# Efficient Catalog Cross-matching Sky Mining Hackastron 2018

Team members:Kristof, Karl & Casey

# Our Goal

## Take the 50 catalogues

- Cross match them (49 comparisons) and create light curves with the unique measurement in each epoch
- Report the row ID for each source in each epoch
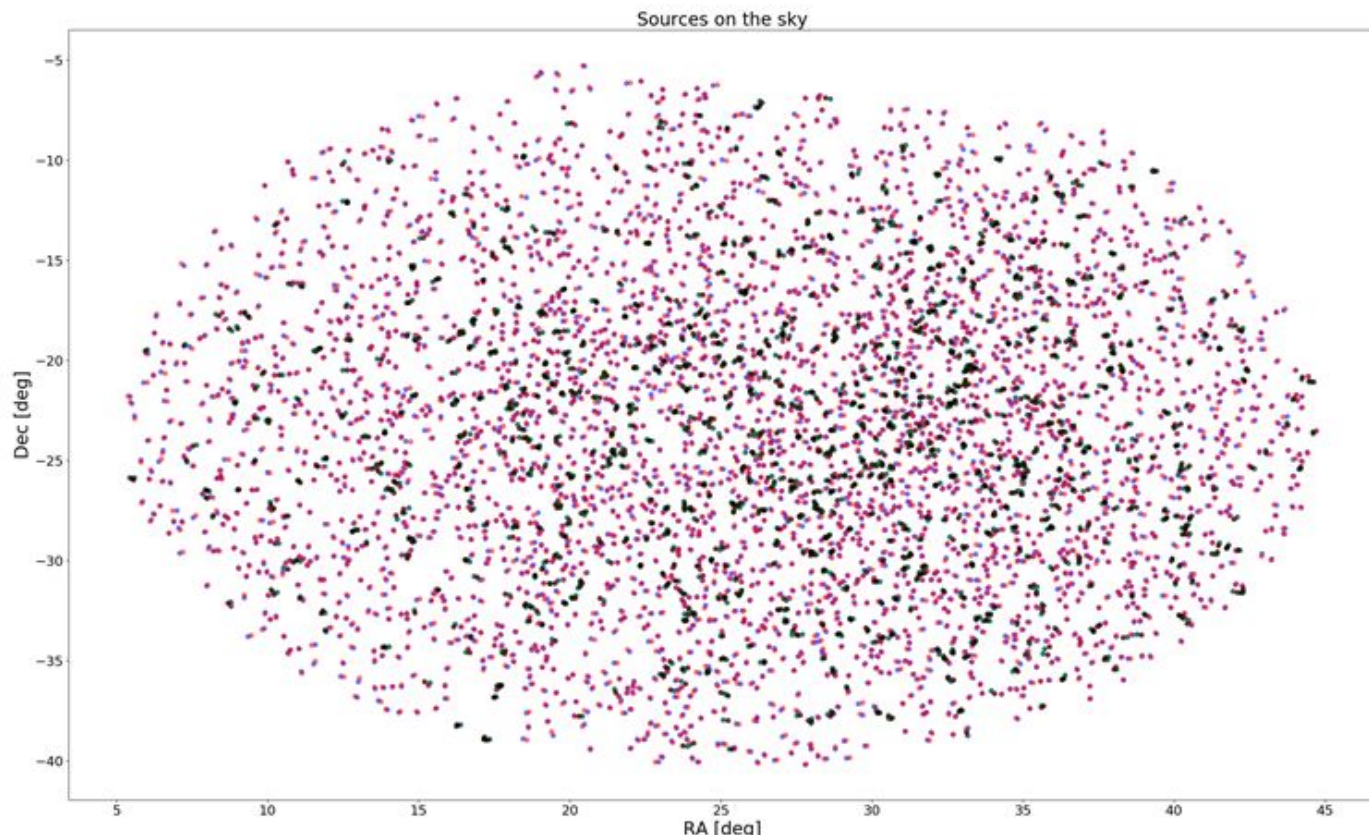- Sort the answer key by epoch00 row ID

# Our Approach

- Simple check: if have 'Clear Neighbour', no further analysis

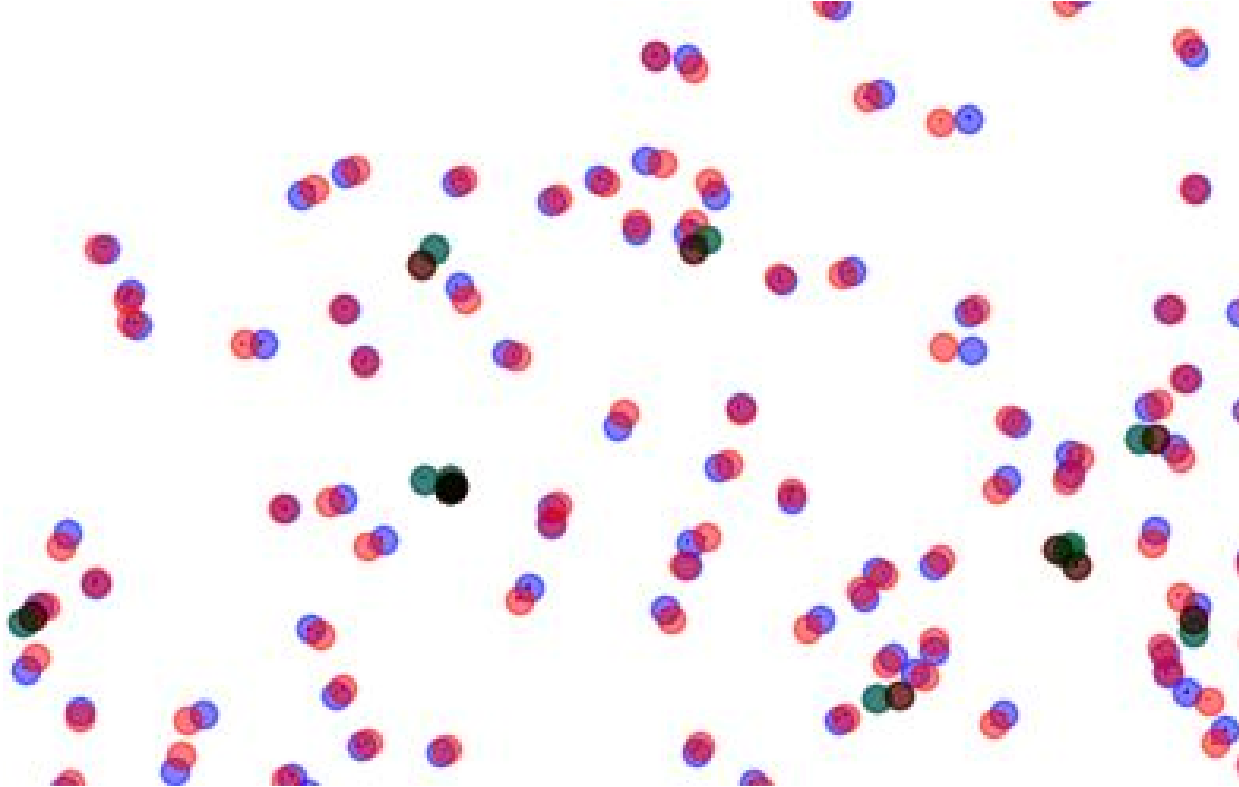## Supplemented with

- Complex check: Hungarian Tinder

# Simple check: 'Clear neighbour', Unoptimised= 77%.

Sources on the sky

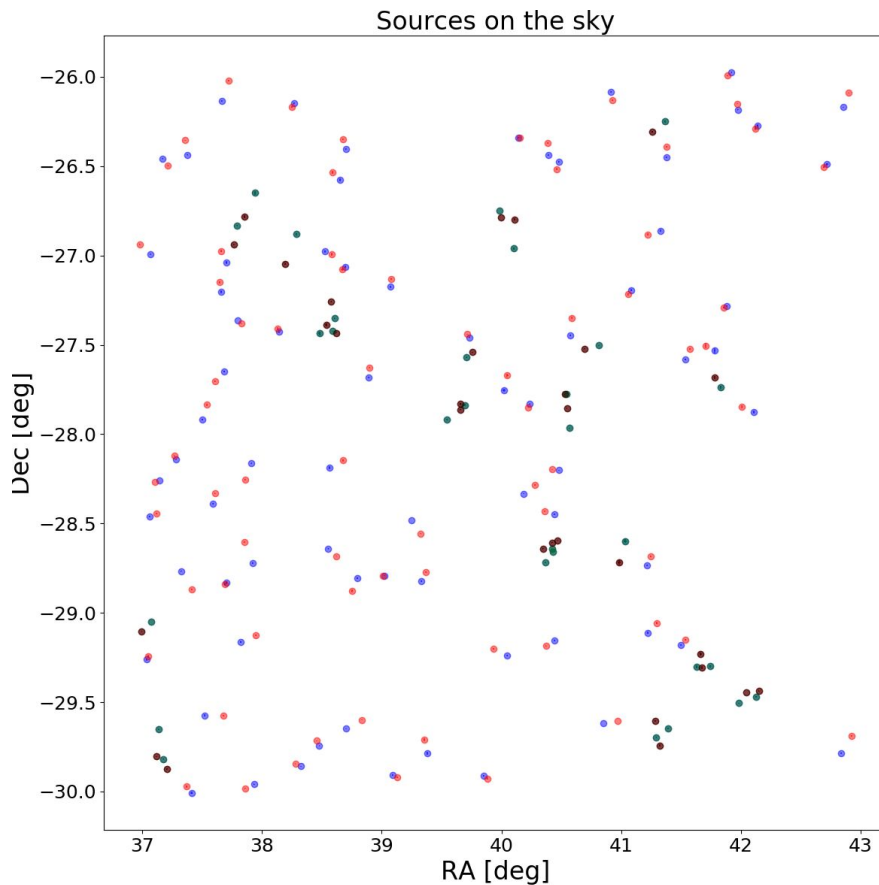1. Neighbr1 vs Nbrs 2 distance ratio
2. Checks for duplicate neighbrs

Ambiguous = green & black.

# Clear neighbour - zoomed in

Ambiguous = green & black.

# Clear neighbour filter on small dataset of 100.



Sources on the sky

Epoch 00 vs Epoch 01 = 72 unambiguous.

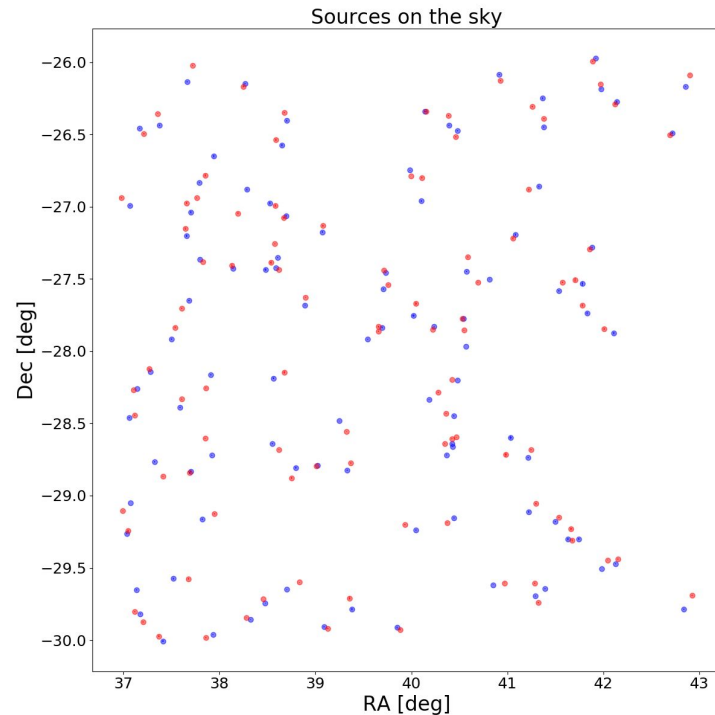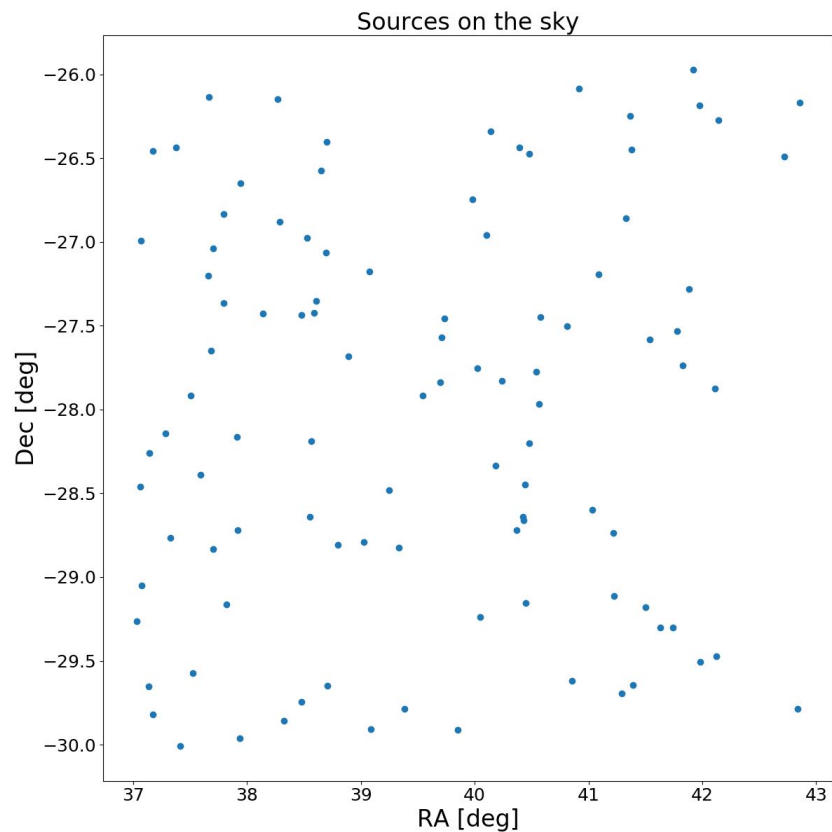Epoch 00 vs 01 vs… vs 48 vs 49 all unambiguous = 16.

6

# Complex Check
## Tinder for galaxy positions: the Hungarian algorithm

- Minimizes the cost for an NxN matching problem
- Need to introduce a metrics, for matching an observed galaxy position to a galaxy in our sky model.
- **The metrics need to be dimensionless!**
- <u>**Our metics:**</u>

  We measure the probability of a point, that drawn from a 3D distribution defined by our sky model, is closer to the distribution center than the observed position.

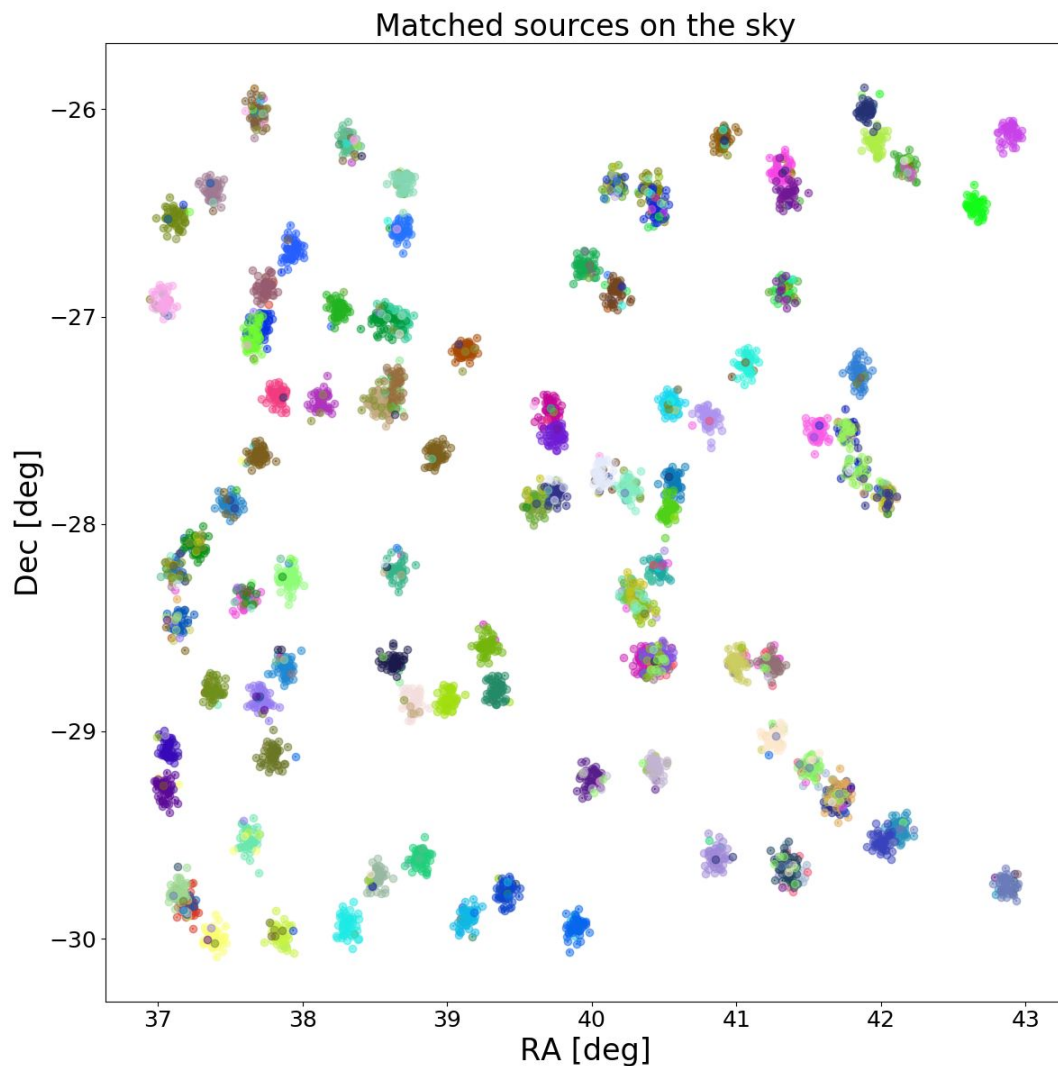  - We assumed Gaussian distribution in each marginal direction

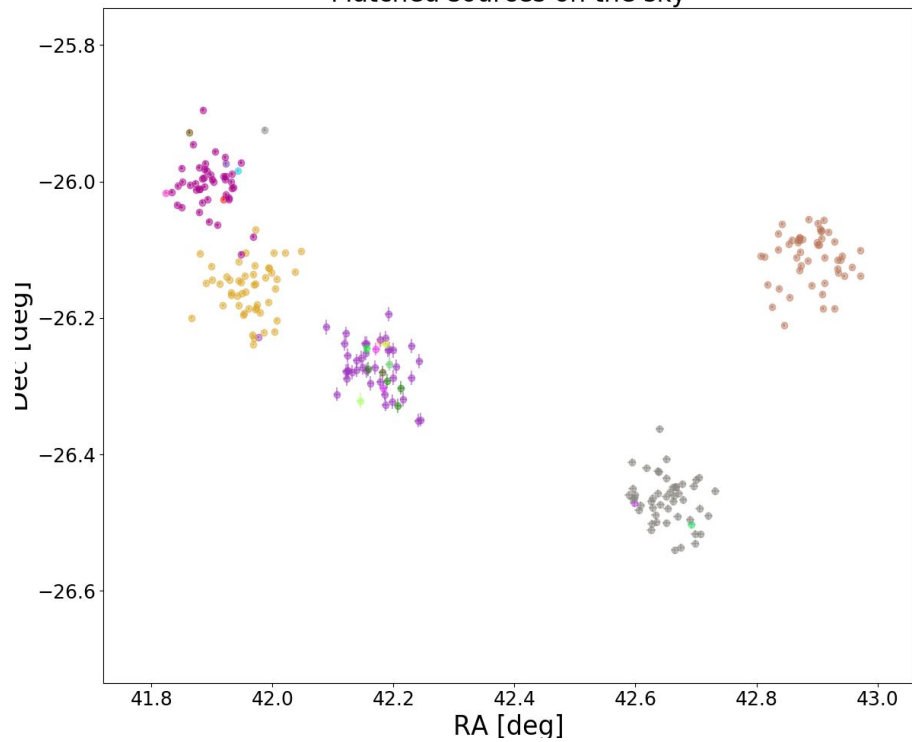# Initial model

Sources on the sky



Sources on the sky



- We set our initial model based on the first observation
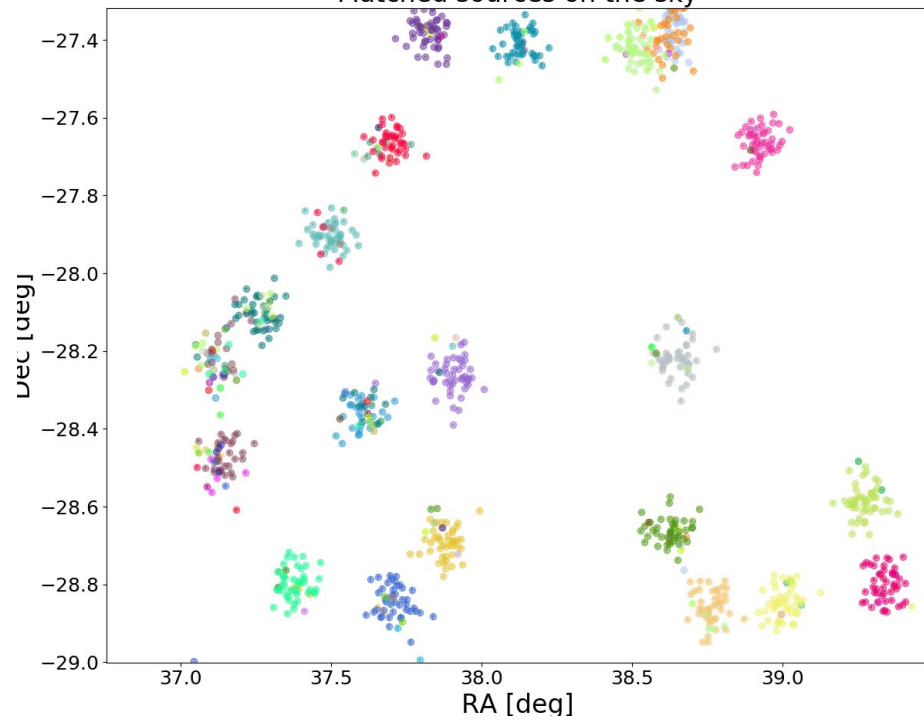- We updated the model in while processing observations

# Final results

- We identified the isolated galaxies, but we had some issues with the cluster galaxies.
- Unfortunately, we couldn't quantify the algorithm preciseness
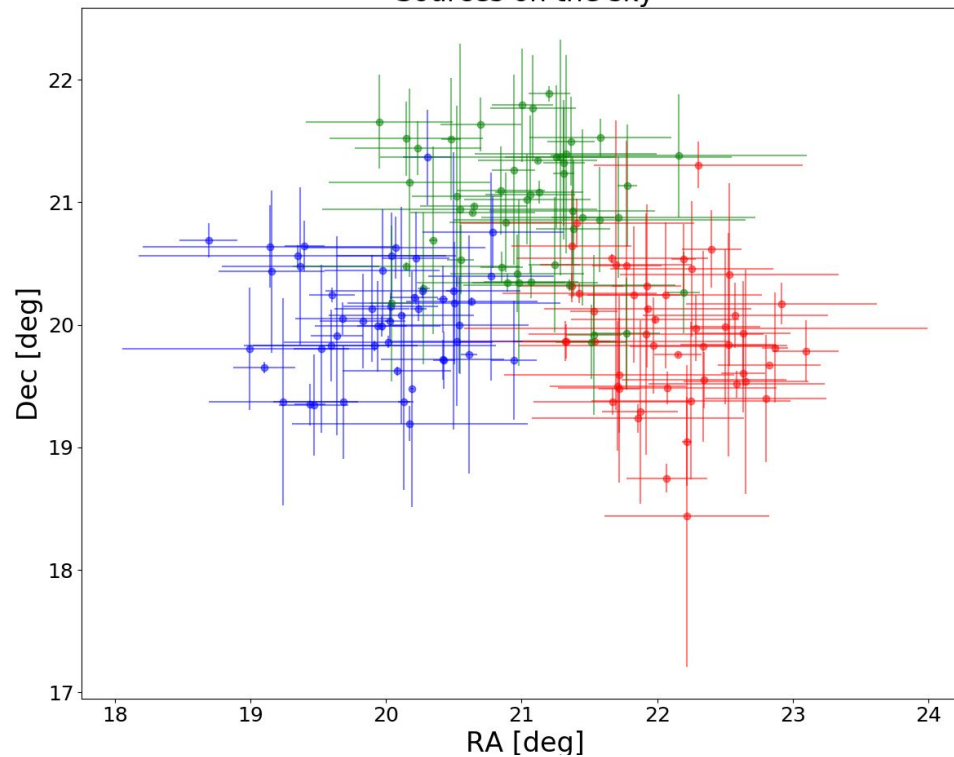


Matched sources on the sky
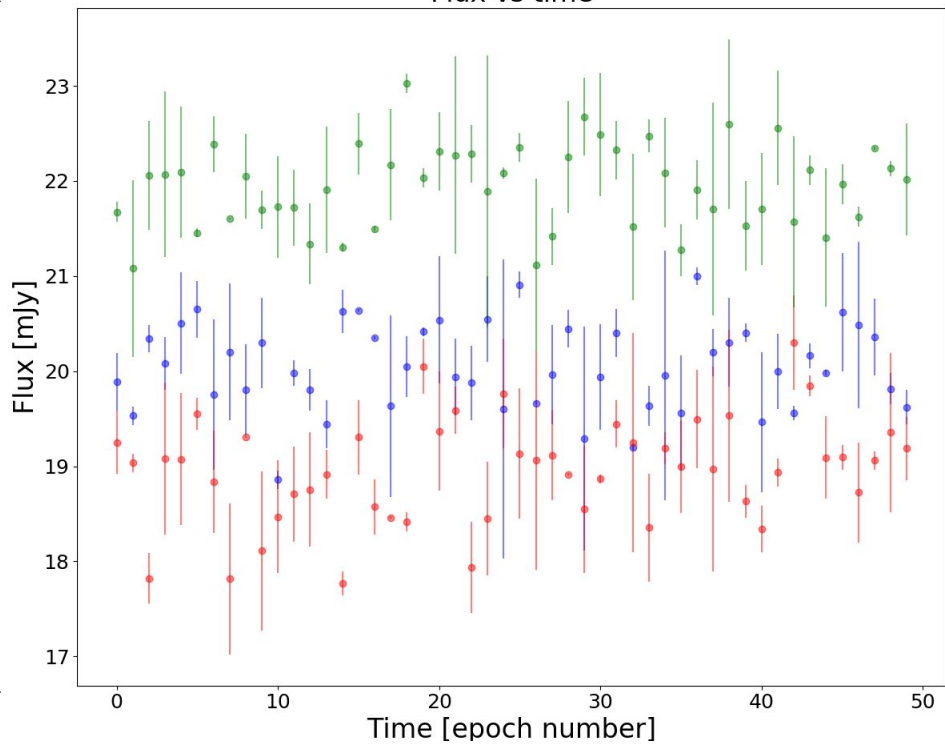
Matched sources on the sky

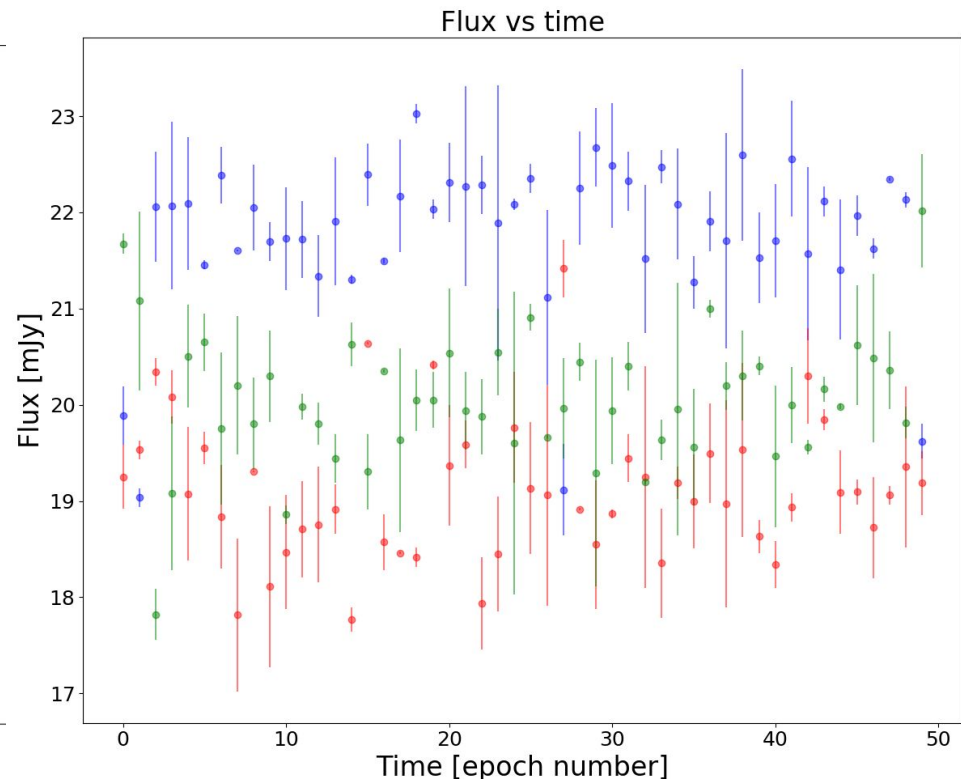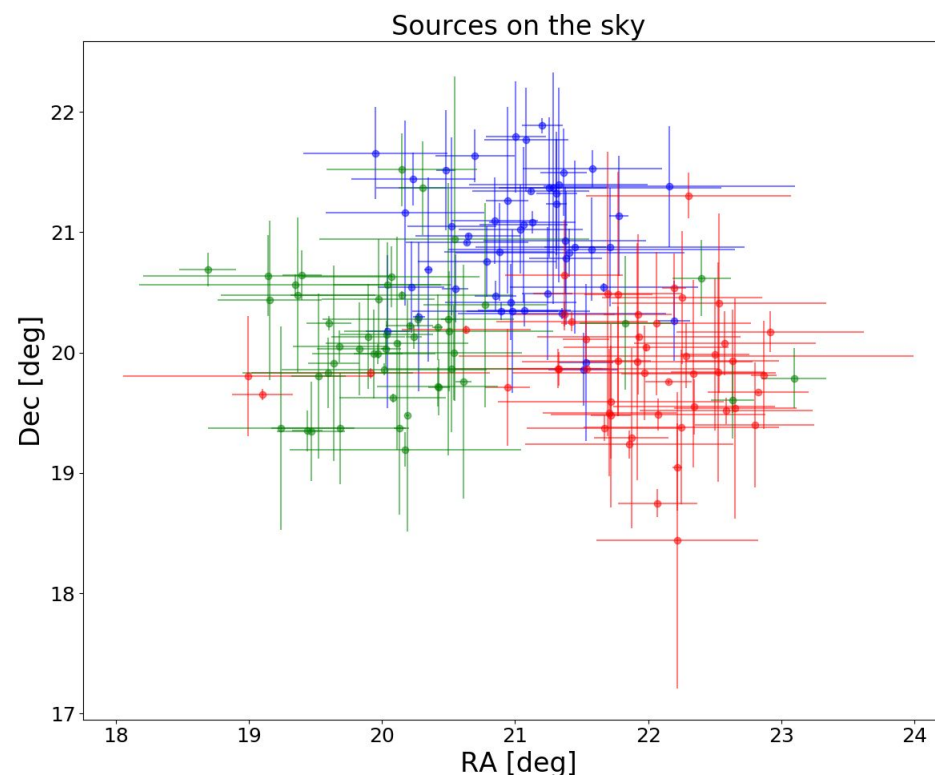# Zoom in simulation, issues & future improvement

- We simulated 3 galaxies close to each other to quantify the algorithm efficiency
- We find that the method worked in ~95% of the data (43 out of 50)
- However the method scales with O(N^2)
- A possible solution would be to use the **Bellmann-Ford algorithm**

Sources on the sky — Flux vs time

# Future Plans

1. Quantify Tinder comparison

2. Optimise Simple
   a. Already checks for false positives; add check for true negative

3. Combine:
   a. For each comparison, combine output of both algorithms to maintain 90%.