

Access to massive catalogues in the Gaia archive: a new paradigm

Jesús Salgado

*Juan González, Raúl Gutiérrez, Juan Carlos Segovia,
Christophe Arviset, Sara Nieto, Bruno Merín*

*Gaia Archive Development Team
ESAC Science Data Center (ESDC)*

Issue/Revision: 1.0

Reference: Gaia Archive

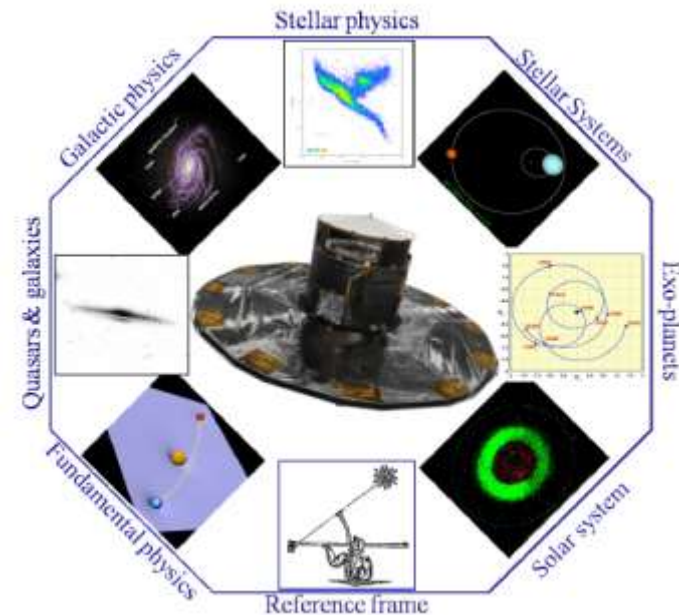
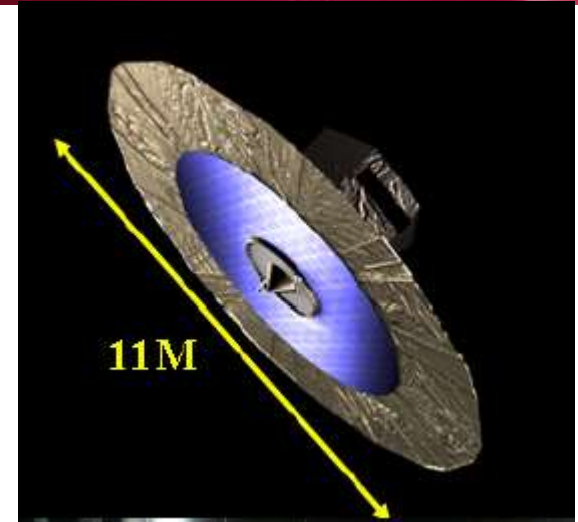
Status: Issued

ESA UNCLASSIFIED - Releasable to the Public

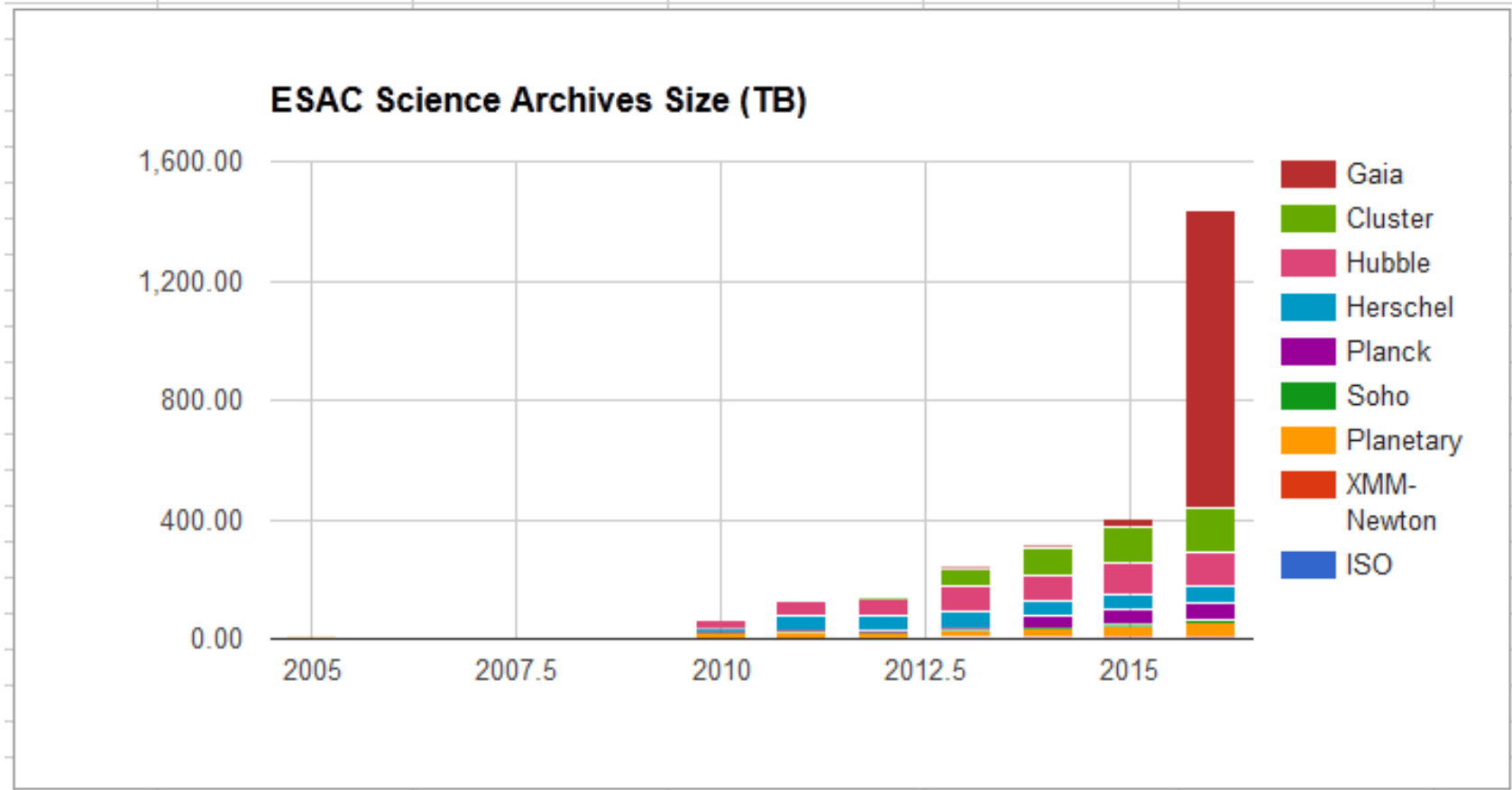
Gaia Satellite Overview



1. ESA Corner Stone 6
2. ESA provided the hardware and launch
 - a. Mass: 2120 kg (payload 743 kg)
 - b. Power: 1631 W (payload 815 W)
3. Launched 19th December 2013.
4. Stereoscopic Census of Galaxy over 5 years
5. Possible extension of 1 year (fuel should be ok)
6. Astrometry $G < 20$ (10^9 sources)
7. $25 \mu\text{arcsec}$ at $V=15$
8. Radial Velocities $G < 16$
9. Spectro Photometry $G < 20$ (millimag)
10. Final catalogue ~ 2022



ESAC Archives Volume evolution



All data stored on hard disks and distributed through Internet

Euclid will add up to ~150 PBs by 2023

Statistics up to June 5

Nominal Mission Data

Type of Data	Amount	
Science telemetry	17 TB	
Astrometry transits	22.5×10^9	225×10^9 images
Photometry transits	22.5×10^9	45×10^9 images
Spectroscopy transits	1.5×10^9	4.5×10^9 spectra
Main Database	44TB	

- $\approx 30GB/day$ – $> \approx 100TB$ total
- with products can be 1PB total data by mission end

Timo Prusti (2015)



Gaia Archive current content



Simulations

• GUMS		
• Milky Way:	2×10^9	rows
• Large Magellanic Cloud:	7.5×10^6	rows
• Small Magellanic Cloud:	1.2×10^6	rows
• Galaxies:	38×10^6	rows
• Quasars:	10^6	rows
• GOG	1.8×10^9	rows

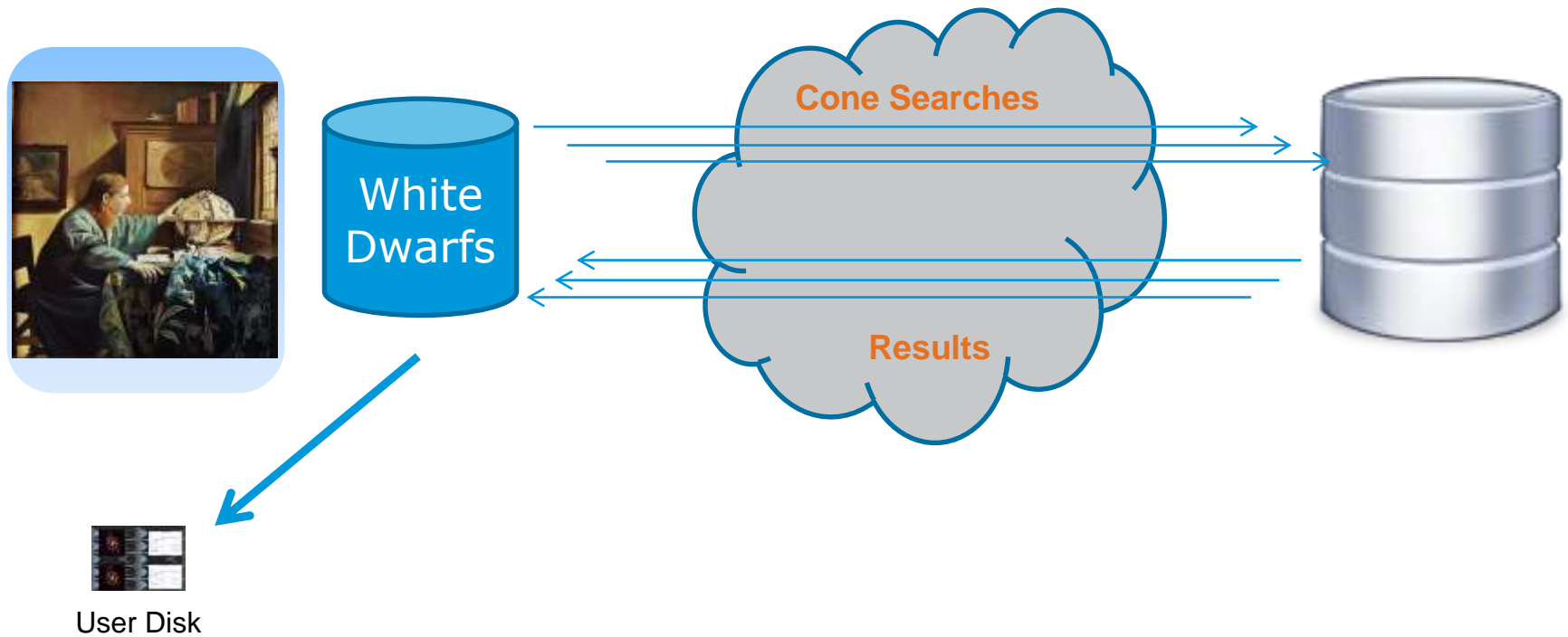
External Catalogues

• IGSL (Initial Gaia Source List)	1.2×10^9	rows
• 2MASS	9.4×10^8	rows
• Tycho2	2.5×10^6	rows
• UCAC4	1.1×10^8	rows

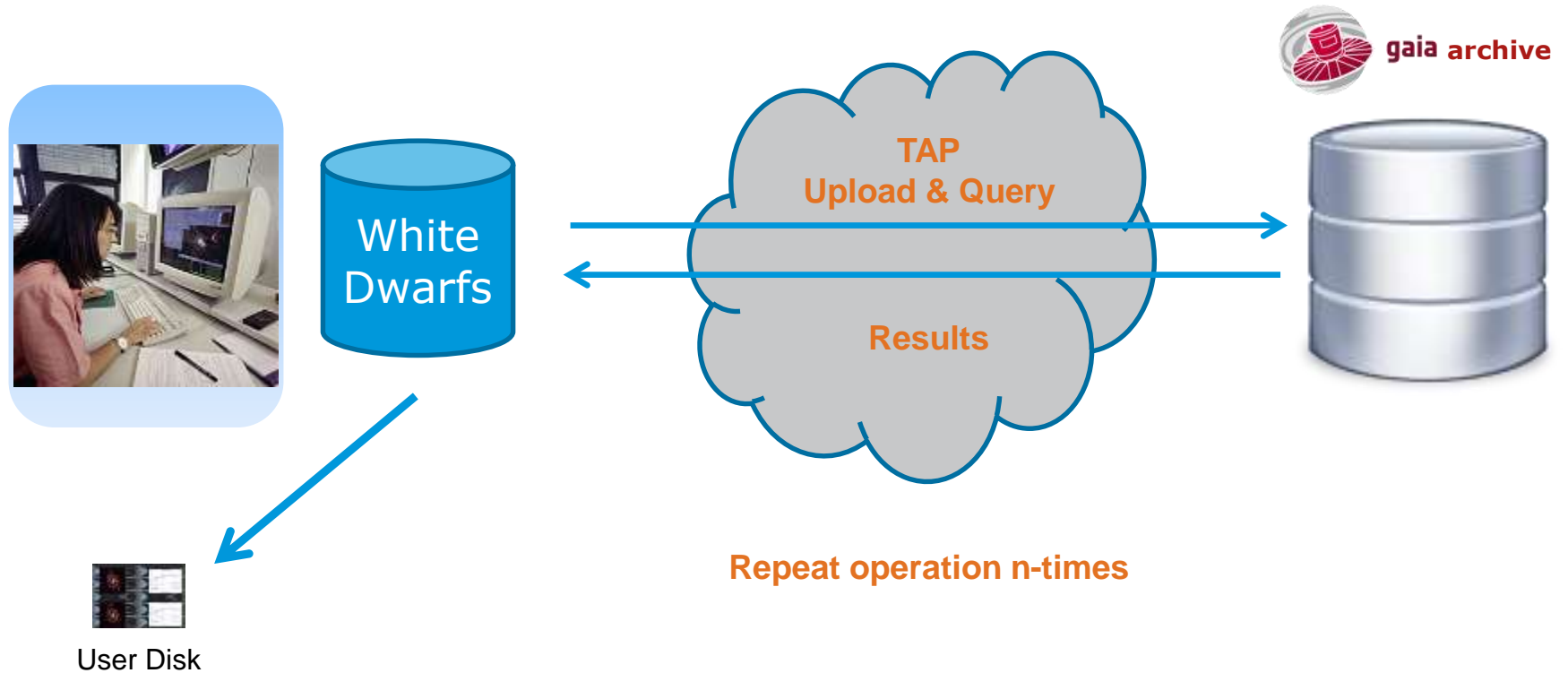
Gaia (not public data)

• TGAS	2.2×10^6	rows
• MDB	$\sim 3 \times 10^9$	rows
• Foreseen, first Gaia catalogue	$> 2 \times 10^9$	rows

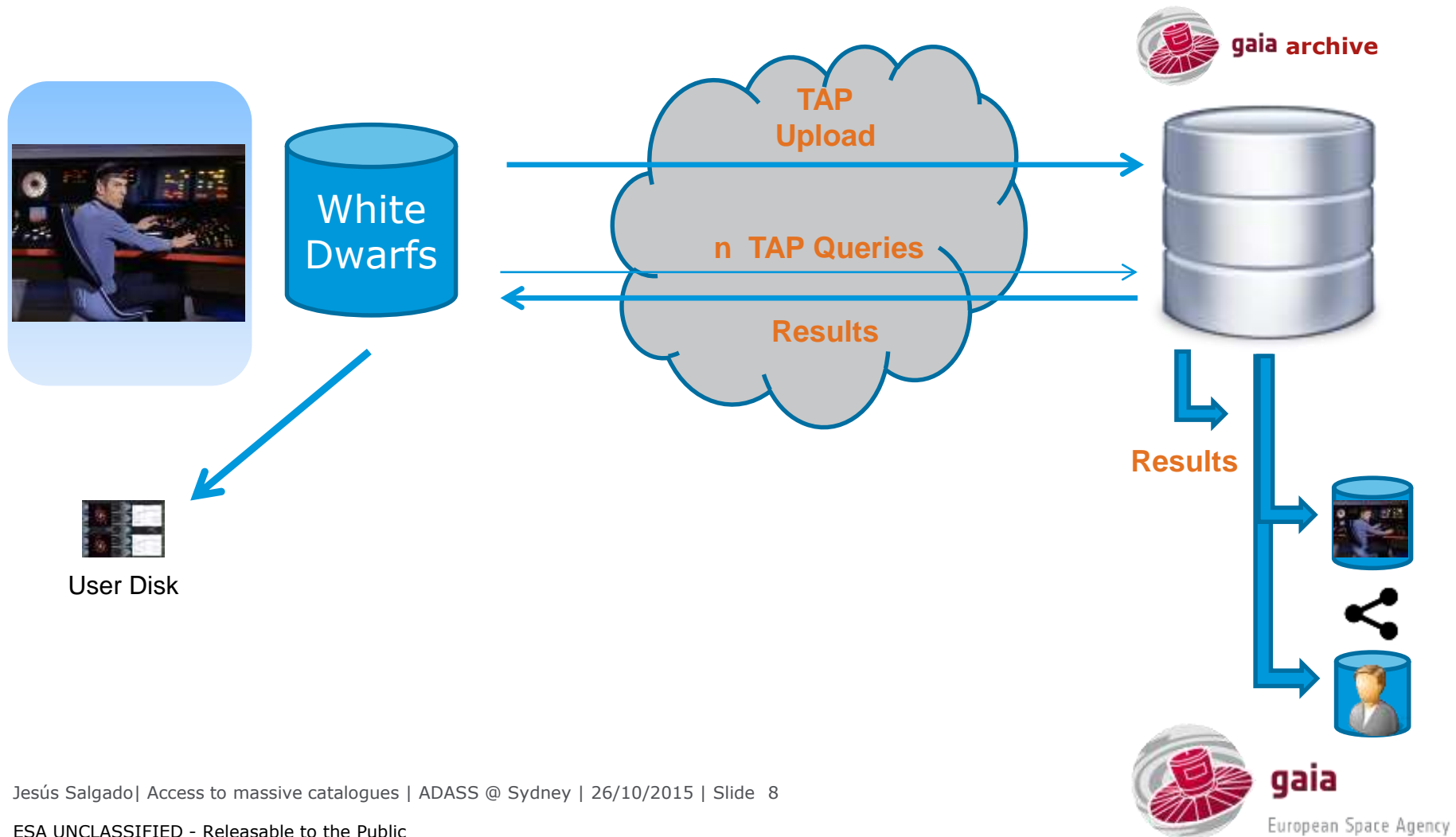
Crossmatch for my objects



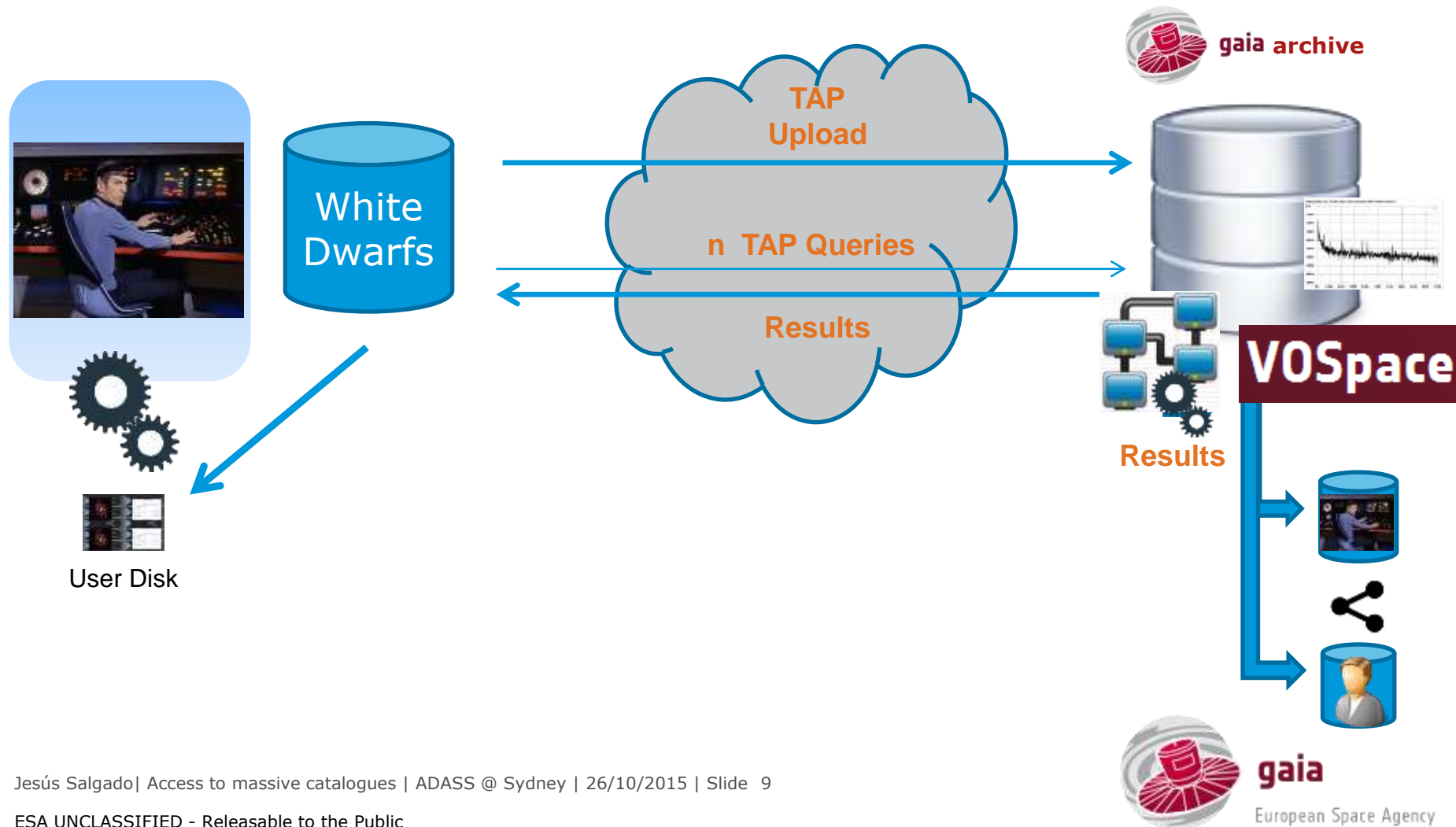
Crossmatch for my objects (II): TAP approach



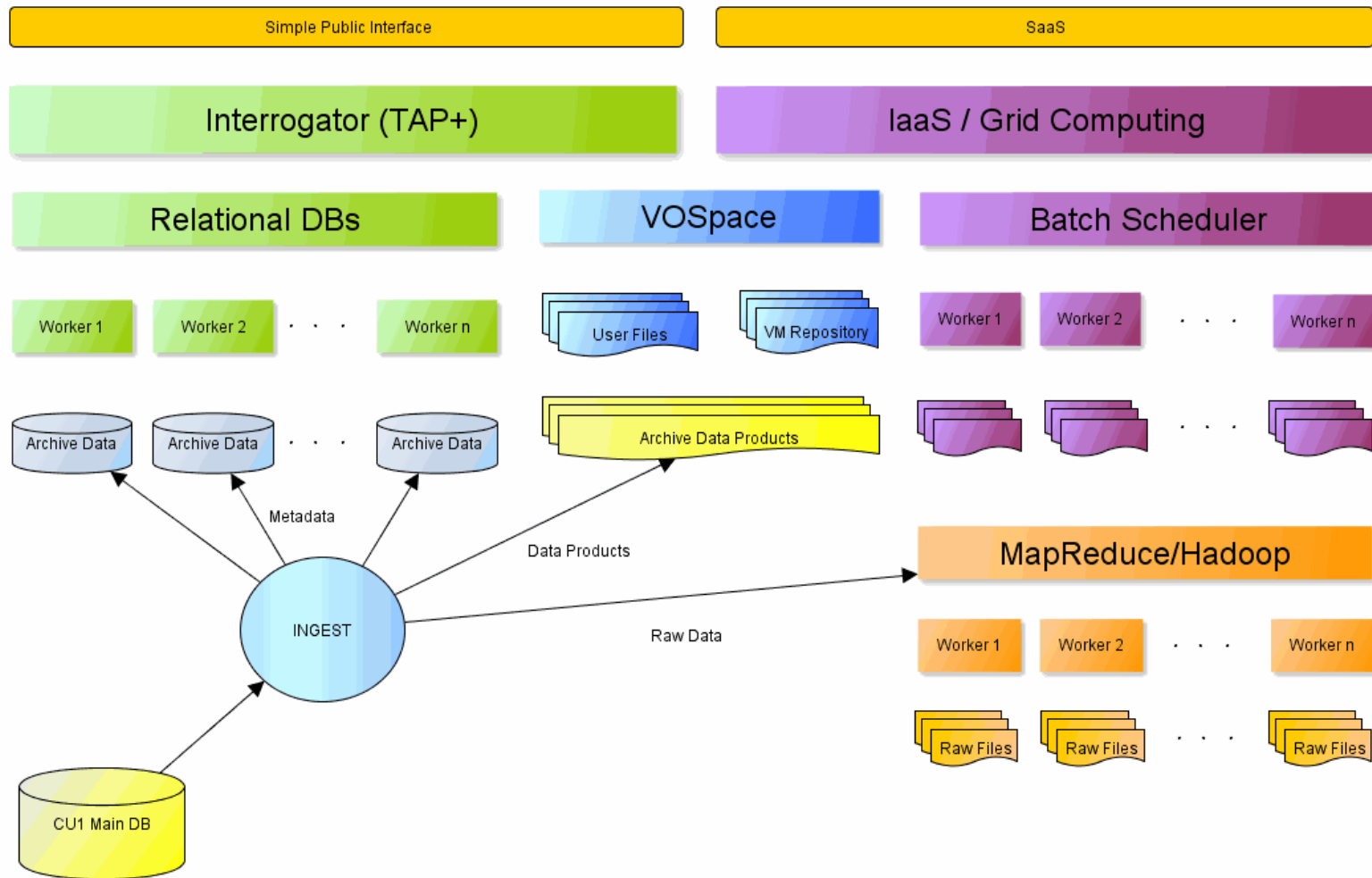
Crossmatch for my objects (III): TAP+ approach & VOspace+



Crossmatch for my objects (III): TAP+ approach & VOSpace+ & SaaS



Gaia Archive Architecture



Gaia Archive: TAP+ Interface



The screenshot displays the Gaia Archive TAP+ interface. At the top, there is a navigation bar with links for HOME, SEARCH, STATISTICS, HELP, DOCUMENTATION, VOSPACE, SHARE, and ADMIN. Below this, there are tabs for Sample Form, ADQL Form, and Query Results. The main area is divided into a left sidebar with a tree view of Gaia tables and a central query editor. The query editor contains the following SQL query:

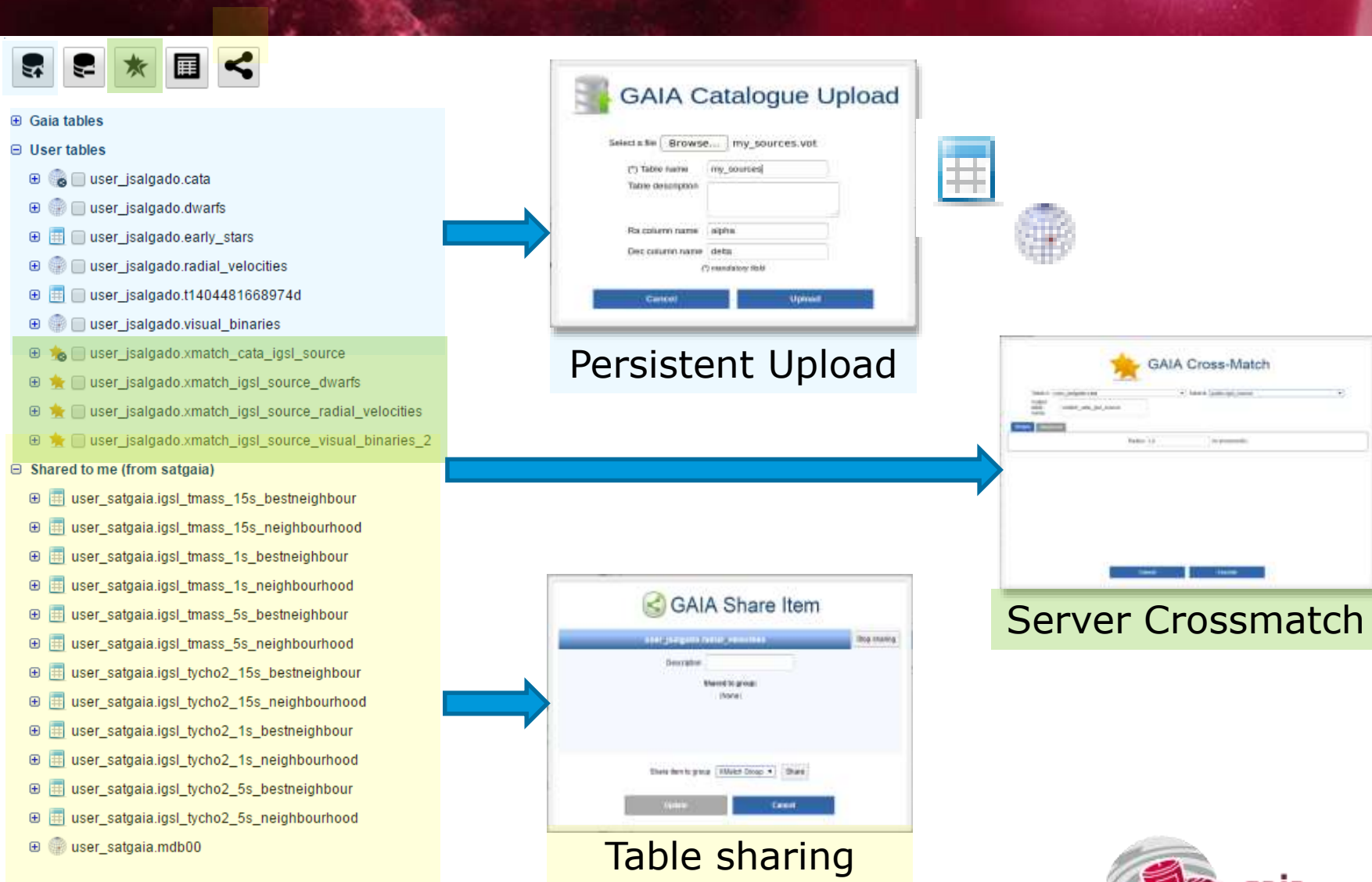
```
SELECT DISTANCE(POINT('ICRS',alpha,delta), POINT('ICRS',266.41683,-29.00781)) AS dist, *
FROM public.gai_cataloguesources
WHERE 1<CONJAINS(POINT('ICRS',alpha,delta),CIRCLE('ICRS',266.41683,-29.00781, 0.00333333)) ORDER BY dist ASC
```

Below the query editor are buttons for 'Reset Form' and 'Submit Query'. The results table below shows a list of jobs with columns for Status, Job ID, Creation date, Num. rows, and Size. The table contains 20 rows of data, with the first 10 rows showing jobs with 4266 rows and 1 MB size, and the last 10 rows showing jobs with 0 rows and 0 KB size.

Status	Job	Creation date	Num. rows	Size
✓	1444318635024D	08-Oct-2015, 17:37:15	4266	1 MB
✓	1444318623040D	08-Oct-2015, 17:37:03	4266	1 MB
✓	1444318596137D	08-Oct-2015, 17:36:36	4266	1 MB
✓	1444318513603D	08-Oct-2015, 17:35:13	4266	1 MB
✓	1444318461116D	08-Oct-2015, 17:34:21	4266	1 MB
✓	1435659073865D	30-Jun-2015, 12:11:14	4266	1 MB
✓	1426505662813D	16-Mar-2015, 12:34:43	1313	666 KB
✓	rmatch_cata_igsl_source	16-Mar-2015, 12:33:04	0	0 KB
✓	1425915487847D	09-Mar-2015, 16:38:08	1830	270 KB
✓	1425915456515D	09-Mar-2015, 16:37:37	1313	666 KB
✓	rmatch_igsl_source_cata	09-Mar-2015, 16:36:23	0	0 KB
✓	1425318251124D	02-Mar-2015, 18:44:11	1313	666 KB
✓	rmatch_igsl_source_cata	02-Mar-2015, 18:42:07	0	0 KB
✓	1424280930031D	18-Feb-2015, 18:35:30	1	0 KB



Gaia Archive: TAP+ Interface



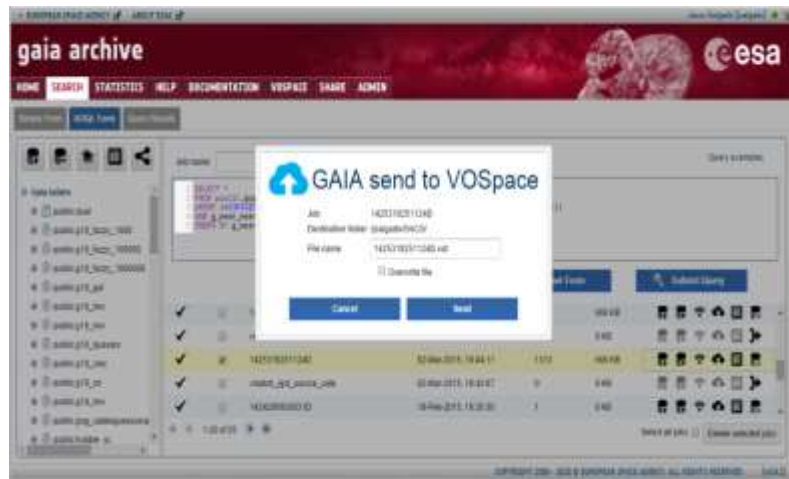
Gaia Archive Crossmatch Examples (20 threads)



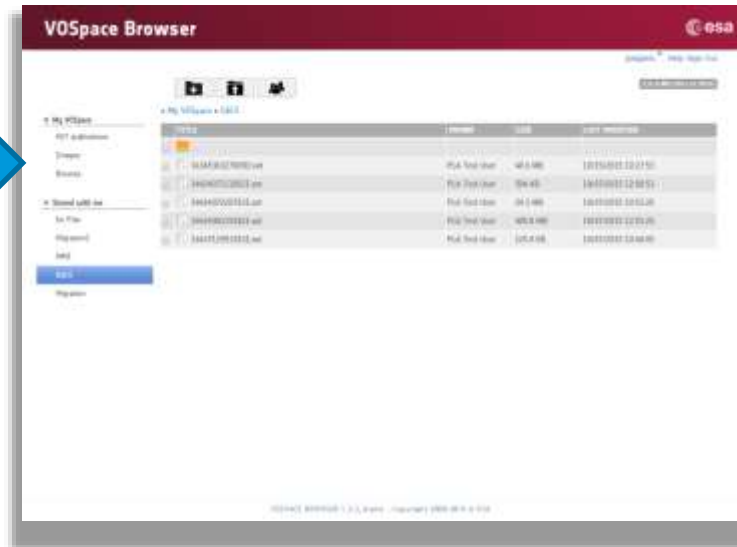
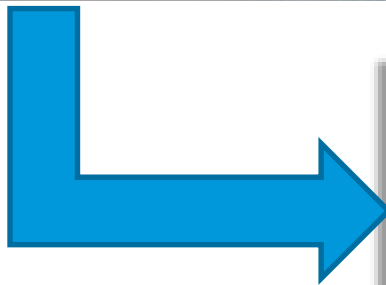
Catalogue 1	Catalogue 2	Radius (arcsec)	# results	Time
Tycho2 2.5x10 ⁶ rows	2MASS PSC 4.7x10 ⁸ rows	1''	2,495,304	49s
Tycho2 2.5x10 ⁶ rows	2MASS PSC 4.7x10 ⁸ rows	5''	2,614,163	116s
Tycho2 2.5x10 ⁶ rows	IGSL 1.2x10 ⁹ rows	1''	2,600,542	46s
Tycho2 2.5x10 ⁶ rows	IGSL 1.2x10 ⁹ rows	5''	2,829,401	55s

Tycho2 vs IGSL crossmatches are even faster than the ones with 2MASS as IGSL is located in the fastest local storage (PCIe), even when IGSL (similar to the final Gaia catalogue) is around 3 times bigger than 2MASS

VOSpace : Virtual storage for collaboration

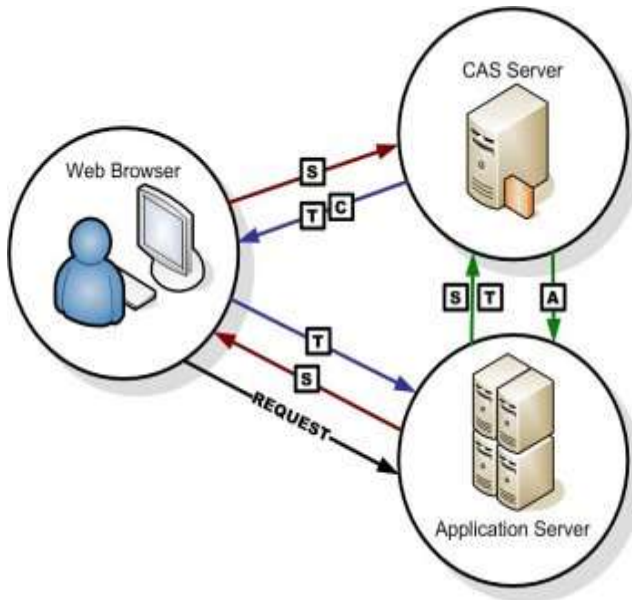


1. Astronomy DropBox
2. VO protocol
3. InterOperable
4. Accessible by VO Applications
5. Accessible by REST
6. Close to data
7. Repository VOSpace IDs?



P082 Sara Nieto





AUTHENTICATION PROCESS

1. CLIENT REQUEST →
2. AUTHENTICATION REDIRECT →
3. TICKET FORWARDING →
4. TICKET VALIDATION →

TRANSMITTED DATA

- T** ICKET
- S** ERVICE ID
- C** OOKIE (CAS SSO)
- A** UTHENTICATED ID

1. Single Sign-on through CAS server (both GACS and VOSpace)
2. Authorization through CAS proxy tickets
3. SAMP simple authentication
4. Full support on SAMP https for Web profile still pending at protocol level



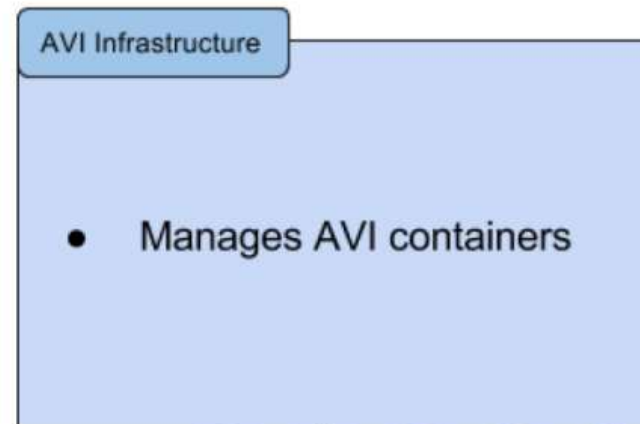
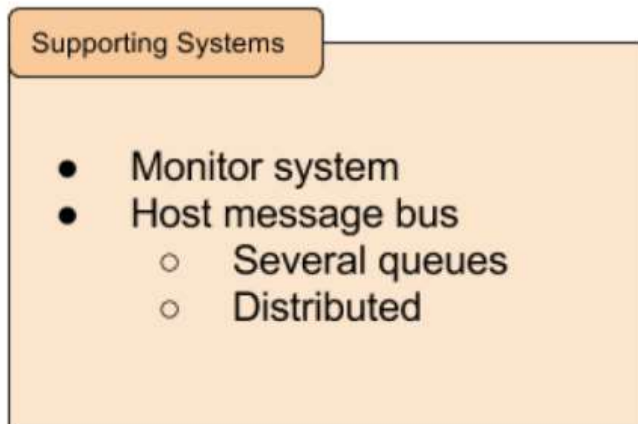
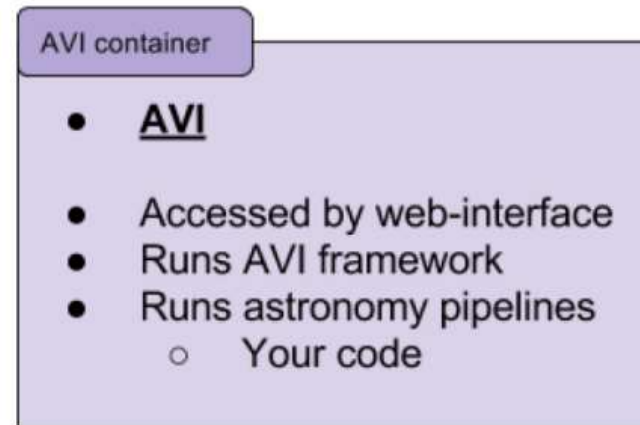
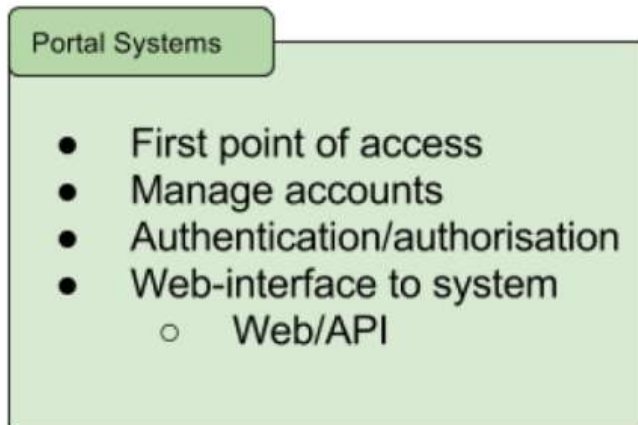
Authenticated SAMP



gaia

European Space Agency

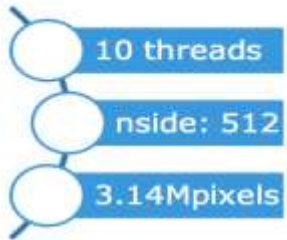
Gaia Added Value Interfaces Portal



O1.4 William O'Mullane
O11.3 Christophe Arviset
P080 Vicente Navarro Ferreruela



Vizualization

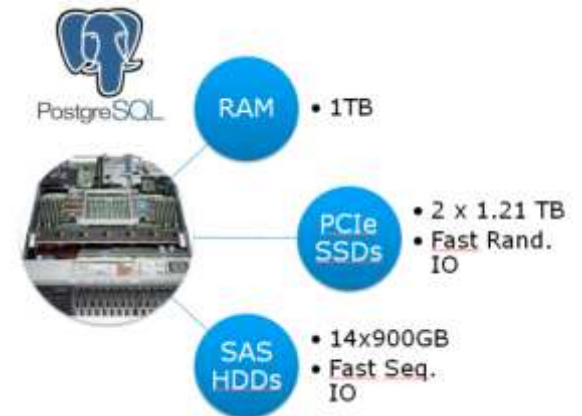


	G10_MW <ul style="list-style-type: none"> • 2.14 billion rows • 1.8 minutes
	IGSL_SOURCE <ul style="list-style-type: none"> • 1.22 billion rows • 65 seconds
	G10_GAL <ul style="list-style-type: none"> • 38 Million rows • 10 seconds

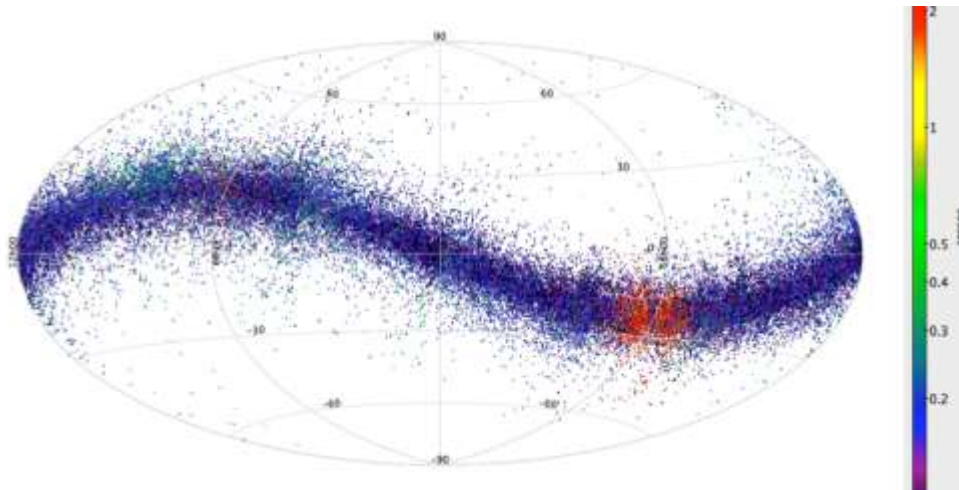


	G10_MW <ul style="list-style-type: none"> • 2.14 billion rows • 14 minutes
	IGSL_SOURCE <ul style="list-style-type: none"> • 1.22 billion rows • 7 minutes
	G10_GAL <ul style="list-style-type: none"> • 38 Million rows • 16 seconds

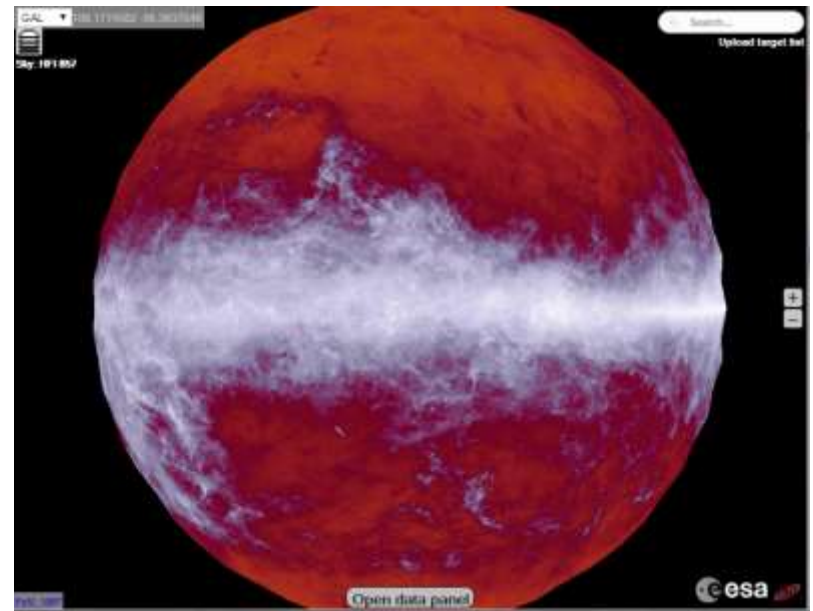
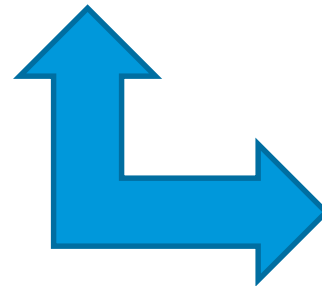
1. Visualization is a need
2. Statistics provide holistic views. Big data techniques
3. Specific visualization packages within Gaia consortium



Asteroids



“Asteroids position knowledge will be increased on a factor of 50 to 100 due to Gaia mission”



Conclusions



1. New Astronomical missions are facing the same challenges about big data of other communities. That implies:
 - a. A new way to offer the data
 - b. A new way to work for scientists
2. Techniques can be reused like:
 - a. Asynchronous jobs, visualization, Map/Reduce
 - b. VLDBs (Very Large Databases)
3. Resources for the community
 - a. User work space *inside* the archive
 - b. Analysis work is done where the data is
4. Security is now an issue
5. First public version coming soon: Summer 2016