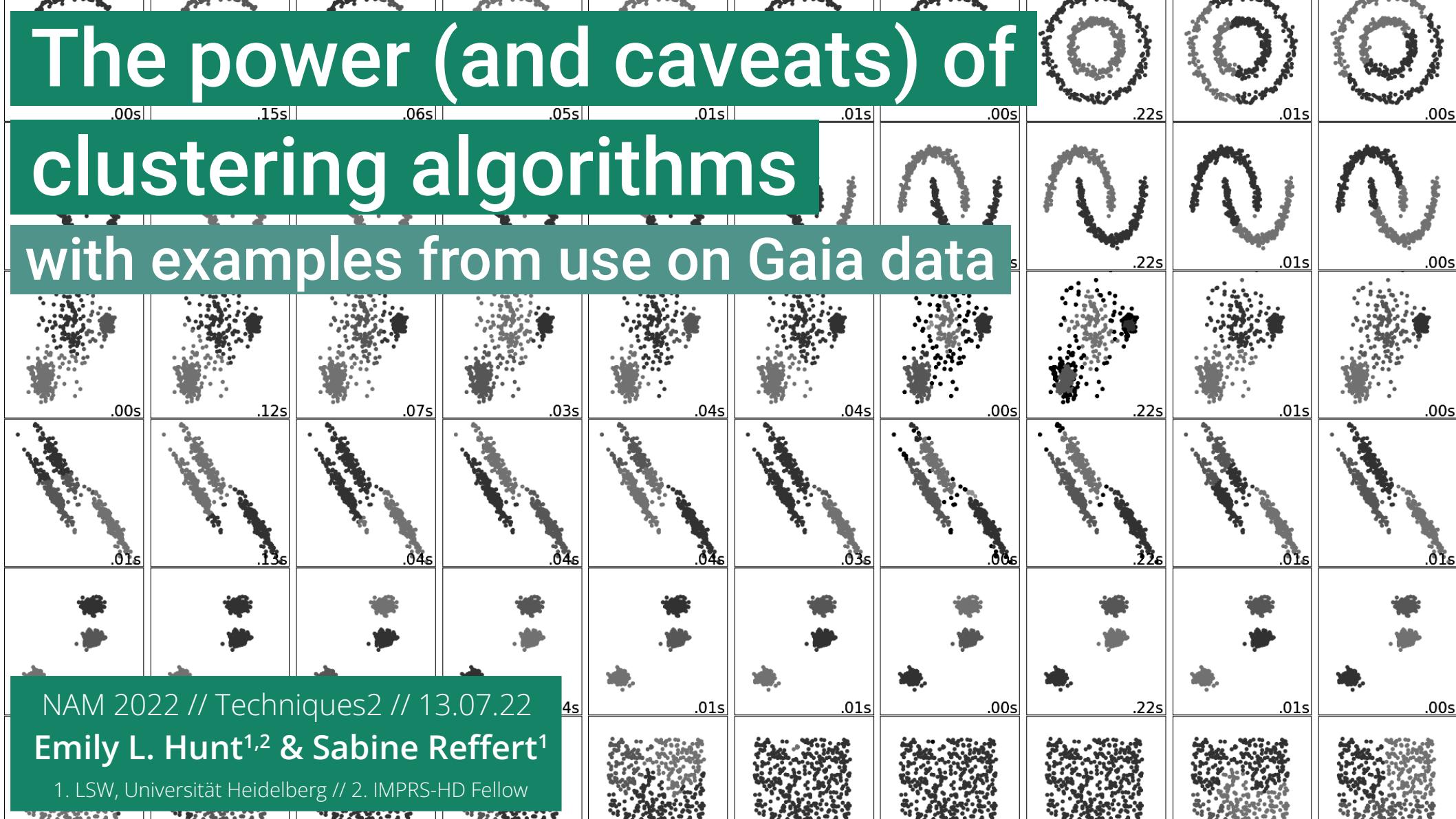


# The power (and caveats) of

## clustering algorithms

with examples from use on Gaia data



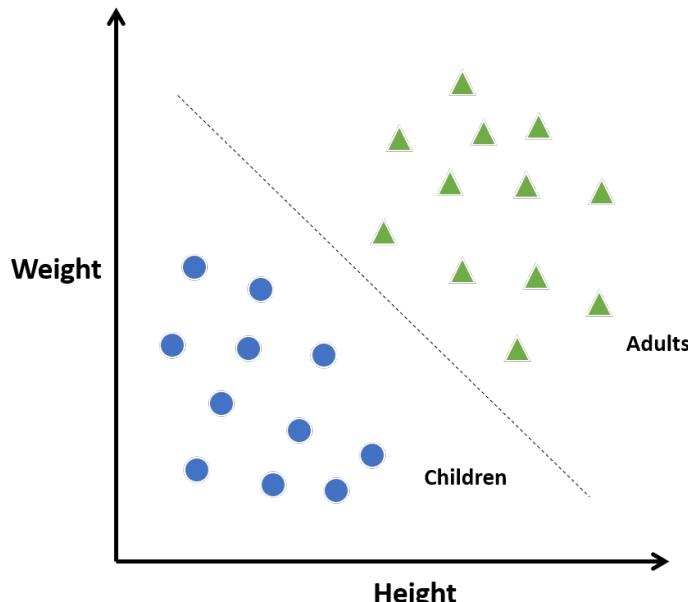
NAM 2022 // Techniques2 // 13.07.22

Emily L. Hunt<sup>1,2</sup> & Sabine Reffert<sup>1</sup>

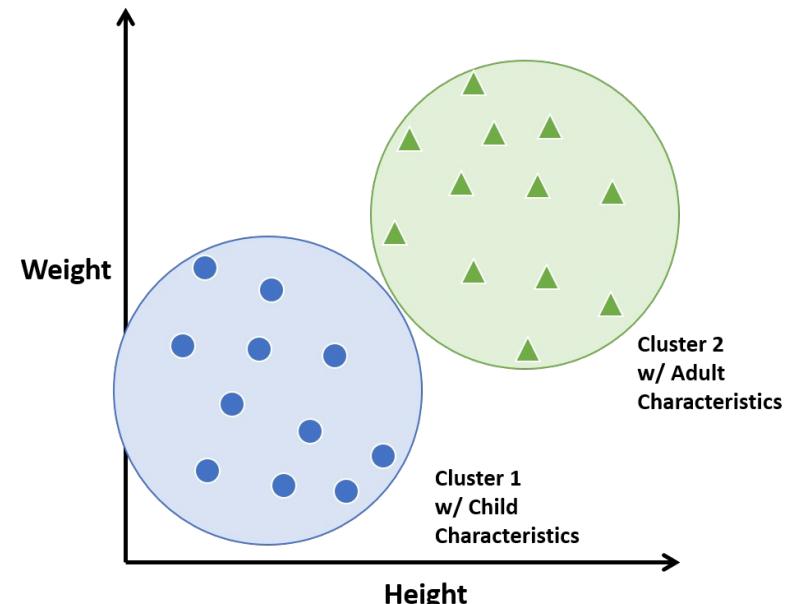
1. LSW, Universität Heidelberg // 2. IMPRS-HD Fellow

# what is *unsupervised* ML?

**Classification** in supervised ML is about making a **decision boundary** between different regions



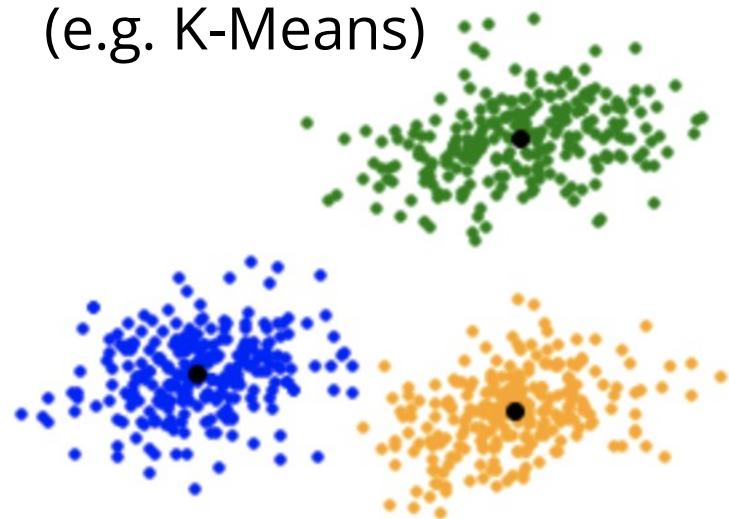
**Clustering algorithms** (a type of unsupervised ML) do this using **properties of the data itself**



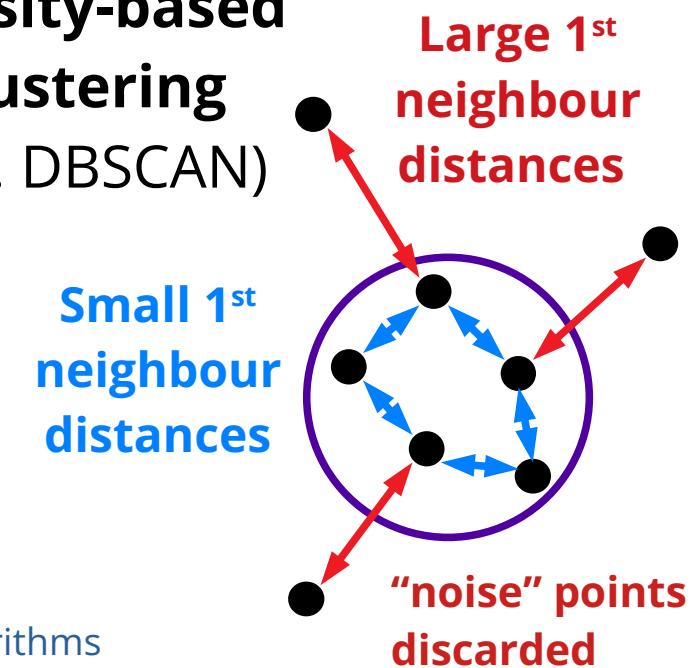
# the many, many algorithms...

There are many different types, each suited to different problems

## Partitioning (e.g. K-Means)



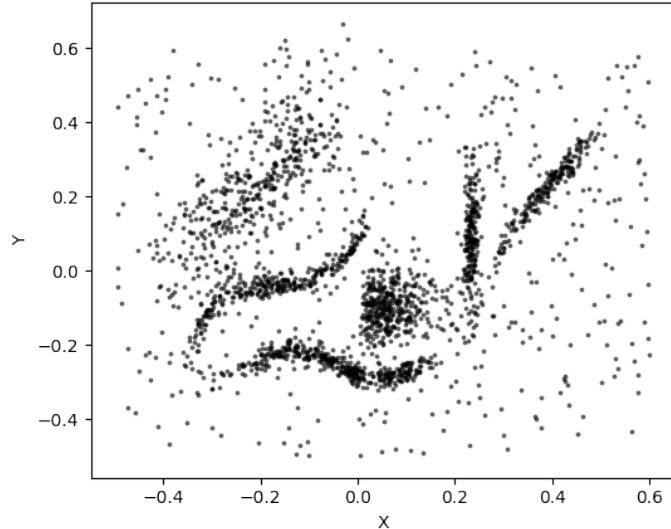
## Density-based clustering (e.g. DBSCAN)



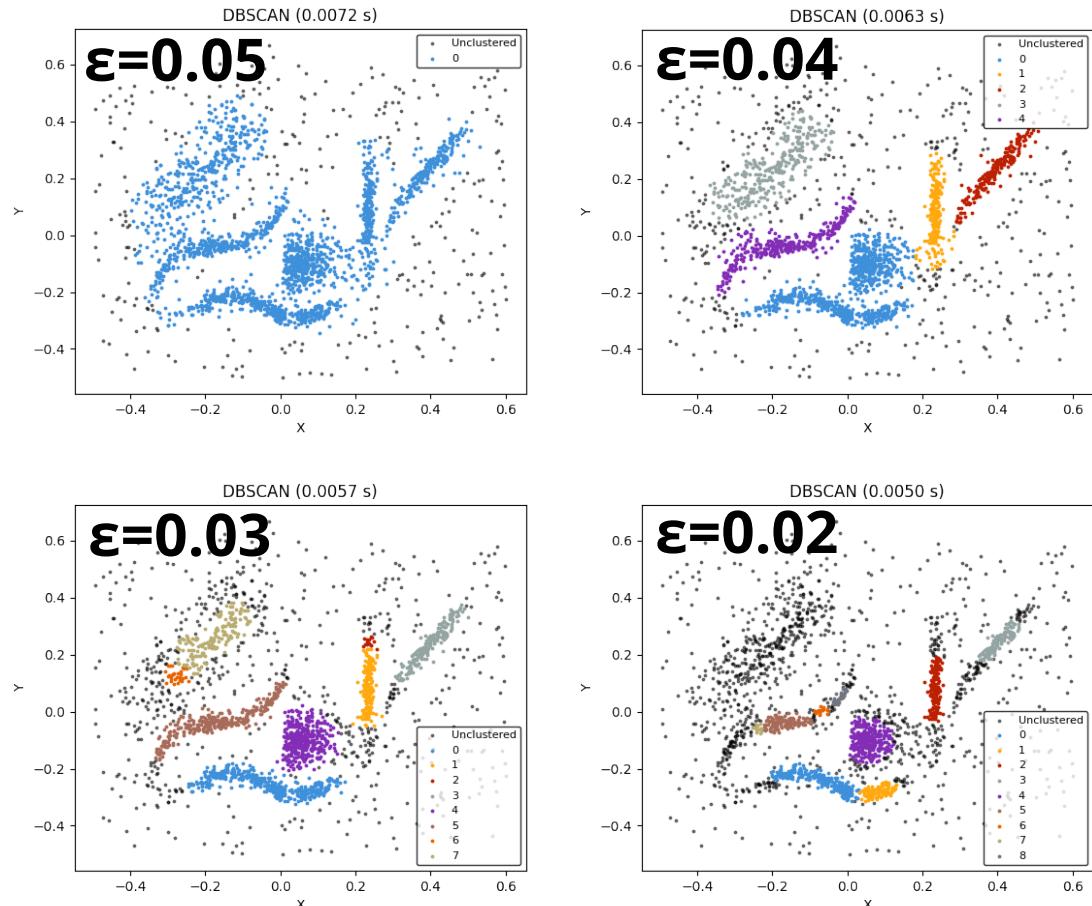
Link: [Overview of various algorithms](#)

# the hard part: setting parameters.

For some basic data...



...what  $\varepsilon$  and  $m_{Pts}$  to use? (this is DBSCAN)



# play around with this in a notebook!

[github.com/emilyhunt/nam\\_2022\\_talk](https://github.com/emilyhunt/nam_2022_talk)  
(also includes these slides)

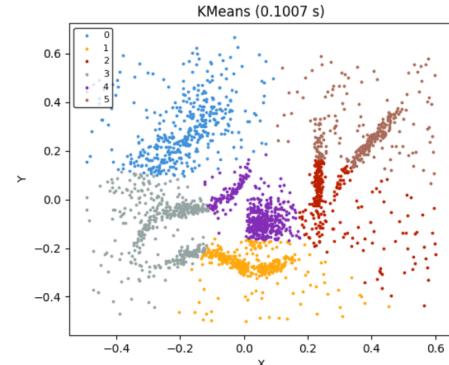
## 4. The clustering algorithms

Let's try various algorithms! Feel free to play around with the parameters.

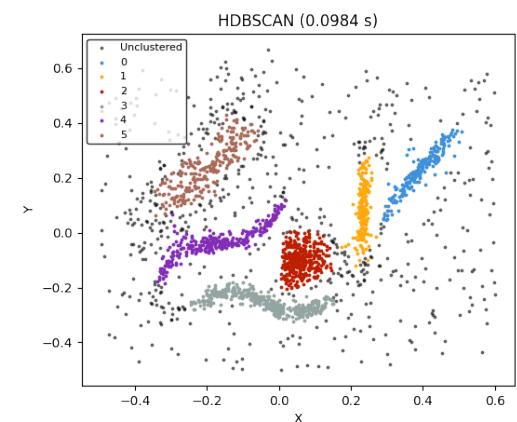
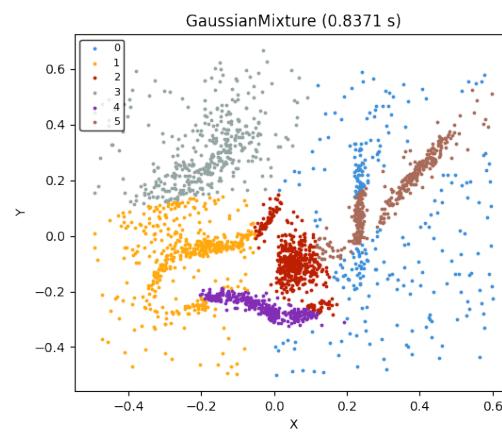
### 4.1. K-Means

The archetypal (and most simple) clustering algorithm...

```
[118]: run_clustering_algorithm(data, sklearn.cluster.KMeans, n_clusters=6, savename='KMeans')
Running clustering algorithm!
clustering took 0.1007 seconds
Plotting results!
```



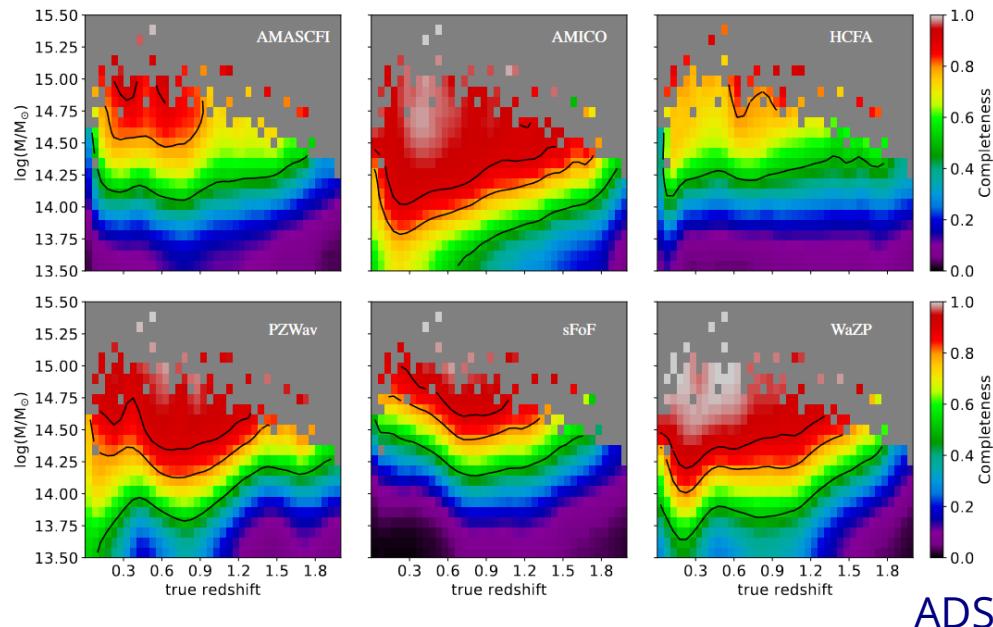
Basic algorithm; some things are roughly separated, but not so well.



# many astro problems come in clusters.

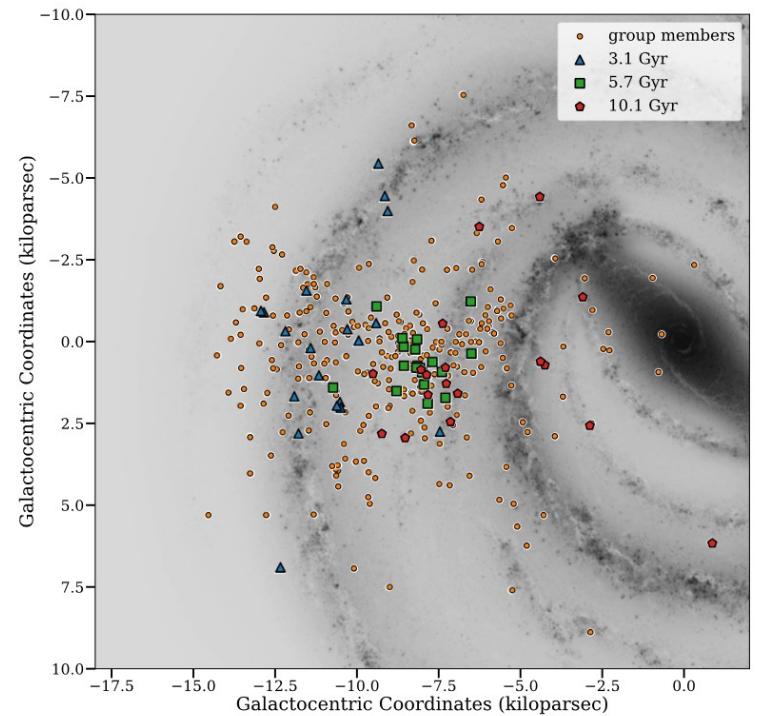
## ***Euclid preparation III. Galaxy cluster detection in the wide photometric survey, performance and algorithm selection***

Euclid Collaboration, R. Adam<sup>1,2,3\*</sup>, M. Vannier<sup>2</sup>, S. Maurogordato<sup>2</sup>, A. Biviano<sup>4</sup>, C. Adami<sup>5</sup>, B. Ascaso<sup>6</sup>,



## **Strong chemical tagging with APOGEE: 21 candidate star clusters that have dissolved across the Milky Way disc**

Natalie Price-Jones<sup>1,2\*</sup>, Jo Bovy<sup>1,2</sup>, Jeremy J. Webb<sup>1</sup>, Carlos Allende Prieto<sup>3,4</sup>,



# our problem: open clusters of stars.

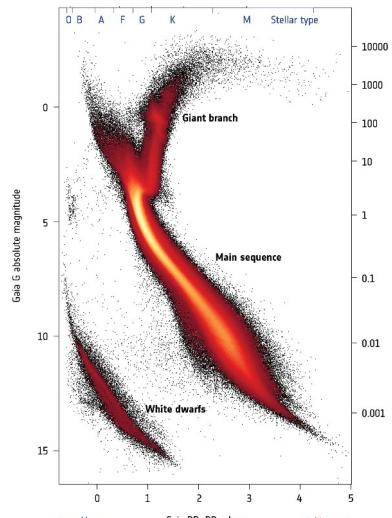
The Pleiades, age: ~100 Myr



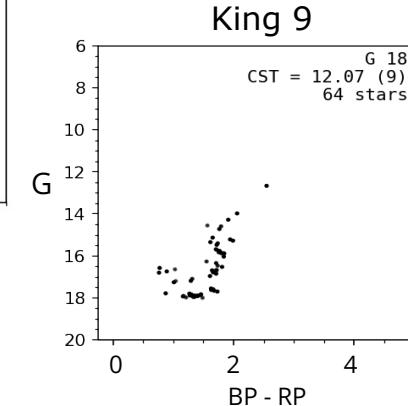
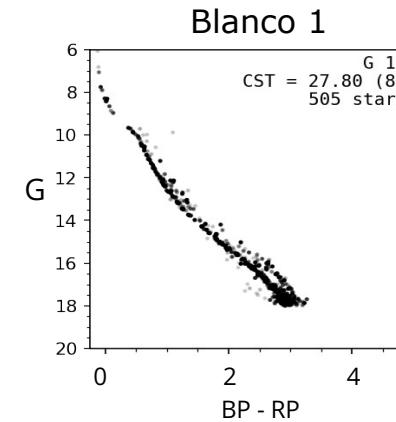
Credit: ESO

Open clusters are **homogenous** and **bound** groups of **coeval** stars

**Highly useful** to stellar & galactic science!

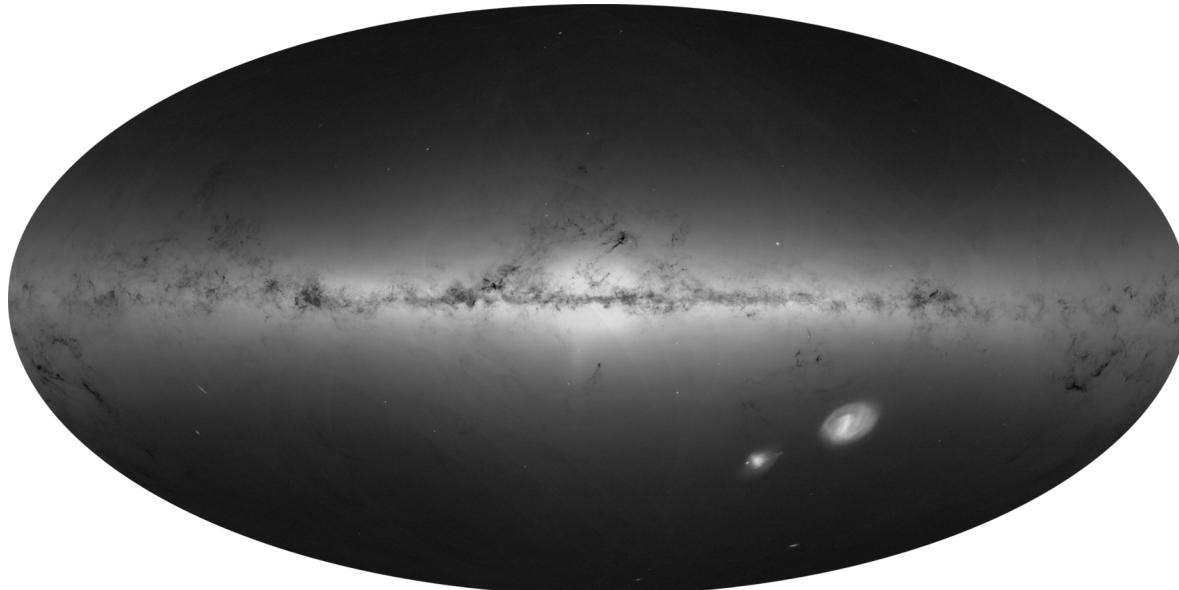


Credit: ESA



# Gaia data is a challenge.

*Gaia* satellite has reliable astrometry for **~1 billion** stars  
=> **~0.1%** of which in **open clusters**

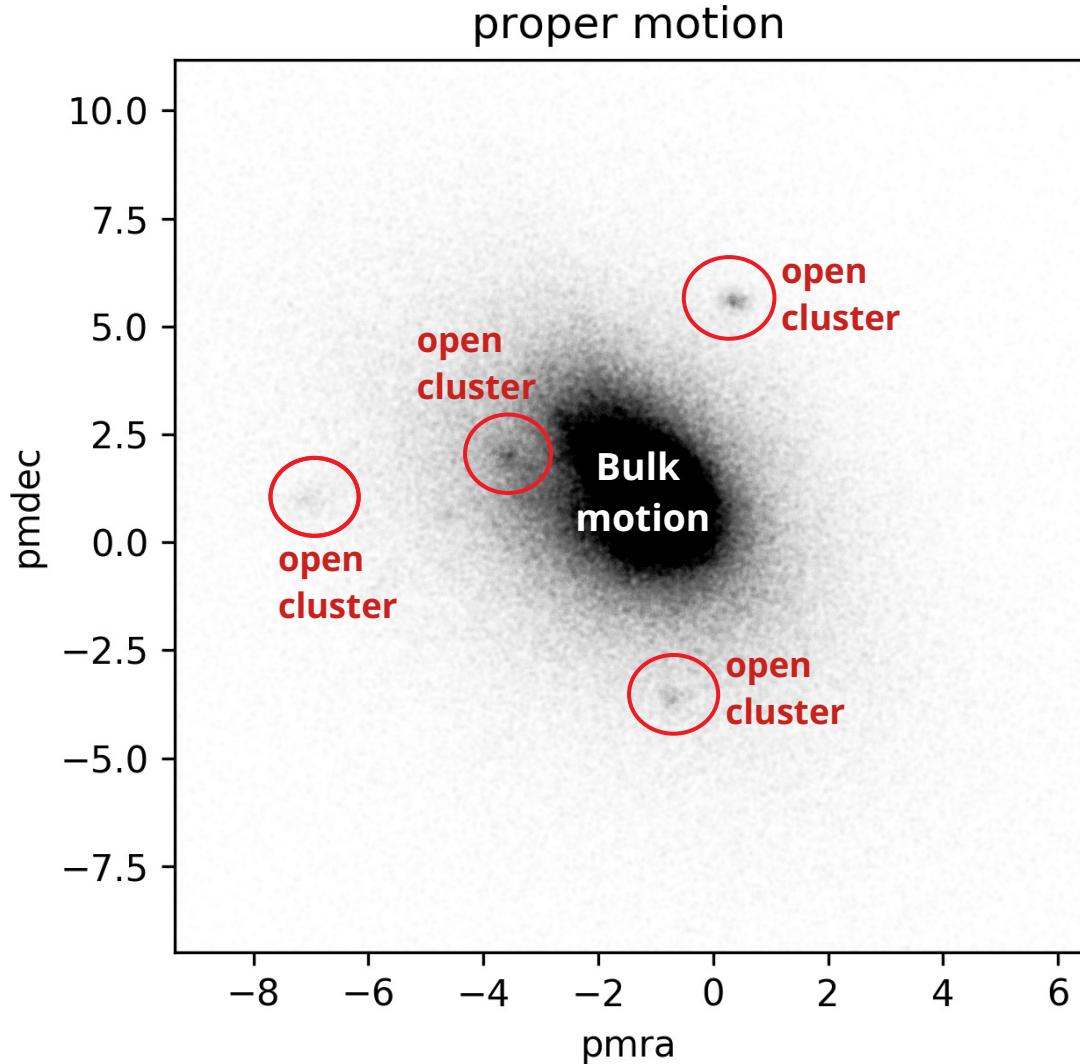


*Gaia* source density skymap. Credit: ESA

*Some challenges:*

- **So many sources...**
- **~99.9%** of which must be discarded!
- Open clusters vary from **~0.1°** to **~10°** in size

**but... Gaia is still  
perfect for open  
clusters!**

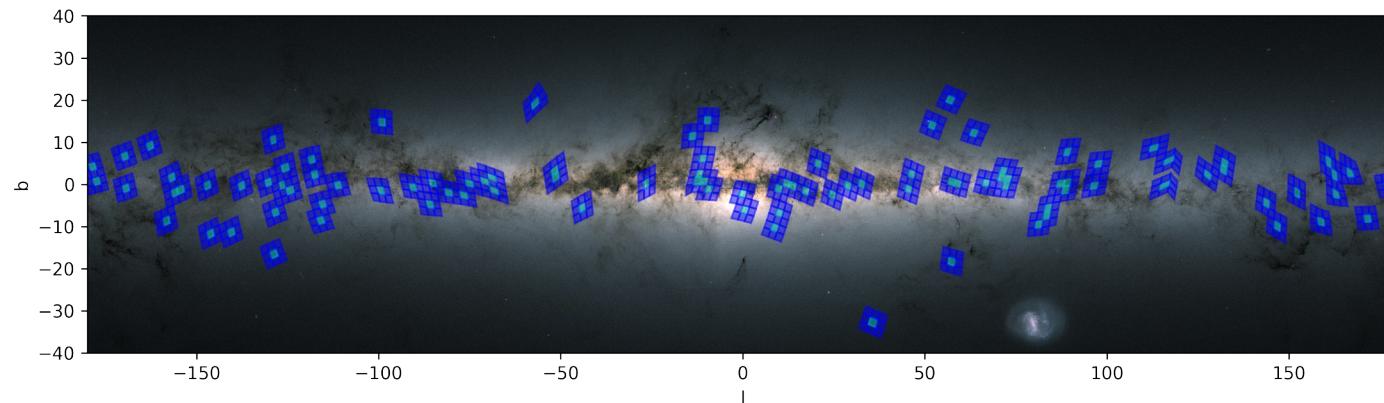


# testing different algorithms.

Hunt & Reffert (2021): search for literature clusters in randomly selected fields

After initial trials of what worked best...

- Algorithms given **rescaled** positions, proper motions and parallaxes (5 dimensions) + corrected for spherical distortions
- Paper used **DBSCAN**, **HDBSCAN** and **Gaussian mixture models**

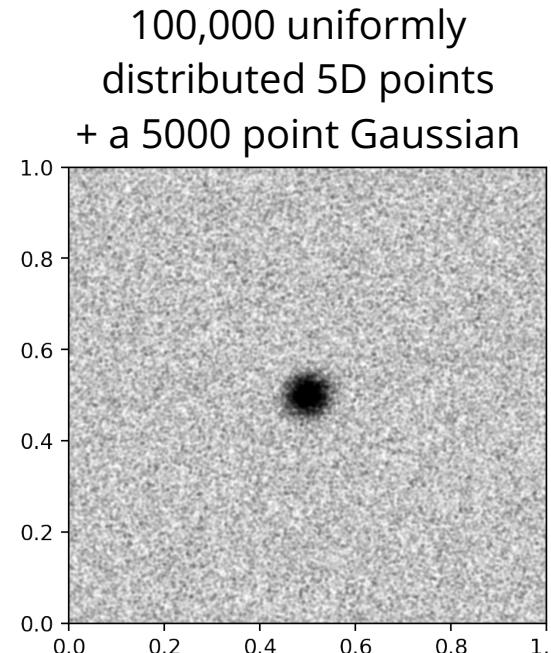


# which one is best?

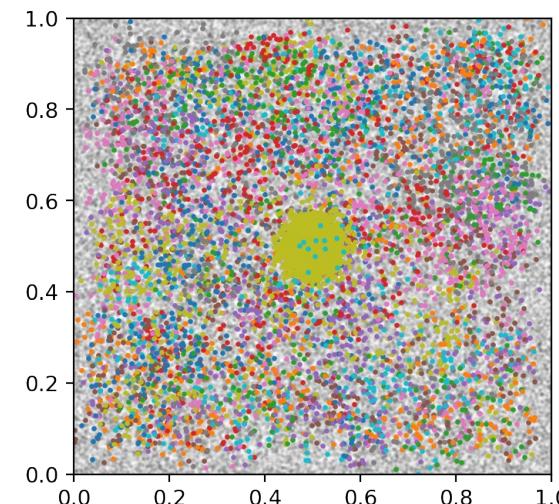
Algorithm	Speed	Sensitivity $TP / (TP + FN)$	Precision $TP / (TP + FP)$	
<b>DBSCAN</b> (Castro-Ginard+18 parameters)	Fast	0.53	1.00	← Main literature approach
<b>DBSCAN</b> (My parameters)	Fast	0.62	0.93	
<b>HDBSCAN</b>	Quite fast	0.82	0.82	← My favourite
<b>Gaussian Mixtures</b>	Slow	0.33	1.00	

# let's talk about false positives.

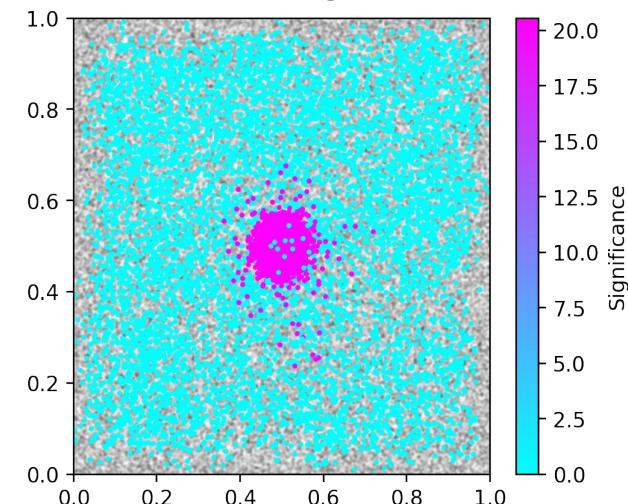
The **most sensitive** clustering methods can also produce **the most false positives**. Care **must** be taken when using "off-the-shelf" methods



**HDBSCAN**,  $m_{clSize} = 20$ ...  
151 clusters?!



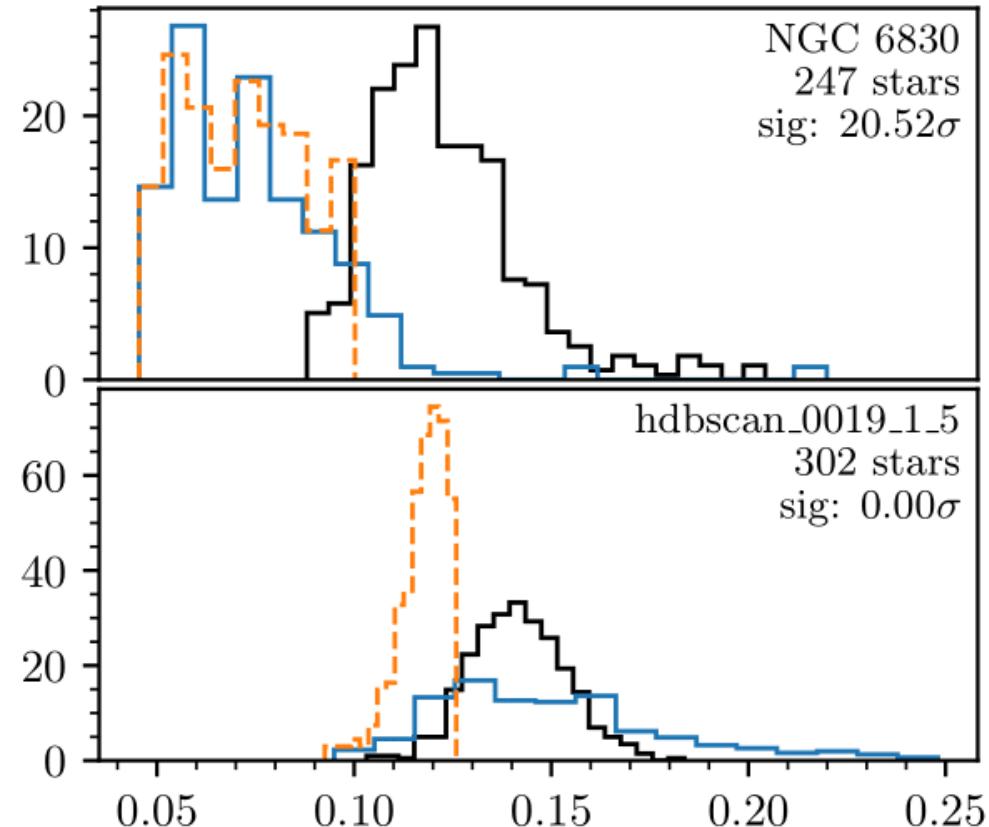
Density test → only the real cluster is significant



# how the significance test works.

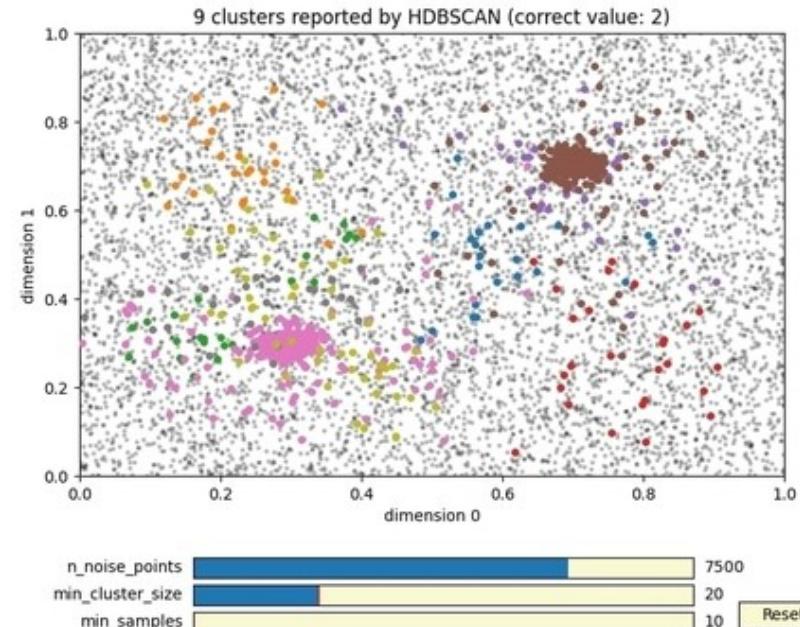
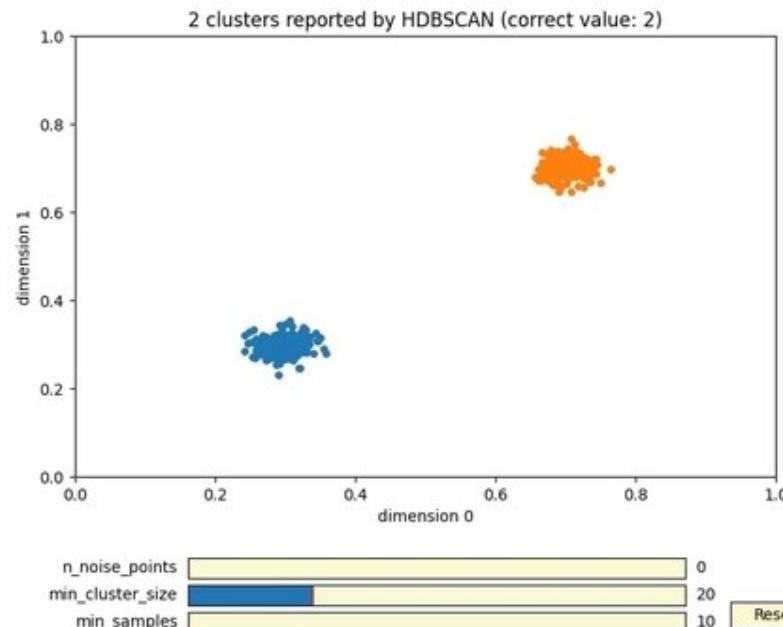
Simple density test compares  
**nearest neighbour  
distances** in cluster with field

Can tell **random chance  
associations** and **real  
clusters** apart



# a way to try this in real time!

I made a script that shows this in real time! ([GitHub link](#))

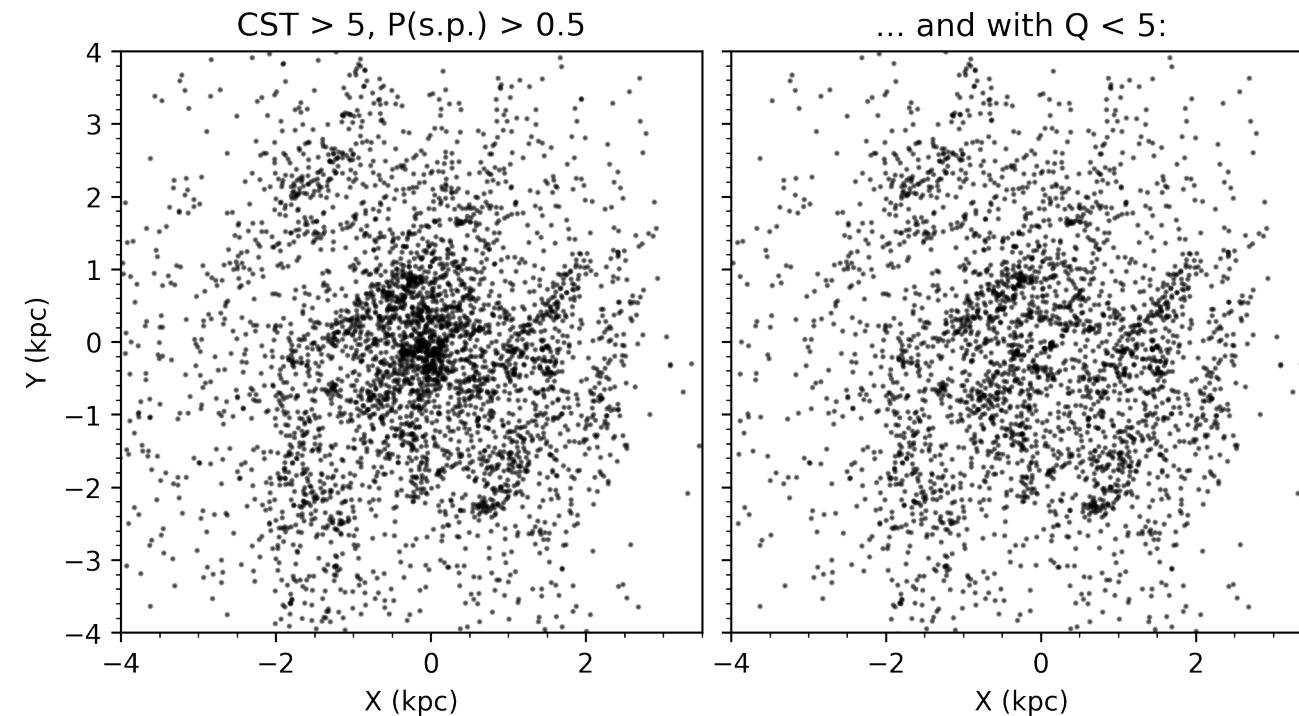


# clustering 729 million stars in Gaia EDR3.

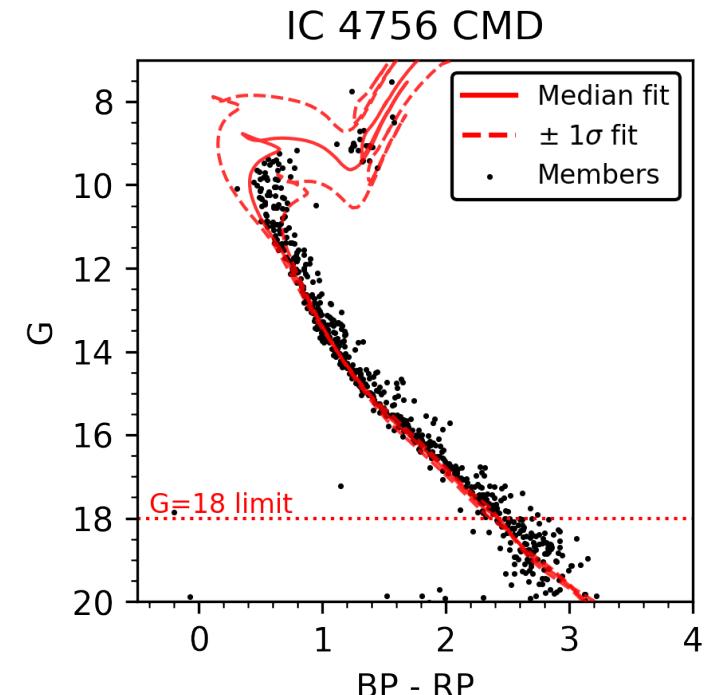
- Ran **HDBSCAN** on largest-ever sample of stars
- Used **HEALPix** tessellation scheme for regions
- About **8 days of wall time** on a powerful machine (actually not bad)

# after much work: results.

Distribution,  $|z| < 500$  pc:



Example CMD:



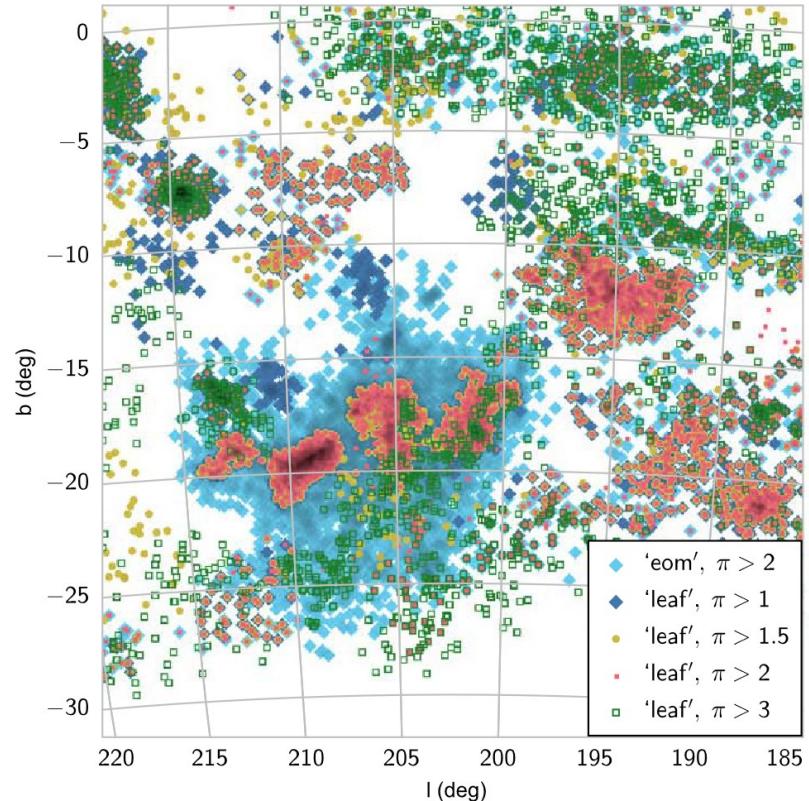
# and a cautionary tale on false positives...

**Kounkel+19,20:** reported **thousands** of groups and ‘strings’ of stars using **HDBSCAN** on Gaia DR2 data

(+ without using a significance test to clean HDBSCAN results)

In EDR3: we re-detect just **18.1%** of their groups (with better data!)

Also echoes **Zucker+22:** “many Kounkel+ group members inconsistent with having common origin”



# where next for clustering algorithms?

- Clustering algorithms are **fantastic** tools for astronomy
- Many off-the-shelf algorithms available (e.g. scikit-learn)
- BUT: they aren't designed for astronomy problems, and can run into weird issues

**In the future:** astro must collaborate with computer scientists, mathematicians etc. to develop algorithms

# datasets to come...

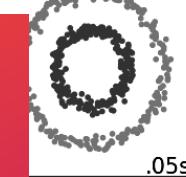
**Gaia DR5** (~2030): 2 billion stars

**Euclid** (late 2020s): ~1 billion galaxies with redshifts

**LSST** (final data ~2035): 17 billion stars, 20 billion galaxies

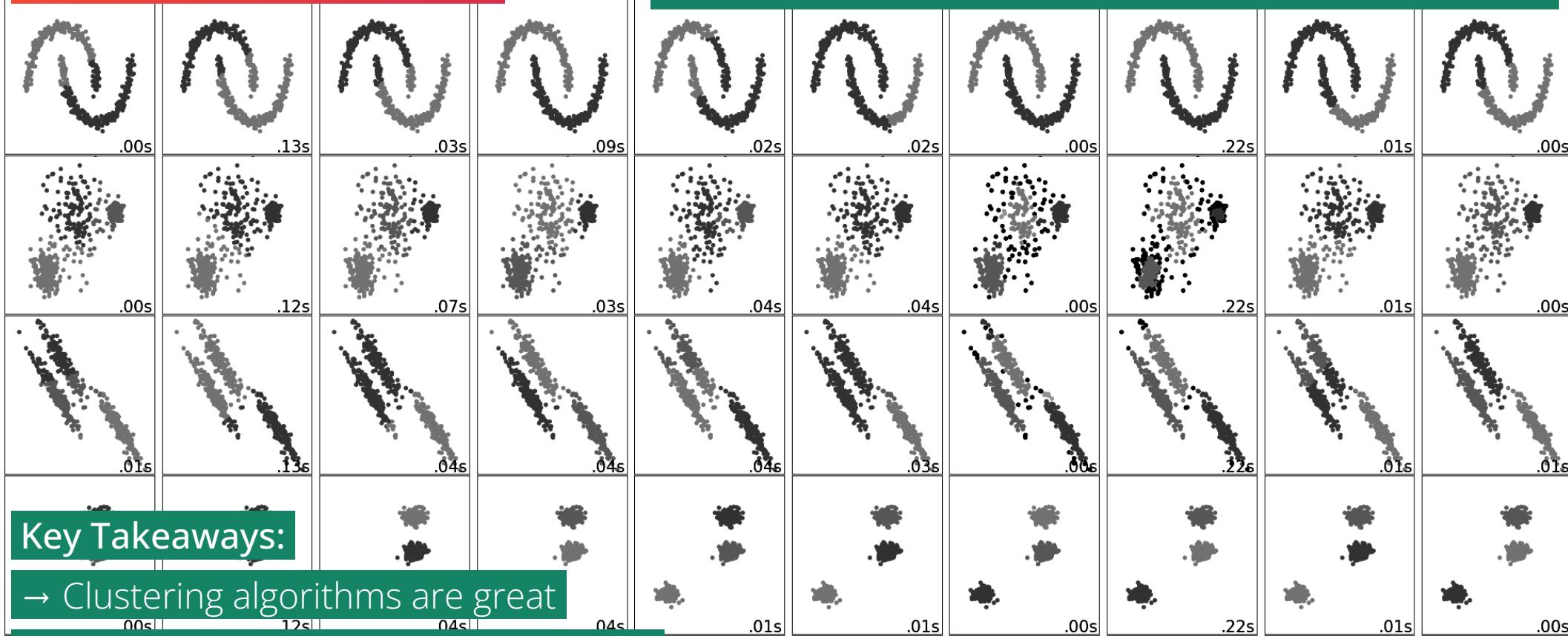
**GaiaNIR** (~2050): 12 billion stars

+ many, many other uses across astronomy for clustering algorithms



Get the slides at...

[github.com/emilyhunt/nam\\_2022\\_talk](https://github.com/emilyhunt/nam_2022_talk)



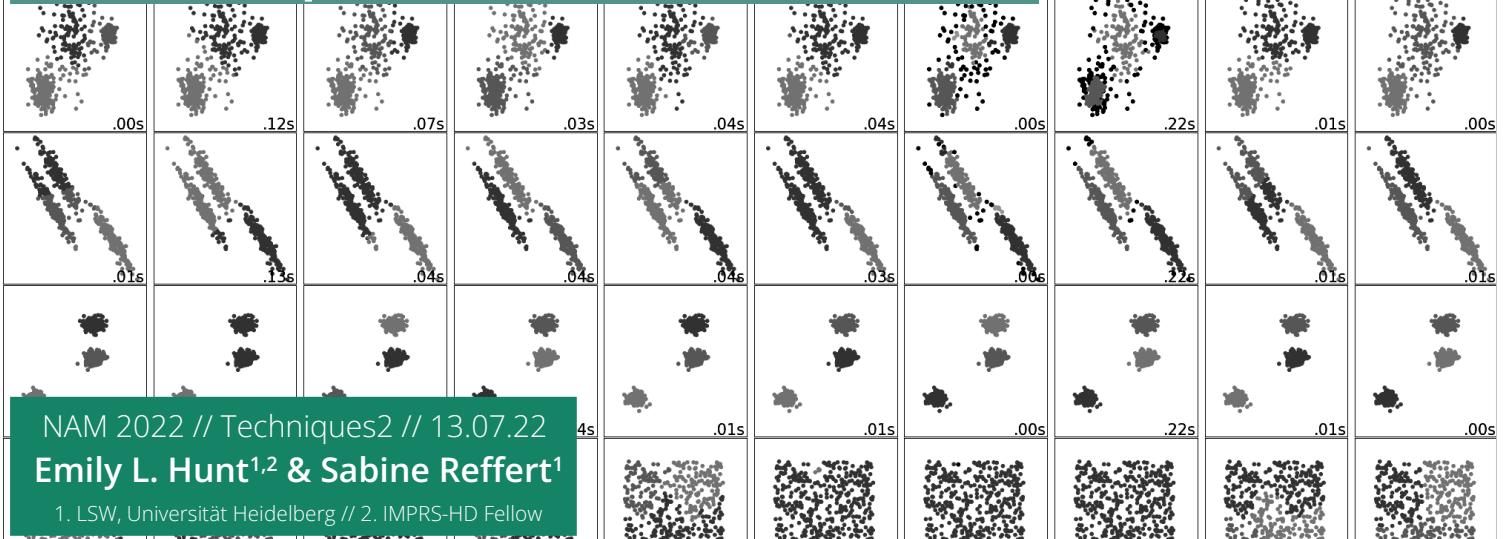
### Key Takeaways:

→ Clustering algorithms are great

→ Many off-the-shelf solutions available

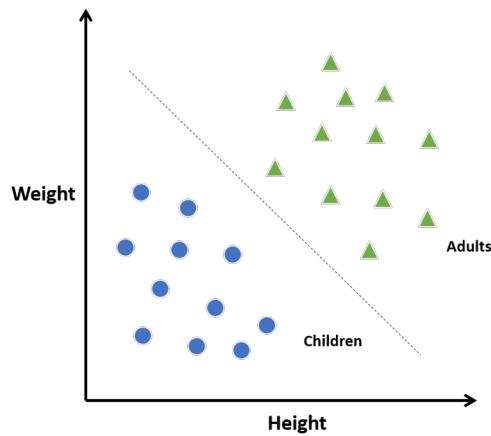
→ BUT: be careful! They can be wrong (and not tell you)

# The power (and caveats) of clustering algorithms with examples from use on Gaia data

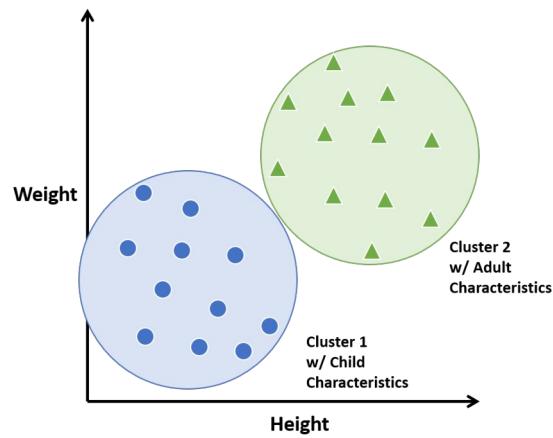


# what is *unsupervised* ML?

**Classification** in supervised ML is about making a **decision boundary** between different regions



**Clustering algorithms** (a type of unsupervised ML) do this using **properties of the data itself**



Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

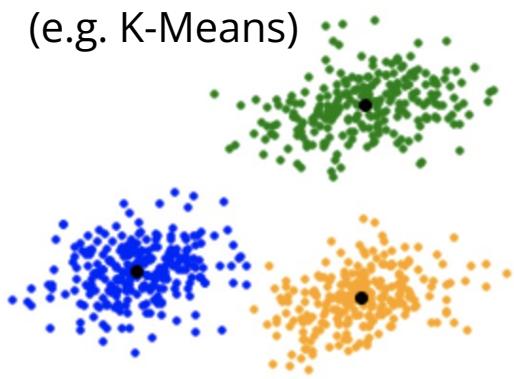
2 of 20

# the many, many algorithms...

There are many different types, each suited to different problems

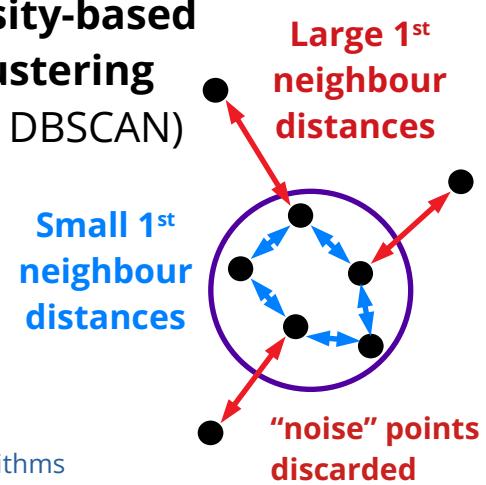
## Partitioning

(e.g. K-Means)



## Density-based clustering

(e.g. DBSCAN)



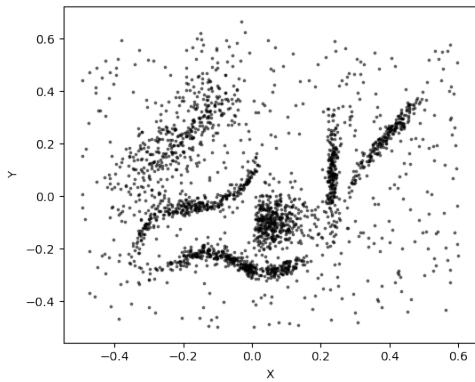
Link: [Overview of various algorithms](#)

Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

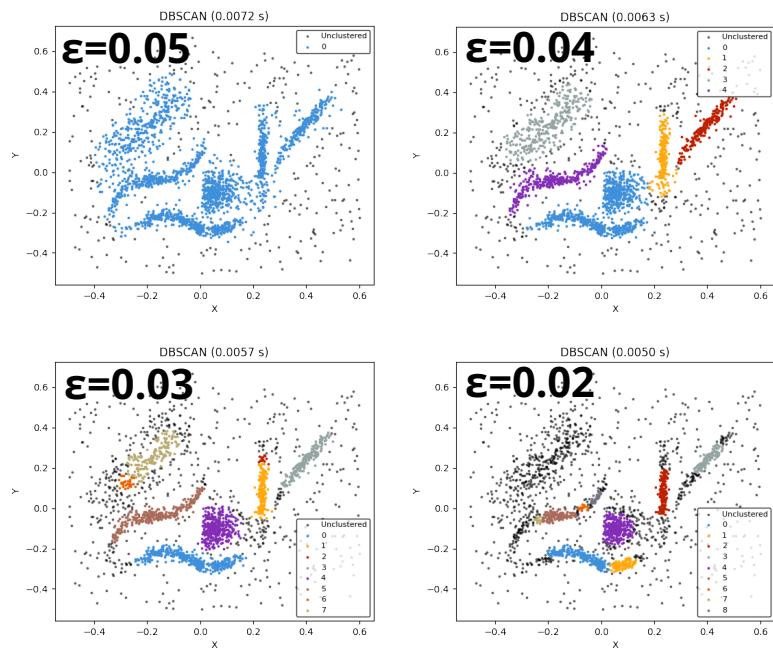
3 of 20

# the hard part: setting parameters.

For some basic data...



...what  $\varepsilon$  and  $m_{pts}$  to use? (this is DBSCAN)



Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

4 of 20

# play around with this in a notebook!

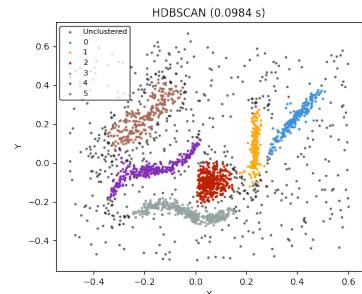
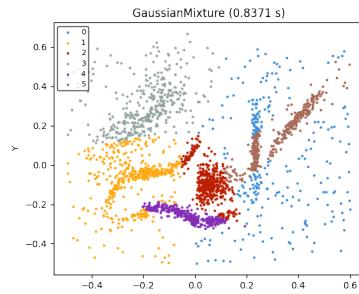
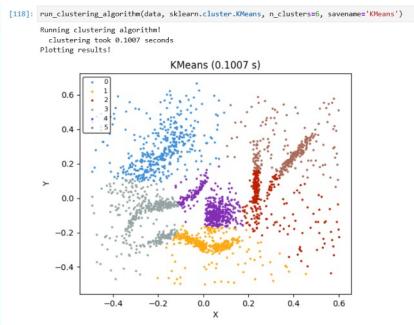
[github.com/emilyhunt/nam\\_2022\\_talk](https://github.com/emilyhunt/nam_2022_talk)  
 (also includes these slides)

## 4. The clustering algorithms

Let's try various algorithms! Feel free to play around with the parameters.

### 4.1. K-Means

The archetypal (and most simple) clustering algorithm...



Basic algorithm: some things are roughly separated, but not so well.

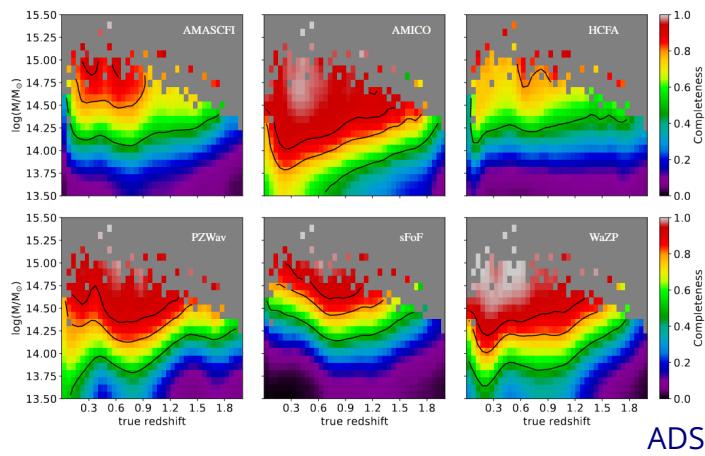
Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

5 of 20

# many astro problems come in clusters.

## *Euclid preparation III. Galaxy cluster detection in the wide photometric survey, performance and algorithm selection*

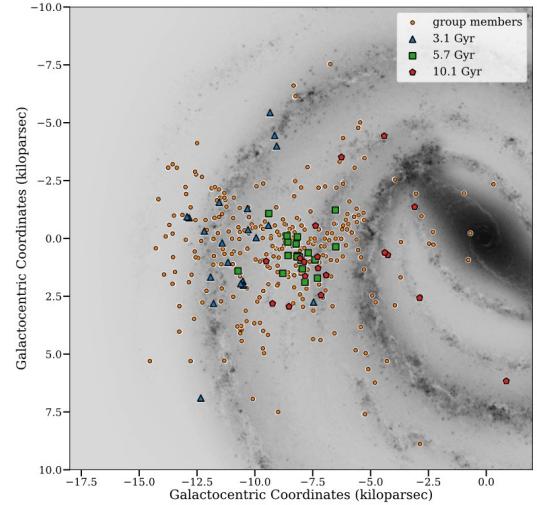
Euclid Collaboration, R. Adam<sup>1,2,3\*</sup>, M. Vannier<sup>2</sup>, S. Maurogordato<sup>2</sup>, A. Biviano<sup>4</sup>, C. Adami<sup>5</sup>, B. Ascaso<sup>6</sup>,



ADS

Strong chemical tagging with APOGEE: 21 candidate star clusters that have dissolved across the Milky Way disc

Natalie Price-Jones<sup>1,2\*</sup>, Jo Bovy<sup>1,2</sup>, Jeremy J. Webb<sup>1</sup>, Carlos Allende Prieto<sup>3,4</sup>,



ADS

Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

6 of 20

# our problem: open clusters of stars.

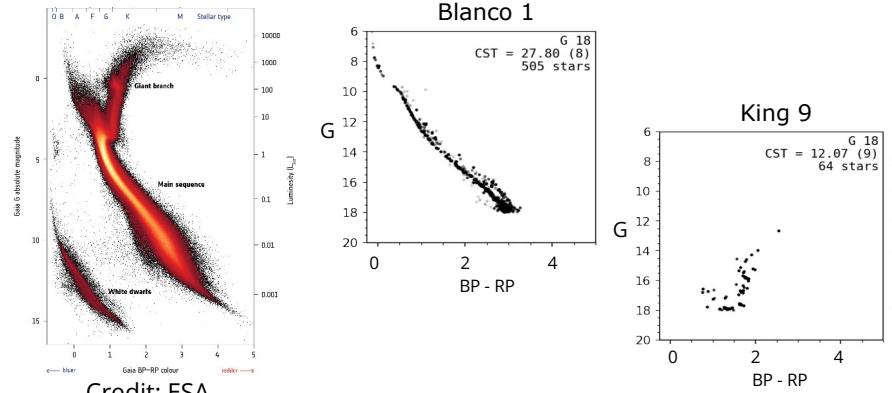
The Pleiades, age: ~100 Myr



Credit: ESO

Open clusters are **homogenous** and **bound** groups of **coeval** stars

**Highly useful** to stellar & galactic science!

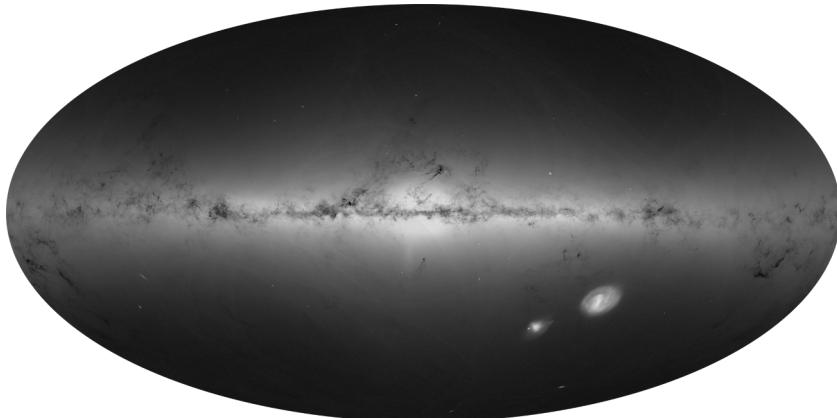


Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

7 of 20

# Gaia data is a challenge.

*Gaia* satellite has reliable astrometry for **~1 billion** stars  
=> **~0.1%** of which in **open clusters**



*Gaia* source density skymap. Credit: ESA

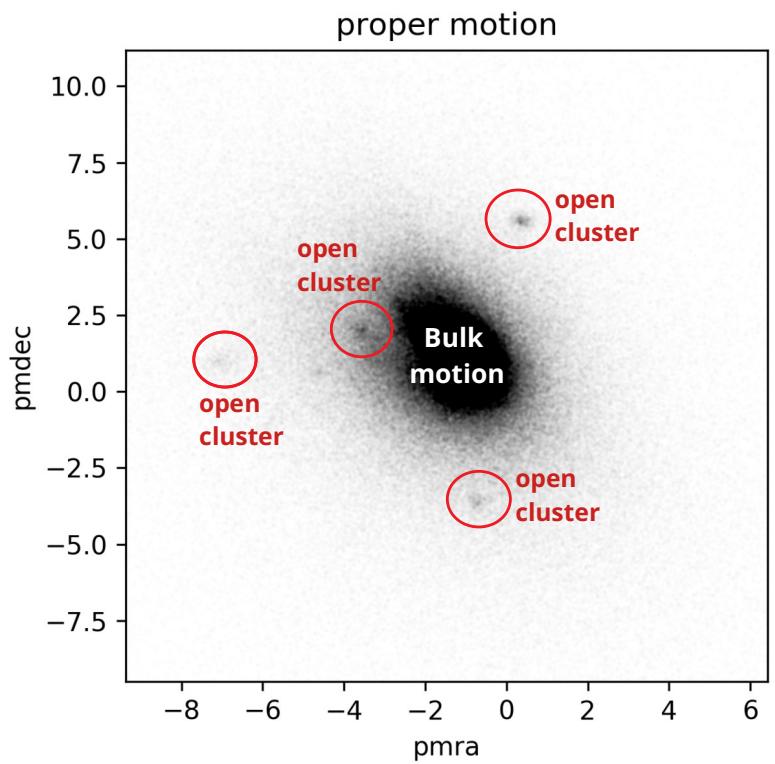
*Some challenges:*

- **So many sources...**
- **~99.9%** of which must be discarded!
- Open clusters vary from **~0.1°** to **~10°** in size

**Emily L. Hunt.** *The power (and caveats) of clustering algorithms.*

**8 of 20**

but... Gaia is still  
*perfect* for open  
clusters!



Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

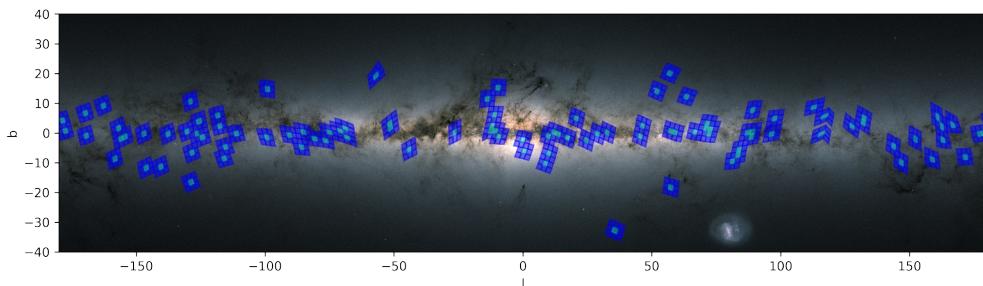
9 of 20

# testing different algorithms.

Hunt & Reffert (2021): search for literature clusters in randomly selected fields

After initial trials of what worked best...

- Algorithms given **rescaled** positions, proper motions and parallaxes (5 dimensions) + corrected for spherical distortions
- Paper used **DBSCAN**, **HDBSCAN** and **Gaussian mixture models**



Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

10 of 20

# which one is best?

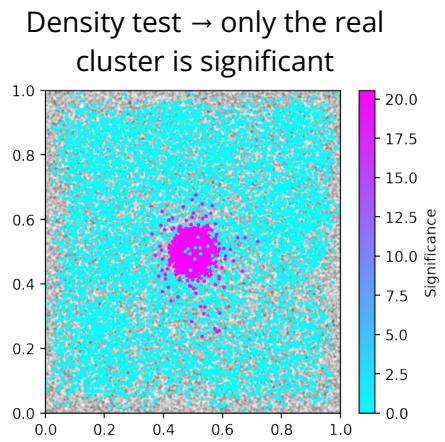
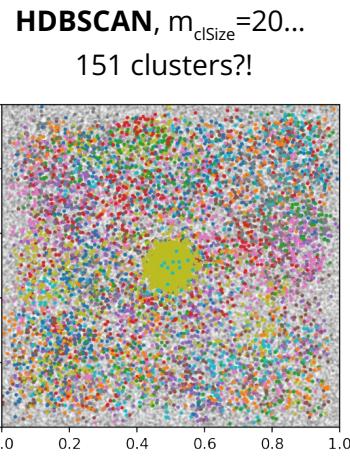
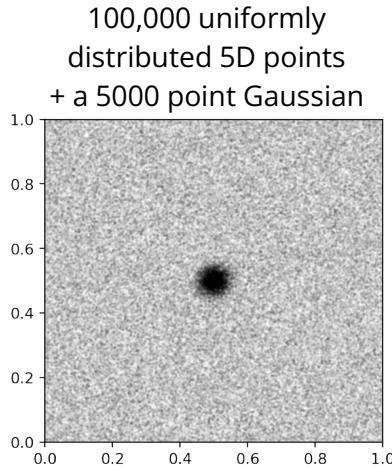
Algorithm	Speed	Sensitivity $TP / (TP + FN)$	Precision $TP / (TP + FP)$	
DBSCAN (Castro-Ginard+18 parameters)	Fast	0.53	1.00	← Main literature approach
DBSCAN (My parameters)	Fast	0.62	0.93	
HDBSCAN	Quite fast	0.82	0.82	← My favourite
Gaussian Mixtures	Slow	0.33	1.00	

Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

11 of 20

# let's talk about false positives.

The **most sensitive** clustering methods can also produce **the most false positives**. Care **must** be taken when using “off-the-shelf” methods



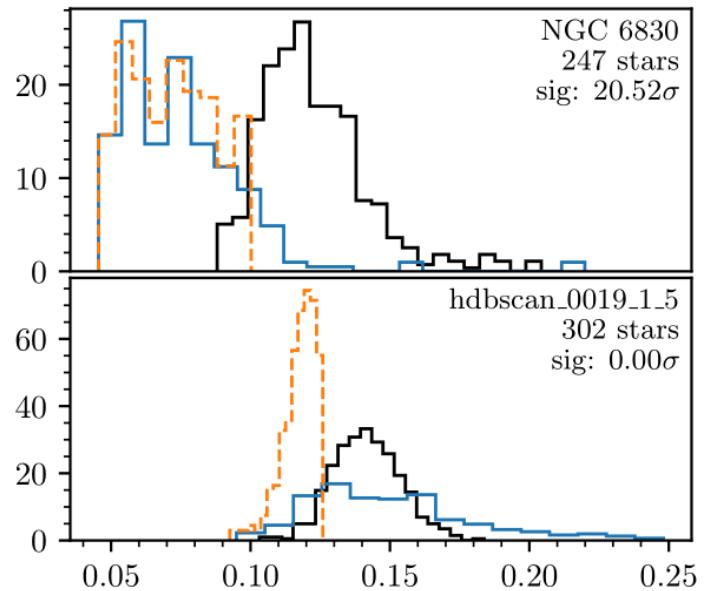
Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

12 of 20

# how the significance test works.

Simple density test compares  
**nearest neighbour distances** in cluster with field

Can tell **random chance associations** and **real clusters** apart

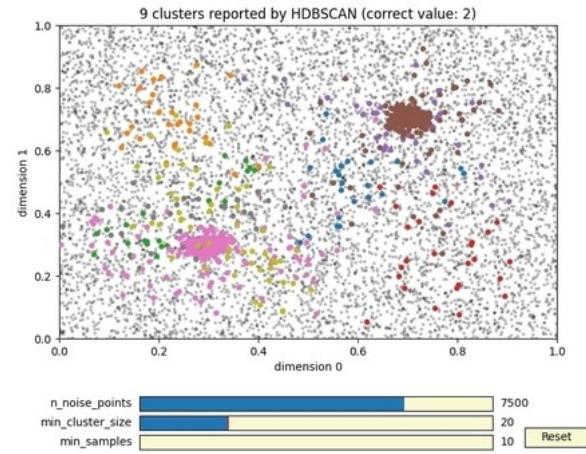
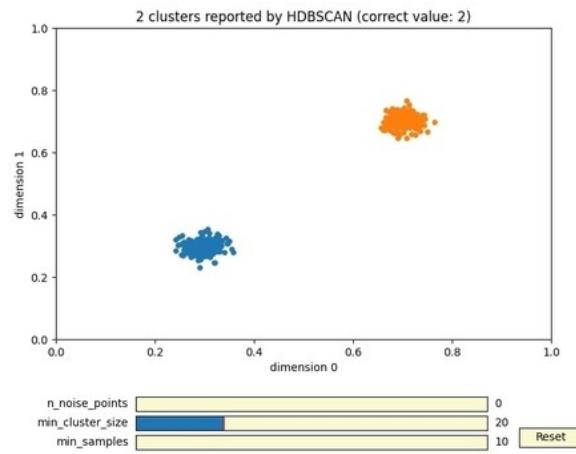


Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

13 of 20

# a way to try this in real time!

I made a script that shows this in real time! ([GitHub link](#))



Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

14 of 20

# clustering 729 million stars in Gaia EDR3.

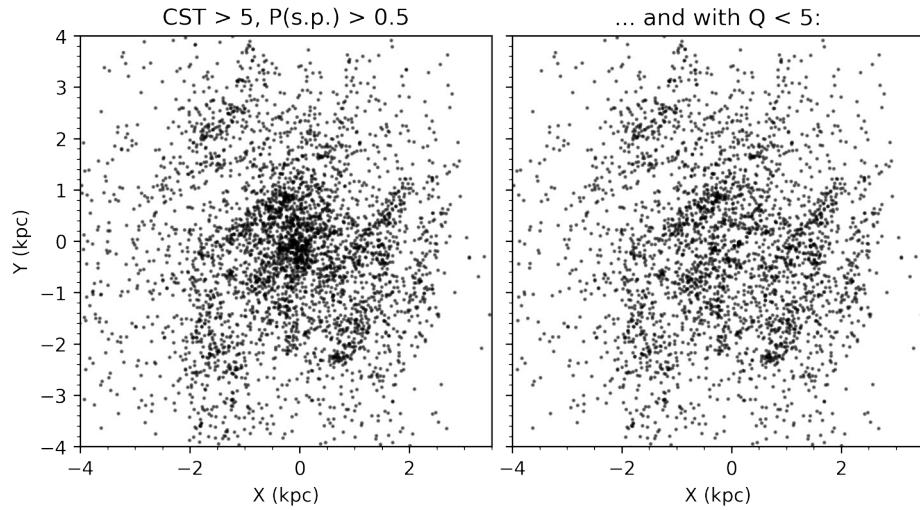
- Ran **HDBSCAN** on largest-ever sample of stars
- Used **HEALPix** tesselation scheme for regions
- About **8 days of wall time** on a powerful machine (actually not bad)

Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

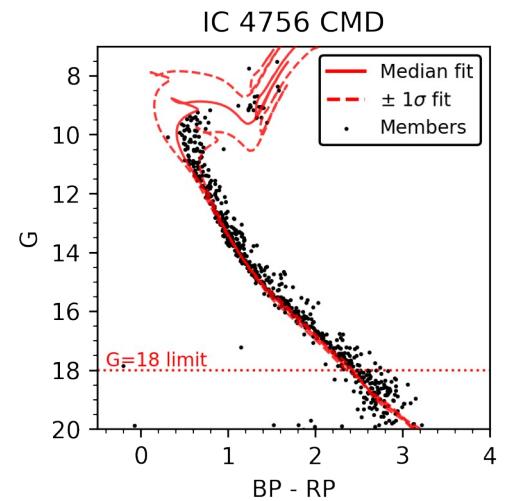
15 of 20

# after much work: results.

Distribution,  $|z| < 500$  pc:



Example CMD:



**Emily L. Hunt.** *The power (and caveats) of clustering algorithms.*

16 of 20

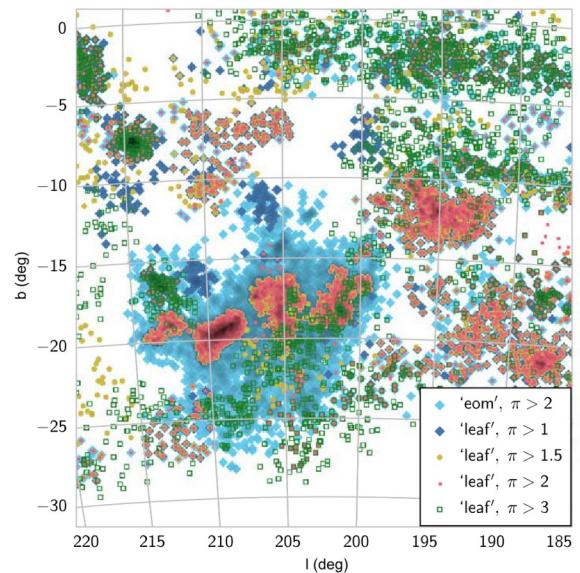
# and a cautionary tale on false positives...

**Kounkel+19,20:** reported **thousands** of groups and ‘strings’ of stars using **HDBSCAN** on Gaia DR2 data

(+ without using a significance test to clean HDBSCAN results)

In EDR3: we re-detect just **18.1%** of their groups (with better data!)

Also echoes **Zucker+22:** “many Kounkel+ group members inconsistent with having common origin”



Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

17 of 20

# where next for clustering algorithms?

- Clustering algorithms are **fantastic** tools for astronomy
- Many off-the-shelf algorithms available (e.g. [scikit-learn](#))
- BUT: they aren't designed for astronomy problems, and can run into weird issues

**In the future:** astro must collaborate with computer scientists, mathematicians etc. to develop algorithms

Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

18 of 20

# datasets to come...

**Gaia DR5** (~2030): 2 billion stars

**Euclid** (late 2020s): ~1 billion galaxies with redshifts

**LSST** (final data ~2035): 17 billion stars, 20 billion galaxies

**GaiaNIR** (~2050): 12 billion stars

+ many, many other uses across astronomy for clustering algorithms

Emily L. Hunt. *The power (and caveats) of clustering algorithms.*

19 of 20

emilydoesastro.com  
@emilydoesastro

Get the slides at...  
[github.com/emilyhunt/nam\\_2022\\_talk](https://github.com/emilyhunt/nam_2022_talk)

