

SKA Science Data Processor

B. Nikolic &
SKA SDP Consortium Team

Astrophysics Group, Cavendish Laboratory, University of Cambridge
<http://www.mrao.cam.ac.uk/~bn204/>

7 May 2013
NRAO Charlottesville



UNIVERSITY OF
CAMBRIDGE

- ▶ SKA “Science Data Processor”
 - ▶ What it is?
 - ▶ The computing and data challenge
 - ▶ Architectures and approaches to solving this challenge
-
- ▶ Similarities of LLST? Differences?

About the Square Kilometre Array

Science Data Processor

Architectures, Technologies

Summary

Introduction

SKA SDP

B. Nikolic / SDP
Team

SKA₁-Mid



SKA₁-Low



SKA₁-Survey



About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary

Introduction

SKA SDP

B. Nikolic / SDP
Team

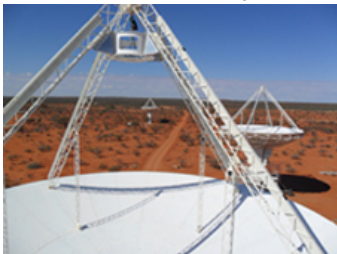
SKA₁-Mid



SKA₁-Low



SKA₁-Survey



HPC Computer



About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary

Introduction

SKA SDP

B. Nikolic / SDP
Team

SKA₁-Mid



SKA₁-Low



SKA₁-Survey



HPC Computers



About the Square
Kilometre Array

Science Data
Processor

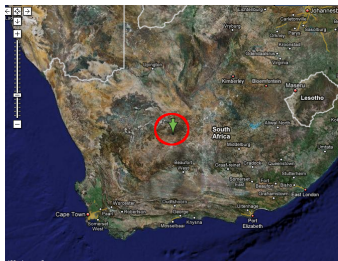
Architectures,
Technologies

Summary

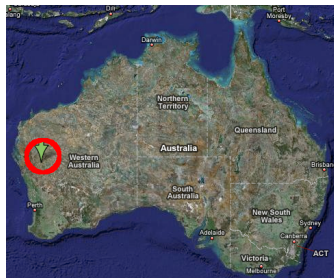
Summary of planned features:

- ▶ $100\times$ better sensitivity than current best telescopes
- ▶ $10^6\times$ sky survey speed compared to current facilities
- ▶ Frequency coverage 50 MHz– \sim 10 GHz
- ▶ Extremely radio quiet sites
- ▶ Staged construction:
 - ▶ Precursors: ASKAP/MeerKAT – under construction
 - ▶ SKA₁: Construction start in 2016-7. Full Ops 2020
 - ▶ SKA₂: Construction start in \sim 2023
- ▶ Will require ‘exascale’ computing to form images and analyse them

Karoo Desert, South Africa



Western Australia



About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary

Within the SKA project, Information and Communications Technology (ICT) is:

- ▶ Major part of the design budget
- ▶ Major part of the construction budget
- ▶ Major part of the operations budget (power, S/W maintenance)
- ▶ Major part of the **risk** 'budget'
 - ▶ Strong interaction with industry, other astronomy projects, other science project
 - ▶ Direct involvement of major existing HPC facilities
 - ▶ Reusing, Road-mapping, Best-practices, Sustainability

About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary

About the Square Kilometre Array

Science Data Processor

Architectures, Technologies

Summary

Wide area dataflow for SKA

SKA SDP

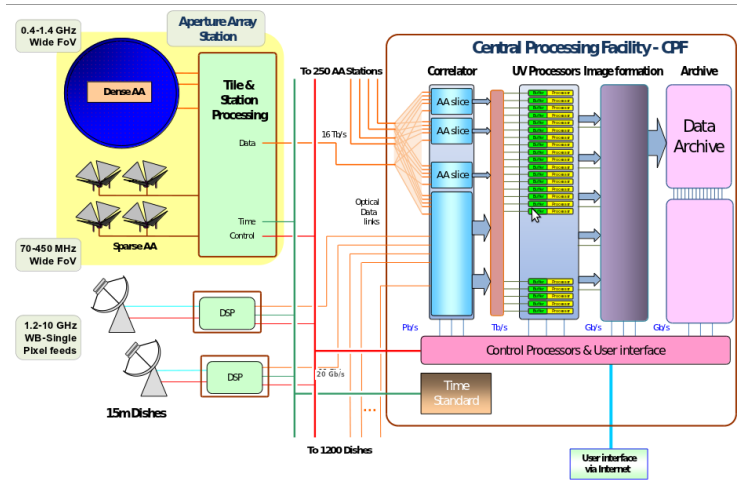
B. Nikolic / SDP
Team

About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary



SKA simplified dataflow

SKA SDP

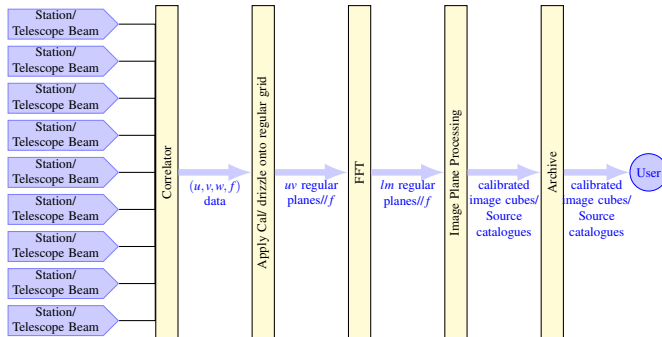
B. Nikolic / SDP
Team

About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary



What drives the challenge?

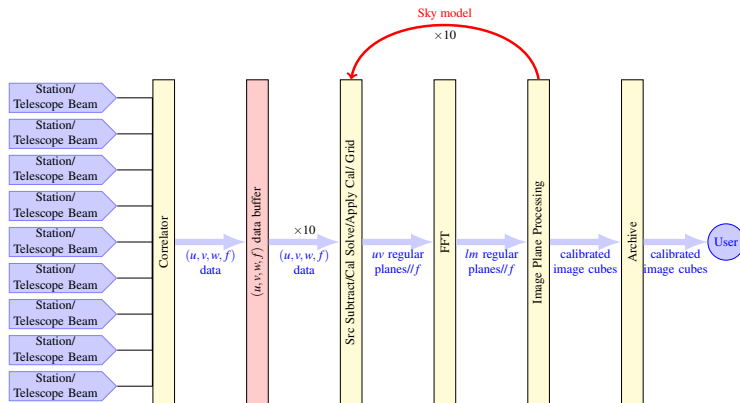
Complex to do

- ▶ Irregular 3d sampling of the signal
⇒ grid & expensive 3d correction
- ▶ Image reconstruction (CLEANing)
⇒ identify and remove 'bright' sources
- ▶ Changing electrical properties of telescope **and the atmosphere**
⇒ self-calibrate



Computationally intensive

- ▶ Very high data rates



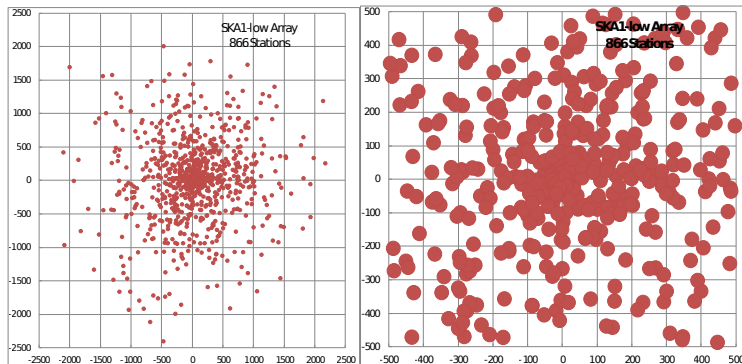
- ▶ # visibilities scales as N^2
- ▶ Very large N :
 - ▶ $N \sim 900$ for SKA₁-Low
 - ▶ $N \sim 250$ for SKA₂-Mid

About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary



SKA₁-Low Core stations configuration. Credit: P. Dewney, SKA1 Baseline Design/RFP

Data rate – multiple beams

SKA SDP

B. Nikolic / SDP
Team

About the Square
Kilometre Array

Science Data
Processor

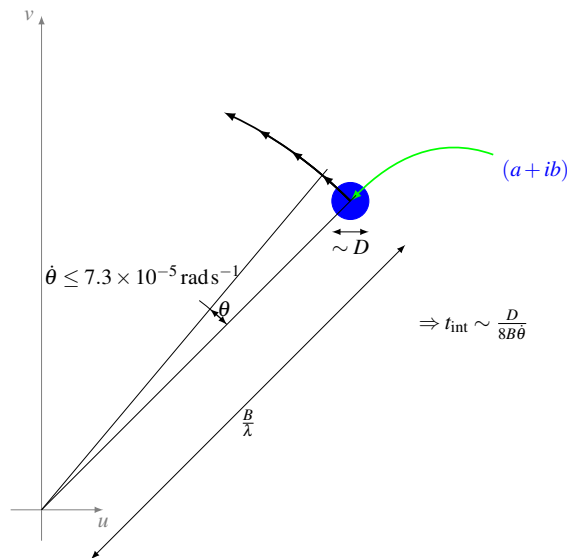
Architectures,
Technologies

Summary

- ▶ # visibilities scales as # beams
- ▶ SKA₁-Survey will have 36 simultaneous beams



Data Rate – time smearing



Data Rate – bandwidth smearing

SKA SDP

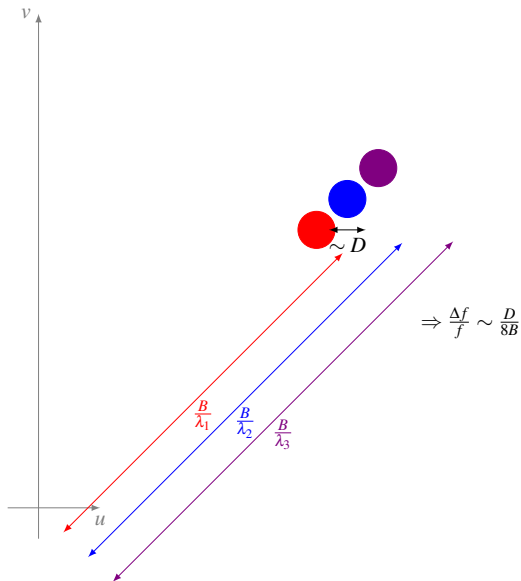
B. Nikolic / SDP
Team

About the Square
Kilometre Array

Science Data
Processor

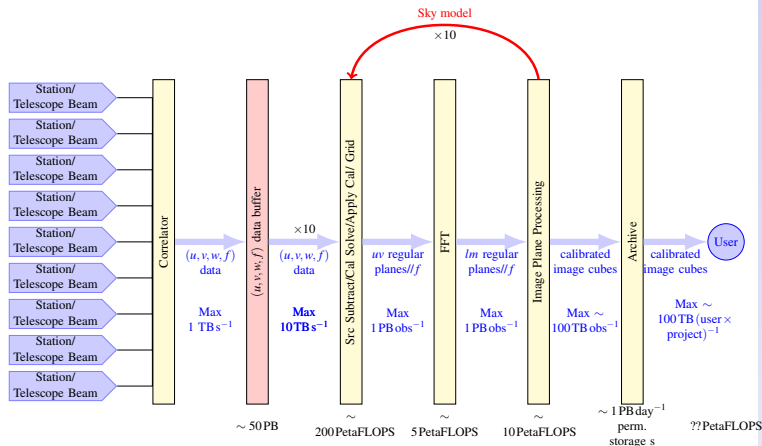
Architectures,
Technologies

Summary



- ▶ Data volume most demanding for 'high resolution spectral line survey' observation:
 - ▶ Roughly $32k \times 32k \times 32k$ cube
 $\implies \sim 200$ TB/obs
 - ▶ Very low re-observation cadence
 - ▶ Potentially useful for cosmology \rightarrow many fields?
- ▶ Other observations typically coarser spatial and/or frequency resolution \rightarrow TB scale datasets
- ▶ However some of the other observations could have high re-observation cadence
- ▶ **Transient Imaging pipeline for selectively retaining data**

Data flow with rates, computes and storage requirements



About the Square Kilometre Array

Science Data Processor

Architectures, Technologies

Summary

Rapid progress, often driven by consumer technologies:

- ▶ > 10 TeraFLOPS/ accelerator card likely in 2017
 - ▶ 100 PetaFLOPS $\implies 10^4$ cards, ~ 2 MW
 - ▶ At 2 cards/per node, $> 50 \text{ TB s}^{-1}$ total inter-node network capacity easily achieved
- ▶ $> 1 \text{ TB s}^{-1}$ total network throughput possible today (600 port Infiniband FDR)
- ▶ Nodes with $2 \times \text{FDR} + 2 \times 40 \text{ GbE}$ already being deployed
 - $\implies 10/\text{GB s}^{-1}$ onto one node
- ▶ 40 kW/rack with water-cooled cabinet doors
 - ▶ 100 PetaFLOPS, 2 MW $\implies 50$ racks

High-level dataflow/possible architecture

SKA SDP

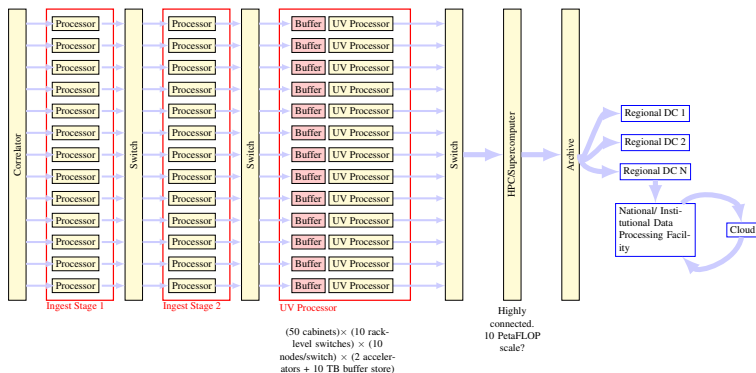
B. Nikolic / SDP
Team

About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary



Data Distribution, Remote Compute & Visualisation

SKA SDP

B. Nikolic / SDP
Team

- ▶ Tiered data distribution allows use of existing/forthcoming national compute facilities
- ▶ Remote (institutional/national/international) visualisation and compute
- ▶ Efficient collaboration for large science teams essential

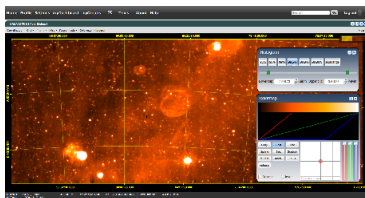
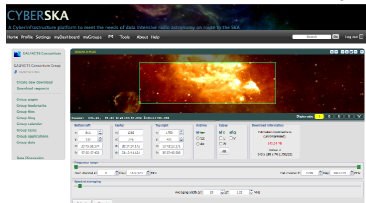
About the Square
Kilometre Array

Science Data
Processor

Architectures,
Technologies

Summary

CyberSKA



- ▶ A large software engineering project: **hybrid** between conventional HPC, 'big data' and streaming processing
- ▶ Need both flexibility and the ability to operate autonomously (without human intervention)
- ▶ Data flow and data-locality critical
- ▶ Few *iterations*, mostly data-parallel
- ▶ Streaming / soft-realtime processing
- ▶ **Evolving underlying hardware**

About the Square Kilometre Array

Science Data Processor

Architectures, Technologies

Summary

- ▶ Very large data rates in SDP:
 - ▶ Internal rates $\sim 1 - 10 \text{ TB s}^{-1}$
 - ▶ Into archive $\sim 0.1 - 1 \text{ PB day}^{-1}$
- ▶ Complex, varied, processing *before* data reduction
- ▶ SDP cost analysed and balanced against overall SKA system
- ▶ Commercial COTS H/W roadmaps look promising
- ▶ Challenges:
 - ▶ S/W design, construction, maintenance
 - ▶ System complexity
 - ▶ H/W & S/W failures
 - ▶ Schedule
 - ▶ Science Analysis on large data-cubes