# DISTRIBUTED DATA PROCESSING USING SPARK IN RADIO ASTRONOMY



Panos Labropoulos
Bright Computing, Inc.

Sarod Yatawatta ASTRON

#### Who are we?





- Panos Labropoulos
  - Works at Bright Computing
  - ► HPC, accelerators, high-speed interconnects
  - ▶ PhD Radio Astronomy, University of Groningen

- Sarod Yatawatta
- Works at ASTRON
- Signal and Image processing
- ▶ PhD Electrical engineering, Drexel





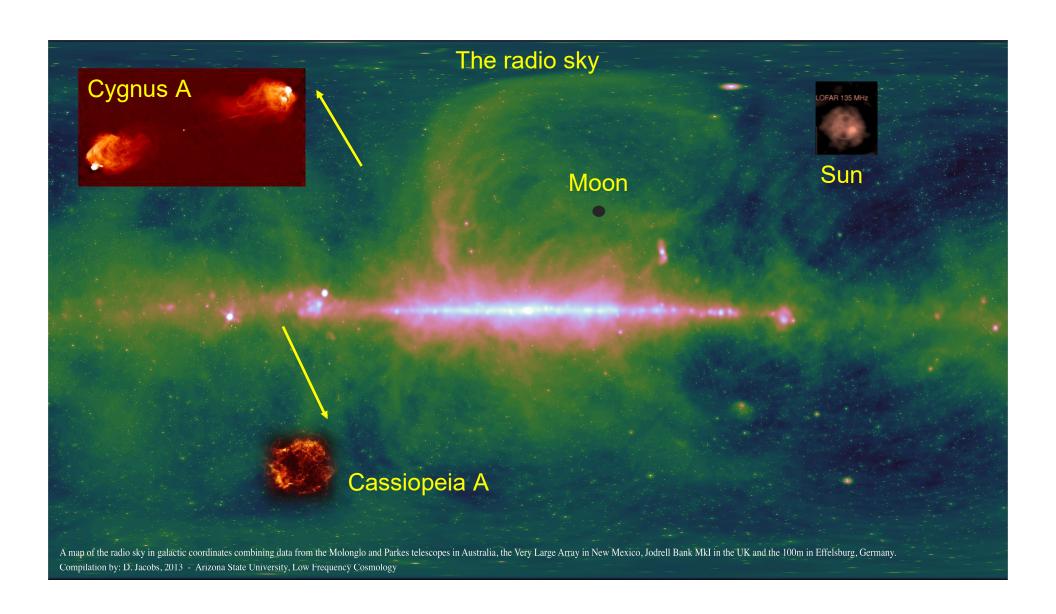


# Outline

- ▶ Basics of radio interferometry
- Scientific motivation
- Why is Spark relevant?
- Two case studies: calibration and imaging

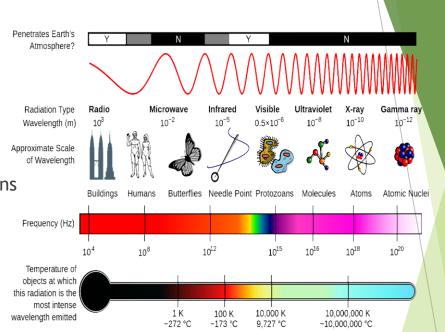




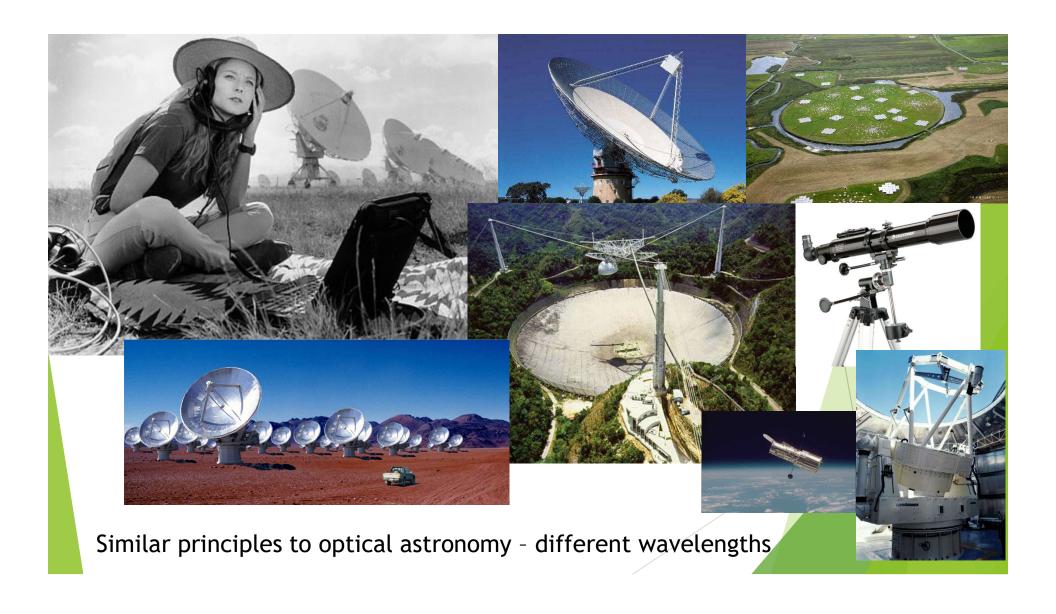


# Why radio?

- Thermal Radiation
- Synchrotron Radiation
  - ► Relativistic e<sup>-</sup> in magnetic fields
- Bremstrahlung
  - ▶ "Breaking Radiation" e-/ion collisions
- Maser
  - Microwave Laser e<sup>-</sup> oscillations in molecular clouds
- Atomic Transitions (emission spectra)
  - ► Hydrogen e<sup>-</sup> spin flip



Tools required: antennas, receivers (e.g. voltage amplifiers, ADCs, computers)

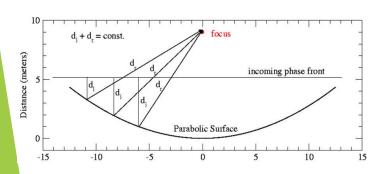


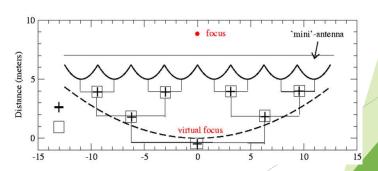
# Basic characteristics of a telescope

#### Angular resolution and sensitivity proportional to size

- ▶ Oldest telescope: human eye
  - $\triangleright$  Human pupil: 0.5 cm diameter, 0.05 deg. Resolution at  $\lambda$ =600 nm (yellow)
- $\triangleright$  Radio telescope: 100m diameter, 1.4 deg. resolution at  $\lambda$ =2m
  - ▶ 2.8 km diameter required for same resolution as human eye
  - ▶ 250 km dimeter required for same resolution as an optical telescope

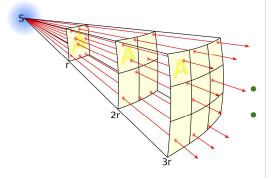
#### Constructing a 250km parabolic dish is simply not feasible





Use smaller antennas and add-up signals using a computer

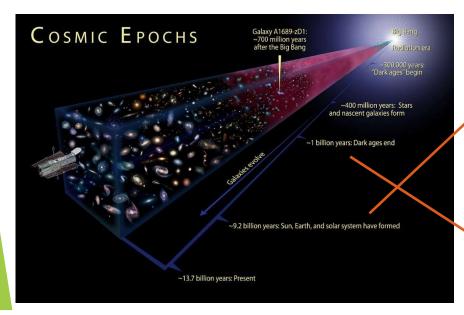
### How many antennas? How big?



#### Intensity follows an inverse square law

Source 3 times further away become 9 times fainter Size of telescope depends on distance at which we wan t to look objects at.

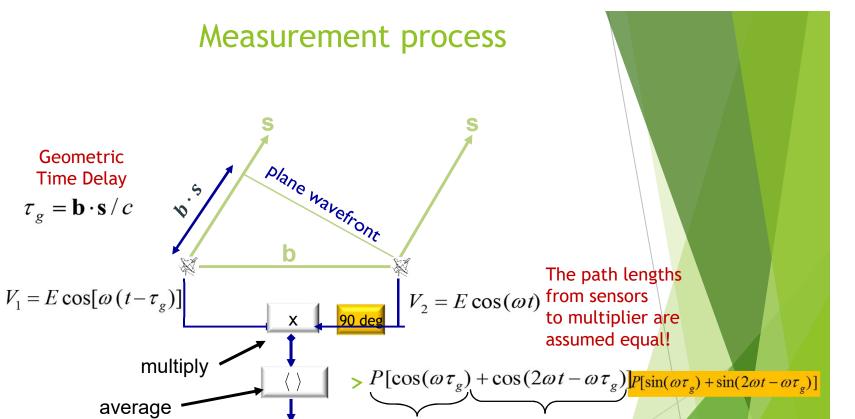
10.000 sq. m. VLA





10x further -> 100x bigger 1 sq. Km





Rapidly varying,

$$V_{v}(\mathbf{b}) = R_{C} - iR_{S} = \iint I_{v}(s)e^{-2\pi\lambda i \,\mathbf{b}\cdot\mathbf{s}} d\Omega$$
 2-D FT of the sky's brightness distribution

 $R_C = P\cos(\omega \tau_g)$  Unchanging

 $R_s = P\sin(\omega \tau_g)$ 

Geometric

Time Delay

 $\tau_g = \mathbf{b} \cdot \mathbf{s} / c$ 

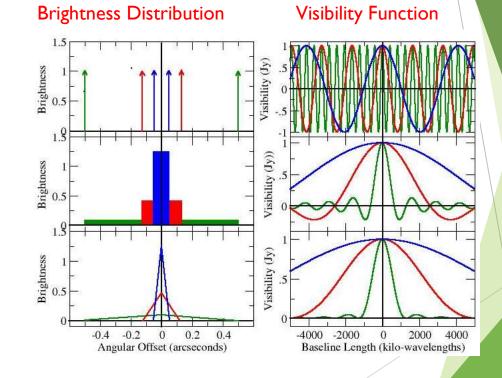
#### Examples of 1-Dimensional Visibilities

Simple pictures are easy to make illustrating I-dimensional visibilities.



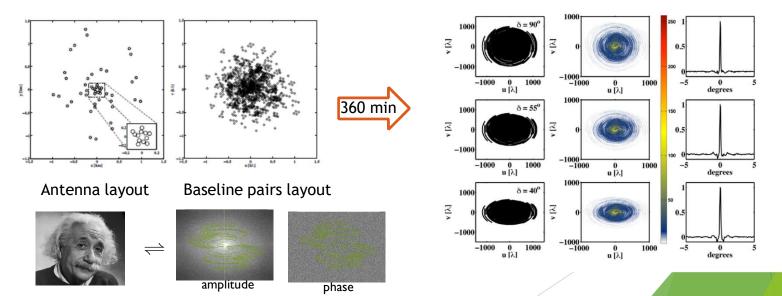






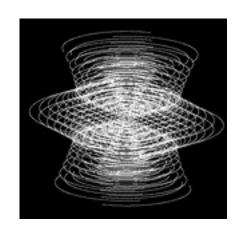
## Earth rotation synthesis

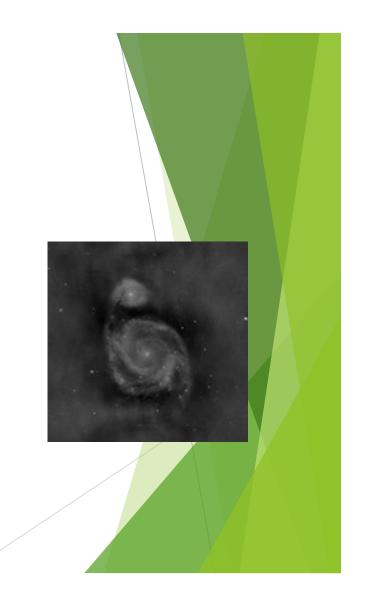
- ▶ An interferometer with N antennas has N (N -1) / 2 "baseline" pairs
- Real interferometers are built on the surface of the earth a rotating platform. From the observer's perspective, sources move across the sky, but from the source's perspective the antenna array is rotating in such a way that the baseline pair points form tracks over time
- u-v tracks have gaps: incomplete sampling of the 2-D transform / missing information



# Example: synthesis image

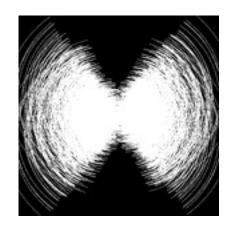






# Example: synthesis image

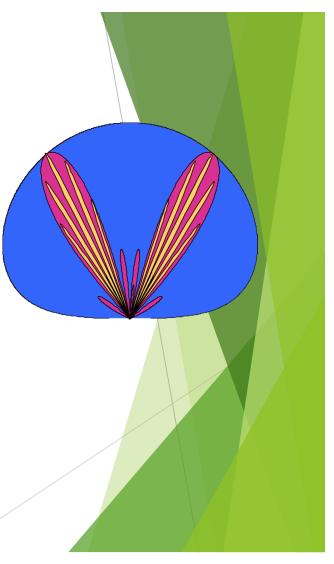


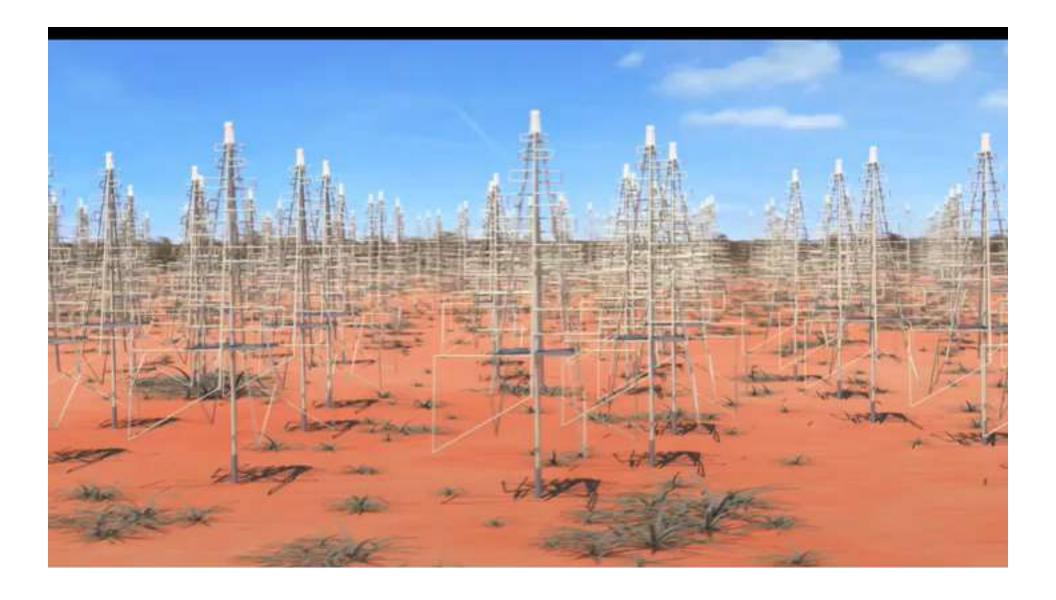






- ▶ 100x more sensitive, 1,000,000x faster imaging
- ▶ Up to 5 sq.km collecting area and 3500km baselines
- ► Sub-arcsecond resolution, large field-of-view
- ▶ 3 types of antennas, 2 locations, a single observatory
- ▶ No mechanical steering for low- and mid-frequency antennas, multiple beams
- ▶ Will address key questions of astronomy, cosmology and physics
- Ambitious IT project. Probably will be the first production exascale system
- ▶ Will be generating several exabytes of processed data per year
- ▶ Builds up on techniques developed in Europe
- Major HPC facility
- ▶ UK, RSA, AUS, NZ, Canada, China, NL, Germany, Italy, India, Sweden collaboration







# Five Key Science Areas for the SKA

Topic	Goals						
Probing the Dark Ages	Map out structure formation using HI from the era of reionization (6 < z < 13)     Probe early star formation using high-z CO     Detect the first light-emitting sources						
Gravity: Pulsars & Black Holes	Precision timing of pulsars to test theories of gravity approaching the strong-field limit (NS-NS, NS-BH binaries, incl Sgr A*)      Millisecond pulsar timing array for detecting longwavelength gravitational waves						
Cosmic Structure	Understand dark energy [e.g. eqn. of state; W(z)]     Understand structure formation and galaxy						
Cosmic Magnetism	Determine the structure and origins of cosmic magnetic fields (in galaxies and in the intergalactic medium) vs. redshift z						
The Cradle of Life	Understand the formation of Earth-like planets     Understand the chemistry of organic molecules and their roles in planet formation and generation of life     Detect signals from ET						

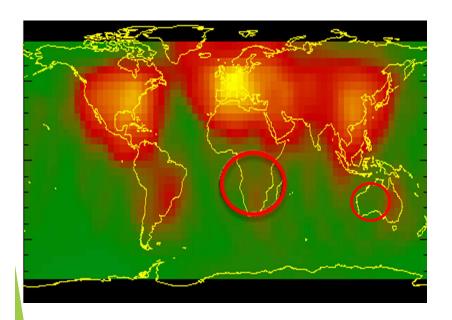
								\						
KSP ID	KSP Description	Frequency Range GHz						FoV	Sens- itivity	Survey Speed	Resn.	Base- line	Dyn. Range Driver	Poln. Driver
		0.1	0.3	1.0	3.0	10	30	deg <sup>2</sup>	m²/K	deg <sup>4</sup> m <sup>4</sup> K <sup>4</sup>	mas	Km		
1	The Dark Ages													
1a <sup>†</sup>	EoR	+								>~3x10 <sup>7</sup>		1	*	**
1b	First Metals					-	+	0.003	15,000		50	125		
1c	First Galaxies & BHs			+		ļ			20,000		10	4500	*	**
2	Galaxy Evolution, Cosmology & Dark Energy													
2a <sup>†</sup>	Dark Energy			+						6x10 <sup>9</sup>		5		
2b <sup>†</sup>	Galaxy Evolution		+	+					20,000	1x10 <sup>9</sup>		10		
2c	Local Cosmic Web		Ш	-						2x10 <sup>7</sup>		0.5		
3	Cosmic Magnetism		П											
3a <sup>†</sup>	Rotation Measure Sky			H						2x10 <sup>8</sup>		10-30		**
3b	Cosmic Web	+	+	+						1x10 <sup>8</sup>		5		**
4	GR using Pulsars & Black Holes													
	Search			-						1×10 <sup>8</sup>		< 1		
4a <sup>†</sup>	Gravitational Waves		1	+	-	-		-	>15,000		1	200		**
4b	BH Spin			+		-		1	10,000					**
4c <sup>†</sup>	Theories of Gravity					-			>15,000		1	200		**
5	Cradle of Life					1								
5a <sup>†</sup>	Protoplanetary Disks						-	<10 <sup>-5</sup>	10,000		2	1000		
5b	Prebiotic Molecules			_		+	+	0.5-1	10,000		100	60		
5c	SETI		Ш	-		-		1						
6	Exploration of the Unknown					1	-	Large	Large	Large				

† Headline science, see Section 3.2

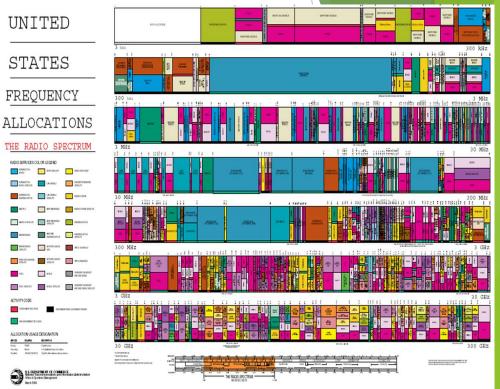
Afterglow Light Pattern 400,000 yrs. Inflation 1st Stars about 400 million yrs. Dark Ages Big Bang Expansion
13.7 billion years Development of Galaxies, Planets, etc. Dark Energy Accelerated Expansion

NASA/WMAP Science Te

#### Site selection: determined by man-made interference



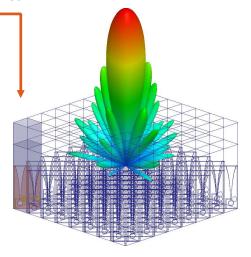
The majority of the radio spectrum has been allocated

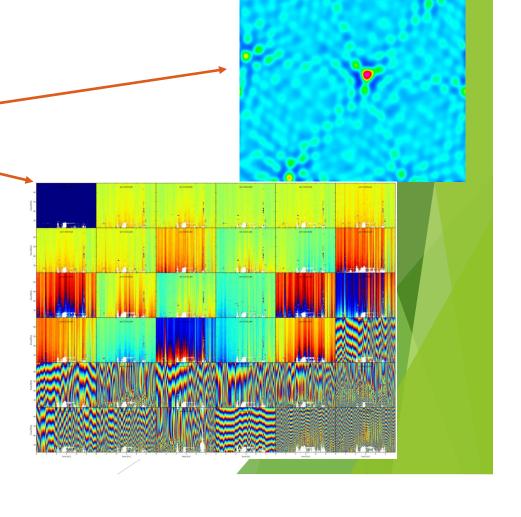


Only available radio-quit sites: south Africa and western Australia (or the back side of the Moon)

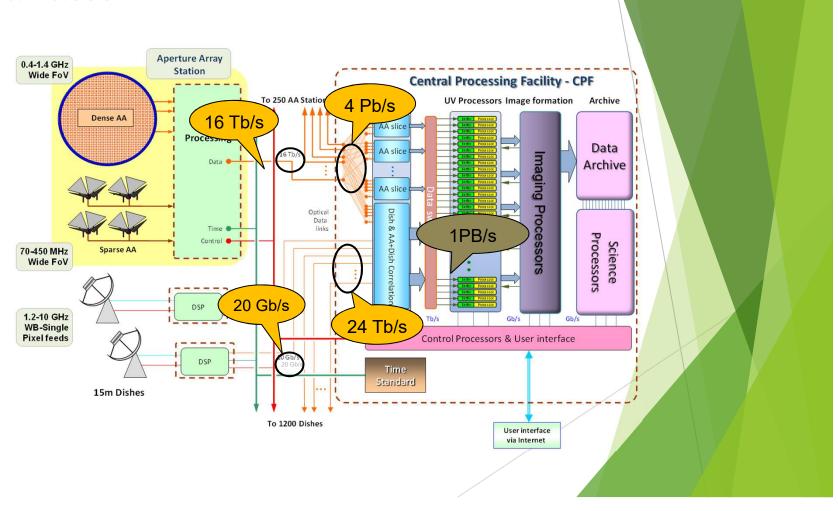


- Astronomical radio signals are very weak
- ▶ Calibration is the name of the game
- Several sources of corruption
  - ▶ Ionospheric Faraday Rotation/phase
  - Antenna/station beam patterns
  - Receiver gain and phase errors
- If not tackled, they will limited the sensitivity and therefore the science

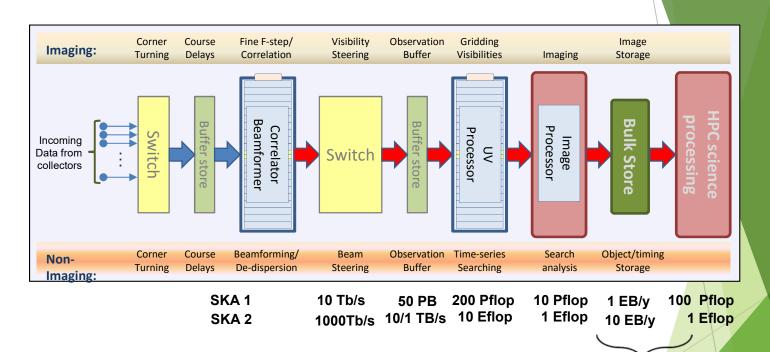




## SKA data rates



# Science data processor pipeline

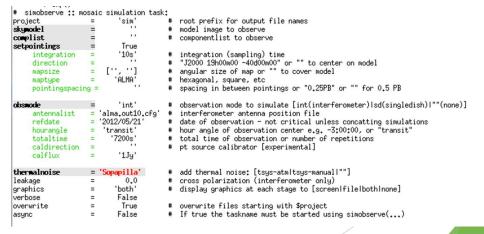


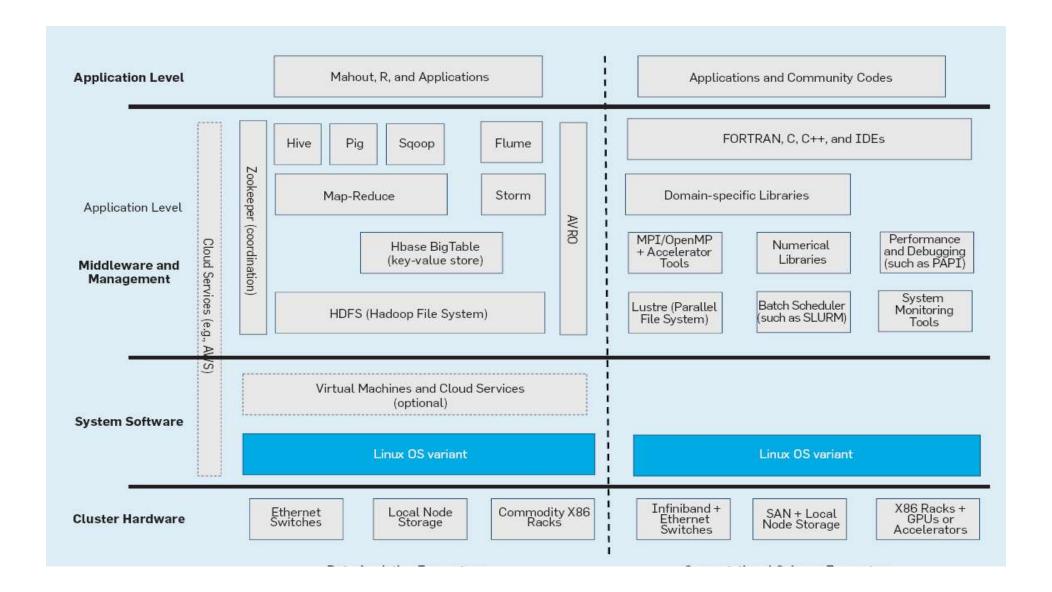
Imaging and calibration determines problem size

30x more computationally intensive (deconvolution, gridding, prediction, FFTs) 90% of power consumption

### Current data reduction scheme

- Software developed in the 80s/90s
- Fortran libraries
- Minimal use of multi-cores/No use of accelerators
- Design for interactive processing of single dataset/Support for scripted batched processing
- ipython/custom REPL
- Will not work for SKA





# What is Spark?

Fast and Expressive Cluster Computing System Compatible with Apache Hadoop

Up to **10** × faster on disk, **100** × in memory

# **Efficient**

- General execution graphs
- ►In-memory storage



# **Usable**

- Rich APIs in Java, Scala, Python
- Interactive shell

# **Key Concepts**

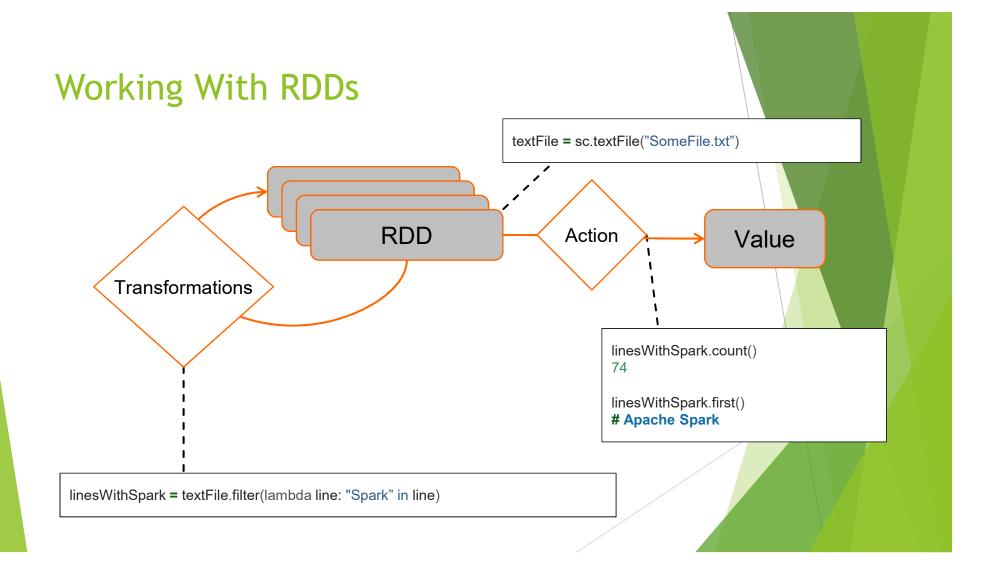
Write programs in terms of transformations distributed datasets

# Resilient Distributed Datasets

- Collections of objects spread across a cluster, stored in RAM or on Disk
- Built through parallel transformations
- Automatically rebuilt on failure

# **Operations**

- Transformations (e.g. map, filter, groupBy)
- Actions (e.g. count, collect, save)



# Example: Log Mining

Load error messages from a log into memo then interactively search for various patter

Cache 1

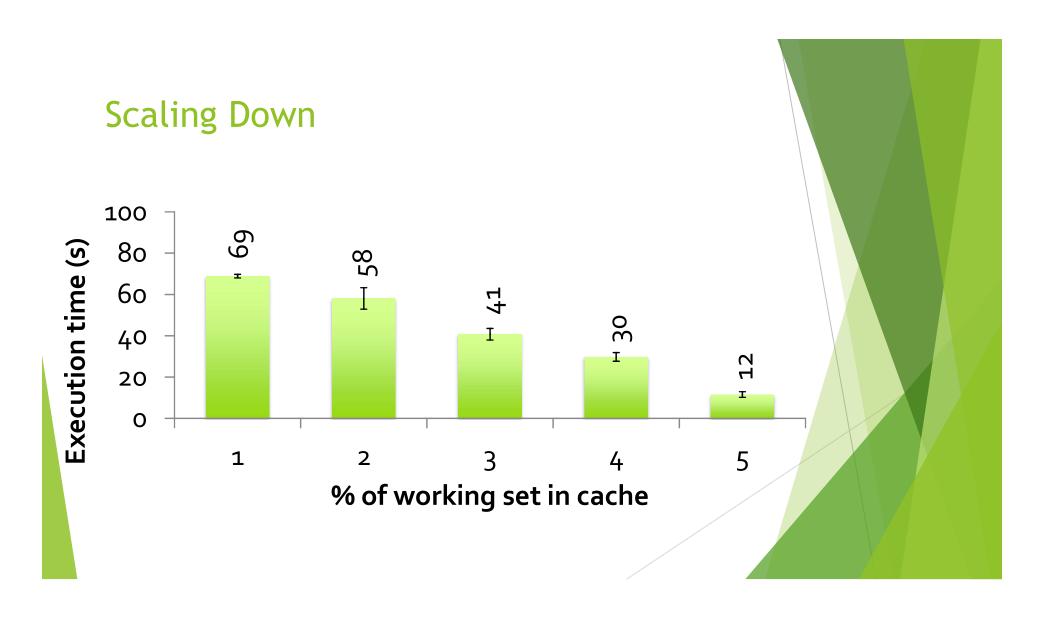
Cache 2

```
Base RDD
              Transformed RDD
lines = srark.textFile("hdfs://...")
                                                                      Worker
                                                             results
errors = lines.filter(lambda s: s.startswith("ERROR"))
                                                                tasks
                                                                      Block 1
messages = errors.map(lambda s: s.split("\t")[2])
                                                      Driver
messages.cache()
                                                     Action
messages.filter(lambda s: "mysql" in s).count()
                                                                     Worker
messages.filter(lambda s: "php" in s).count()
                                                     Cache 3
                                                   Worker
         Full-text search of Wikipedia

    60GB on 20 EC2 machine
```

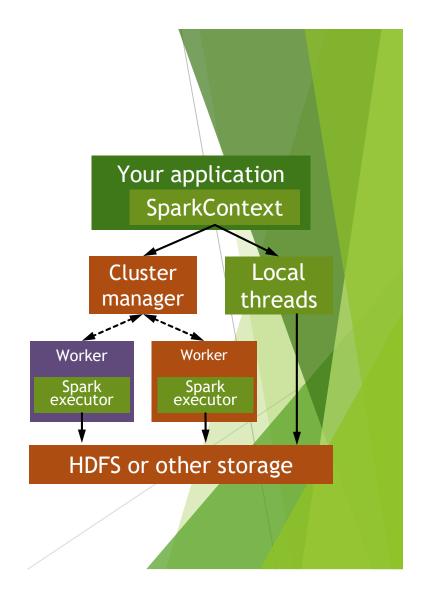
Block 3

- 0.5 sec vs. 20s for on-disk



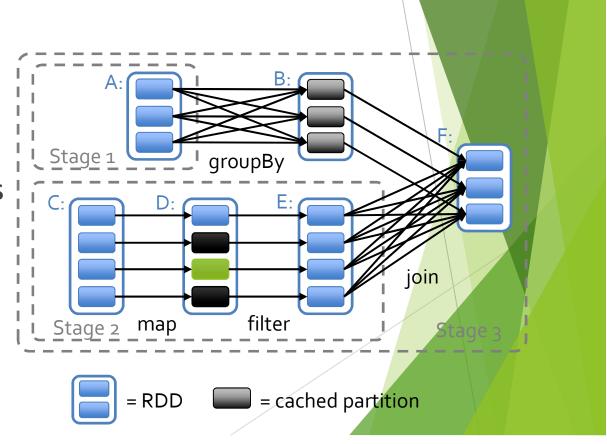
# **Software Components**

- Spark runs as a library in your program (1 instance per app)
- Runs tasks locally or on cluster
  - ▶ Mesos, YARN or standalone mode
- Accesses storage systems via Hadoop InputFormat API
  - ► Can use HBase, HDFS, S3, ...



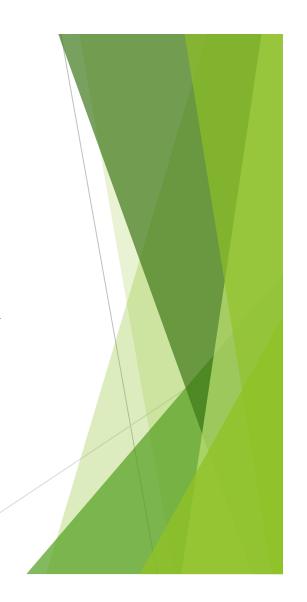
## Task Scheduler

- General task graphs
- Automatically pipelines functions
- Data locality aware
- Partitioning aware to avoid shuffles



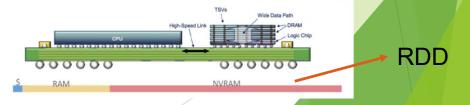
## **Advanced Features**

- ▶ Controllable partitioning
  - ▶ Speed up joins against a dataset
- ► Controllable storage formats
  - ▶ Keep data serialized for efficiency, replicate to multiple nodes, cache on disk
- ► Shared variables: broadcasts, accumulators
- See online docs for details!

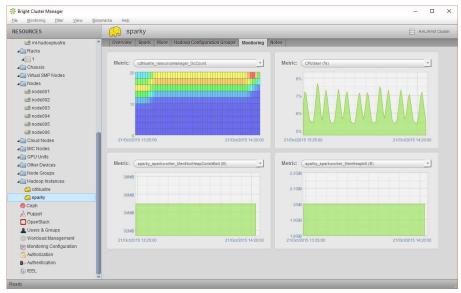


## Why Spark

- Power consumption will be the most important differentiator -minimize movement
  - ▶ supercomputer in the middle of the desert
  - ▶ data locality automatic work placement
  - ► In-memory processing
- Iterative algorithms / Pipelines that enable piggybacking
- New algorithms that have modest inter-process communication
  - ▶ Not all features of MPI are required -> Simple collectives: broadcast, aggregation
  - ▶ Can MPI handle the increased pararellism available on exascale systems?
- Fault tolerance
  - ► Current MTBF: 2 out of 600 nodes per week -> 300 nodes per week for SKA
  - ▶ MPI applications need to be designed for fault-tolerance/checkpointing
- Straggler mitigation
- Friendly APIs

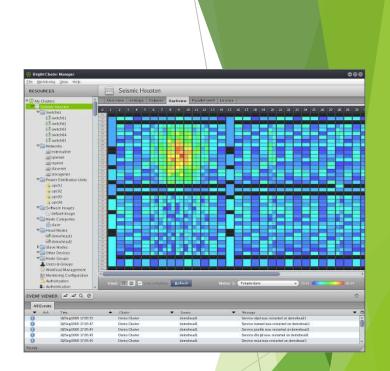


## Fault Tolerance





- ► HDFS/Lustre/Ceph support
- ▶ Early warning via e-mail, SMS etc and automate actions
- ▶ Single plane of glass: monitor nodes, switches, PDUs, multiple clusters
- ► Hadoop/Spark-specific metrics and healthchecks
- ► Simple Python/REST/C++ APIs

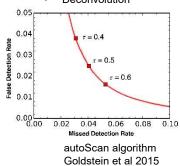


## How is data processed

- o Datasets contain multiple sources of information
- · Ionosphere, lightning, cosmic rays, ET?, interference, compact and extended astronomical sources
- · Different algorithms for each
- o Multidimensional data: spatial, time, frequency, polarization
- Pipeline should be able to chain different algorithms in order to extract as much information as possible with minimal data movement

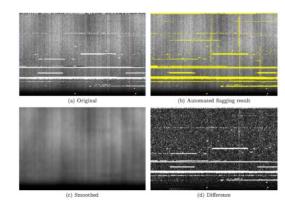
#### Types of processing:

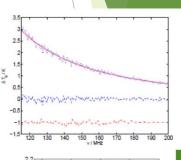
- Signal detection/extraction
- Object classification
- AGN spectral classification
- Time series
- Clustering
- LSS
- Class discovery
- RFI mitigation
- Transient detection
- Deconvolution

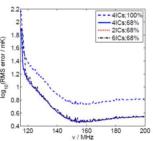


#### Types of algorithms:

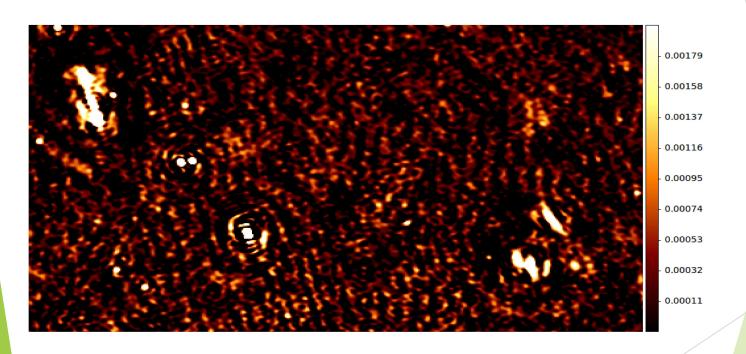
- ANNs (this talk)
- Decision trees
- SVM
- · Nearest neighbor
- Expectation maximization
- Non-linear optimization

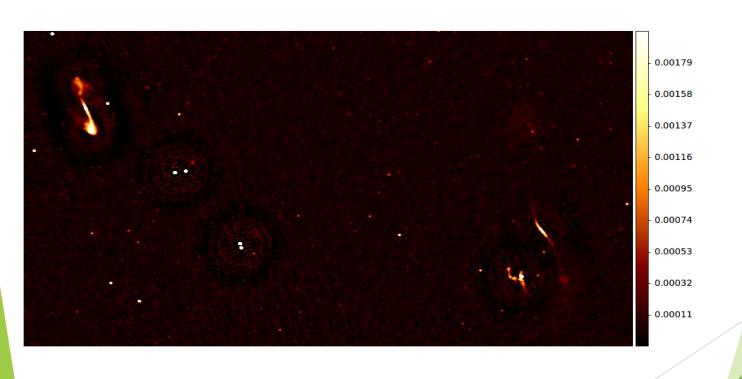




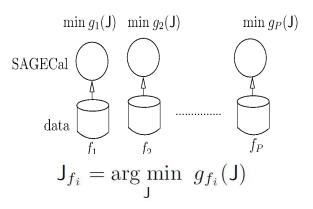


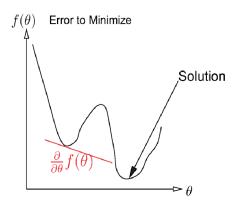
Chapman et al 2012



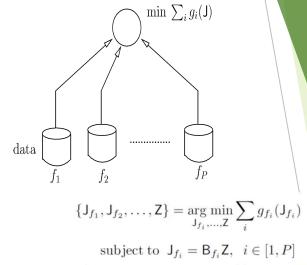


Data distributed over a network: all existing calibration algorithms work independently



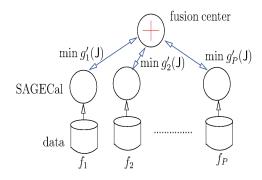


We want a unified solution, collecting information from all data



Where  $B_{f_i}Z$  is a smoothness constraint along time, space and frequency

Does not work in practice: data does not fit in memory, model not accurate enough to parametrize



$$L(\mathbf{J}_{f_1},\ldots,\mathbf{Z},\mathbf{Y}_{f_1},\ldots) = \sum_i g_{f_i}(\mathbf{J}_{f_i}) + \|\mathbf{Y}_{f_i}^H(\mathbf{J}_{f_i} - \mathbf{B}_{f_i}\mathbf{Z})\| + \frac{\rho}{2}\|\mathbf{J}_{f_i} - \mathbf{B}_{f_i}\mathbf{Z}\|^2$$

- Consensus optimization: exploit smoothness of systematic errors to find a unified solution
- Use of a fusion center is not essential, but easier
- Communication overhead is little, more important is robustness and fault tolerance (can throw away lost data).
- Large rho makes problem convex

#### 1. INITIALIZE

- Read data as Breeze BlockMatrix
- P sub-problems, each with its own copy of Z,Y
- 2. REPEAT until convergence or N\_max\_iter
  - 1. Each worker solves:

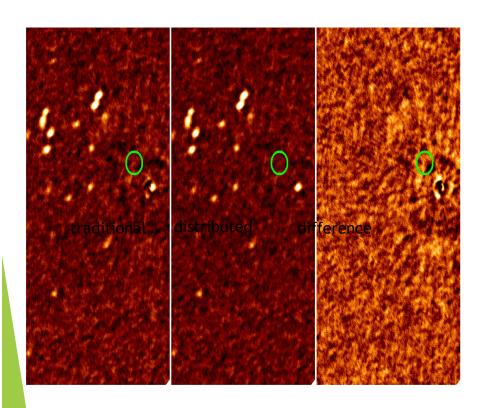
$$(\mathsf{J}_{f_i})^{n+1} = \arg\min_{\mathsf{J}} \, L_i \, (\mathsf{J}, (\mathsf{Z})^n, (\mathsf{Y}_{f_i})^n)$$

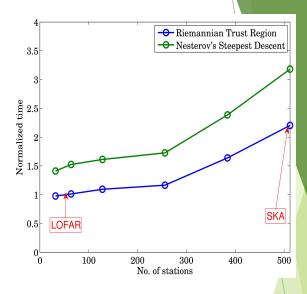
2. Globally calculate average:

$$(\mathsf{Z})^{n+1} = \arg\min_{\mathsf{Z}} \sum_{i} L_{i} \left( (\mathsf{J}_{f_{i}})^{n+1}, \mathsf{Z}, (\mathsf{Y}_{f_{i}})^{n} \right)$$

- 3. Broadcast Z to P sub-problems Zb = sc.broadcast(Z)
- 4. Locally calculate Lagrange multiplier

$$(Y_{f_i})^{n+1} = (Y_{f_i})^n + \rho((J_{f_i})^{n+1} - B_{f_i}(Z)^{n+1})$$





- Almost linear caling with number of antennas
- Broadcasted data is 10x less (4KNP) than actual data  $(4TP^{\frac{N(N-1)}{2}})$
- Can be further reduced with:
- Use frequency multiplexing
- Broadcast data only to neighbors

#### Case study: Imaging

#### 1. INITIALIZE

- a residual image set to the Fourier Transforr of the visibilities
- var row: Int = \_, var col: Int = \_
- (row , col, flux) MAX= tuple(l: Int, m: int)
- a Clean Component list to empty

#### 2. WHILE ( img.getNoise < THERMAL\_NOISE )

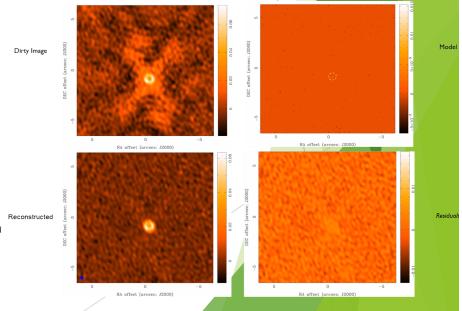
- 1. (row, col, flux) = img.getMax
- 2. resimg = resimg g \* PSF.imshift(l,m,size)
- 3. CCList.++ = (row, col, flux)
  - linear scaling
  - fits in memory
  - easy to implement
  - gridding/FFT can be offloaded to GPU via JNI
  - Does not work well with extended sources, still useful for catalogs a nd calibration

Poor man's Compressed Sensing (matching pursuit)

#### Variations:

- •Clarke: subtract multiple components in on iteration (~ StOMP)
- Cotton-Schwab: Work in measurement space instead of image space

In the absence of noise, CLEAN is equivalent to a least squares fit of sinusoids to visibilities



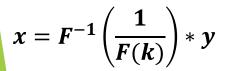
#### **Imaging**

The "dirty" image is the convolution of the true underlying image with the PSF, corrupted by Gaussian random noise

$$y = k * x + n$$

Convolution Theorem

$$F(y) = F(k) F(x)$$



This pseudo-inverse cannot be computed in the presenc e of corruptions

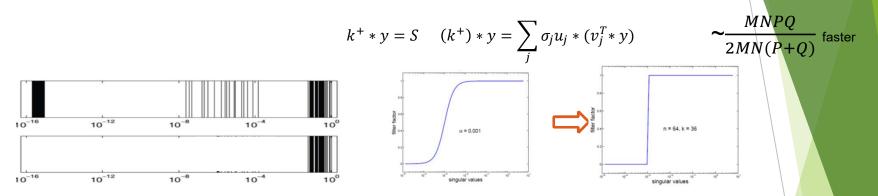
$$x = F^{-1}\left(\frac{1}{F(k)}\right) * y \qquad x = F^{-1}\left(\frac{1}{F(k)}\frac{|F(k)|^2}{|F(k)|^2 + SNR^{-1}}\right) * y$$

Tikhonov Regularization High noise strong → regularizatio

$$x = \operatorname{argmin} ||y - x * k|| + R(x) + \cdots$$

#### **Imaging**

The previous kernel can be transformed as a sum of 1D filters:



Goal: use the following feed-forward, multi-layer convolutional neural network and train it to minimize

$$h_{l} = \begin{cases} \hat{y}, l = 0 \\ \tanh(W_{l} * h_{l-1} + b_{l-1}), l \in \{1,2\} \\ W_{3} * h_{2} \end{cases}$$

Initialize W from  $k^+$  or from random, uniform distribution 43

Astrophysical sources are complication, but most radio sources have simple shapes, and occasionally irregular:





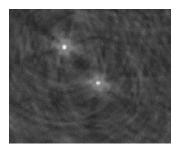




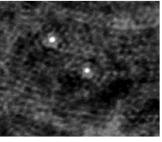
Use high SNR images from the VLA FIRST survey supplemented with simulated sources (elliptical Gaussians, points sources, narrow gaussians) to construct a training pair set of 10.000 64x64 images



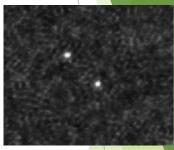
uncorrupted



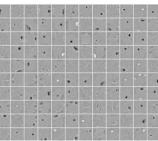
convolved

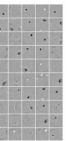


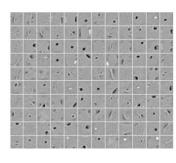
noise added

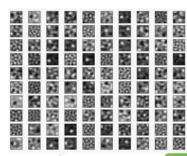


recovered









#### Challenges

Spark is feature already but some useful components are not yet available



- JNI
- Analytical implementations needed
- ▶ Hard to debug and argue about performance sometimes
- Application specific formats
- Accelerators such as GPUs are almost mandatory JNI

### **Conclusions**

- New generation interferometers will enable us
- The availability of HPC (exascale) resources is a prerequisite
- ▶ Energy cost will be the decisive factor when it comes to design
- New hardware architectures require new algorithmic approaches
- Is Spark the answer? To early to answer but maybe
  - Distributed algorithms not only address computational issues, but are more robust and accurate as well
  - Consensus calibration is the best way forward
  - Spark offers many of the tools need to construct radio astronomical pipelines