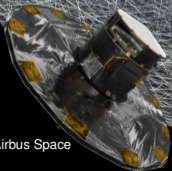# Gaia: the two (well, almost three) billion star surveyor

Anthony Brown

Leiden Observatory, Leiden University
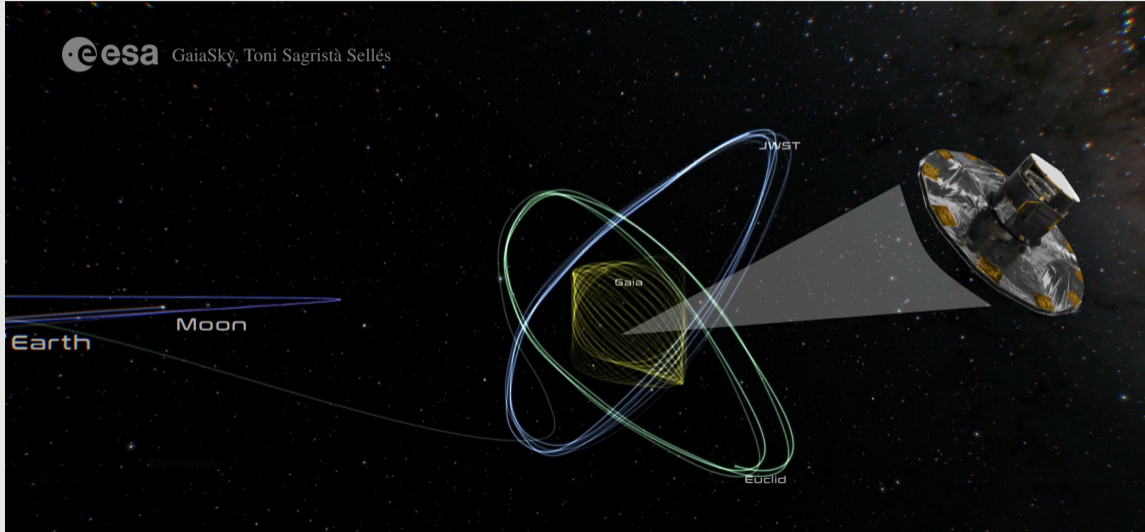
brown@strw.leidenuniv.nl

Airbus Space

ESA/Gaia/DPAC

# Gaia collects fundamental astronomical data
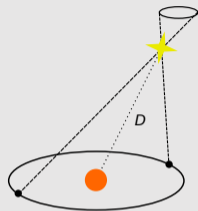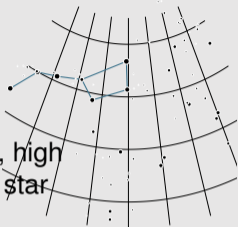


Parallaxes and proper motions

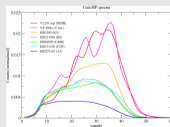All-sky, complete, high accuracy star atlas

Astrometric, photometric, spectroscopic, radial velocity time series

Astrophysical properties

Brown AGA. 2021
*Annu. Rev. Astron. Astrophys.* 59:59–115

# Gaia data collection and processing

**a** Spin · P · Γ · F · Combined field of view

**b** Fiducial observation line (×2) · CCD on sky · AC · AL

**c** CCD pixel stream · μ · κ

Brown AGA. 2021
*Annu. Rev. Astron. Astrophys.* 59:59–115

Unpack telemetry → Project pixel stream on sky → Group observations by source → Astrometric, photometric, spectroscopic processing → Source astrophysical characterization → Gaia data release

Iterate

# How big is Gaia data?

| CURRENT DATE AND TIME | 2024-01-21T19:04:49 (TCB) |
|---|---|
| **MISSION STATUS** | |
| Satellite distance from Earth (in km) | 1,386,884 |
| Number of days having passed since 25 July 2014 | 3467 |
| Number of days in mission extension | 1650 |
| **OPERATIONS DATA (collected since 2014/07/25)** | |
| Volume of science data collected (in GB) | 127,376 |
| Number of object transits through the focal plane | 241,170,660,648 |
| Number of astrometric CCD measurements | 2,377,253,654,952 |
| Number of photometric CCD measurements | 478,159,154,700 |
| Number of spectroscopic CCD measurements | 47,043,571,347 |
| Number of object transits through the RVS instrument | 15,799,464,121 |

# HPC, numerical algorithms, machine learning, AI(?)

- Computing is done on large clusters in six data processing centres
- Data centres communicate via central hub: 'Main data base'
- Data releases are extracted from MDB contents



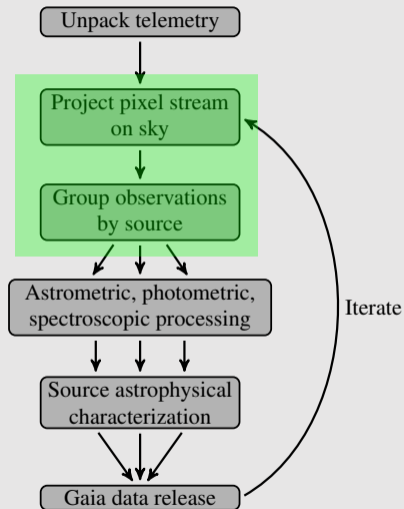MareNostrum 4 @ Barcelona Supercomputing Center



Unpack telemetry

↓

Project pixel stream on sky

↓

Group observations by source

↓

Astrometric, photometric, spectroscopic processing

↓

Source astrophysical characterization

↓

Gaia data release

Iterate

# HPC, numerical algorithms, machine learning, AI(?)



- Sophisticated clustering algorithm to associate observations with sources
- Needs to account for conflicting matches, fast moving sources, variable sources, spurious sources



Unpack telemetry

Project pixel stream on sky

Group observations by source

Astrometric, photometric, spectroscopic processing

Source astrophysical characterization

Gaia data release

Iterate

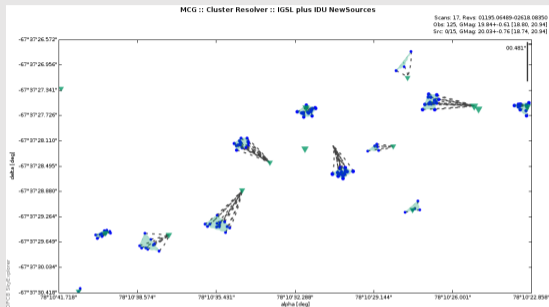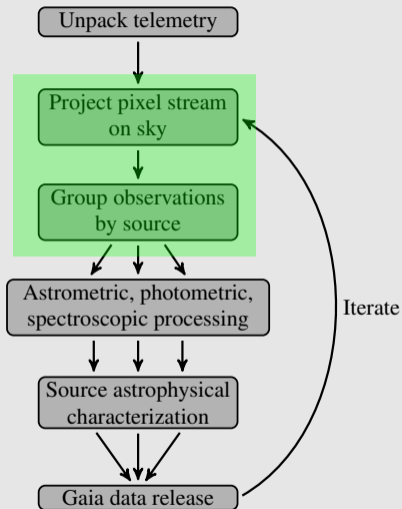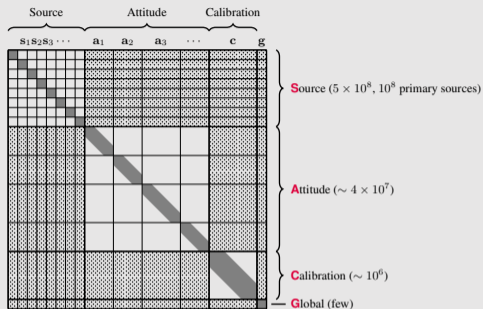# HPC, numerical algorithms, machine learning, AI(?)
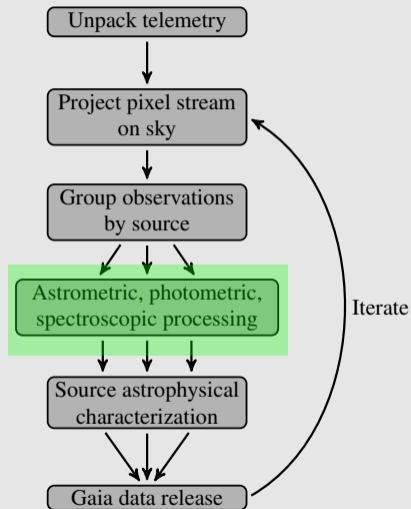


- Sophisticated clustering algorithm to associate observations with sources
- Needs to account for conflicting matches, fast moving sources, variable sources, spurious sources
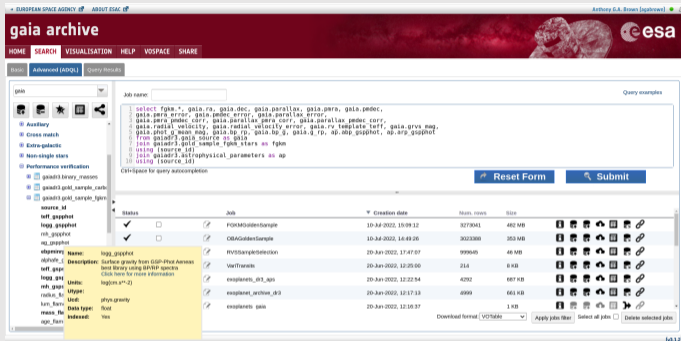


Unpack telemetry

Project pixel stream on sky

Group observations by source

Astrometric, photometric, spectroscopic processing

Source astrophysical characterization

Gaia data release

Iterate

- Large systems of linear equations; iterative least squares solution
  - calibration terms lead to non-sparse matrices
- Also classical astronomical data processing methods (e.g., cross correlation of spectra with templates to derive radial velocity)
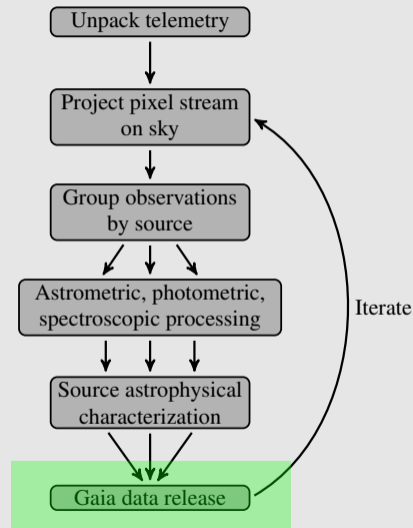
# HPC, numerical algorithms, machine learning, AI(?)
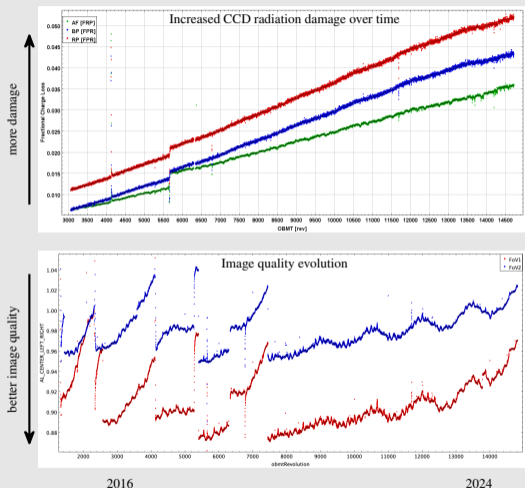


Credit: Eyer et al. (2018)
Adapted from: Eyer & Mowlavi (2008)

- Machine learning algorithms: classification, clustering, outlier detection, period searching in time series, ...
- Inference algorithms: MCMC, likelihood optimization, ...

# HPC, numerical algorithms, machine learning, AI(?)



- Data access through ADQL, TAP, Python modules, bulk download
- Extensive documentation, tutorials, code examples
- ESA + Several partner data centres



Unpack telemetry

↓

Project pixel stream on sky

↓

Group observations by source

↓

Astrometric, photometric, spectroscopic processing

↓

Source astrophysical characterization

↓

Gaia data release

Iterate

# Specifics of the Gaia data processing

- Of order $3 \times 10^{12}$ individual observations
  - ▶ they all count!
- Observations are connected over large angles on the sky
- Data access: spatially or temporally grouped, time series per source
- All instruments are self-calibrated
  - ▶ Complex calibration models with millions of parameters
- Instrument characteristics evolve over time
- Division of data processing tasks and physical location was partly political
- Most of this is not unique to Gaia...



Increased CCD radiation damage over time



Image quality evolution

2016

2024

# Gaia and Big Data: the practical challenges

- Develop and maintain a well-documented data model
- Coordination across DPAC units (JIRA, e-mail, many telecons. . . )
- Organizing data ahead of a processing run
- Data transfer between processing centres
- Develop/ test/ validate new or improved pipeline code
- Validation of the data products
  - are they scientifically correct?
  - no missing or out of bounds values?
  - what to do if we find a problem: fix, or document?
  - . . .
- Transfer of data from the main data base to the public archive
  - make sure data mapping is done carefully and without mistakes
- Documentation for a data release