# Summary

- SKA Scope: Is limited!

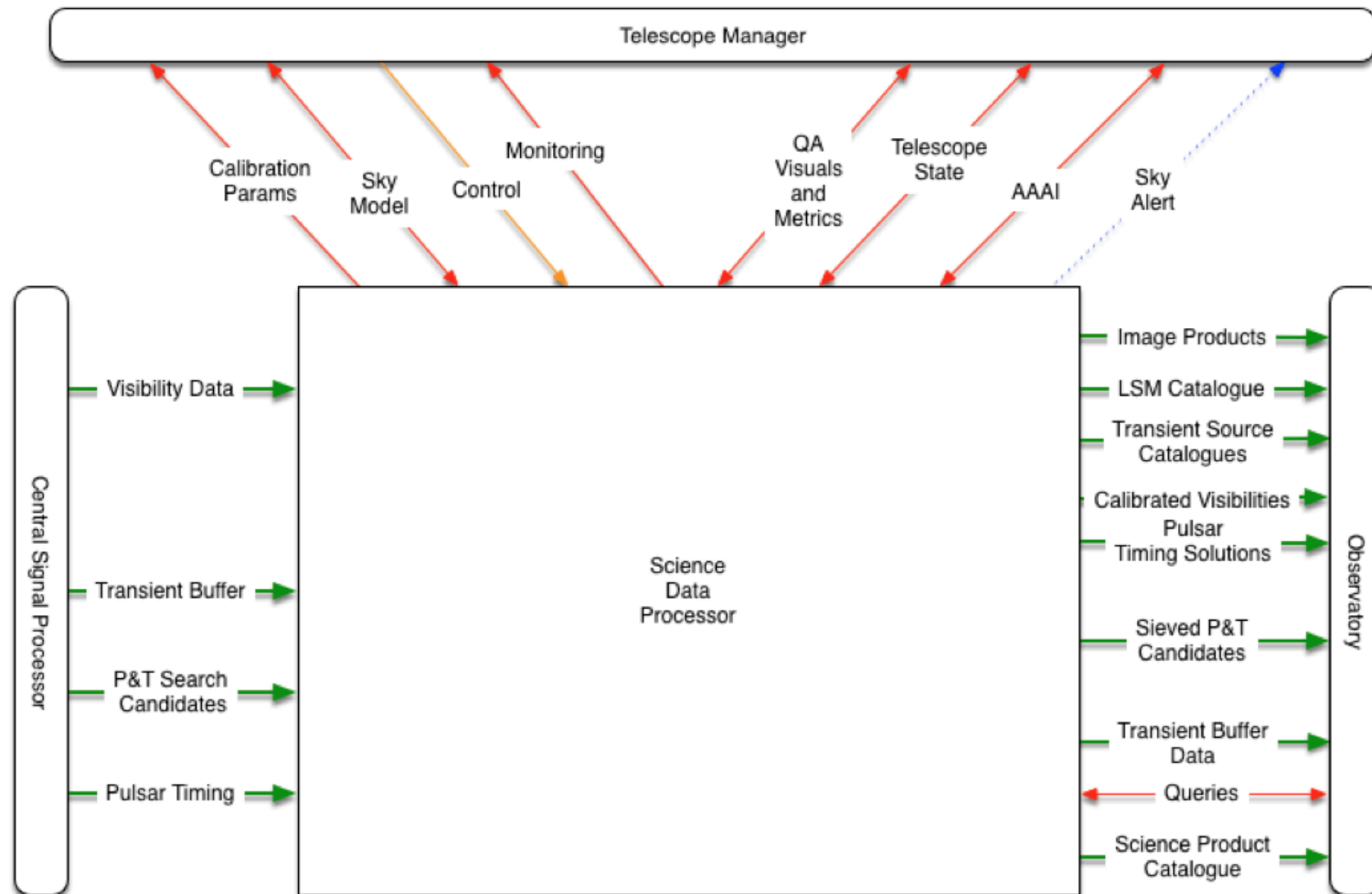- SKA will deliver only from a subset of standard **data products**.

- Paradigm shift: The SDP compute is a resource, in an identical way to the "on-sky" time.

- No option for reprocessing data in current design.

- But the board has adopted regional centres…

# SDP Visibility processing



RCAL

Real Time Calibration

"real-time"; latency seconds to minutes

Ingest

Ingest

FastImg

Fast Imaging

Double Buffer: Process 6 hours of data, in 6 hours, whilst collecting and performing the real-time processing on the next 6-hours' worth.

ICAL

Selfcalibration

DPrepA, DPrepB, DPrepC

Final Imaging

"off-line"; latency 12hours since start of observation
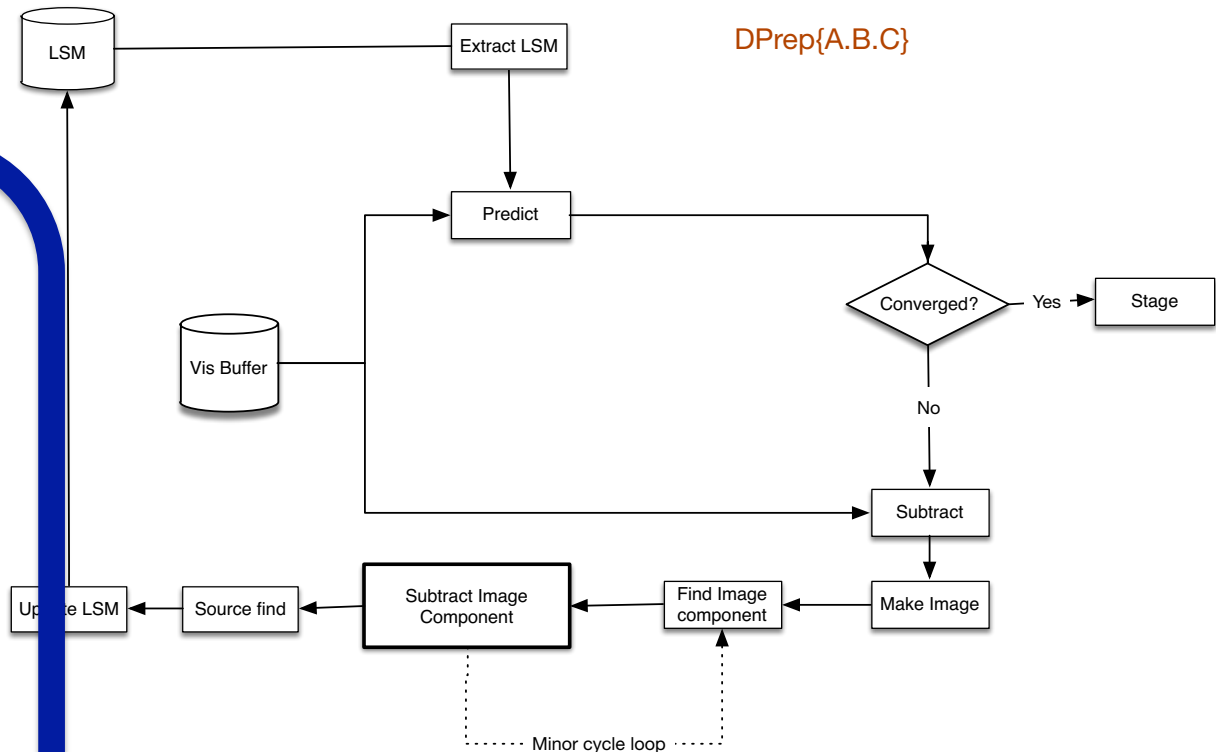
UNIVERSITY OF CAMBRIDGE

# Image products

- **Continuum Taylor Term** Images

- Coarsely channelised "continuum" full Stoke's images (~500 ~~channels~~

1. **Residual image (continuum & spectral line)**

2. **Clean component image (continuum and spectral line)**

3. **Spectral line cube** (continuum subtracted)

4. **Calibrated gridded visibilities**, at the spatial and frequency resolution required by the experiment.

5. **Accumulated Weights** for each uv cell in each grid.

- Representative Point Spread Function for observation

DPrep{A.B.C}

LSM → Extract LSM → Predict

Predict → Converged? → Yes → Stage

Converged? → No → Subtract

Vis Buffer

Subtract → Make Image → Find Image component → Subtract Image Component → Source find → Update LSM

Minor cycle loop

Auxiliary products multiply up the required archive size (by a factor of 5-10) but they are necessary to enable science product generation
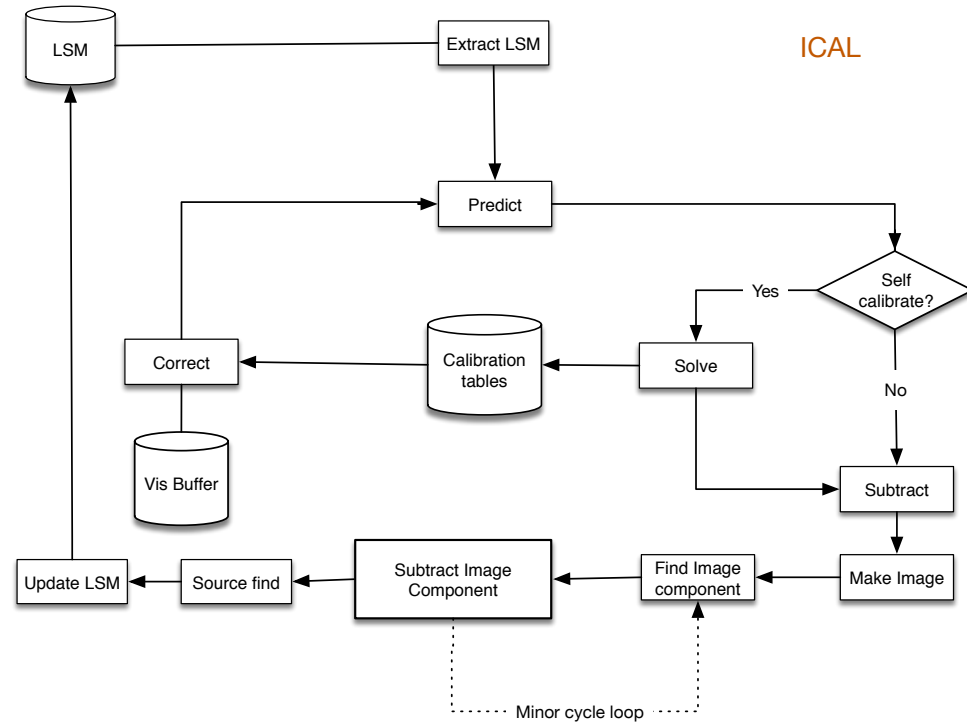
UNIVERSITY OF CAMBRIDGE

# SDP Visibility processing: Parametric model

| Pipeline: ICAL, Self-calibration |
|---|
| Extract LSM; Update LSM; Notify LSM <br> … (not yet modelled) |
| Predict via DFT <br> $C_{\text{Predict}} = (1 + N_{\text{SelfCal}})N_{\text{major}}32N_{\text{source}} \times R_{\text{vis,full FoV}}$ |
| Solve <br> $C_{\text{Solve}} = (1 + N_{\text{SelfCal}})N_{\text{beam}}N_{\text{pol}}\left(\frac{1}{t_{\text{ICAL,G}}} + \frac{f_{sol}^{B}}{t_{\text{ICAL,B}}}\right) \times 48N_{\text{it}}N_{\text{a}}^{2}$ |
| Subtract <br> $C_{\text{Subtract}} = 16 \times ((1 + N_{\text{SelfCal}})N_{\text{major}}) \times R_{\text{vis,full FoV}}$ |
| Correct <br> $C_{correct} = 224 \times (1 + N_{\text{SelfCal}})N_{\text{major}} \times R_{\text{vis,full FoV}}$ |
| Weight visibilities <br> $C_{\text{Weight}} = 16(1 + N_{\text{SelfCal}})N_{\text{major}}R_{\text{vis,full FoV}}$ |
| Phase rotation <br> $C_{\text{phrot,total}} = 2 \times (1 + N_{\text{SelfCal}})N_{\text{major}}N_{\text{facet}}^{2} \times 28R_{\text{vis,full FoV}}$ |
| Coalesce <br> $C_{\text{Coalesce}} = 8(1 + N_{\text{SelfCal}})N_{\text{major}}N_{\text{facet}}^{2}R_{\text{vis,full FoV}}$ |
| Grid <br> $C_{\text{grid}} = (1 + N_{\text{SelfCal}})N_{\text{major}}N_{\text{facet}}^{2}N_{\text{mm}} \times \sum_{\text{i,j}}^{\text{baselines}} R_{\text{vis,facet FoV}}(B_{\text{i,j}})8(N_{\text{pix,kernel}}(B_{\text{i,j}}))^{2}$ |
| De-grid <br> $C_{\text{de-grid}} = C_{grid}$ |
| Gridding kernel update <br> $C_{\text{Kernels}} = \frac{(1+N_{\text{SelfCal}})N_{\text{major}}N_{\text{facet}}^{2}N_{\text{pol}}N_{\text{beam}}}{t_{\text{kernel,update}}}(\sum_{\text{i,j}}^{\text{baselines}} 5N_{\text{cvff}}^{2}\log_{2}(N_{\text{cvff}}^{2})N_{\text{f,kernel}})$ |
| FFT <br> $C_{\text{FFT}} = \left(\frac{(1+N_{\text{SelfCal}})N_{\text{major}}N_{\text{facet}}^{2}N_{\text{pol}}N_{\text{beam}}(N_{\text{subbands}}N_{\text{Tt}})}{t_{\text{snap}}}\right)2.5N_{\text{pix,facet}}^{2}\log_{2}(N_{\text{pix,facet}}^{2})$ |
| IFFT <br> $C_{\text{IFFT}} = C_{\text{FFT}}$ |
| Re-project <br> $C_{\text{Reproj}} = \frac{(1+N_{\text{SelfCal}})N_{\text{major}}r_{\text{facet}}^{2}N_{\text{pol}}N_{\text{beam}}(N_{\text{subbands}}N_{\text{Tt}})}{t_{\text{snap}}}50N_{\text{pix}}^{2}$ |
| Image spectral fitting <br> $C_{\text{spectralfit}} = 2\frac{(1+N_{\text{SelfCal}})N_{\text{major}}N_{\text{pol}}N_{\text{beam}}(N_{\text{subbands}}N_{\text{Tt}})N_{\text{pix}}^{2}}{t_{\text{obs}}}$ |
| Identify Component <br> $C_{\text{Identifycomponent}} = \frac{2(1+N_{\text{SelfCal}})N_{\text{major}}N_{\text{minor}}N_{\text{beam}}N_{\text{pol}}N_{\text{subbands}}N_{\text{Tt}}N_{\text{pix}}^{2}}{t_{\text{obs}}}$ |
| Subtract Image component <br> $C_{\text{Subtractimagecomponent}} = \frac{2(1+N_{\text{SelfCal}})N_{\text{major}}N_{\text{minor}}N_{\text{pol}}N_{\text{beam}}(N_{\text{subbands}}N_{\text{Tt}}N_{\text{patch,pix}}^{2})}{t_{\text{obs}}}$ |
| Source Find <br> Insignificant |

- We continue to develop a parametric model for the SDP compute.

- From this we estimate the required compute rate for each pipeline.



UNIVERSITY OF CAMBRIDGE

# SDP will be cost constrained

- There is always more processing one *could* do on the data
- Limited budget puts constraints on both the hardware built and the use made of it (power)
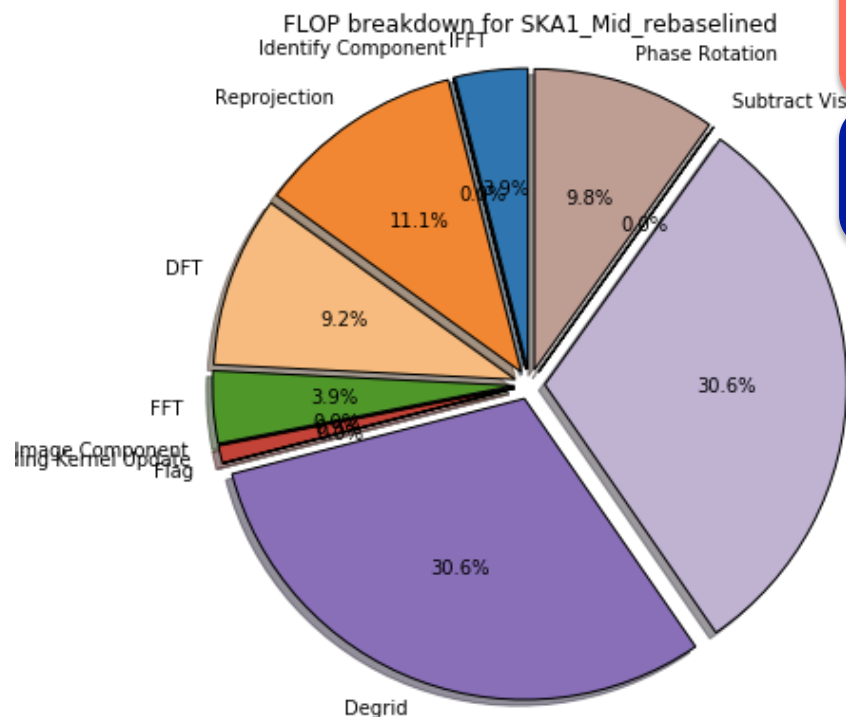- SKAO have provided example power caps to consider

# Power caps: 4MW, 10MW

# SDP will be cost constrained

SKA1 SDPs are ~200-300 PFLOPS systems (25% efficiency)



FLOP breakdown for SKA1_Mid_rebaselined

|  | LOW (50-350M Hz) | MID Band 1 (350-1050 MHz) | MID Band 2 (950-3050 MHz) | MID Band 5 (4.6 to 9.6 GHz) |
|---|---|---|---|---|
| DD CAL (not in iPython) | 18.3 | 17.4 | 17.4 | 17.4 |
| ICAL | 4.9 | 9.5 | 7.5 | 6.3 |
| DPrep A+B | 4.8 | 10.8 | 9.2 | 6.8 |
| DPrep C | 12.0 | 30.4 | 23.0 | 17.4 |
| Fasting | | | 3.0 | 2.5 |
| Sustained Compute Load Total (PFLOPS) | 41.5 PFLOPS | 72.1 PFLOPS | 60.2 PFLOPS | 50.5 PFLOPS |
| Actual P | 4.4 | 7.0 | 5.8 | 4.9 |
| Apparent power, with PUE and power factor (MVA)** | 5.8 | 9.9 | 8.3 | 6.9 |
| Hardware CAPEX Estimate (M€)*** | 57 | 110 | 92 | 77 |

Need 6MVA for LOW, 10VA for MID c.f.
4MVA for LOW, 10MVA for MID

UNIVERSITY OF CAMBRIDGE

# SDP will be a one way street

- Power caps are at the right sort of level to enable processing, according to our model

However, there's no "slack" here

- No re-processing of data

- Each 6 hour observation is passed through the pipelines only once

- Each 6 hour observation has data products pushed through to the archive about 12 hour after the observation started.

# High Priority Science Objectives

- SKAO has developed a list of HPSO experiments – programmes targeting specific scientific goals and taking long periods (~5000-16000 hours) of telescope time.
- Draft schedule for these taking 5-15 years to complete
- Just an example – *NOT* to be taken as "real" time allocation

| Sci Goal | SWG | Science Objective | SKA1 Compone | Hours |
|---|---|---|---|---|
| 1 | CD/EoR | EoR - I. Imaging | LOW | 5000 hrs |
| 2 | CD/EoR | EoR - II. Power spectrum | LOW | 10000 hrs |
| 4 | Pulsars | Reveal pulsar population | LOW+MID | 12750+3200 hrs |
| 5 | Pulsars | High precision timing | LOW+MID | 4300+3200 hrs |
| 13 | HI | Resolved HI out to z~0.8 | MID | 5000 hrs |
| 14 | HI | ISM in the nearby Universe. | MID | 2000 hrs |
| 15 | HI | ISM in our Galaxy | MID | 12600 hrs |
| 18 | Transients | Fast Radio Bursts | MID | 10000 hrs |
| 22 | Cradle of Life | Map dust grain growth | MID | 6000 hrs |
| 27 | Magnetism | All-Sky magnetic fields | MID | 10000 hrs |
| 32 | Cosmology | Gravity on super-horizon scales. | MID | 10000 hrs |
| 33 | Cosmology | Non-Gaussianity and the matter dipole. | MID | 10000 hrs |
| 37+38 | Continuum | Star formation history of the Universe | MID | 16000 hrs |

100,000 hours = 11 years

We can use these to generate example SDP use cases and archive growth rates.
Also could enable load balancing if we relax latency requirement of off-line processing.

# High Priority Science Objectives

| HPSO | Total (PFLOPS) | Hours of telescope time | Fraction of time |
|---|---|---|---|
| U.HPSO-4a Pulsar Search MID SPF1 | ~0 | 800 | 0.01 |
| U.HPSO-4b Pulsar Search MID SPF2 | ~0 | 2400 | 0.04 |
| U.HPSO-5a Pulsar Timing MID SPF2 | ~0 | 1600 | 0.02 |
| U.HPSO-5b Pulsar Timing MID SPF3 | ~0 | 1600 | 0.02 |
| U.HPSO-13  Hi Kinematics and Morphology | 25.6 | 5000 | 0.07 |
| U.HPSO-14 Hi MID | 32.7 | 2000 | 0.03 |
| U.HPSO-15 Studies of the ISM in our Galaxy | 26.2 | 12600 | 0.19 |
| U.HPSO-18 Transients MID | ~0 | 10000 | 0.15 |
| U.HPSO-22 Cradle of Life MID Band 5 | 25.4 | 6000 | 0.09 |
| U.HPSO-27 All Sky Magnetism | 26.3 | 10000 | 0.15 |
| U.HPSO-37a Continuum Survey MID band 2 | 28.1 | 2000 | 0.03 |
| U.HPSO-37b Continuum Survey MID band 2 (deep) | 28.1 | 2000 | 0.03 |
| U.HPSO-37c  Continuum survey, band 2 wide | 28.1 | 10000 | 0.15 |
| U.HPSO-38a Continuum Survey MID band 5 | 26.1 | 1000 | 0.01 |
| U.HPSO-38b Continuum Survey MID band 5 | 26.1 | 1000 | 0.01 |
| Weighted average FLOPS value for MID HPSOs | | | 20 PFLOPS |
| Approximate AVERAGE Apparent power requirement[2] | | | ~2.7 MVA |

Average FLOPS values 3.5x lower for MID than maximal case.

Could still build system for maximal case but not power it up all the time

or

Relax latency requirement and save both capital cost and power cost

# Data rates: be careful what you wish for!

- SKA images ~30k pixels on a side for MID, ~8k for LOW
  - (Low images have 20x fewer pixels than MID from simple $B_{max}/D$ ratio)
- Each instrument is capable of 65k channels, 4 stokes parameters in output images
- Auxiliary products give 7x increase in data size c.f. single image cube
- 4 Bytes per pixel (32 bit)
- Data creation rate scales as number of pixels (2D), number of channels and inverse of observation length
- If observations are too short (or too many channels), data rate going out is larger than visibility data rate coming in!

*So much for " data reduction "*

**UNIVERSITY OF CAMBRIDGE**

# HPSO analysis of archive growth rate

MID:

Weighted average archive growth rate for all experiments is
                                        9 Gbits/s

Much lower than the once per 6 hours maximal rate of
                                        3 Tbits/s*          (or 30x 100Gbits/s)


25% of time is "allocated" to Non-imaging experiments
Many imaging experiments are continuum experiments
The spectral line experiments (so far) do not need all 65k channels

*Yes, the observant of you will have calculated that the incoming data rate for both SKA instruments is around 40G visibilities per second, or ~3.8 Tbits/s.
But an ultra deep blind spectral line (Hi) survey at full spatial resolution is on the list of "things which might kill us"
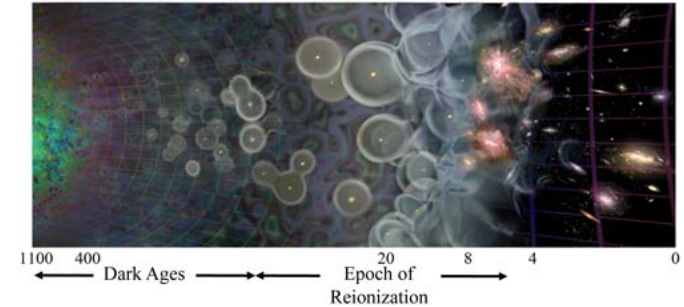
# HPSO analysis of archive growth rate



LOW:

EoR Science is expected to need access to ungridded visibility data.

Even if we average this (post-calibration) in frequency (<10kHz) and time (<10s) we estimate that each of the 3 EoR experiments will generate data rate of averaged visibilities is 200Gbits/s.

There are 3 EoR experiments, each 5000 hours – total for all three is 1.4 Exabytes!

SKAO estimate only 11% (1000 hrs pa) of time likely "good" enough for EoR… (so combined EoR experiments take 15 years)

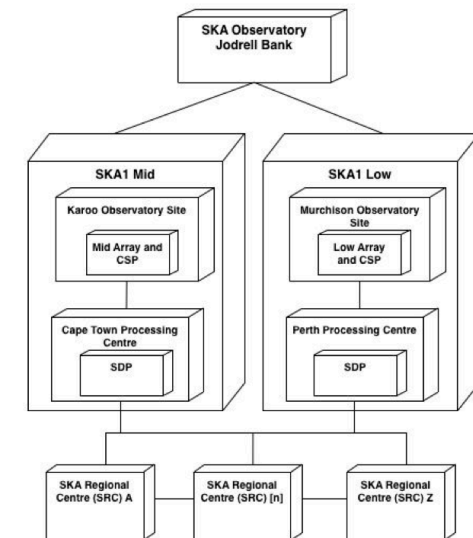*She might look cute, but she's going to eat a hell of a lot of peanuts!*

# Summary

- SKA Scope: Is limited!

- SKA will deliver only from a subset of standard data products.

- Paradigm shift: The SDP compute is a resource, in an identical way to the "on-sky" time.

- No option for reprocessing data in current design.

- But the board has adopted regional centres

- Collaborations with UK science teams are needed to foster better understanding of requirements of SDP and SKA regional centres.





*Don't be afraid to talk about the elephants!*

END