# Dataspex

# big data management solutions for astronomy

ADASS 2019
Groningen, The Netherlands
6 – 10 October 2019

Bart Scheers, Arjen de Rijke

## Spin-off from CWI's Database Architectures Group

- Bart Scheers
- Arjen de Rijke
- MonetDB
  The open-source column-store relational database system
- CWI Inc.
  The Dutch national research institute for mathematics and computer science

Tomorrow's telescope will be a "total telescope", where the overall design integrates telescope hardware with computer hardware and software; all observations will be inspected, stored and queried very quickly for scientific analysis.
The science is in the data and smarter management of the data will raise your science to the next level.

Dataspex is specialised in managing big data streams on modern hardware and in finding the right solutions for your system.
We can assist you with the design, software development, implementation or improvements of such systems.

## Modern Telescope Back-End System

### New instruments

Radio:
- LOFAR – MeerKAT
- ASKAP – SKA

Optical:
- MeerLICHT – BlackGEM – LSST

Common challenges and overlapping strategies:
- Automated pipeline between telescope and database
- Inspect data streams for transient and variable events
- Database is wealthy laboratory for doing statistics and complementary science

Integrate DB technologies to reach science goals:
- Full-source database for complementary science
- Move algorithms and statistics inside database engine
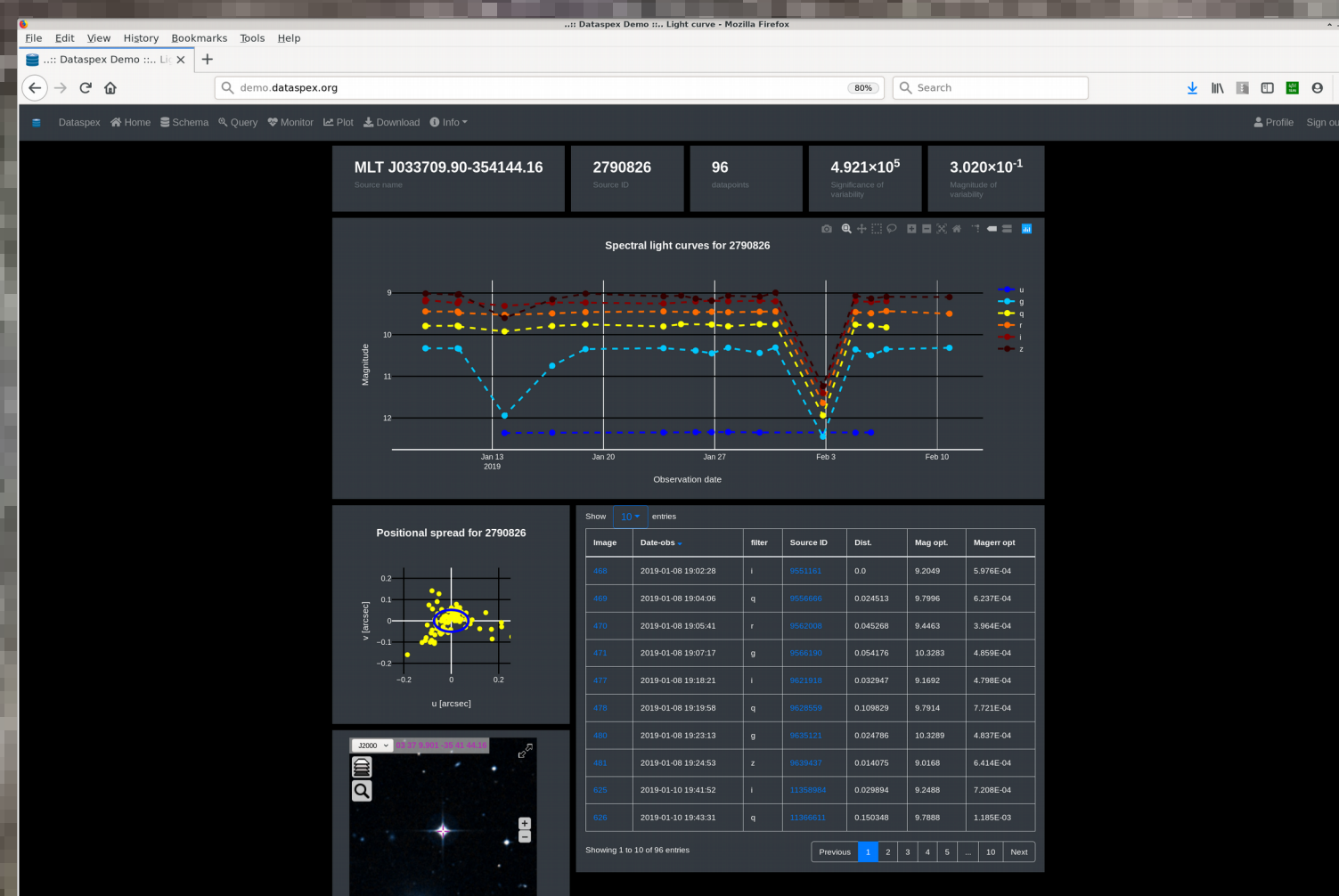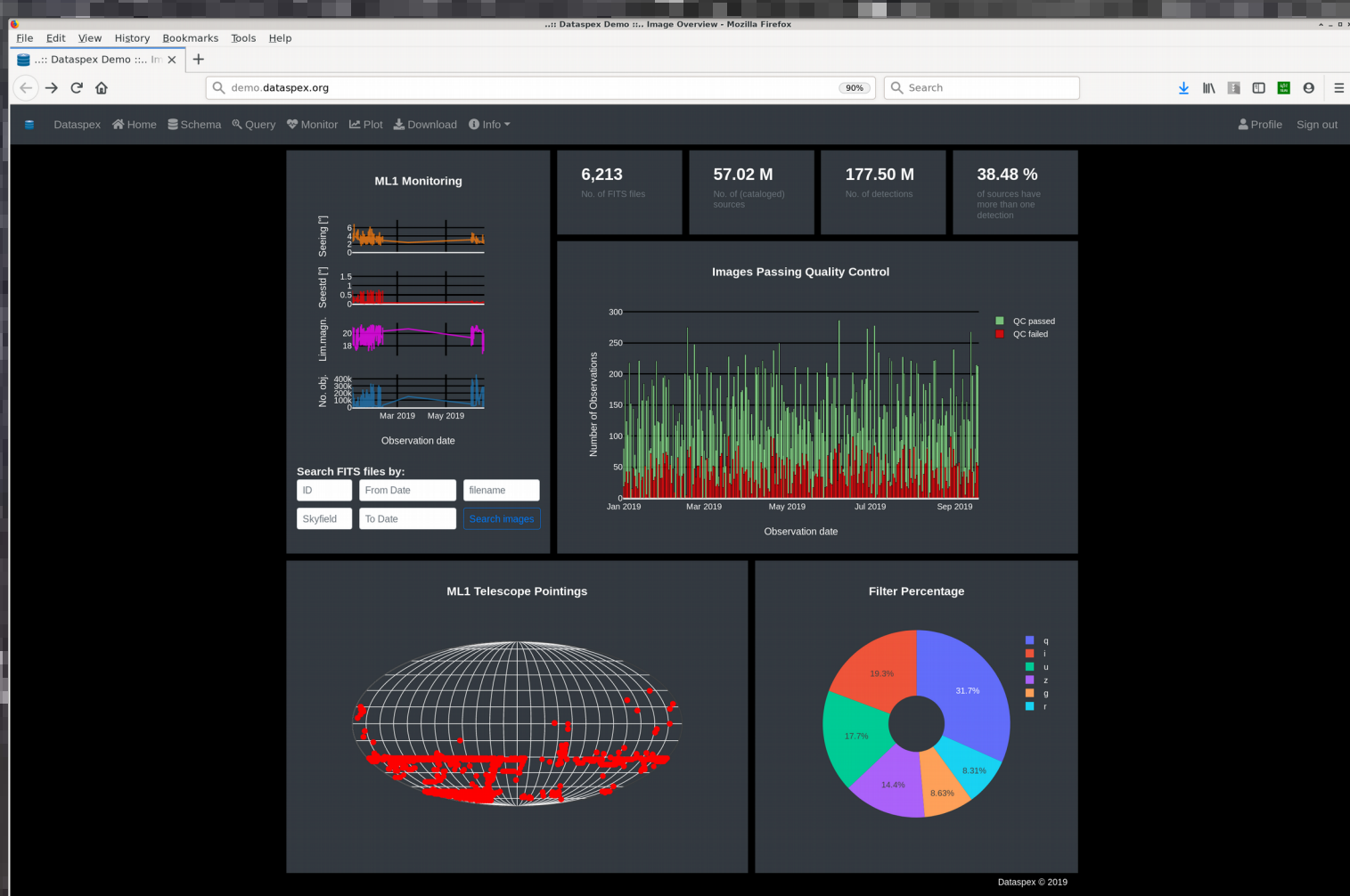- Interactive visualisation

### High-cadence Astronomy

High-cadence astronomy is a relatively new field in observational astronomy. Advances in hardware and software technology have made it possible to stream large volumes of observational data over fast links to clusters of computers that process the data in one or more automated pipelines for scientific analysis.

The time available to do real-time analysis is limited by the cadence of the instrument. Therefore, additional and complementary scientific data analyses are forced to shift to non-real time environments. Here, all data accumulates over time and the growth may vary in the range of 0.1–100 PB/yr.

### Data Processing Pipelines

We have centered our astronomical data pipeline around the open-source database system of MonetDB. Processing an image's source list – i.e. storing, cross-matching, cataloguing 100,000s of sources - involves about 100 SQL queries. While the pipeline keeps on processing the high-cadance data, the database size keeps on growing dynamically.
In order to keep the processing constant in time and independent from the database size, we partition tables in multiple dimensions and distribute them over multiple nodes.
This allows to parallelise pipeline runs and scale up to even process data from the largest (arrays of) telescopes.

## Cloud Solutions

### The main pillars for running Big Data Pipelines

- Partitioning
- Distributing
- Parallelising

### Partitioning

Partitioning source tables according to their sky location along the great circle of their declination serves as a natural index. Secondary partitioning dimension is the time domain.

### Distributing

In a Cloud environment we can easisly distribute the partitioned source tables without replication to specific storage and processing nodes. The isolated parts are ideal for parallel processing.

### Parallelising

The data processing pipeline can be split to process only dedicated portions of the data that fit to the chosen distributed partitions.
In this way the pipeline can run multiple instances in parallel.

## Long-Term Production

### DevOps & CI

Dataspex combines its hardware and software expertise with DevOps to manage the big data streams from modern telescopes.
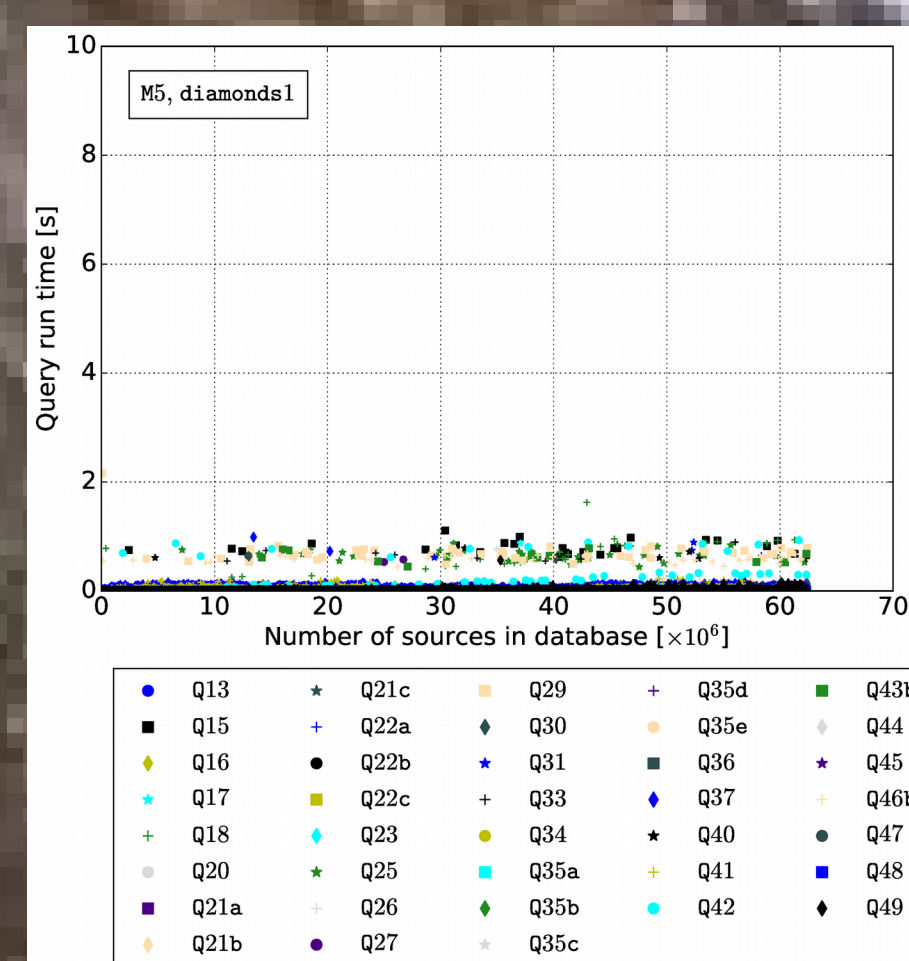
All code repositories are configured to automatically build, test and deploy applications based on strict Continuous Integration configurations.

### Pipeline & Application Code

Pipeline and application code is written for the longer term.
Therefore, Dataspex is commited to the following practices for code quality:

- Reusability of code
- Tests
- Automated deployment
- Microservice architecture
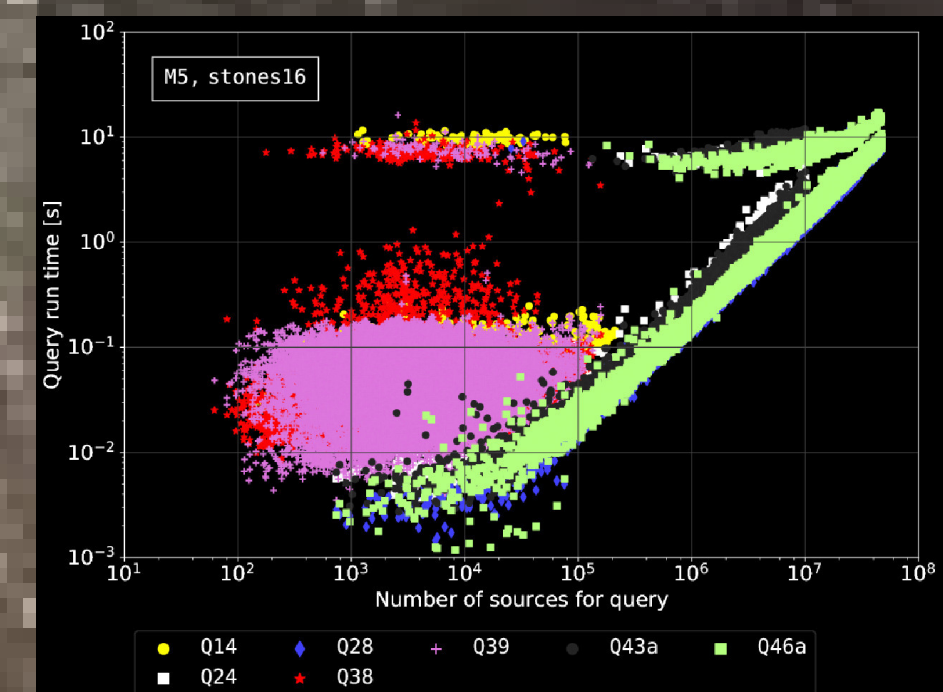- Continuous performance monitoring

## Cost Reductions, Efficiency Improvements

### Constant Query Run Times

Most pipeline queries are independent of the database size. The graph on the left clearly shows the constant-time complexity O(1). This means that the processing is scalable for the longer runs, where the database size continues to grow.

### Linear Query Run Times

For some pipeline queries there is a linear relationship between the number of sources that have to be processed and the query run times (see the graph on the left). Based on this ratio we can reliably estimate the characteristics of the data partitions.

This allows us to distribute the data over multiple partitions, optimising the processing by running several pipelines in parallel.

### Distinction

Dataspex is a young company, but its founders have more than 20 years of experience in the research fields of astronomy and computer science.
Its customers – research institutes, universities, consortia – benefit from a highly skilled team, that implements specific "big data" solutions much more efficiently than in the traditional way.
Distinguishing features are:

- Expertise
- Efficiency
- Code robustness
- Stable production software

## Info

bartscheers@dataspex.nl

https://www.dataspex.nl