       |
| Channels                   | 31  | 16384  | 30720                     |
| Velocity resolution (km/s) | 40  | 3.5  | 3.5                       |
| Instantaneous z coverage   | 0 <z<0.004< td=""><td>0<z<0.193< td=""><td>0<z<0.5< td=""></z<0.5<></td></z<0.193<></td></z<0.004<> | 0 <z<0.193< td=""><td>0<z<0.5< td=""></z<0.5<></td></z<0.193<> | 0 <z<0.5< td=""></z<0.5<> |



#### CHILES





#### ASKAP

#### Murchison Radio-Astronomy Observatory (MRO)

S26° 42' 15", E116° 39' 32"

MURCHISON SHIRE AREA: 1 NL POPULATION: 114 DENSITY: 2.7 mPeople km<sup>-2</sup>





#### ASKAP



# ASKAP is complicated

- 36 antennas
- 36 PAFs

- Data from PAFs



188 sensing elements per PAF  $36 \times 36 \times 2 \times 300 = 777600$  beam formers

: 100 terabits / sec Visibility data to disk : 2.3 gigabytes /sec

> 500,000 monitor points



#### **People have been doing (HI) Surveys for decades.**

- Things were new at the time
- Computers were not up to spec
- Working with a new instrument is always challenging





## Are the challenges of today bigger ?

#### **Are we in a different Era?**







## The challenges of today are bigger !

- RFI Environment is much more challenging
- Going from small (~10MHz) to large (>300MHz) bandwidth adds many complications
- Data volumes are unprecedented
- Data processing requirements are unprecedented
- Workload is larger than humanly possible
- We have more stringent requirements as new surveys go larger/wider/deeper/higher
- We go from ~1 instrument/survey per decade to N where N > 7

#### Are we in a different Era?







#### **RFI environment**:

In the Netherlands from the occasional tv tower to  $\sim 17$  million people with  $\sim 17$  million mobile phones, being stuck in electric cars, while getting information from 7 positioning satellites simultaneously.

#### **Challenge: Radio Frequency Interference**







## L-Band Spectra, 2012 June B-config



log10(Amp) (Jy)

#### **Challenge: Radio Frequency Interference**





## L-Band Spectra, 2012 June B-config



#### **Challenge: Radio Frequency Interference**

#### Freq (MHz)













Baseline Length [km]

#### **Challenge: Radio Frequency Interference**







#### It is hard to find good flagging parameters







1230

(**2** HW) 1250

1230

1280

1230

(ZHW) 1250

S.

Lequer 1500 1520

1280



#### **Reference antenna based RFI tracking**

- Many instruments are limited by the GPS bands.
- RFI is the hot potato of many surveys
- Current methods are masking effected data rather than taking out the bad data.
- Use a reference antenna to track the RFI source and subtract the power from observations
- Potentially powerful method, especially for a telescope with a phased array feed
- Needs a serious commitment and hence investment to make it work









#### **Data volumes are uprecedented**







#### **Data volumes are uprecedented**





#### **Commissioning archive use**



#### Data volumes are uprecedented

A. Hotan

![](_page_17_Picture_5.jpeg)

![](_page_18_Picture_0.jpeg)

- ASKAP is the first telescope that will not be able to hold on to the raw visibilities
- We will require real time processing
- We will require very robust pipelines
- We need high quality pipelines that we trust

## Will we ever reach a state where single-pass processing will be sufficient?

#### No more data storage

#### **ASKAP** data rates

![](_page_18_Figure_8.jpeg)

![](_page_18_Picture_10.jpeg)

![](_page_19_Picture_0.jpeg)

- - Low-level RFI and weak continuum sources
  - Poorly removed sidelobes and other systematic calibration errors

![](_page_19_Figure_5.jpeg)

Visibility storage for deep spectral line surveys using uv grids

## **Visibility storage using UV-grids**

By performing one-pass data processing we lose the calibration versatility

Deep surveys are extremely sensitive to small systematic uncentairtinies

Kristóf Rozgonyi

![](_page_19_Picture_14.jpeg)

![](_page_20_Picture_0.jpeg)

- However visibilities evaluated to a uniform grid for FFT
- Compressed uv data is naturally available because grids are sparse
- By storing some cells several times, we can hold on to baseline information
- ► Time is present on the grid

## **UV grid storage**

There are ways to compress visibility for storage, e.g average by time, baseline...etc. (e.g Wijnholds, S. J. et al 2018, arxiv: 1802.09321)

![](_page_20_Figure_9.jpeg)

uv gridding using convolutional resampling

Kristóf Rozgonyi

![](_page_20_Picture_14.jpeg)

![](_page_21_Picture_0.jpeg)

- However mosaicking reduces the required storage the most, we lose a lot of – even the beam – information

|                  | ASKAP – core<br>(2km max baseline) |                 | ASKAP – full<br>(6km max baseline) |                 |
|------------------|------------------------------------|-----------------|------------------------------------|-----------------|
|                  | $T_{obs} = 0.5$ h                  | $T_{obs} = 8 h$ | $T_{obs} = 0.5$ h                  | $T_{obs} = 8 h$ |
| Quick & Dirty    | 64.14 PB                           | 4.01 PB         | 577.26 PB                          | 36.08 PB        |
| Precise          | 81.85 PB                           | 5.12 PB         | 736.69 PB                          | 46.04 PB        |
| Imaging pipeline | 434.52 PB                          | 27.16 PB        | 1.84 EB                            | 117.87 PB       |
| Mosaicking       | 34.58 PB                           | 2.16 PB         | 311.25 PB                          | 19.45 PB        |
| Sparse uv grid   | 10.00 PB                           | 4.54 PB         | 17.36 PB                           | 9.33 PB         |
| #1, 8h mask      | 69.75 PB                           | 4.36 PB         | 143.24 PB                          | 8.95 PB         |
| #16, 0.5h mask   | 8.63 PB                            | 4.36 PB         | 15.15 PB                           | 8.95 PB         |
| On the fly       | 7.97 PB                            | 3.63 PB         | 13.81 PB                           | 7.46 PB         |

Would enable new processing methodologies with improved: RFI mitigation, continuum subtraction, cleaning, resolution

Visibility storage for deep spectral line surveys using uv grids

## **DINGO storage requirements**

Holding on to the >2km baselines don't increase data rate dramatically on the grids  $\Rightarrow$  method scales well with additional long baselines

Kristóf Rozgonyi

![](_page_21_Picture_11.jpeg)

![](_page_22_Picture_0.jpeg)

#### With the increased amount of data, the processing is also more demanding

![](_page_22_Picture_2.jpeg)

#### **Data processing requirements**

![](_page_22_Picture_4.jpeg)

![](_page_23_Picture_0.jpeg)

![](_page_23_Picture_1.jpeg)

## Imaging

![](_page_23_Picture_3.jpeg)

![](_page_24_Picture_0.jpeg)

![](_page_24_Figure_1.jpeg)

## Imaging

![](_page_24_Picture_3.jpeg)

![](_page_25_Picture_0.jpeg)

![](_page_25_Picture_1.jpeg)

## Imaging

![](_page_25_Picture_3.jpeg)

![](_page_26_Picture_0.jpeg)

Conventional Cluster (pleiades) 5 nodes each node has 2x Intel Xeon X5650 2.66GHz CPUs (6 cores / 12 HTs) with 64-192 GB of RAM

Super computer (MAGNUS) Cray XC40 - 24 cores per node 2.6GHz Intel Xeon E5-2690V3 64GB per Node 35,712 cores available 3PB of storage #58 in the world

AWS

Whatever we wanted r3.xlarge 16 cores 122GB Ram

## **Computing efforts**

![](_page_26_Figure_6.jpeg)

![](_page_26_Picture_7.jpeg)

![](_page_27_Picture_0.jpeg)

![](_page_27_Picture_1.jpeg)

|            | On demand        | Spot Price           |
|------------|------------------|----------------------|
| r3.4xlarge | \$1.68           | \$0.20               |
| r3.2xlarge | \$0.840          | \$0.09               |
| m3.xlarge  | \$0.392          | \$0.04               |
| m3.medium  | \$0. <b>0</b> 98 | \$0.0 <mark>1</mark> |

#### Spot Instance Pricing History

![](_page_27_Figure_4.jpeg)

## **AWS pricing**

![](_page_27_Picture_6.jpeg)

![](_page_27_Picture_8.jpeg)

![](_page_28_Picture_0.jpeg)

|                      | AWS        | Magnus<br>(HPC)      | Pleiades      |
|----------------------|------------|----------------------|---------------|
| Completion<br>Time   | 96hr       | 110hr                | 1,060hr (est) |
| Capital Costs        | AUD\$0     | AUD\$12,000,000      | AUD\$50,000   |
| Operational<br>Costs | AUD\$2,000 | AUD\$3,240<br>(free) |               |
| Control              | Root       | Limited              | Root          |
| Usability            | Complex    | Good                 | Good          |

## **Computing efforts**

![](_page_28_Picture_3.jpeg)

![](_page_29_Picture_0.jpeg)

- Department Cluster
  - Not very satisfactory
- HPC
  - Very fast
  - You can't have root access
  - Installing new software is done by the admin's
  - In WA it is effectively free
- Cloud
  - You can do what you like (a good and a bad thing)
  - EBS volumes are slow
  - Directed attached SSDs are very quick
  - You pay for what you use... And if you forget to turn it off you are still paying

## **Compute conclusion**

![](_page_29_Figure_15.jpeg)

**Cloud computing can provide an excellent alternative, however will require a change in mindset** 

![](_page_29_Picture_18.jpeg)

![](_page_30_Picture_0.jpeg)

![](_page_30_Figure_1.jpeg)

## We have more stringent requirements

- Low level continuum artefacts after combining data
- Not visible in individual observations
- Need better continuum model and go back to raw data to subtract.

Note: we may have deleted the raw visibilities at this stage ....

![](_page_30_Picture_7.jpeg)

![](_page_30_Picture_8.jpeg)

![](_page_31_Picture_0.jpeg)

![](_page_31_Picture_1.jpeg)

![](_page_31_Figure_2.jpeg)

#### **Complicated workflows**

![](_page_31_Picture_4.jpeg)

![](_page_32_Picture_0.jpeg)

| Level 0               | Level 1          | Level 2                     |
|-----------------------|------------------|-----------------------------|
| Reproducible results  | Documented       | Smart use of resources      |
| Easy to build/install | Portable         | Scalable                    |
| Easy to use           | Good diagnostics | Cache intermediate results? |
|                       | Configurable     |                             |
|                       |                  |                             |

## The ideal pipeline

S. Makhathini

![](_page_32_Picture_4.jpeg)

![](_page_33_Picture_0.jpeg)

![](_page_33_Figure_1.jpeg)

### **Typical data reduction**

S. Makhathini

![](_page_33_Picture_4.jpeg)

![](_page_33_Figure_5.jpeg)

![](_page_34_Picture_0.jpeg)

![](_page_34_Figure_1.jpeg)

#### Workflows have to be modular

keep data in memory

### **Typical data reduction**

S. Makhathini

To make this effective, the input and output format of all modules has to be similar, so you can

![](_page_34_Picture_7.jpeg)

![](_page_35_Picture_0.jpeg)

The Data Activated Liu Graph Engine (DALiuGE) developed by ICRAR is an execution framework for processing large astronomical datasets at a scale required by the SKA1.

DALiuGE provides:

- an interface for expressing complex data reduction pipelines consisting of both data sets and algorithmic components; and
- an implementation run-time to execute pipelines on distributed resources.

#### **Chen Wu**

DALiuGE

![](_page_35_Figure_7.jpeg)

![](_page_35_Picture_8.jpeg)

![](_page_36_Picture_0.jpeg)

#### DALiuGE execution

![](_page_36_Picture_2.jpeg)

**Minimise:** 

![](_page_36_Figure_4.jpeg)

**Chen Wu** 

#### We require workflows and graph engines

![](_page_36_Picture_7.jpeg)

![](_page_37_Picture_0.jpeg)

Cray XC30 Series Supercomputer

472 Compute Nodes:

- 2 x 3.0 GHz Intel Xeon CPUs
- 10 Cores per CPU
- 64 GB DDR3-1866
- Cray Aries Interconnect
- Cray Dragonfly Topology
- 200 TeraFLOPS
- **1.4 Petabytes Lustre Data Storage**

**ASKAPsoft Data Pipelines** 

#### What will be the model for the SKA?

#### **ASKAP: telescope + computer**

![](_page_37_Picture_13.jpeg)

![](_page_37_Picture_14.jpeg)

![](_page_37_Picture_15.jpeg)

![](_page_38_Picture_0.jpeg)

- Projects are not run by individuals but large (international teams)
- project management and team interaction is very important
- Should we use redmine / confluence / zoom / slack / pbworks / email / jira / git / svn to share information
- Should we use all (e.g. ASKAP) to make sure no information gets lost, at the risk of losing track overall
- Human pipeline vs Automated pipeline
- Do you spend your time working on the data or working on the pipeline
- We have to accept that the results will not be as good as theoretically possible
- Loss of expertise

The unwritten rule of a PhD is that you put in a large amount of hard work (i.e. flagging data, eyeballing sources etc.) The one condition here is that at least it is possible to flag or eyeball this data within the duration of a PhD. This condition is changing which is good and bad

## Workload larger than humanly possible

![](_page_38_Picture_12.jpeg)

![](_page_38_Figure_13.jpeg)

![](_page_39_Picture_0.jpeg)

![](_page_39_Picture_1.jpeg)

#### Are we going to many or too many ...?

#### Are we doing too much ?: new telescopes, new software, new projects etc ...

- Most people are involved in many large projects, realistically you can only fully commit to one There is a risk in using new software with a new telescope:
- CHILES statement: We have found this problem, it is either a bug or feature of CASA
- ASKAP: in early days it was very difficult to verify results.
- If you use new software, make sure you have motivated developers in your team
- Keep things transparent, software should be modular and not a black box (e.g. CASA clean task)
- We need a collaborative approach to solve the same problems (e.g. SoFiA)

#### We go from ~1 instrument/survey per decade to many

- **JVLA:**
- ASKAP:
- APERTIF:
- MeerKAT:
- **FAST:**
- •

![](_page_39_Picture_19.jpeg)

![](_page_40_Picture_0.jpeg)

![](_page_40_Picture_1.jpeg)

- SKA precursor surveys
- International collaborative project

![](_page_40_Figure_5.jpeg)

#### **SoFiA - the Source Finding Application**

HI source finding and parameterisation pipeline for extragalactic

#### Originally initiated by the WALLABY team

T. Westmeier (chair), P. Serra, B. Koribalski, B. Winkel, R. Jurek, H. Courtois, A. Popping, N. Giese, L. Flöer, L. Staveley-Smith, T. van der Hulst, M. Meyer

![](_page_40_Picture_11.jpeg)

![](_page_41_Picture_0.jpeg)

#### ★ Model galaxies

- Created in GIPSY using galmod
- ► Wide range of
  - sizes
  - fluxes
  - inclinations
  - rotation velocities
- Placed on a regular grid for efficiency
  - 600 galaxies
  - cube size: 800 × 800 pixels 800 channels
  - 1.9 GB of data

![](_page_41_Picture_12.jpeg)

### **ASKAP Early Science Data**

![](_page_41_Figure_14.jpeg)

Model galaxies

![](_page_41_Picture_16.jpeg)

Model galaxies + ASKAP noise

#### T. Westmeier

![](_page_41_Picture_19.jpeg)

![](_page_42_Picture_0.jpeg)

| ★ SoFiA settings                                  | 1.0         |
|---|-------------|
| ► S+C finder with 3σ threshold                    | 0.8-        |
| Reliability estimation enabled                    | Хец 0.6 -   |
| ★ Results   | านรู้ 0.4 - |
| About 63,000 detections                           | 0.2         |
| Half of these are negative                        | 0.0-        |
| 276 sources remain after<br>reliability filtering | -0.2-1.0    |

High reliability thanks to reliability filter!

#### **ASKAP Early Science Data**

![](_page_42_Figure_4.jpeg)

T. Westmeier

![](_page_43_Picture_0.jpeg)

#### **★** Completeness

![](_page_43_Figure_2.jpeg)

# Peak SNR:100% above $2.6\sigma$ detections at $\gtrsim 0.5\sigma$

#### **ASKAP Early Science Data**

Integrated SNR: 100% above  $5\sigma$ detections at  $\gtrsim 2\sigma$ 

![](_page_44_Picture_0.jpeg)

#### ★ Data set II

- ► ASKAP-12
- ► HI data of the NGC 7232 group
- WALLABY early science
- Mosaic of multiple PAF beams
- ► 12.6 GB in size
- Thanks to Karen Lee-Waddell

![](_page_44_Picture_8.jpeg)

*Wide-field ASKAP L-band Legacy All-sky Blind surveY* (Koribalski & Staveley-Smith)

#### **ASKAP Early Science Data**

![](_page_44_Picture_11.jpeg)

ASKAP, Boolardy, Western Australia

T. Westmeier

![](_page_44_Picture_14.jpeg)

![](_page_45_Picture_0.jpeg)

## ★ Challenges for SoFiA

- Strong noise variation across field
- Residual sidelobes from shallow deconvolution
- RFI stripes from insufficient flagging
- Other imaging artefacts

![](_page_45_Picture_6.jpeg)

- ► A lot of these problems with ASKAPsoft are now under control or being resolved
- Nevertheless a good test data set to throw at SoFiA

#### **ASKAP Early Science Data**

![](_page_45_Picture_10.jpeg)

![](_page_46_Picture_0.jpeg)

#### **★** Issues and solutions

- Low detection threshold would pick up too many artefacts
- Use  $5\sigma$  threshold

![](_page_46_Figure_5.jpeg)

![](_page_46_Picture_7.jpeg)

![](_page_47_Picture_0.jpeg)

#### **★** Issues and solutions

- those of artefacts

![](_page_47_Figure_5.jpeg)

#### **ASKAP Early Science Data**

![](_page_47_Picture_8.jpeg)

![](_page_48_Picture_0.jpeg)

#### ★ Result

- Detailed HI map of NGC 7232 group
- Lee-Waddell et al. (in prep.)

![](_page_48_Picture_4.jpeg)

#### **ASKAP Early Science Data**

![](_page_48_Figure_6.jpeg)

T. Westmeier

![](_page_49_Picture_0.jpeg)

★ Source finding on ASKAP early science data using SoFiA works

#### ★ Clean data free of artefacts

- Low threshold of  $3\sigma$  possible
- ► Completeness of 100% at  $SNR_{peak} \gtrsim 2.6$  and  $SNR_{int} \gtrsim 5.0$
- ► Reliability near 100%
- ★ Data affected by artefacts
  - Higher threshold of  $5\sigma$  required
  - Limited completeness and reliability
  - Negative artefacts render reliability filter useless

### **Source Finding**

![](_page_49_Figure_12.jpeg)

#### BEST CASE **SCENARIO**

![](_page_49_Figure_15.jpeg)

#### **WORST** CASE **SCENARIO**

T. Westmeier

![](_page_49_Picture_19.jpeg)

![](_page_50_Picture_0.jpeg)

#### AP's Masters thesis (with Thijs) ~ year

![](_page_50_Figure_2.jpeg)

#### LGG 334 then and now

SoFiA (with Thijs) ~20 seconds

![](_page_50_Picture_5.jpeg)

![](_page_51_Picture_0.jpeg)

Did you solve the problem? "No, but I am very positive. Otherwise why coming to work"

![](_page_52_Picture_0.jpeg)

![](_page_53_Picture_0.jpeg)

- Can we afford to delete raw visibilities ?
- Will we ever reach a state where we can do single pass processing ?
- **Is it still realistic to do data reduction by hand ?**
- Should we put effort in decreasing data volume rather than more processing ?
- Who should be working on or take responsibility for RFI mitigation ?
- Have you implemented data verification and quality control ?
- How do we avoid reinventing the wheel for every telescope ?
- **•** Do you have enough software engineers in your team ?

## **Questions (to you)**