

# A.I. for Astrophysics: Masterclass

Dr. Rob Lyon

Edge Hill University

[robert.lyon@edgehill.ac.uk](mailto:robert.lyon@edgehill.ac.uk)

 @scienceguyrob



# Welcome

In the next few hours we'll learn how to apply ML to a real-world astrophysics problem. If you've never done machine learning before don't worry. I'll cover everything you need to know including:

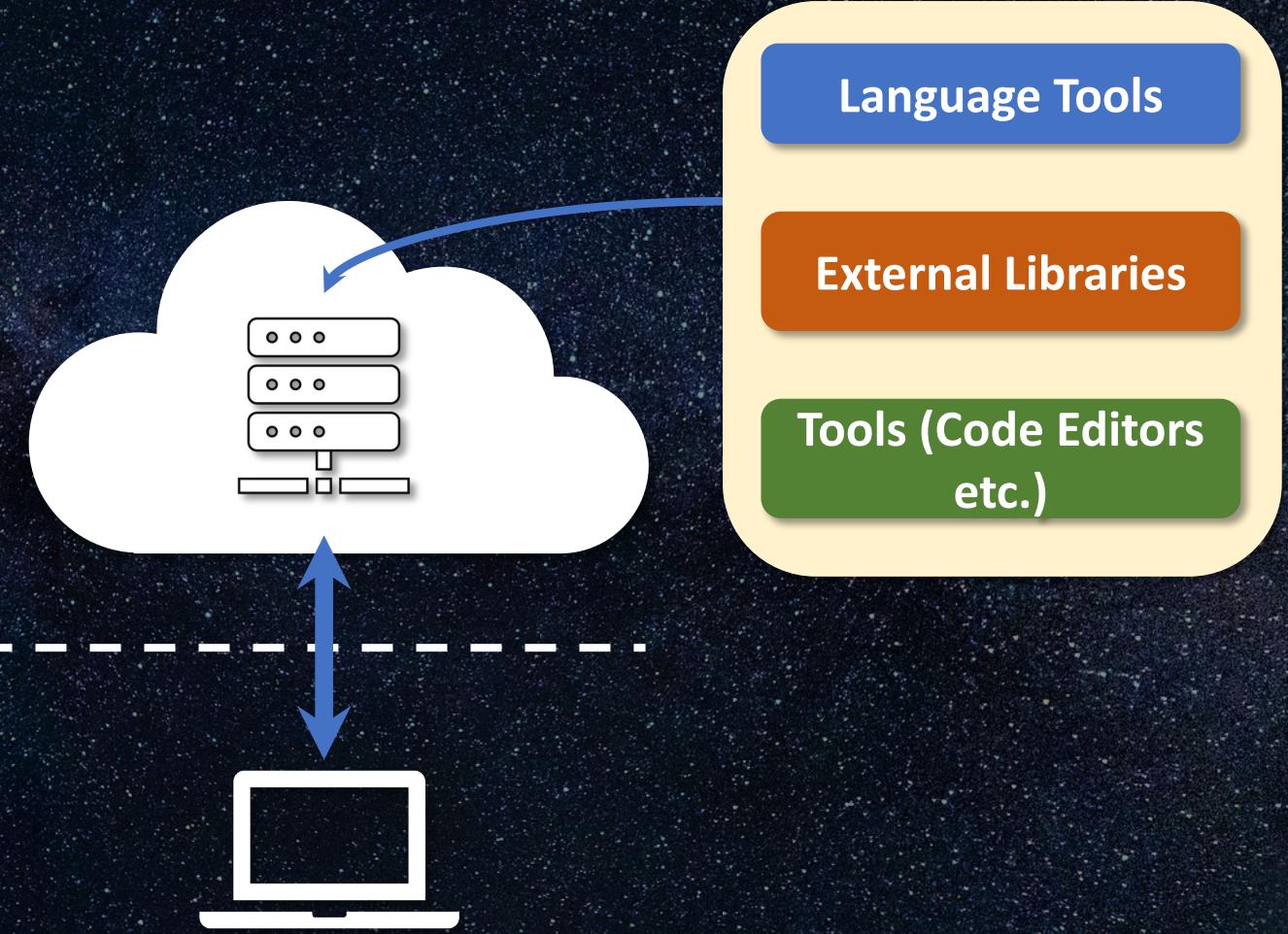
- Machine learning basics.
- Python basics.
- The astronomy background you'll need to understand the challenge.

If you have questions at any point – just ask. We'll move through things slowly. We want this masterclass to be enjoyable!

# Masterclass Environment

# Google Colab

- A cloud environment containing the tools need to write & execute Python programs.
- You'll need is a google account to use, e.g. a Gmail Account and the *Chrome Web Browser*.
- When you login to the Collaboratory, it creates a computer just for you.
- Inside you can create & execute Python code.
- Before proceeding, please create a Google account if you don't already have one.

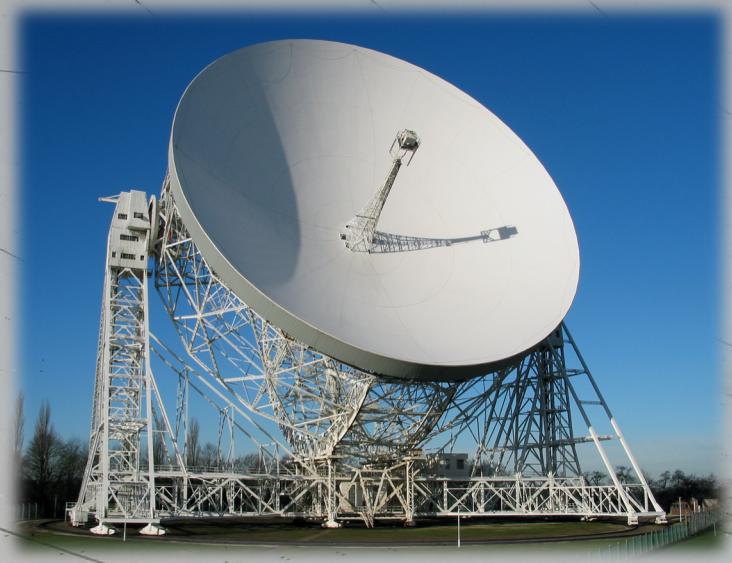


# Accessing the Masterclass Resource

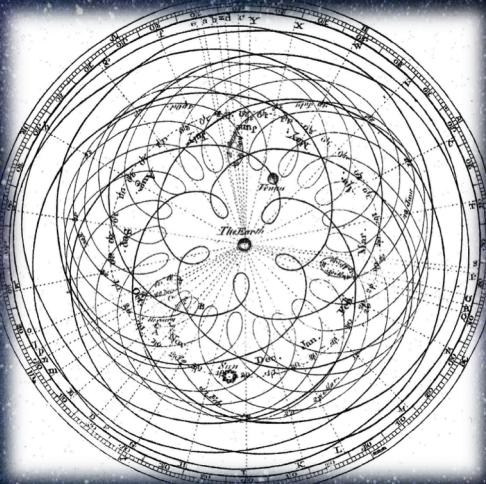
1. Open Google Chrome.
  2. Navigate to:  
<https://github.com/scienceguyrob/AI4AstroMasterclass>
  3. Click the file: “A. I. for Astrophysics Masterclass – Classifying Pulsars.ipynb”
  4. There is a button with the text “Raw” on it – click it. This will load the resource into your browser.
  5. Save the content to your own machine – left click and choose the “Save as” option that appears on the popup menu. Keep the .ipynb file extension.
  6. Open Google Colab: <https://colab.research.google.com/>
  7. Click the File menu, click “Upload Notebook”. Choose the file you saved on your machine.

```
cells": [
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": "# A.I. for Astrophysics Masterclass - Classifying Pulsars\n\n# A Jupyter notebook containing a machine learning tutorial for Pulsar classification. The notebook was written to support a masterclass delivered at Edge Hill University on 28/10/2019.\n\n[DOI](https://zenodo.org/badge/218981425.svg)](https://zenodo.org/badge/latestdoi/218981425)\n\n## Author(s)\n\nDr. Robert Lyon (mailto:robert.lyon@edgehill.ac.uk)\n\nLecturer([https://www.edgehill.ac.uk/computerscience/people/academic-staff/robert-lyon] @ Edge Hill University, [Department of Computer Science](https://www.edgehill.ac.uk/computerscience/people/academic-staff/robert-lyon)).\nAlso see [www.scienceguyrob.com](www.scienceguyrob.com).\n\nIf using/editing this resource, please cite it appropriately using the DOI above.\n\n# License\n\nRemember if you find this resource useful, you can let me know via social media [<img src="https://github.com/scienceguyrob/AI4AstroMasterclass30]>](https://twitter.com/scienceguyrob?lang=en)\n\n## Table of Contents\n\n1. [Introduction](#Introduction)\n  1.1 [Who is this for?](#who)\n  1.2 [Aims](#aims)\n  1.3 [Tools](#tools)\n  1.4 [Using this Resource](#using)\n\n
```

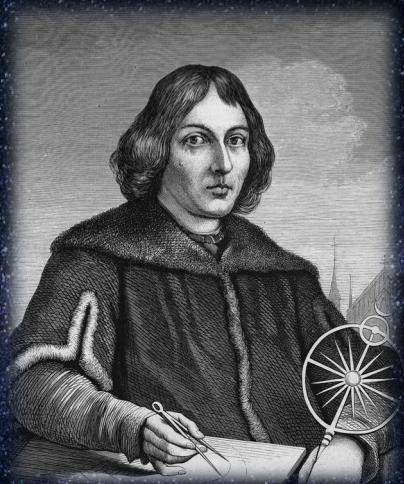
# Astrophysics Background



# Some History



Geocentric Model



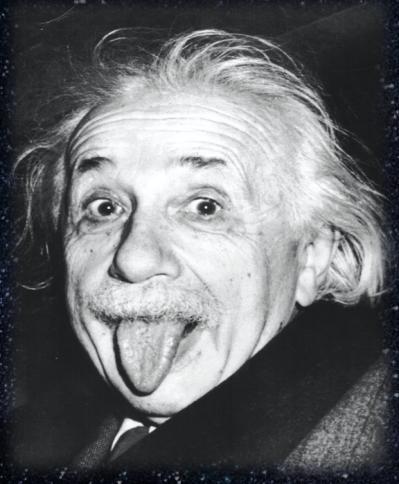
Heliocentric Model



Laws of Motion



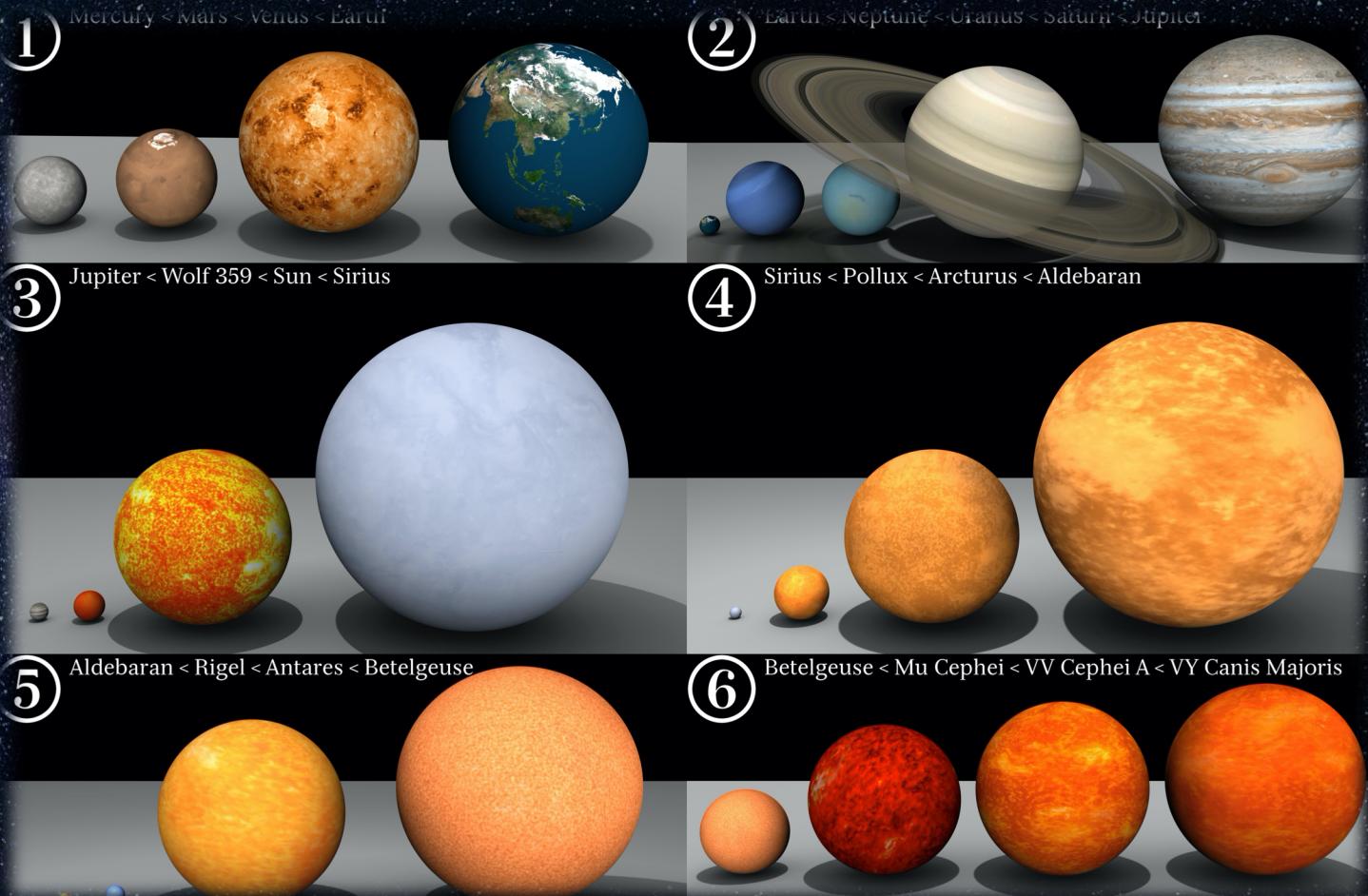
Gravity



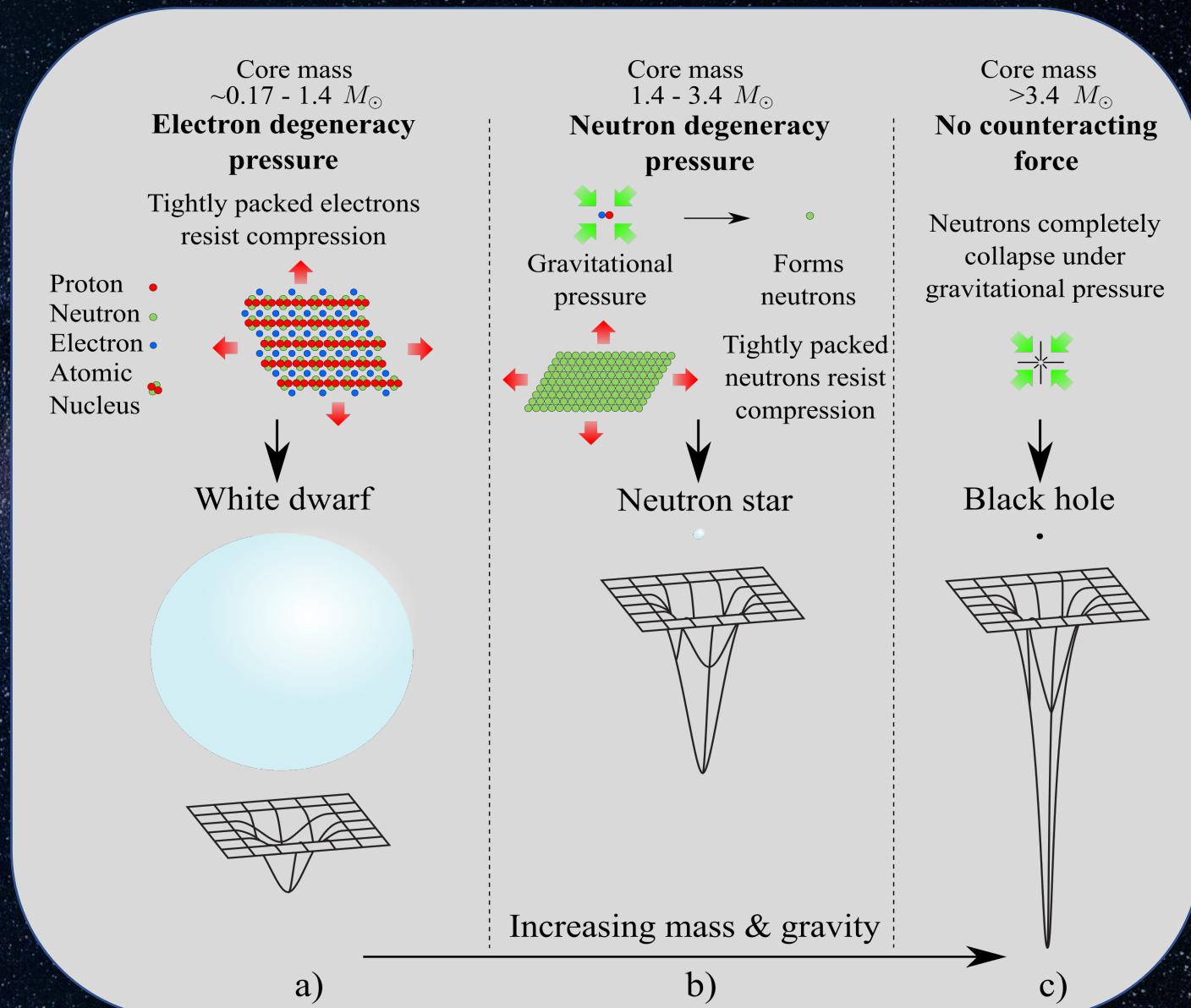
Relativity

Time

# Natural Labs?

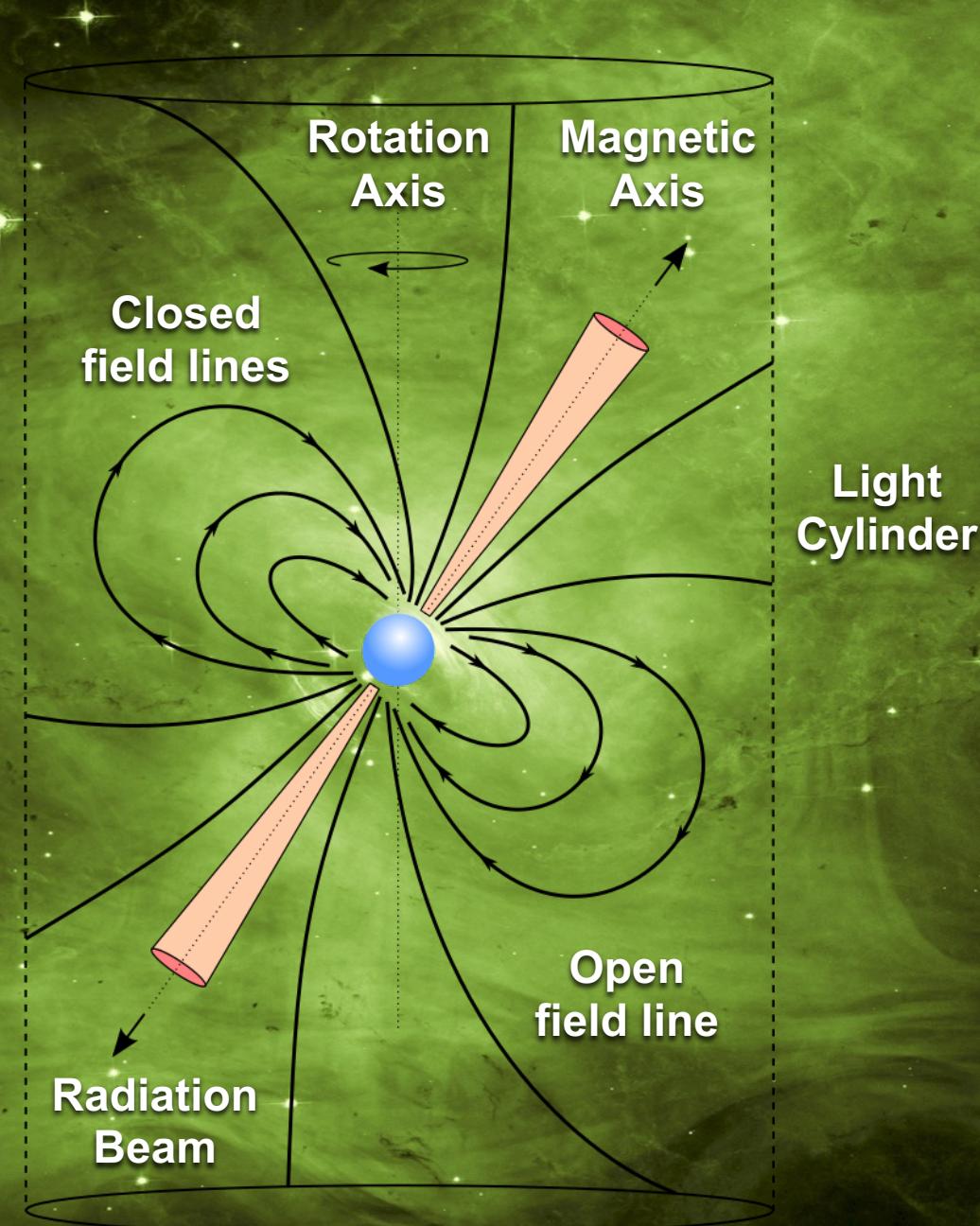


# Natural Labs?



# Radio Pulsars

- Stellar remnant
- Very dense
- ~20 km diameter
- Produce radio emission
- Very useful for science

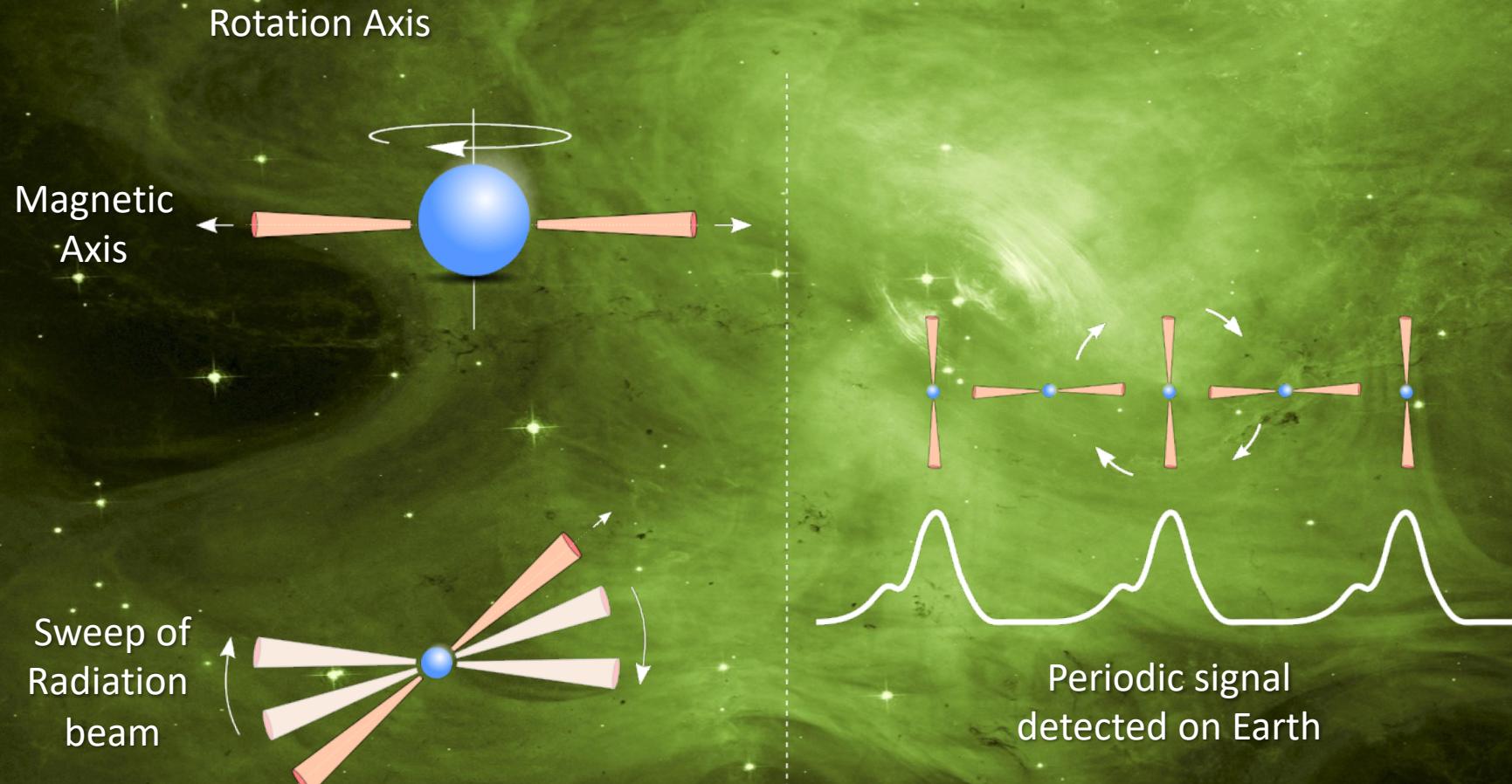


# Pulsar Birth

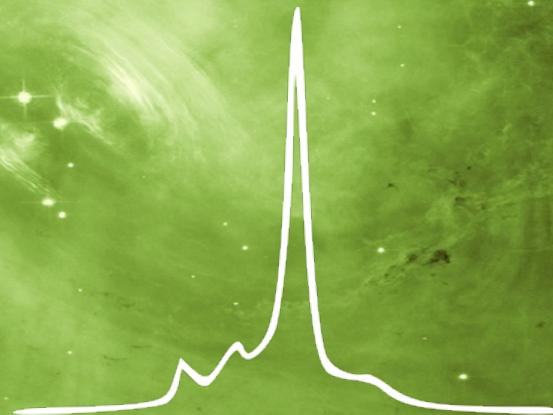
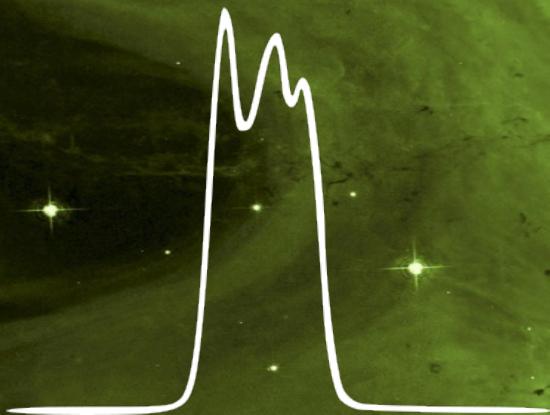


Credit: John Rowe Animations

# Pulsar Signals



# Pulse Profiles



Complex Pulse Profiles

# Search Tools



# Other Signals

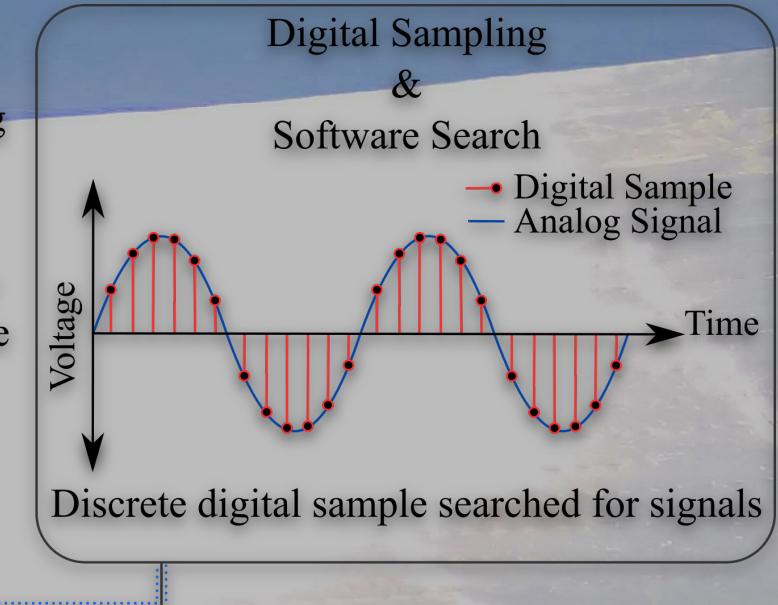
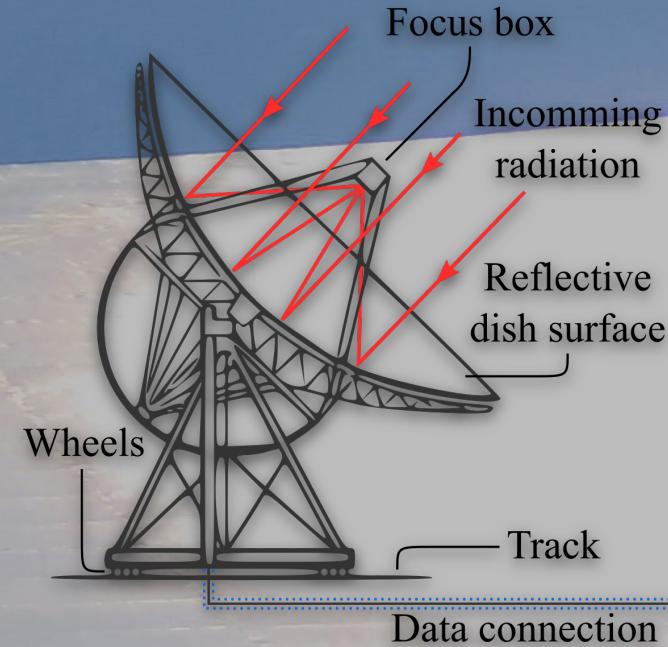
The figure shows a red line graph on a grid background. The x-axis is labeled with the number '10' at the far right. The y-axis has two horizontal dashed lines representing the upper and lower bounds of the data. The data points, represented by small red dots, form a series of peaks and troughs. A smooth red curve is drawn through these points, representing a fitted model. The residuals, which are the vertical distances between the data points and the fitted curve, are plotted as small vertical red lines extending above and below the main curve. These residuals show a clear periodic pattern, indicating a pulsar signal.

- Pulsar Timing
  - SETI Searches
  - Studies of Cosmic Magnetism
  - Studies of stellar evolution
  - Probing the Interstellar Medium
  - Radio Imaging
  - ...

Credit: Big Ear Radio Observatory and North American AstroPhysical Observatory (NAAPO).

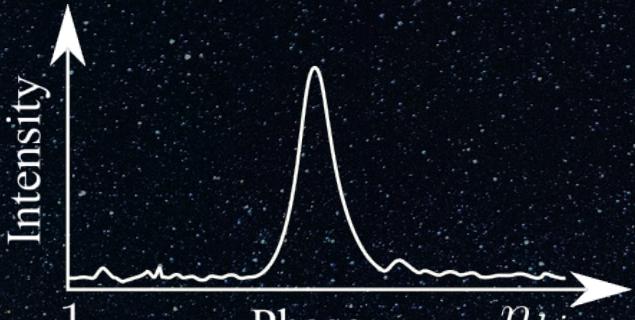
# Data Capture

- Analog to digital samples
- Complex search pipeline applied
- Steps:
  - RFI mitigation
  - Dedisperion
  - FFT
  - Harmonic summing
  - Detection
  - Sifting



# Data Products

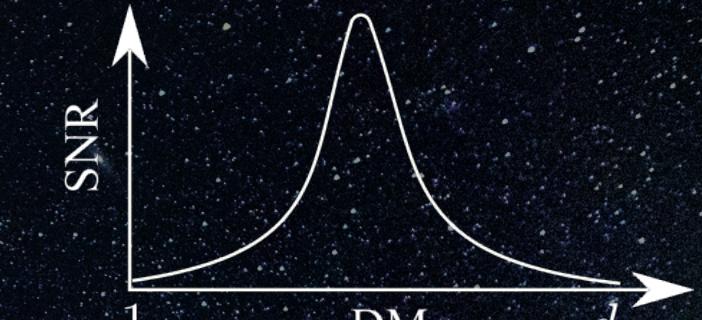
- A candidate file
- Describes a signal detection
- Covers time and frequency space
- At least one per detection (duplicates!)
- Classification appears simple - things do get fuzzy!



$$P = \begin{bmatrix} \square & \square & \square & \dots & \square & \square \end{bmatrix}$$

$\uparrow \quad \uparrow \quad \uparrow$   
 $p_1 \quad p_i \quad p_n$

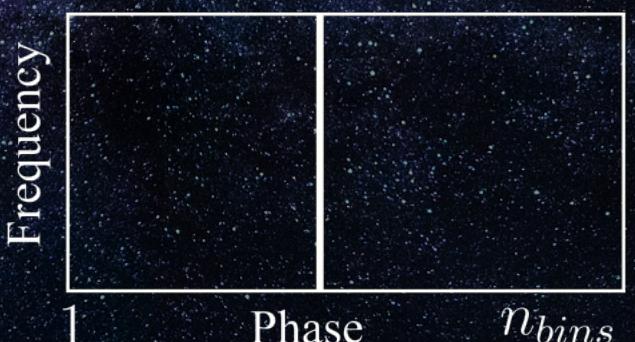
a) Integrated profile



$$D = \begin{bmatrix} \square & \square & \square & \dots & \square & \square \end{bmatrix}$$

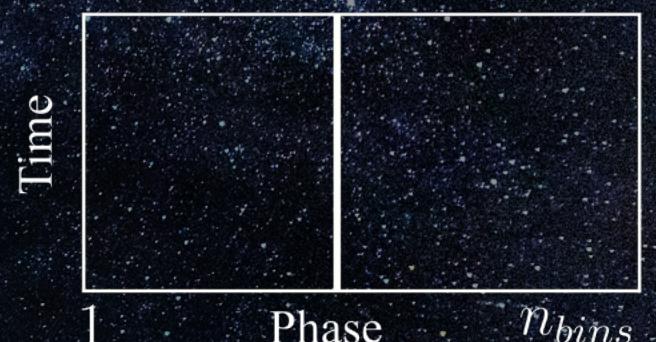
$\uparrow \quad \uparrow \quad \uparrow$   
 $d_1 \quad d_i \quad d_n$

b) DM-SNR curve



$$Sb = \begin{bmatrix} \square & \square & \square & \dots & \square & \square \\ \square & \square & \square & \dots & \square & \square \\ \vdots & \ddots & \vdots & & \vdots & \vdots \\ \square & \square & \square & \dots & \square & \square \end{bmatrix}$$

c) Sub-band matrix



$$Si = \begin{bmatrix} \square & \square & \square & \dots & \square & \square \\ \square & \square & \square & \dots & \square & \square \\ \vdots & \ddots & \vdots & & \vdots & \vdots \\ \square & \square & \square & \dots & \square & \square \end{bmatrix}$$

d) Sub-int matrix

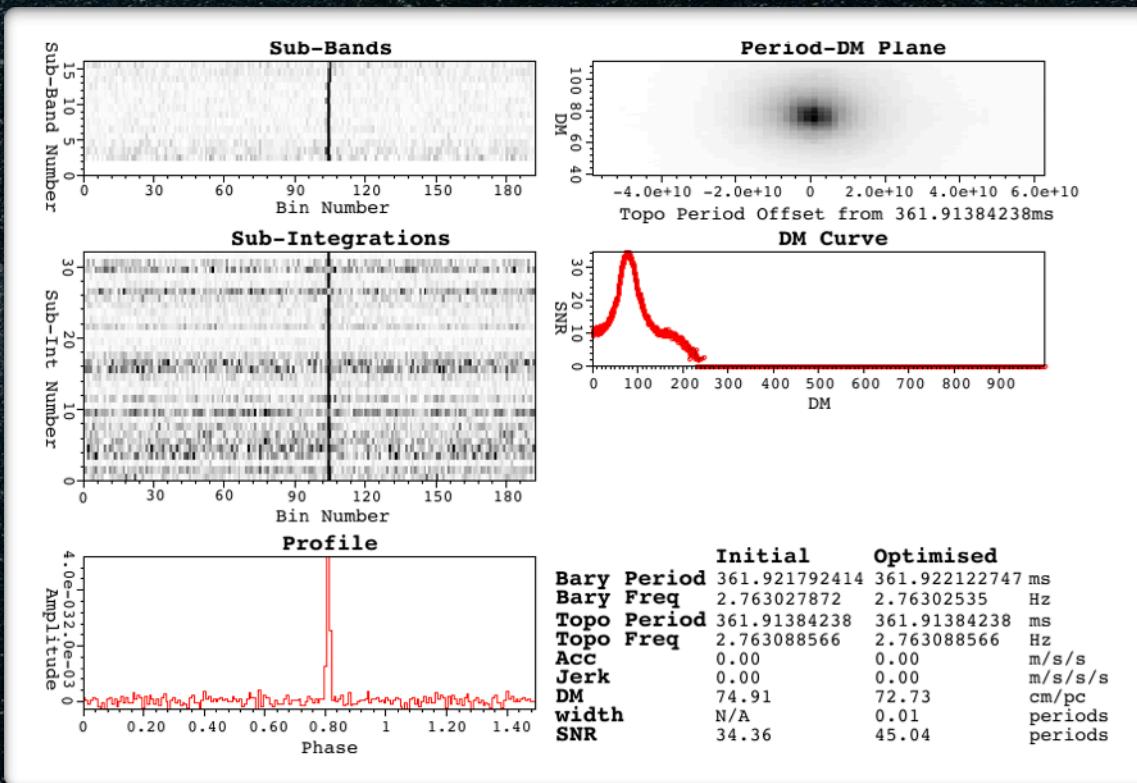
# What to look for?

- **Clear pulse with a defined peak.**
- **Evidence that the signal persists in time, i.e. an indication of a periodic source.**
- **Evidence that the signal persists in frequency - pulsars are broadband emitters.**
- **A DM value greater than zero.**
- **Evidence of other effects, such as scintillation?**



# Candidate Examples (1)

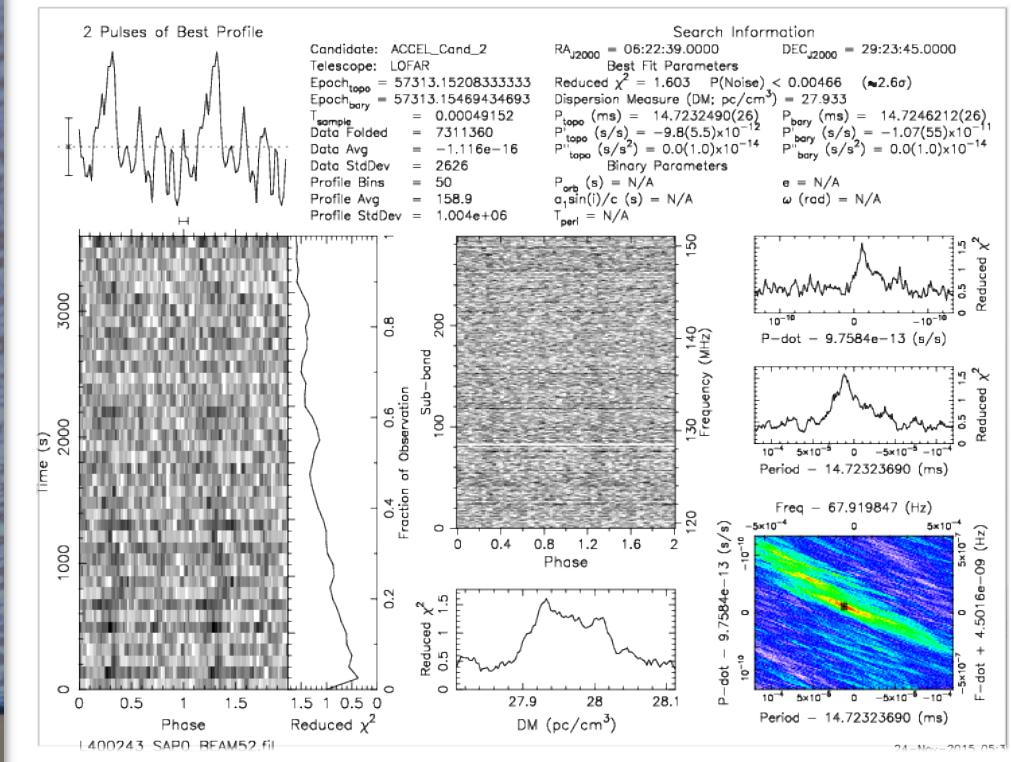
- **Defined peak**
- **Clear DM-SNR peak**
- **Consistent in time**
- **Persistent in frequency**
- **Obvious pulsar!**



Credit: HTRU Collaboration.

# Candidate Examples (2)

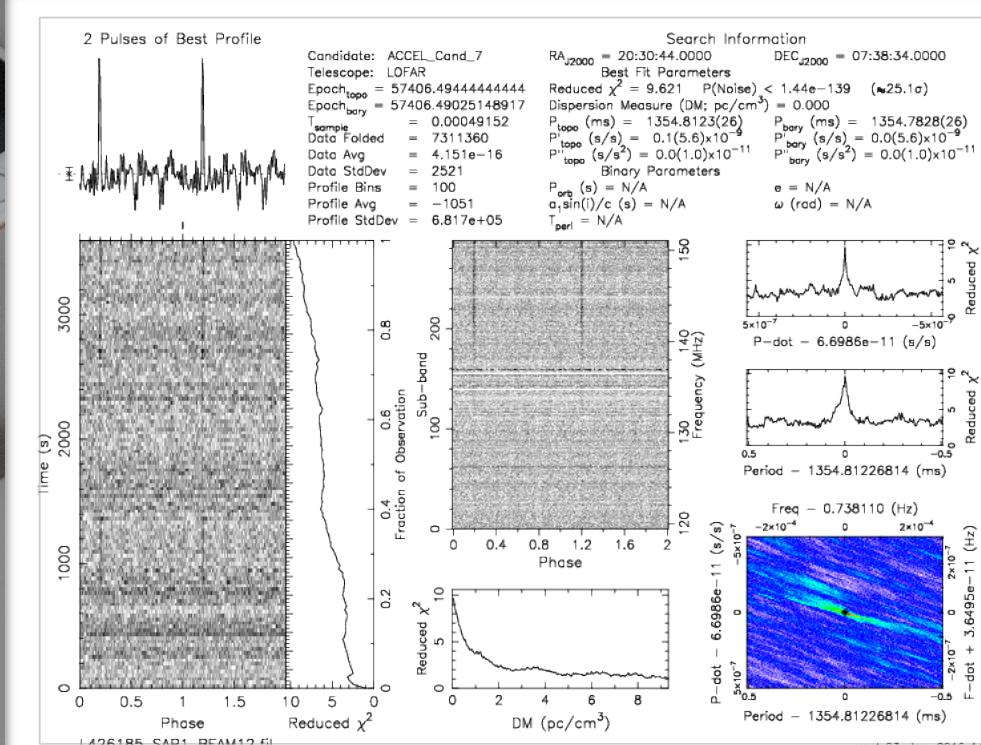
- Clearly defined peak?
- Some persistence in time?
- A DM value > zero
- Labelled as noise



Credit: LOTAAS Collaboration (Chia Min Tan et. al.).

# Candidate Examples (3)

- Clearly defined peak?
- Some persistence in time?
- Some persistence in frequency?
- RFI



Credit: LOTAAS Collaboration (Chia Min Tan et. al.).

# Things go wrong in practice



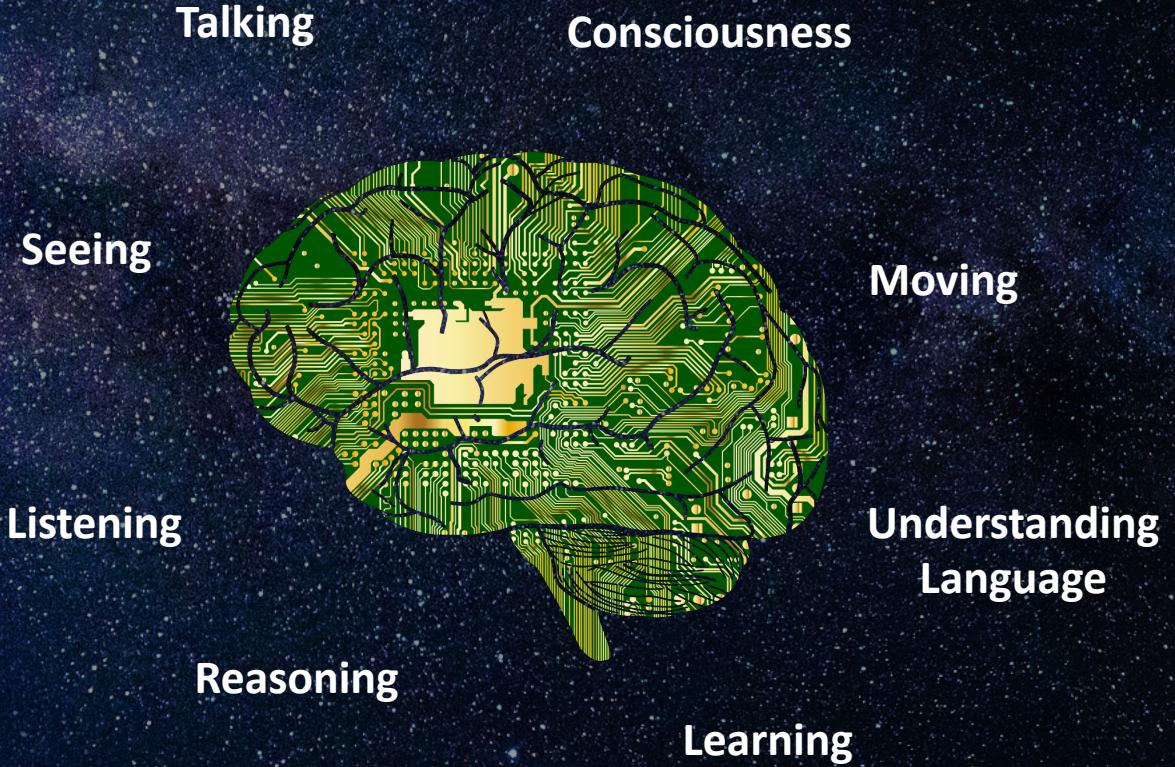
- Not all candidates so easily categorised.
- Huge volume of candidates usually returned by telescopes.
- This makes the problem particularly hard to solve.

# Imbalanced Data too!

- Candidate ratio is 1 real pulsar to 10,000 non-pulsar examples (not a worse case).
- Given that class distributions drift over time, we have a changing problem too.

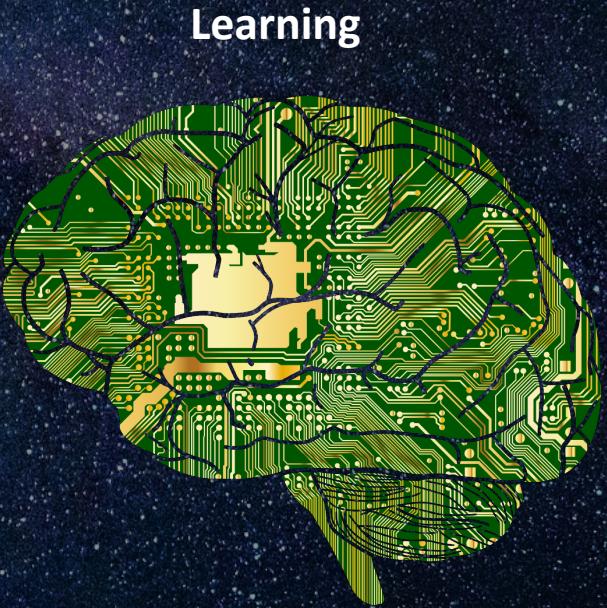
Tackling this with A.I.

# Artificial Intelligence



- Artificial Intelligence (A.I.) is a field of study concerned with reproducing/replicating human intelligence.
- A.I. is an umbrella term – it is comprised of many sub-fields of study.
- Here we'll talk about replicating the human capacity for learning.
- This field is known more generally as “Machine Learning”.

# Machine Learning



- Machine learning (ML) is a fascinating field. It combines insights from statistics, logic, psychology and neuroscience to build automated systems capable of “learning” by themselves.
- ML isn’t concerned with replicating exactly how humans learn, after all, we’re flawed!
- Instead it focuses on developing ways to make optimal decisions/predictions using available information.
- During this module you’ll acquire an understanding of ML that will allow you to apply it moving forward.

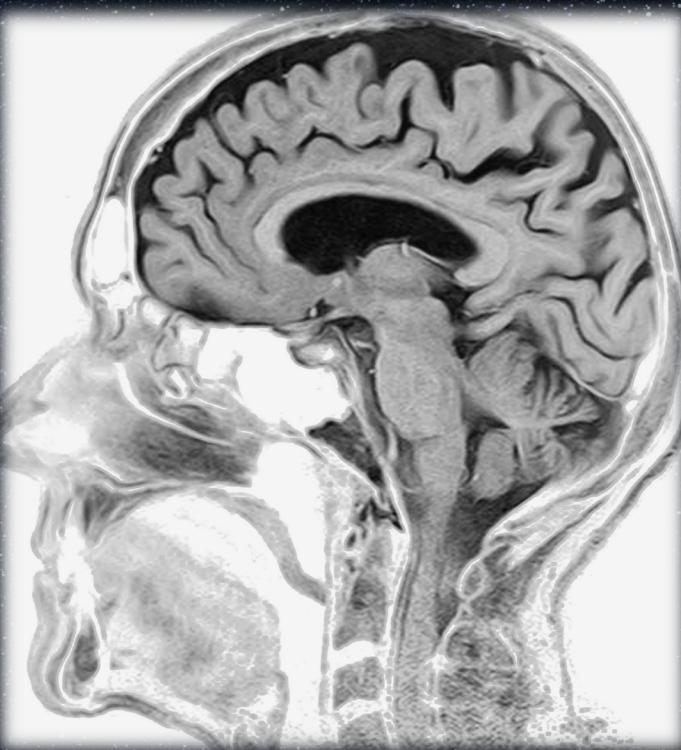
# Example Applications

Searching for rare stars in astronomy data

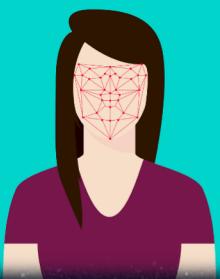


Image Credit: SARAO

Identifying regions of disease in medical images



Climate Modelling



Facial Recognition



Voice Recognition

Uncredited images obtained from <https://pixabay.com> (Creative Commons License, no attribution required).

# Human Decision Making

- Using available information / past experience we try to make optimal decisions.

Data (Sight,  
sound, taste,  
touch, smell)



# Quantifying Experience

- Experience is describable and quantifiable as data.
- You already intuitively understand what data is, but to understand ML, we must explain the terminology used in the field.

Variables / Attributes / Features				Ground Truth	
examples	Mass (Kg)	Height (cm)	Legs	Colour	Class / Label
	4.2	25.4	4	Black	Cat
	18	43.5	4	Brown	Dog
	0.1	5.1	4	Grey/Black	Mouse
	...	...	...	...	...

Table 1. Characteristics of animals we've seen

As the ground truth is known, this is called a labelled dataset.

# Quantifying Experience

**Example  
notation**

$x_i$

	Mass (Kg)	Height (cm)	Legs	Colour
$x_1$	4.2	25.4	4	Black
$x_2$	18	43.5	4	Brown
$x_3$	0.1	5.1	4	Grey/Black

Class / Label
$y$
Cat
Dog
Mouse

**Example  
notation**

$y_1$

$y_2$

$y_3$

$y_i$

Table 1. Characteristics of animals we've seen

**Variables / Attributes / Features**

**Ground Truth**

**Array**

$$x_1 = \{4.2, 25.4, 4, Black\}$$

$$y_1 = \{Cat\}$$

$$x_2 = \{18, 43.5, 4, Brown\}$$

$$y_2 = \{Dog\}$$

$$x_3 = \{0.1, 5.1, 4, Grey/Black\}$$

$$y_3 = \{Mouse\}$$

# Quantifying Experience

Variables / Attributes / Features	Ground Truth
$x_1 = \{4.2, 25.4, 4, Black\}$	$y_1 = \{Cat\}$
$x_2 = \{18, 43.5, 4, Brown\}$	$y_2 = \{Dog\}$
$x_3 = \{0.1, 5.1, 4, Grey/Black\}$	$y_3 = \{Mouse\}$

$$E = \{(x_1, y_1), (x_2, y_1), \dots, (x_n, y_1)\}$$

- Experience is comprised of pairs of feature values and ground truth labels.
- The pairs are actually called “tuples”.

# Sometimes Knowledge is Incomplete

- Sometimes we don't know the ground truth. In such cases, we have to collect information, and provide it for ourselves.
- This can be costly and time consuming.

Variables / Attributes / Features				Ground Truth	
	Ext. Temp (°C)	Int. Temp (°C)	Time (24hr)	Consumption (kWh)	Class / Label $y$
$x_1$	4.2	17.1	15:23	7500	?
$x_2$	33.4	27.8	18:01	7000	?
$x_3$	10.8	19.1	20:36	450	Deactivated
...	...	...	...	...	...

Table 2. Status of a heating control system

- If the ground truth is unknown, this is called an “unlabelled” dataset.
- Datasets may also be “partially-labelled”.

# Feature Values

- Features can be categorical or numerical.
- Usually we “discretise” categorical features as this makes them easier for many algorithms to work with. This is either done automatically using tools, or manually.

Mass (Kg)	Height (cm)	Legs	Colour	Class / Label $y$
0.229	25.4	4	1	1
1.0	43.5	4	2	2
0	5.1	4	4	3
...	...	...	...	...

Table 1. Characteristics of animals we've seen

- We may also “normalise” our data from time to time.

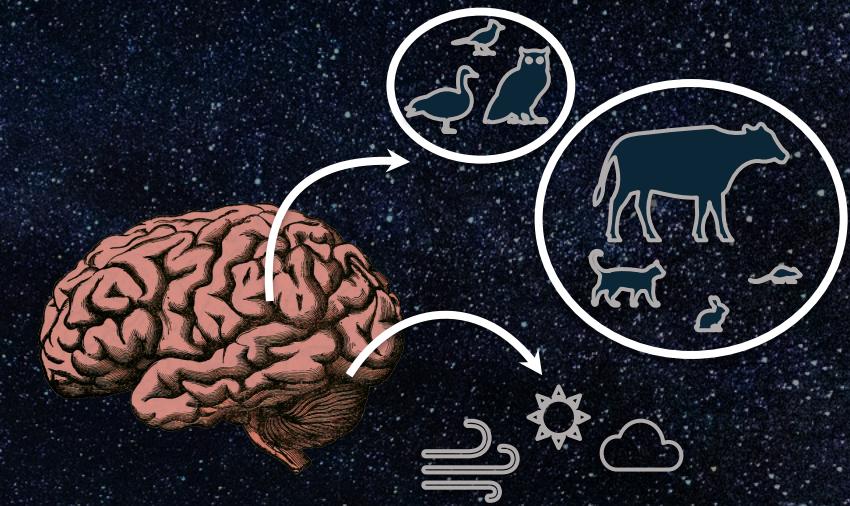
# Human Decision Making

**Features (Variables)**

$$x_i = \{ \text{Height, Mass, ..., Age} \}$$

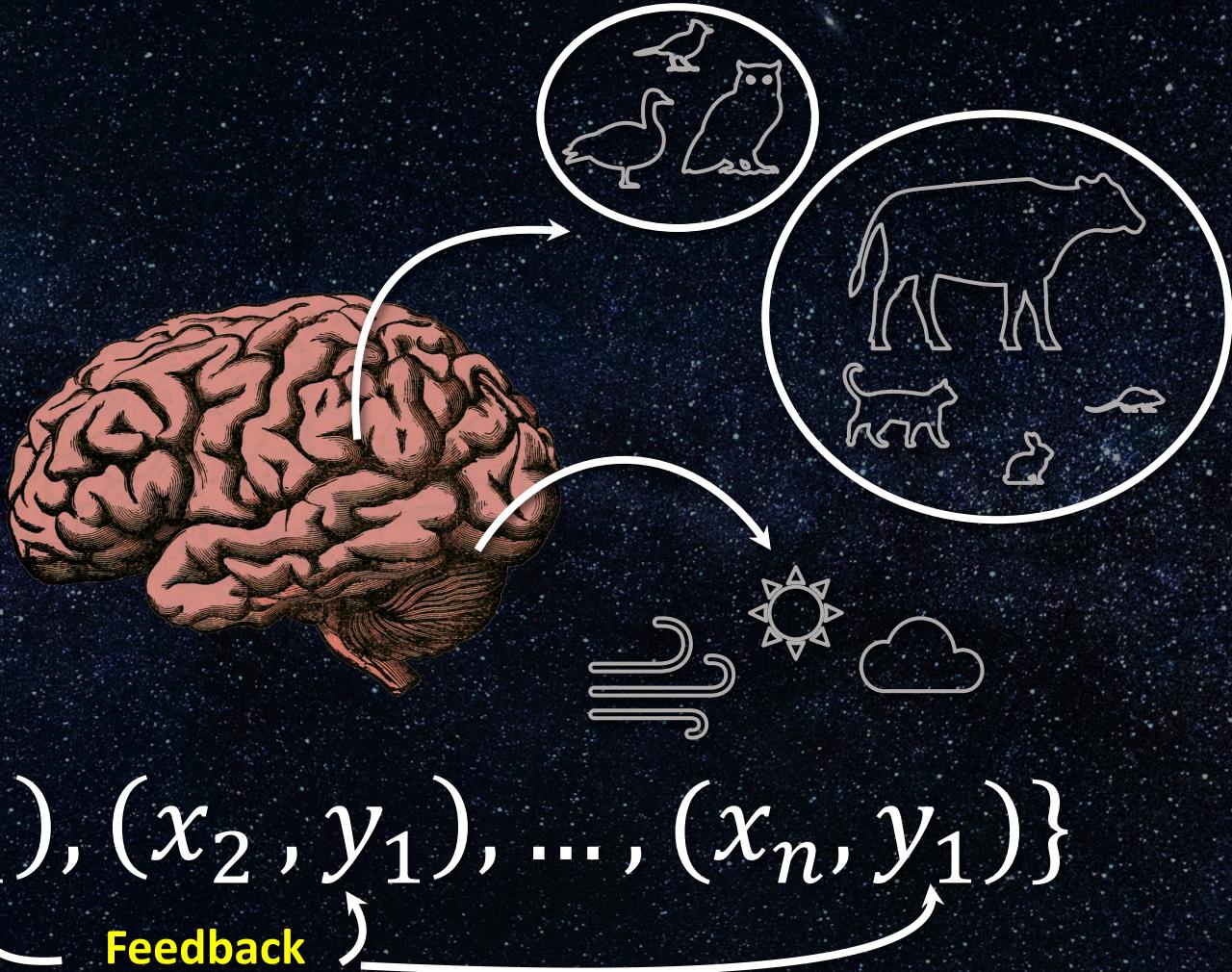
**Labels ( Feedback )**

$$y_i = \{ \text{Cat, Dog, ..., Bird} \}$$



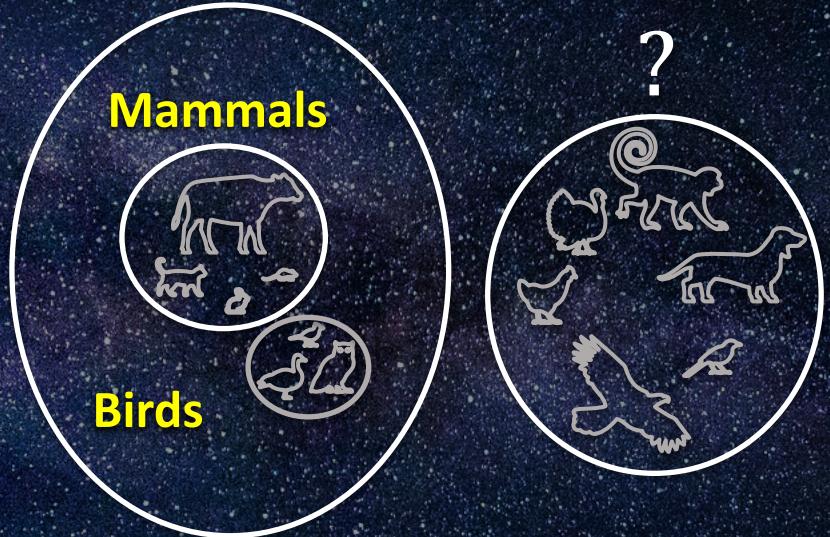
# Human Decision Making

- Humans capable of accurately,
  - Clustering
  - Classifying
  - Predicting
- Made possible via a feedback cycle.



# Approaches to Learning

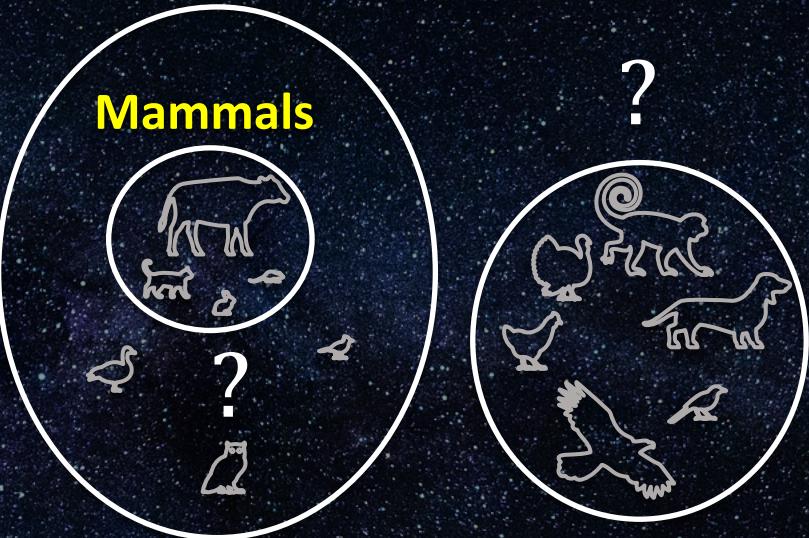
## Supervised Learning



*E* Contains labels for mammals & birds, goal is to correctly group unseen animals into mammal & bird classes.

Labeled data

## Semi-supervised Learning



*E* Contains examples labels for mammals only, goal is to correctly group unseen animals into mammal & bird classes.

Partially-Labeled data

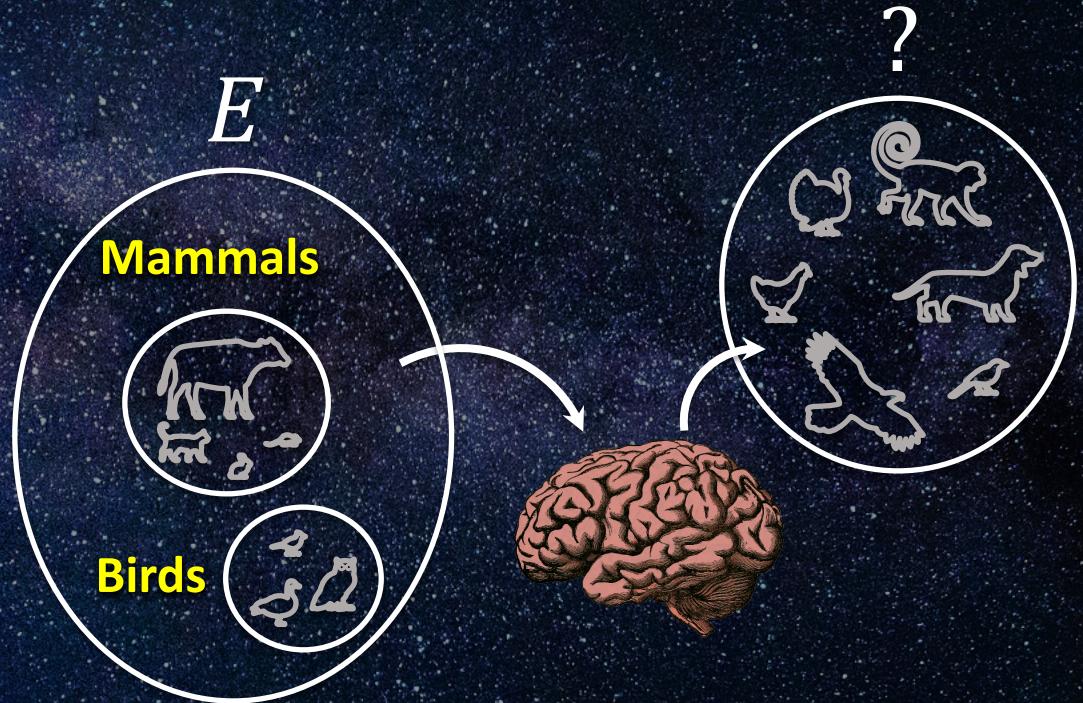
## Unsupervised Learning



*E* Contains no labels. Instead we use the data to self describe class separation.

Unlabeled data

# Classification



- These problems are “binary” classification tasks.
- Called binary, as there are two options (mammal or bird).
- Such classification tasks are everywhere, and we do them every day.
- To complete them we make predictions.
- Classification problems in the real-world are usually far more complex, i.e. more than 2 potential classes to predict.
- These are known as multi-class problems. Let's try one.

# Example

Feline



Lutrinae  
(e.g. Otter)



Canine



?



Mongoose



Viverridae  
(e.g. Civet)



Image Credit: Ran Kirlian - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=4164754>  
Uncredited images obtained from <https://pixabay.com> (Creative Commons License, no attribution required).

# Could you generalise?

- In the last slide, we saw an animal you may not be familiar with.
- Even if you have extensive knowledge of the animal kingdom, you may not have known the animal was a Fossa.
- Perhaps this animal was not described in the dataset,  $E$ .
- Even if you did not know it was a Fossa, you could still tell it was a type of mammal.
- This application of your knowledge, from past experience, illustrates your ability to generalise beyond known facts.
- Humans are gifted generalisers - but they are susceptible to two problems:

Overfitting

Underfitting

# Making Optimal Decisions

- Many have tried to understand how to make optimal decisions.
- We know we should use available evidence at all times.
- Humans are not always so thorough - we do make bad decisions.
- We are biased decision makers - we often use instinct and personal experience to decide.

