



Gaia Spark analysis platform

Progress update

D Morris July 2021





Original GDAF system
Zeppelin & Spark deployed on physical hardware

Initial goal for our project

- replicate GDAF in the cloud
- DONE



Subsequent goals

- scalable deployment
- portable deployment
- more users
- more data
 - DR3
 - DR4



Gaia Data Analytics Framework (GDAF)
description

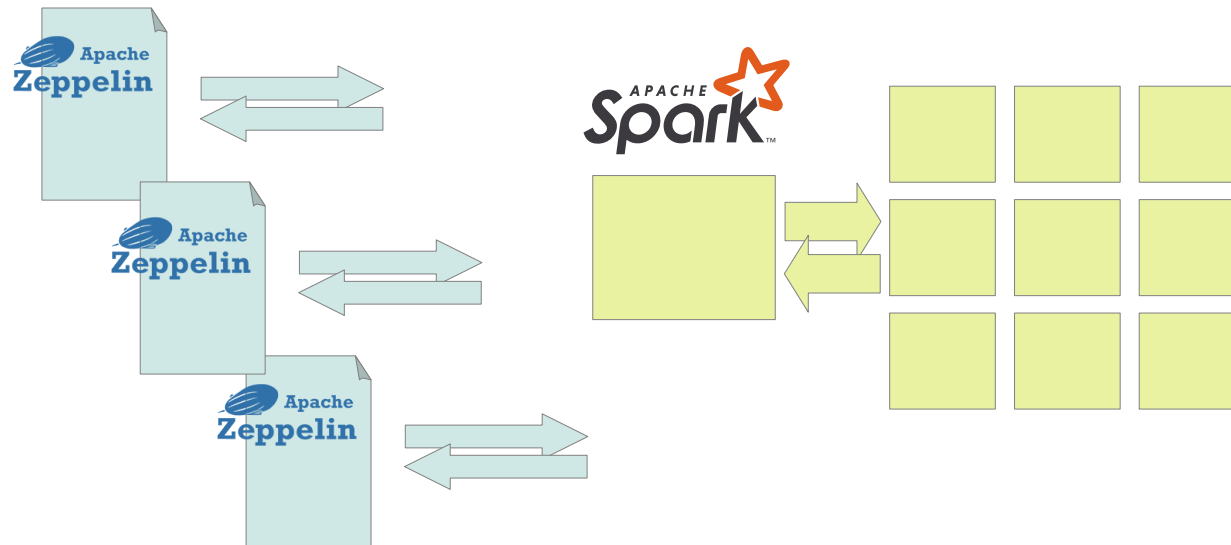




"known unknowns"

Hadoop/Yarn

- Spark cluster deployed on static resources
- Zeppelin notebooks all interact with the same Spark cluster



- Automated with Ansible



99% automated

- create-all
- delete-all

3 deployments

- dev
- test
- live



- In place and working
- eDR3 with neighbors

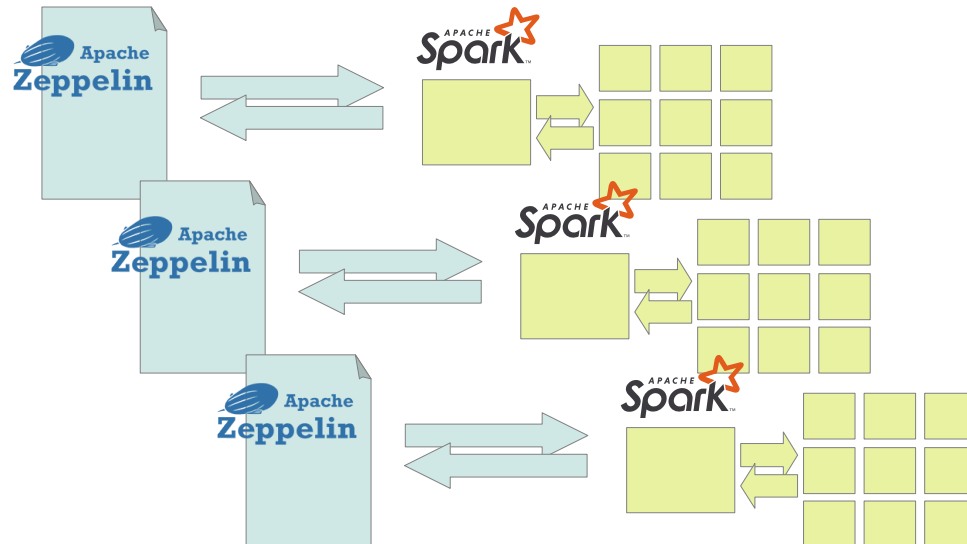




kubernetes

“unknown unknowns”

- Spark cluster on demand
- Notebooks launch their own Spark cluster



- Automated with Helm



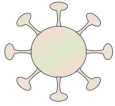
99% automated

- create-all
- delete-all

3 deployments

- dev
- test
- live

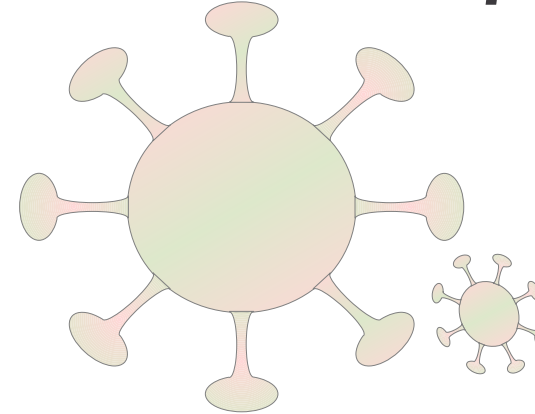
- Experimental in 2020/21
- Not ready for production



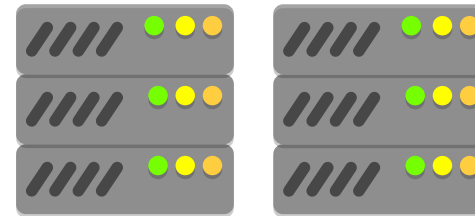
COVID-19



- Cloud based notebook service
- Remote working, ideal application



- Limited access to the data center
- **Huge thank you to the system admins who kept the system running**
- Impact on hardware deployment schedules
- 2021/2022 resources and upgrades will be delayed



UNIVERSITY OF
CAMBRIDGE

HPC cloud



Parquet

<https://parquet.apache.org/>

Apache Parquet columnar storage format

- Gaia eDR3 sources - 561Gbytes
- 2MASS PSC 37G bytes
- 2MASS PSC Gaia eDR3 best neighbors 60G bytes
- Pan-STARRS MeanObjectView 270G bytes
- Pan-STARRS Gaia eDR3 best neighbors 163G bytes
- ALLWISE 341G bytes
- ALLWISE Gaia eDR3 best neighbors 177G bytes



Cross match using best
neighbor tables

Familiar SQL based
JOIN syntax

```
SELECT
    gaia.source_id,
    gaia.ra, gaia.dec,
    ps1.g_mean_psf_mag AS ps1_g,
    ps1.r_mean_psf_mag AS ps1_r
FROM
    gaia_source AS gaia
INNER JOIN
    gaia_source_ps1_best_neighbours AS ps1
ON
    gaia.source_id = ps1.source_id
```

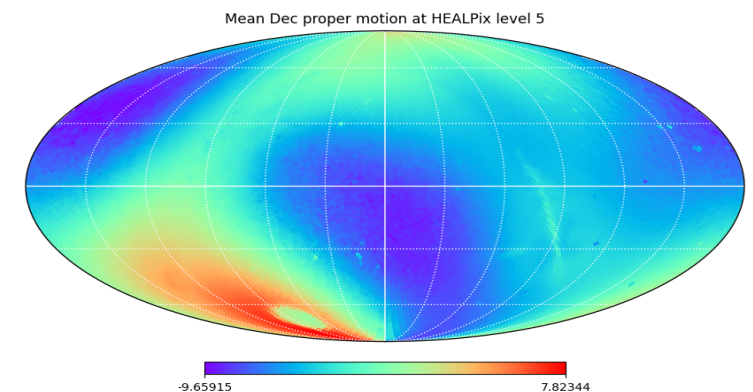
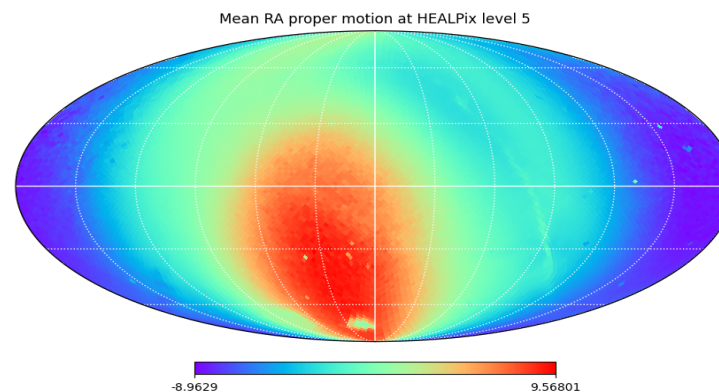


HEALPIX partitioning

Parquet files partitioned based on HEALPIX value embedded in Gaia source_id

Placing adjacent sources in the same file reduces shuffle between Spark workers

```
SELECT
  floor(source_id / 562949953421312) AS hpx5,
  COUNT(*) AS n, AVG(pmra), AVG(pmdec)
FROM
  gaia_source
GROUP BY
  hpx5
```



Mean proper motions over the sky – 1min 28sec to calculate and plot

Machine learning application

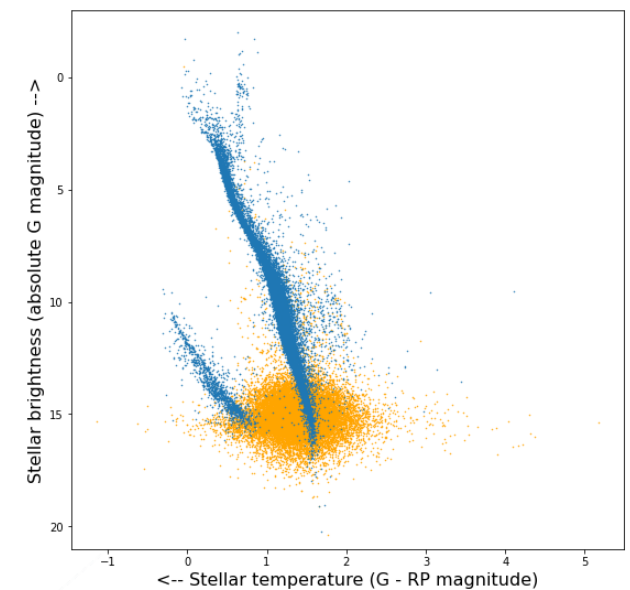
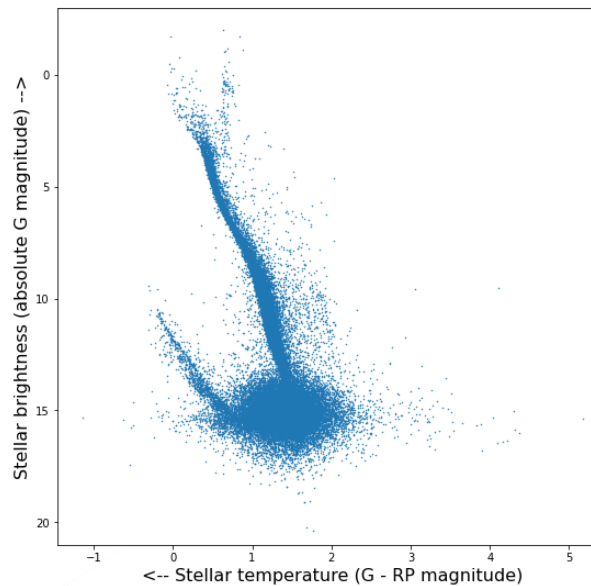
Based on the Gaia EDR3 performance verification *"The Gaia Catalogue of Nearby Stars"* (Smart et al. 2021).

Training a supervised Random Forrest to classify astrometric solutions as 'good' or 'bad'.

SparkSQL queries to generate the training and validation data.

4min to train the classifier

25sec to classify 1,724,028
sources and plot the results





Network filesystem

CephFS via Openstack Manila
to store the science data



Issues with concurrent access and IO bandwidth

ML RandomForest example :

- First pass with empty cache, 4min to train the classifier
- Second pass with populated cache, 24min to train the classifier

Pattern of disc access changes depending on the cache contents.

- Empty cache generates sequential requests for large blocks
- Populated cache generates random access requests for small fragments

Current CephFS system does not handle this pattern well

Solving this may take time – see COVID-19 and hardware deployment





Beta testing

Call for people interested in being beta testers

Booked sessions with tech support

Q4 2021

Outline of the science case you would like to explore

Outline of the technology or algorithm you are interested in using

Contact a member of the development team

- Nigel Hambly <nch@roe.ac.uk>
- Dave Morris <dmr@roe.ac.uk>
- Stelios Voutsinas <stv@roe.ac.uk>

Prizes for the most likely to break the system

