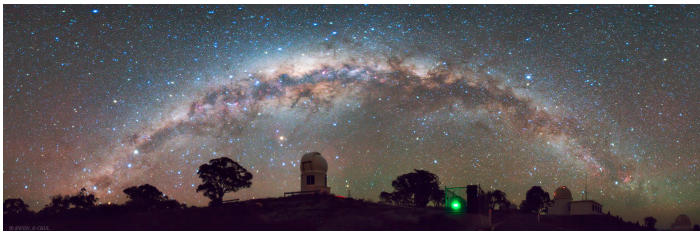# The QUBRICS database for machine learning: architecture and performance.

## Giorgio Calderone (INAF-OATs)

*QUBRICS collaboration:*
*Konstantina Boutsia, Stefano Cristiani, Guido Cupani,*
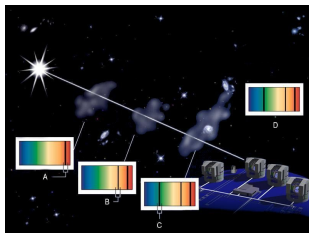*Andrea Grazian, Francesco Guarneri, Luciano Nicastro, Matteo Porru, ...*

## The QUBRICS Project
*QUasars as BRIght beacons for Cosmology in the Southern hemisphere*



**Purpose:**

- Collect photometric datasets in the Southern Hemisphere;
- Use machine learning to select new, bright, high-$z$ ($z > 2.5$) QSOs candidates;
- Spectroscopic follow-up, confirm classification and redshift;
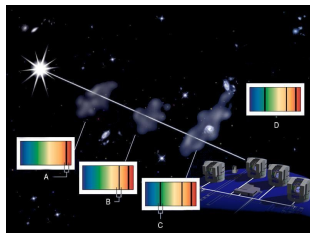- $\Rightarrow$ exploit acquired knowledge!

## The QUBRICS Project
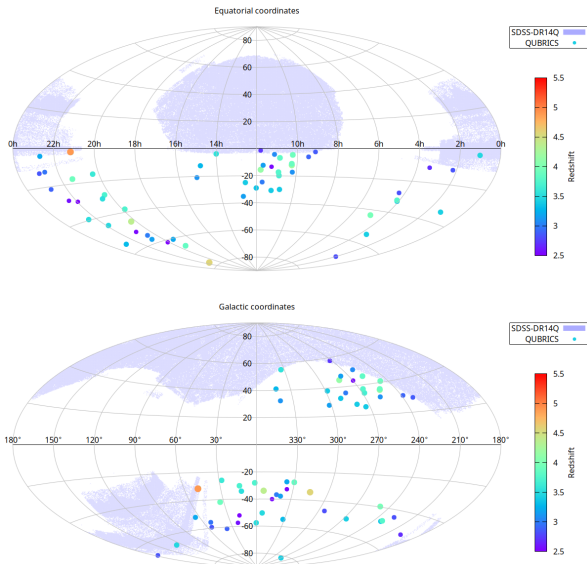*QUasars as BRIght beacons for Cosmology in the Southern hemisphere*



### Purpose:

- Collect photometric datasets in the Southern Hemisphere;
- Use machine learning to select new, bright, high-$z$ ($z > 2.5$) QSOs candidates;
- Spectroscopic follow-up, confirm classification and redshift;
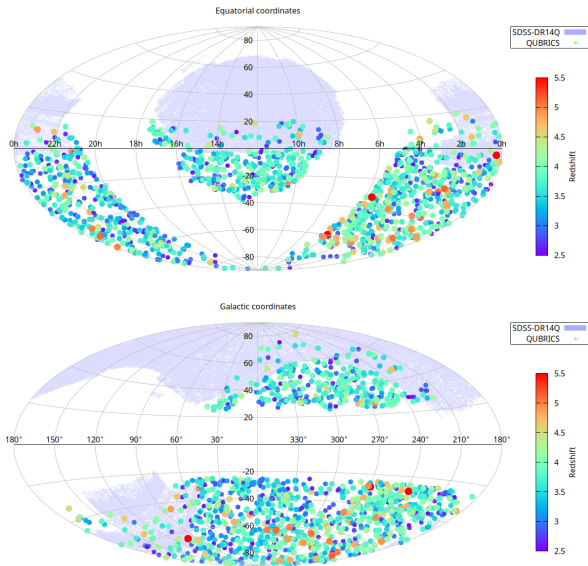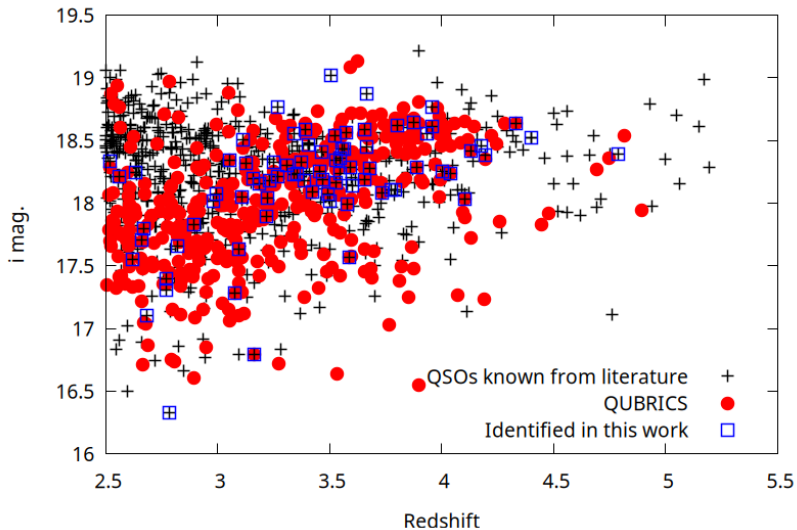- $\Rightarrow$ exploit acquired knowledge!

### Papers:

Calderone+19, Boutsia+20, Boutsia+21, Guarneri+21, Cupani+21, Grazian+21, Guarneri+22, Cristiani+23, Grazian+23, Calderone+24, Grazian+24, More in preparation...

# QUBRICS QSOs ($z > 2.5$) in 2019



Equatorial coordinates

Galactic coordinates
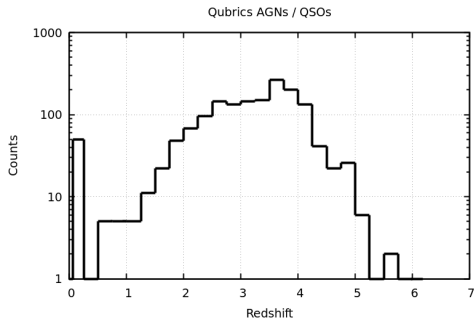
# QUBRICS QSOs ($z > 2.5$) in 2025

# QUBRICS QSOs vs mag. in *i* band
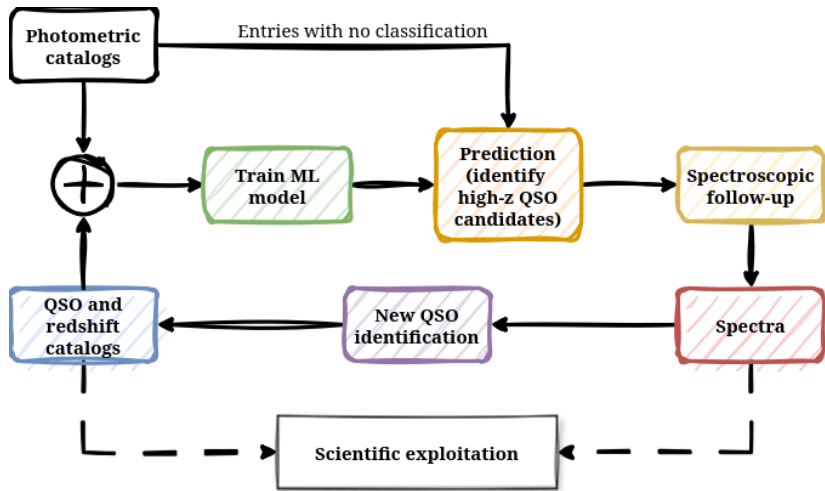


Calderone et al., 2024

## QUBRICS statistics in 2025

**Observations:**

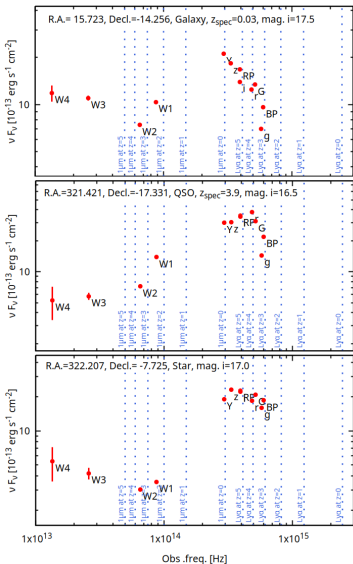- Candidates observed: 2019;
  - good quality: 1764;
    - Stars: 76;
    - Galaxies: 38;
  - bad quality: 255.

- QSOs: 1585 (93%);
  - $z$>2.0: 1438 (84%);
  - $z$>2.5: 1274 (75%);
  - $z$>3.0: 992 (58%);
  - $z$>4.0: 228 (13%)
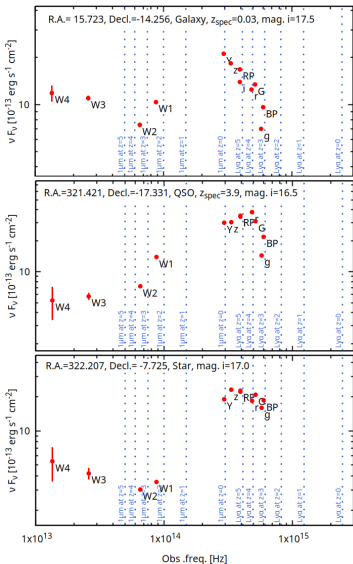- Max. redshift: 5.768.



Qubrics AGNs / QSOs

# QUBRICS self-feeding
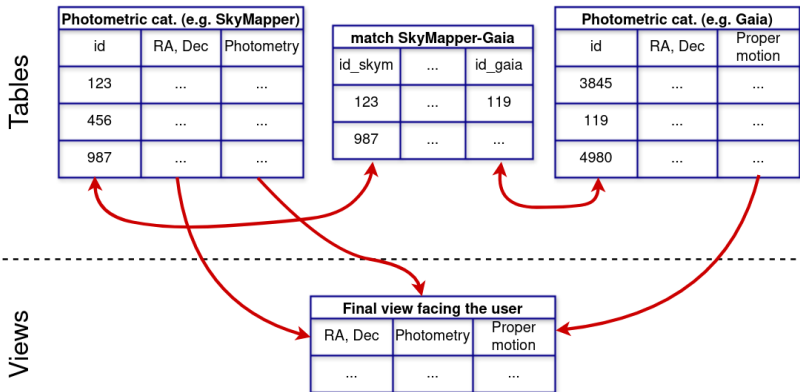
# How data looks like?

# How data looks like?



- Disk usage $\sim$ 4TB;
- Single workstation with 8 CPU, 64GB RAM;
- MariaDB database;
- Less than 10 persons involved (+occasional contributors)
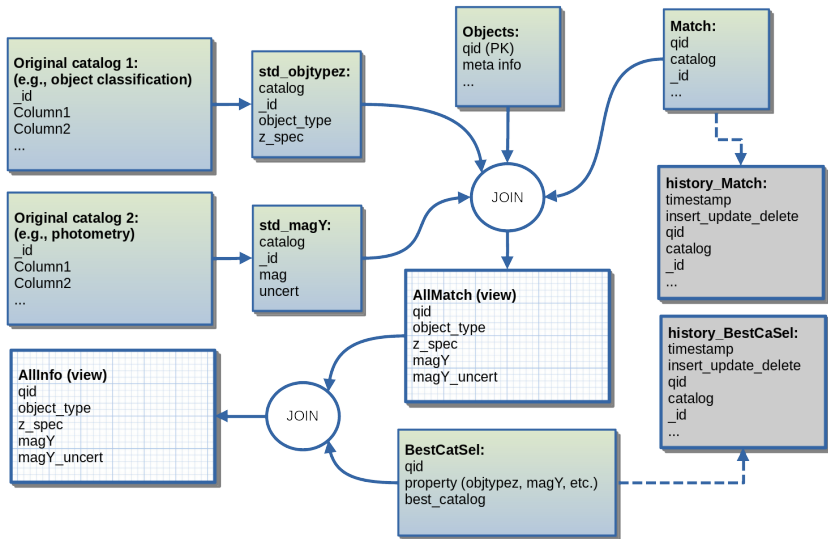- $\Rightarrow$ **"Small data" project**.

## Relational DB, inner joins



- Original catalogs are not modified;
- "Matching" tables contain matching entries among two tables;

- User make queries on a "view";
- **It is always possible to trace back entries to their original catalogs;**

# QUBRICS Database

## QUBRICS Database

### DB Content:

- Photometric: Gaia DR3, PanSTARRS1 DR2, DES DR2, SkyMapper DR4, AllWise and CatWISE, SDSS DR16Q, etc. ;
- Several QSO and inactive galaxies catalogs;
- Stars are identified using Gaia's parallax and proper motion measurements;

## QUBRICS Database

### DB Content:

- Photometric: Gaia DR3, PanSTARRS1 DR2, DES DR2, SkyMapper DR4, AllWise and CatWISE, SDSS DR16Q, etc. ;
- Several QSO and inactive galaxies catalogs;
- Stars are identified using Gaia's parallax and proper motion measurements;

### Performance

- Select data for a single QSO: $\sim 1$ ms;
- Query entire table containing coordinates, classifications and redshifts ($8 \times 10^5$ rows) of known QSOs: $\sim 14$ s;
- Query on photometric catalogs is significantly slower ($\sim 10^8$ rows): pre-matched and stored as separate table;

## QUBRICS Database

### DB Content:

- Photometric: Gaia DR3, PanSTARRS1 DR2, DES DR2, SkyMapper DR4, AllWise and CatWISE, SDSS DR16Q, etc. ;
- Several QSO and inactive galaxies catalogs;
- Stars are identified using Gaia's parallax and proper motion measurements;

### Performance

- Select data for a single QSO: $\sim$ 1 ms;
- Query entire table containing coordinates, classifications and redshifts ($8 \times 10^5$ rows) of known QSOs: $\sim$ 14 s;
- Query on photometric catalogs is significantly slower ($\sim 10^8$ rows): pre-matched and stored as separate table;

### Features:

- All tables have indices on their primary keys, as well as on coordinates: access time do not depends on table size;
- Main table modification histories are recorded;
- DB coherence ensured by triggers.

# QUBRICS on TOCats

- Easy visualization of catalogs on the sky;
- Uses multi-depth indexing to **quickly** access catalogs with $\sim 10^8$ entries;
- Quick access to online services and private repository of spectra;
- $\Rightarrow$ talk by L. Nicastro.

## QUBRICS on TOCats

## Summary

- The QUBRICS project already discovered more than 1000 QSO at $z > 2.5$ (and more than 200 at $z > 4$) at $Y \lesssim 19.5$;
  - Scientific exploitation: luminosity function, cosmological re-ionization, Sandage test, etc.;
- Several machine learning methods adopted: CCA, PRF, XGBoost;
- Dedicated method to deal with severely imbalanced datasets (Calderone+24);

- QUBRICS database is a key part in the project;
- Storage for machine learning project *may* be smaller than for instrumentation...
  - ...and typically has a clear "structure" (feature columns) $\Rightarrow$ relational DB is the perfect tool!
- By using pre-matched tables (e.g. on coordinates), indices, views, etc. we built a highly responsive DB able to support all project activities;
- TOCats is the perfect companion for our DB, to quickly visualize literature data and prioritize observations;

## Summary

- The QUBRICS project already discovered more than 1000 QSO at $z > 2.5$ (and more than 200 at $z > 4$) at $Y \lesssim 19.5$;
  - Scientific exploitation: luminosity function, cosmological re-ionization, Sandage test, etc.;
- Several machine learning methods adopted: CCA, PRF, XGBoost;
- Dedicated method to deal with severely imbalanced datasets (Calderone+24);

---

- QUBRICS database is a key part in the project;
- Storage for machine learning project *may* be smaller than for instrumentation...
  - ...and typically has a clear "structure" (feature columns) $\Rightarrow$ relational DB is the perfect tool!
- By using pre-matched tables (e.g. on coordinates), indices, views, etc. we built a highly responsive DB able to support all project activities;
- TOCats is the perfect companion for our DB, to quickly visualize literature data and prioritize observations;

## Summary

Future challenges:

- Add further photometric data to the NIR (e.g. VISTA, Euclid, etc.);
- Optimize ML selection;
- Probe redshifts up to $\sim$ 5.5 - 6, with completeness $\gtrsim$ 90%;

QUBRICS activities range from technical topics, such as DB management and machine learning, to the science exploitation.

**If you're interested in learning/collaborating, feel free to reach us. Students are very welcome! ;-)**