

LOFAR-IT HPC Activities

A distributed e-Infrastructure for data analysis

Ugo Becciani – Giuliano Taffoni

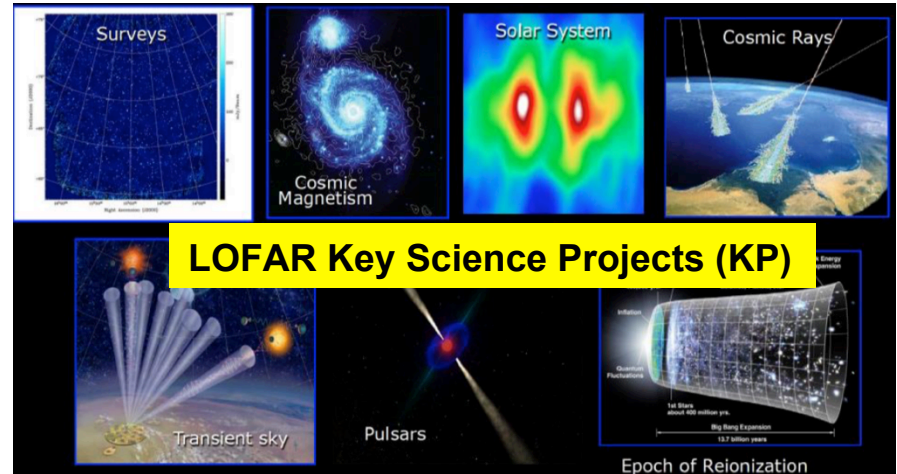


- The LOFAR.IT organization
- Italian participation to LOFAR
- The LOFAR Computing Model
- LOFAR-IT e-infrastructure:
 - Data analysis Requirements
 - Design and setup of the Data Analysis Infrastructure



The International LOFAR Telescope

ILT consists of an interferometric array of dipole antenna stations, and it is distributed throughout 9 EU Countries: NL, Germany, France, Italy, Poland, UK, Sweden, Ireland, Latvia.



Board: **Gianfranco Brunetti (Coordinator INAF-IRA)** **Ugo Becciani (INAF-OA Catania)** Segretario, **Federica Govoni (INAF-OA Cagliari, UTGII)**, **Francesco Massaro (UniTo)**, **Jader Monari (INAF-IRA)**, **Roberto Scaramella (INAF-OA Roma)**

Science Advisory Committee: **Andrea Ferrara (Chair)**, Matteo Murgia, Mauro Messerotti, Grazia Umana, Gianni Bernardi, Ettore Carretti, Isabella Prandoni, Laura Pentericci, Marta Burgay, Rossella Cassano, Andrea Chiavassa (UniTO)

Technological joint WG ASTRON-INAF: Established on March 2018 . Led by Astron. Primary objective: joint development of RCU for LOFAR 2.0 and eventually LBA2.0.

Data WG: **Giuliano Taffoni (INAF-OA Trieste)** - Chair, Alessandro Costa (INAF-OA Catania), Francesco Bedosti (INAF-IRA), Cristina Knapic (INAF-OA Trieste), Manuela Magliocchetti (INAF-IAPS Roma), Annalisa Bonafede (UniBo, Associata INAF IRA)

Main activities of the Consortium

- Build a LOFAR 2.0 station in Medicina (2020-2022)
- Build a LOFAR data analysis infrastructure
- Implement a technological and scientific collaboration with ASTRON
- Develop a community that is able to work with LOFAR data (for science and technology)
- Participation of Italian community to Key Projects (surveys in particular)

Participation to the KPs, LOFAR guarantee time

Close 30 Aug 2018

10 Proposal

Summary and people Involved

- KP SURVEYS => 5 + 17 + 1 TD new tot = 18
- KP EOR => 0 + 4 + 0 TD
- KP TRANSIENT => 0 + 5 + 2 TD
- KP MAG => 1 + 1 + 1 TD
- KP SUN => 0 + 6 + 0 TD

Observation Time obtained : Italian Proposals

The following Single-Cycle proposals were reviewed

Cycle 11

Four Proposal approved
by ILT TAC

Observation Time
obtained : **55 hours**

LC11_001
LC11_002
LC11_004
LC11_005
LC11_006
LC11_007
LC11_008
LC11_009
LC11_010
LC11_011
LC11_012
LC11_013
LC11_014
LC11_015
LC11_016
LC11_017
LC11_018
LC11_019
LC11_020
LC11_021

P. Zucca	72.0	46.0
A. Rowlinson	44.0	17.4
S. Purser	32.0	32.0
P. Leto	13.1	13.1
H. Vedantham	32.0	0.0
O. Wucknitz	9.0	9.0
A. Merloni	104.0	0.0
V. Cuciti	33.3	25.1
F. De Gasperin	67.0	67.0
A. Offringa	16.1	16.1
G. Bruni	25.0	0.0
V. Heesen	25.0	25.0
A. Ignesti	8.7	8.7
S. Mandal	50.0	17.0
F. Bempong-Manful	41.7	41.7
M. Murgia	8.0	8.0
M. Knapp	32.0	16.0
H. Intema	8.4	8.4
J. Callingham	174.0	174.0
J. Broderick	12.0	8.0

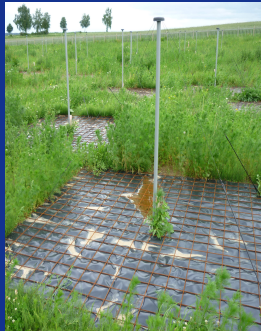
- 53 Stations (24 core (Exloo), 16 remote, 13 International)
- A LOFAR core station consists of 96 Low Band Antennas (LBAs), operating from 10 to 90MHz and 48 High Band Antenna (HBA) tiles that cover the frequency range from 110 to 250 MHz
- Remote stations in the Netherlands have the same number of HBA tiles, and LBAs
- International stations provide a single cluster of 96 HBA tiles and 96 LBAs (6 station in Germany, 3 in Poland, 1 in France, Ireland and UK)



LOFAR - Data flow

Central Processor system

Stations



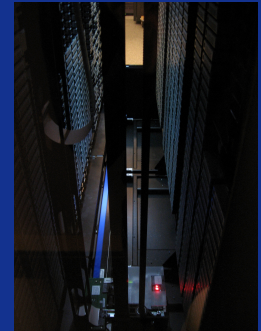
Realtime system



Offline processing



Long Term Archive



240 Gbit/s

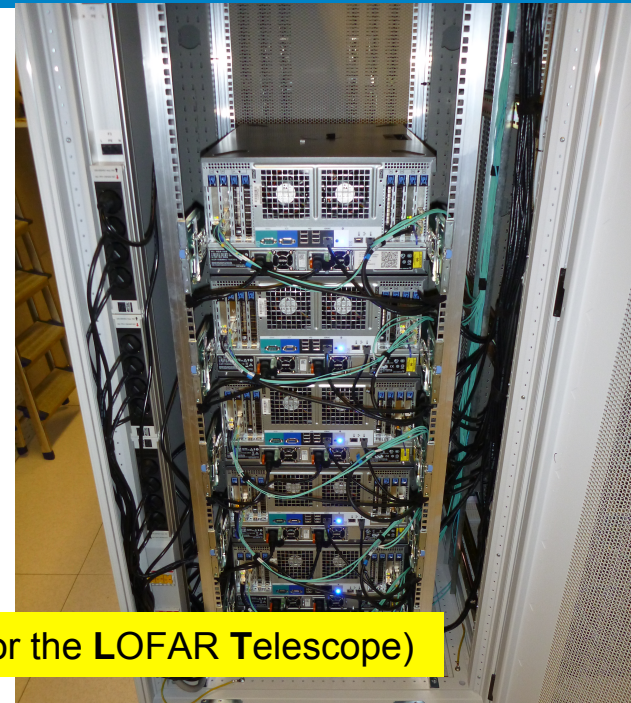
110 Gbit/s

2 Gbit/s

LOFAR COMPUTING MODEL

Central Processor System: The Correlator

- 9 Dell T620 nodes
 - Dual Intel Xeon E5-26xx
 - 2 Nvidia Tesla K10
 - 2 Dual port 10GbE
 - 2 FDR Infiniband HCA



COBALT (**C**orrelator and **B**eamformer Application for the **L**OFAR Telescope)

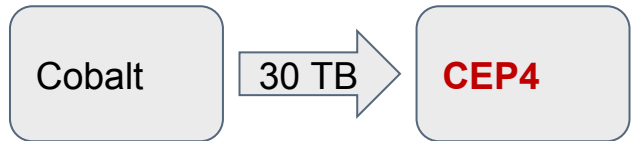
Both the F-stage (Fourier transform) and the X-stage (cross-correlation) are implemented in GPUs

LOFAR COMPUTING MODEL

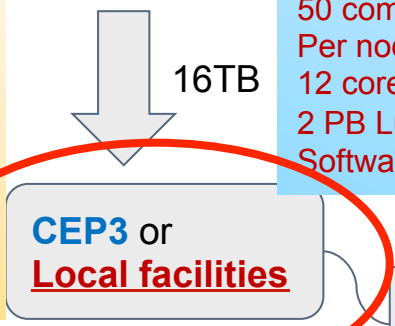
Central Processor System: Post Processing cluster

Two central processing (CEP) clusters in Groningen (i.e. near the correlator). Pipelines use locally-developed generic framework.

Distributed system built using a co-design approach (we know the algorithms and we design the HW)



CEP3: time allocation to PIs based on proposal to do post-processing: 20 nodes
 Per node: 20 cores (2 x Xeon e5 2660v2) 128 GB memory
 2 x 10Gbps Ethernet interface
 22 TB space
 Standard LOFAR software

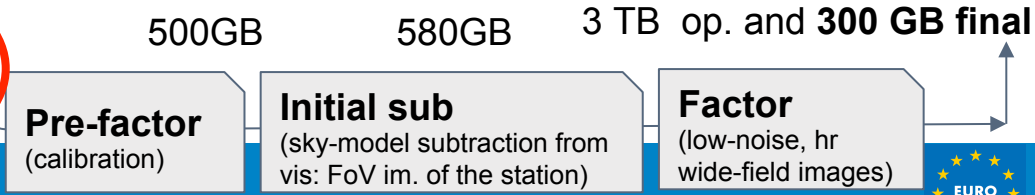


CEP4

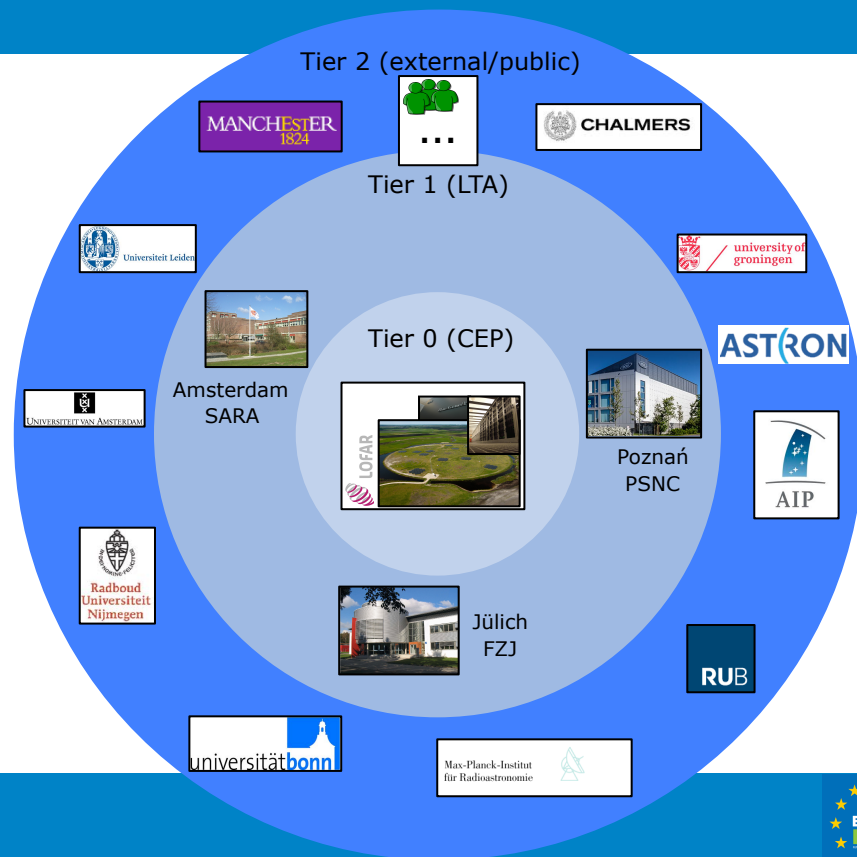
Running **observatory pipelines**: reduction pipelines are used to further process the data into the relevant scientific data products depending on the specific type of observation

Strictly limited to the Radio Observatory -> ingest into LTA
 50 compute nodes (+ 4 with GPU; not used in production).

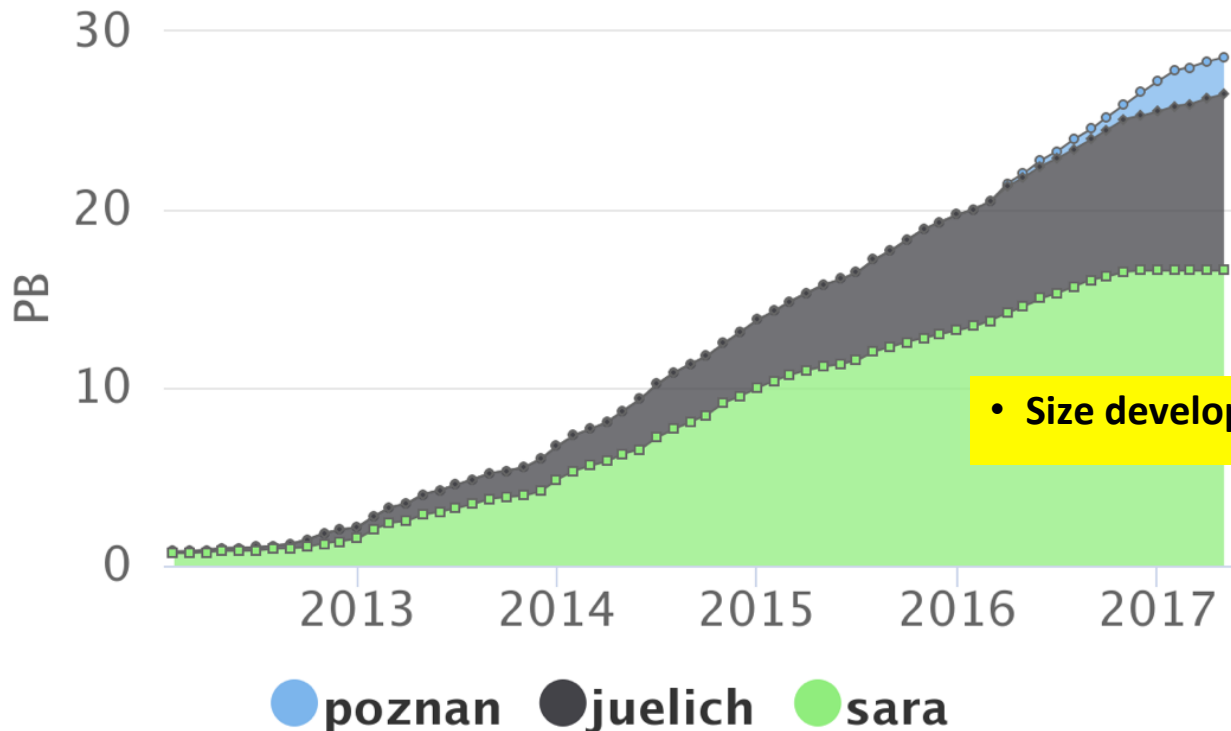
Per node:
 12 cores (Intel Xeon E5-2680v3 2.5 GHz), 256 GB memory
 2 PB LustreFS storage
 Software stack deployed using Docker.



- **7 PB/yr**
- **10 Gbit/s connectivity between Tier 0 and each Tier 1 location**
- **27 PB**
- **~1.5 Bn files**



LTA Storage Site Usage Trend



• Size development in time

Working Group WG-F03-01

UTGII “SKA Precursors and Pathfinders a bassa frequenza” (F03)

UNITO

3 FAT node on OCCAM
4 x Intel Xeon (12 core)
RAM 768 GB DDR4, 1 SSD 800GB, 1 HDD
2TB, 2x10Gb

OATs

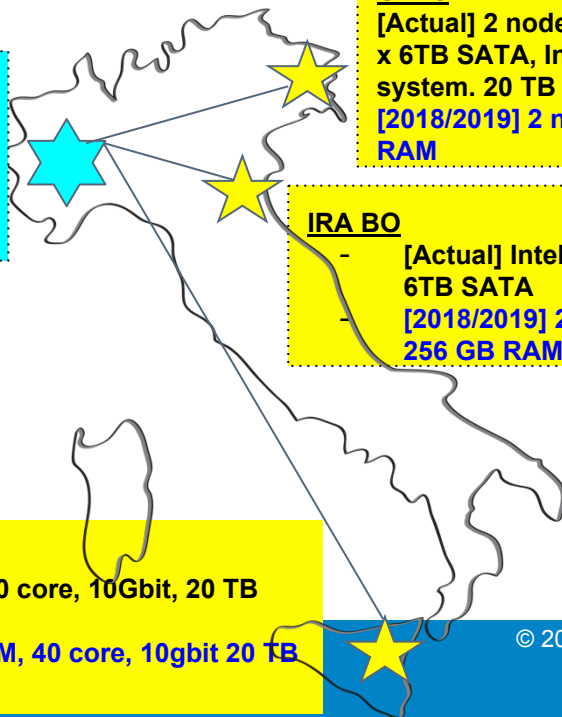
[Actual] 2 nodes: Intel Xeon, 512 GB RAM DDR3 6
x 6TB SATA, Infiniband ConnectX®-3 - Hotcat
system. 20 TB
[2018/2019] 2 nodes: 2 socket (40 core) 256 GB
RAM

IRA BO

- [Actual] Intel Xeon 512 GB RAM DDR3 6 x
6TB SATA
- [2018/2019] 2 nodes: 2 socket (40 core)
256 GB RAM, 10gbit ethernet

OACT

- [Actual] 2 nodes: 256GB RAM, 40 core, 10Gbit, 20 TB
Storage
- [2018/2019] 2 nodes: 256 GB RAM, 40 core, 10gbit 20 TB
storage Bee-GFS.



Alessandro Costa

Gianmarco Maggio

Sara Bertocco

Luca Tornatore

Eva Sciacca

Fabio Vitello

Simone Riggi

Cristina Knapic

Francesco Bedosti

Annalisa Bonafede

Manuela Magliocchetti

Andrea Botteon

...and Giuliano Taffoni and Ugo Becciani

...learn from ILT experience, then....

Direction Independent Pipeline: does the first calibration of LOFAR data

phase 1.a: flag - calibration - transfer of solutions to the target and initial calibration of the target - averaging. Computing time: 3 - 4 days, RAM (core) at least 64GB.

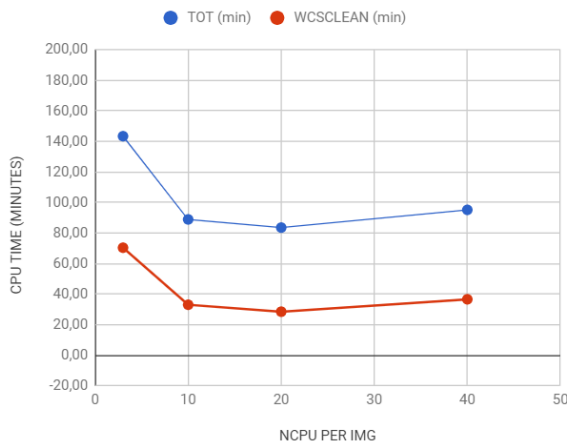
phase 1.b: high and low resolution images and models for auto-calibration. Computing time: 4 -5 days, RAM (core) at least 64GB.

Direction Dependent Pipeline produces low-noise, high-resolution wide-field images from LOFAR HBA data. Computing time: 1 month

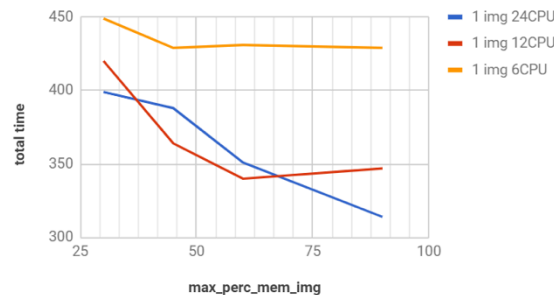
- **Latest versions** of pipelines for the analysis of the **LOFAR HBA observations** in the continuous produced by the **Survey Key Project** and **ASTRON**.
- Pipelines available within **Singularity containers**
- Tests are a **collective work** of the data WG & friends performed at:
 - OACT: MUP node (24 core ht, 64 GB RAM) and LOFAR node (40 core, 256 GB RAM)
 - OATS: HOTCAT node (40 core, 256 GB RAM) with parallel File System
- LoTSS Data tested

- **INPUT Data**: 4 bands (2 MHz each)
- **Relevant input parameters**
 - **CPU Cores**: max_cpus_per_img, max_dppp_threads
 - **Memory**: max_percent_mem_per_img, max_imagers_per_node

INITSUB CPU time vs NCPU



total time vs max_perc_mem_img

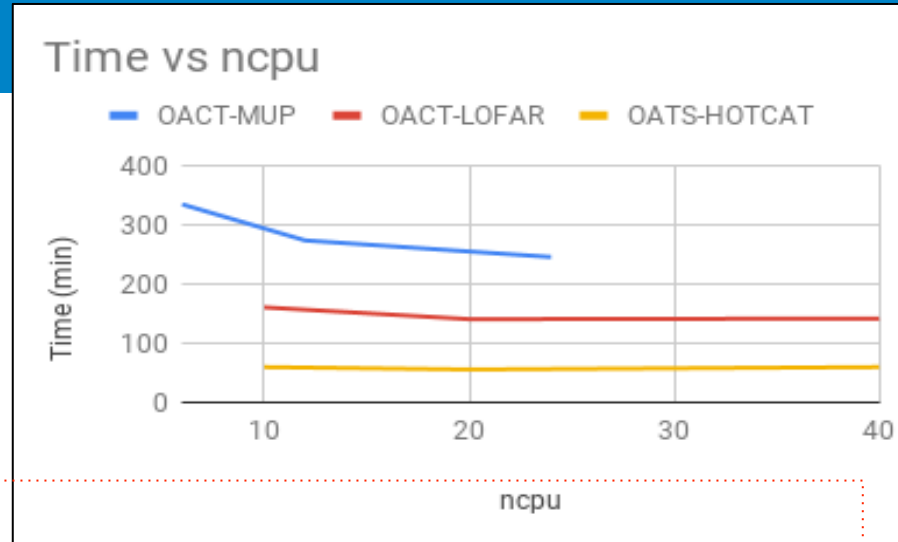


Conclusions:

the pipeline can run also on node with RAM < 256GB.

Recommend node with 3-4 GB per core and 6-10 cores per image.

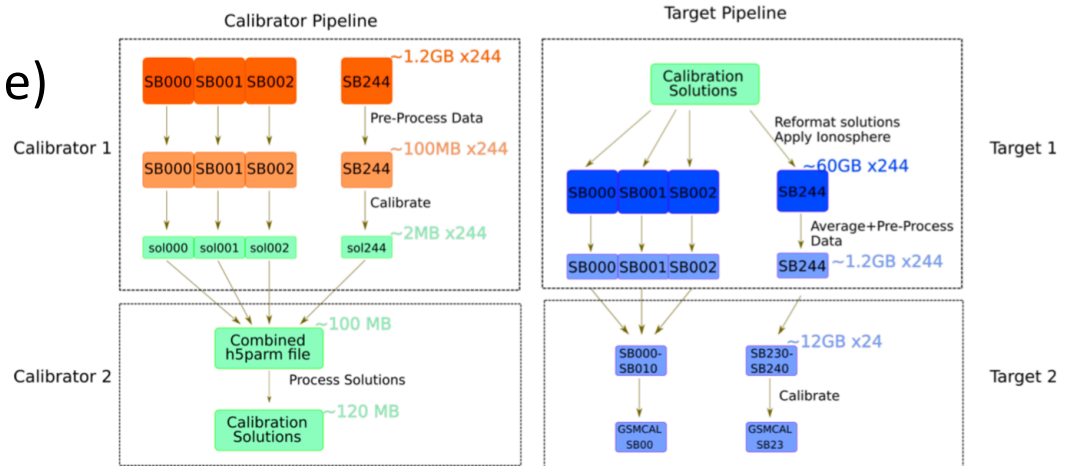
- **INPUT Data:** 6 bands (2 MHz)
- **Relevant input parameters:**
 - **CPU Cores:** ncpu
 - **Memory:** wsclean_fmем



Conclusions: the pipeline does not scale with CPU cores/RAM memory but we have noticed it is very I/O time dependent therefore a parallel File Systems should be employed (e.g. BeeGFS as in OATS node). Recommend nodes with at least 10 cores and 256 GB RAM.

- Embarrassing parallel computation
- Multithreaded (on single node)
- Extremely IO demanding

Machev et al 2018



- HTC cluster
- **Parallel filesystem** optimized for high throughput (~4GB/s)
- Monitoring and Operations (telgraf+grafana)
- Containers (docker/singularity)
- **ObjectStorage** (new HPC approach to storage) testing
- Identify or train **“support” groups** for Italian Astronomers and for HW/SW contribution to ILT



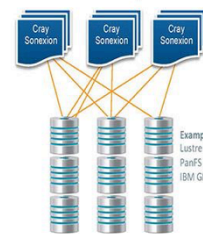
Examples:
NetApp WAFL
Sun ZFS
WinNT
HDS BlueArc

Distributed File System



Examples:
Isilon OneFS
NetApp ONTAP
Gluster

Parallel File System

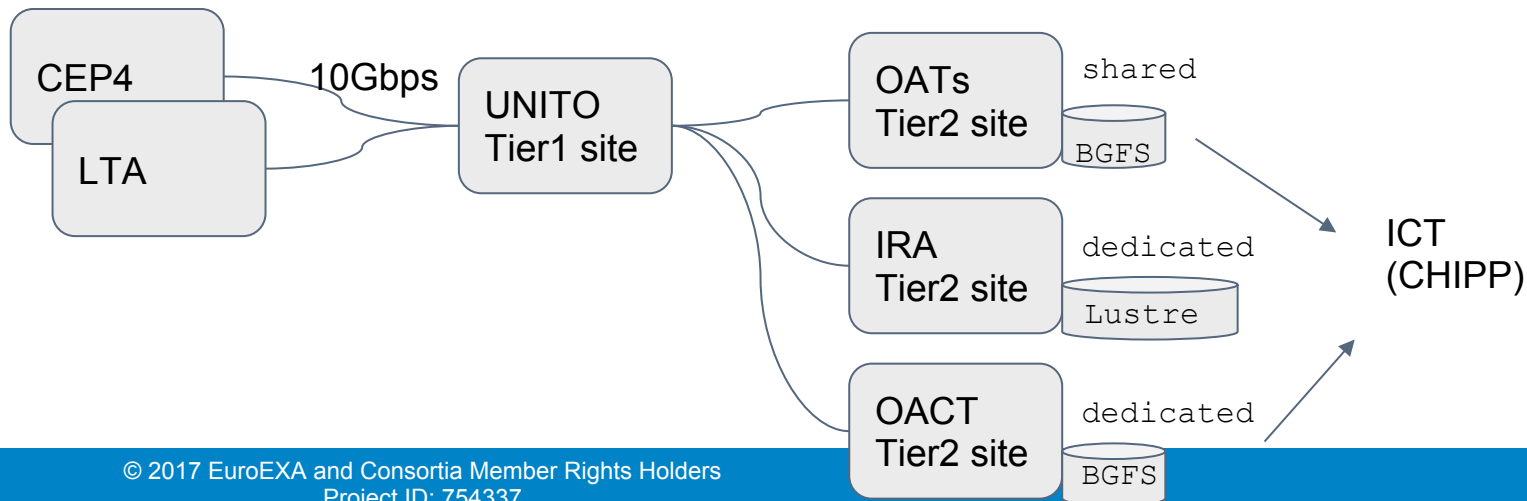


Examples:
Lustre
ParsFS
IBM GPFS



It is not only HW but also **people**

We want to grow our knowledge to the LOFAR software infrastructure and to contribute to its **development and optimization.**



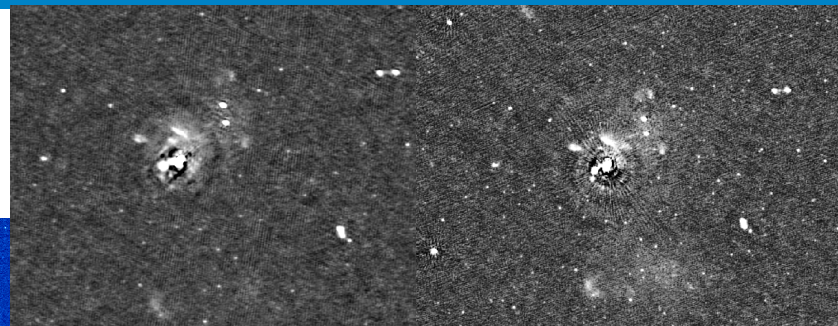
FACTOR



killMS+DDFacet

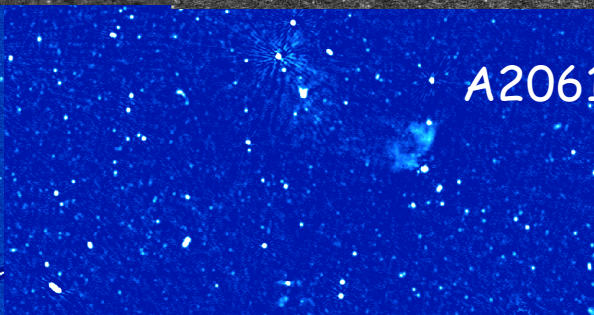
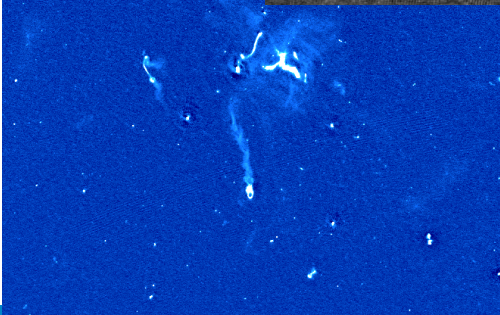
Initially at IRA node
now installed and
operating also at the
TS and CT nodes

A. Botteon @ IRA



A2255

Analysis for the Survey KP (DDFacet)



A2061

Monthly Notices

of the
ROYAL ASTRONOMICAL SOCIETY

MNRAS 478, 885–898 (2018)

Advance Access publication 2018 May 1



doi:10.1093/mnras/sty110

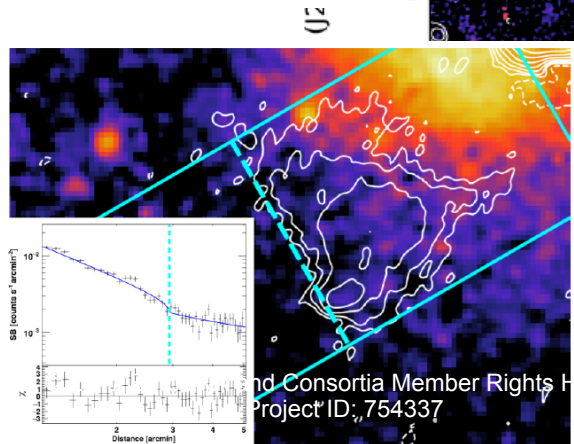
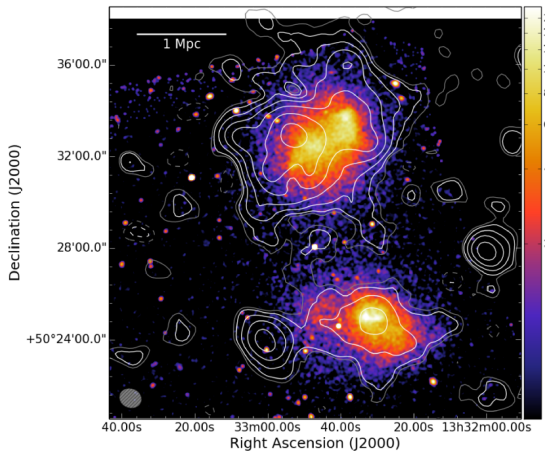
Astronomy & Astrophysics manuscript no. a781

November 21, 2018

©ESO 2018

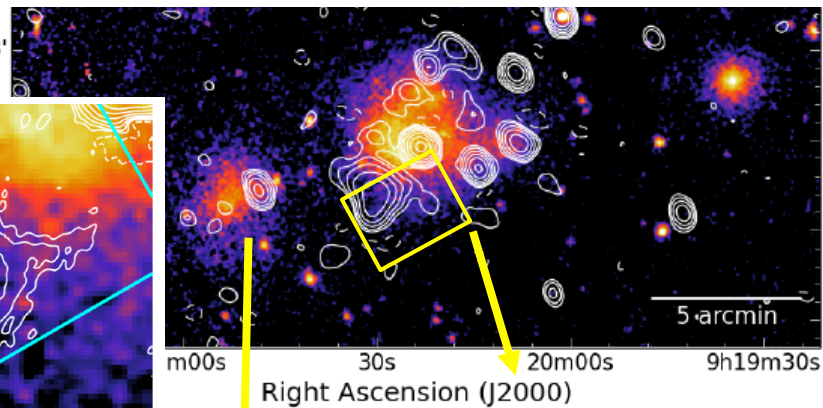
LOFAR discovery of a double radio halo system in Abell 1758 and radio/X-ray study of the cluster pair

A. Botteon,^{1,2}★ T. W. Shimwell,^{3,4} A. Bonafede,^{1,2,5} D. Dallacasa,^{1,2} G. Brunetti,² S. Mandal,⁴ R. J. van Weeren,⁴ M. Brüggen,⁵ R. Cassano,² F. de Gasperin,⁴ D. N. Hoang,⁴ M. Hoeft,⁶ H. J. A. Röttgering,⁴ F. Savini,⁵ G. J. White,^{7,8} A. Wilber⁵ and T. Venturi²



The spectacular cluster chain Abell 781 as observed with LOFAR, GMRT, and XMM-Newton

A. Botteon^{1,2}, T. W. Shimwell^{3,4}, A. Bonafede^{1,2,5}, D. Dallacasa^{1,2}, F. Gastaldello⁶, D. Eckert⁷, G. Brunetti², T. Venturi², R. J. van Weeren⁴, S. Mandal⁴, M. Brüggen⁵, R. Cassano², F. de Gasperin^{4,5}, A. Drabent⁸, C. Dumba^{8,9}, H. T. Intema⁴, D. N. Hoang⁴, D. Rafferty⁵, H. J. A. Röttgering⁴, F. Savini⁵, A. Shulevski¹⁰, A. Stroe¹¹ and A. Wilber⁵



and Consortia Member Rights Holders
Project ID: 754337

- New platforms able to **ExaScale** (1000 times more the sustained performance of actual TIER-0 Facilities)
- 1 Billion Euro EU investment: EuroHPC and EPI
- **HW/SW Co-Design**

- Scientist must **change their algorithms** and approach to computing resources (new system software, massive use of accelerators, massive code optimization ...)
- We must **change the way we program** our codes

- Already some tests and investments for SKA on UK
- INAF is the only A&A Institute in EU Involved in Exascale prototyping
 - LOFAR data reduction and analysis

EuroEXA: towards a platform for exascale in Europe (H2020 20Meuro)

→ hw-sw codesign approach

→ port LOFAR pipelines on Arm

→ Test and verification of system SW that enables the use of exascale capabilities (e.g. OmpSs, GPI etc.)

→ BigData application for EuroEXA (parallel filesystem and software development)

New tech skills

Role in ILT

Improve our capacity to satisfy the needs of scientific community

SKA

Alessandro Costa

Cristina Knapic

Gianmarco Maggio

Francesco Bedosti

Sara Bertocco

Annalisa Bonafede

Luca Tornatore

Manuela Magliocchetti

Eva Sciacca

Andrea Botteon

Fabio Vitello

Simone Riggi

...and Giuliano Taffoni and Ugo Becciani