

ANNA SCAIFE

JODRELL BANK CENTRE FOR ASTROPHYSICS

 @RADASTRAT

REAL-TIME CLASSIFICATION AT SKA- SCALE FOR TIME DOMAIN ASTROPHYSICS

WITH CONTENT FROM ROB LYON

 @SCIENCEGUYROB

ANNA SCAIFE

JODRELL BANK CENTRE FOR ASTROPHYSICS

 @RADASTRAT

REAL-TIME CLASSIFICATION AT SKA- SCALE FOR TIME DOMAIN ASTROPHYSICS

WITH CONTENT FROM ROB LYON

 @SCIENCEGUYROB

GRAVITATIONAL WAVES

- RIPPLES IN SPACE-TIME
- USE A NETWORK OF PULSARS TO DETECT THEM
- CAUSED BY INTER-ACTING SUPER-MASSIVE BLACK HOLES



GRAVITATIONAL WAVES

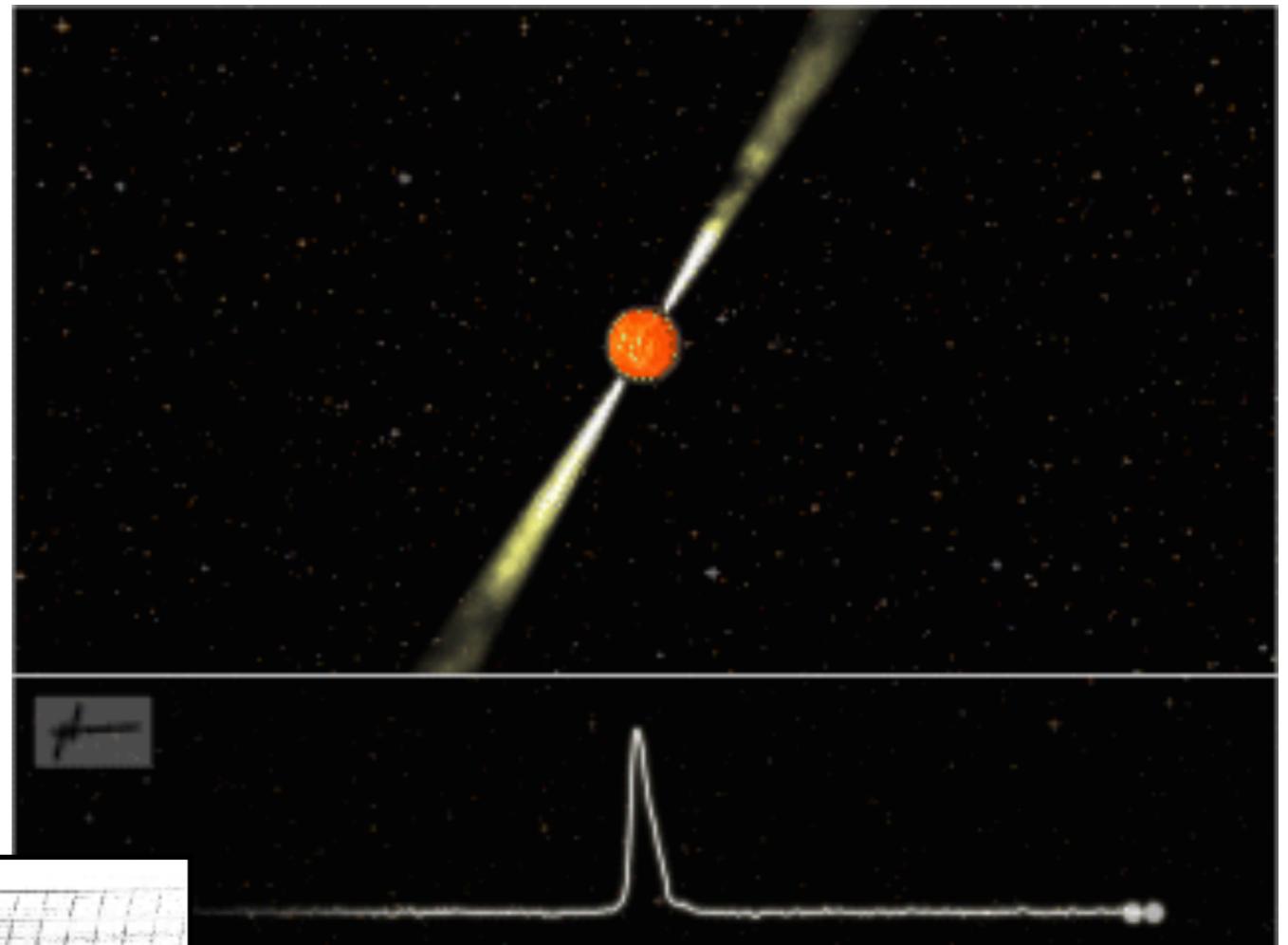
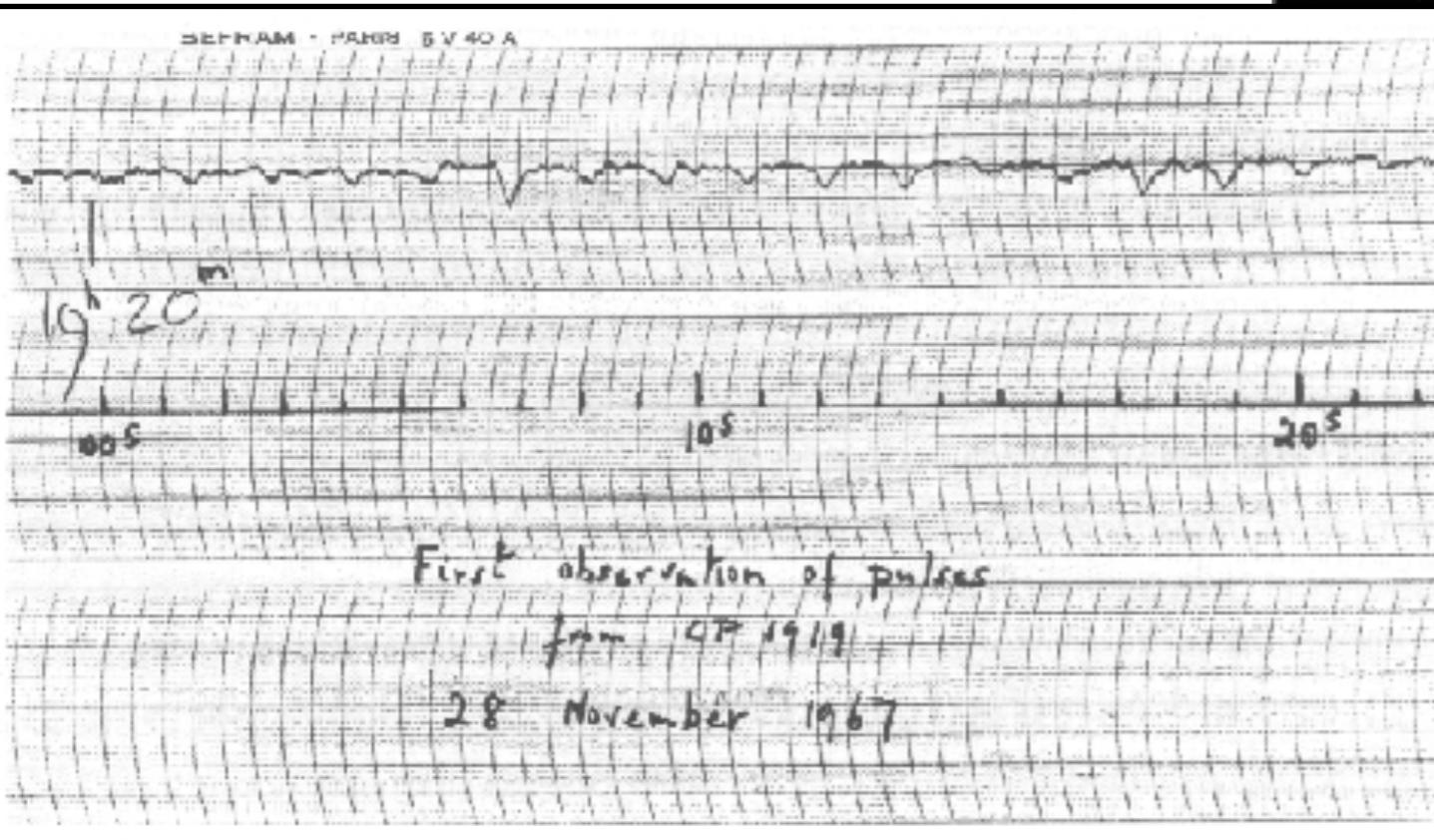
- RIPPLES IN SPACE-TIME
- USE A NETWORK OF PULSARS TO DETECT THEM
- CAUSED BY INTER-ACTING SUPER-MASSIVE BLACK HOLES



WHAT IS A PULSAR?

PSR B0329+54: This pulsar is a typical, normal pulsar, rotating with a period of 0.714519 seconds, i.e. close to 1.40 rotations/sec.

<http://www.jb.man.ac.uk/distance/frontiers/pulsars/section1.html>

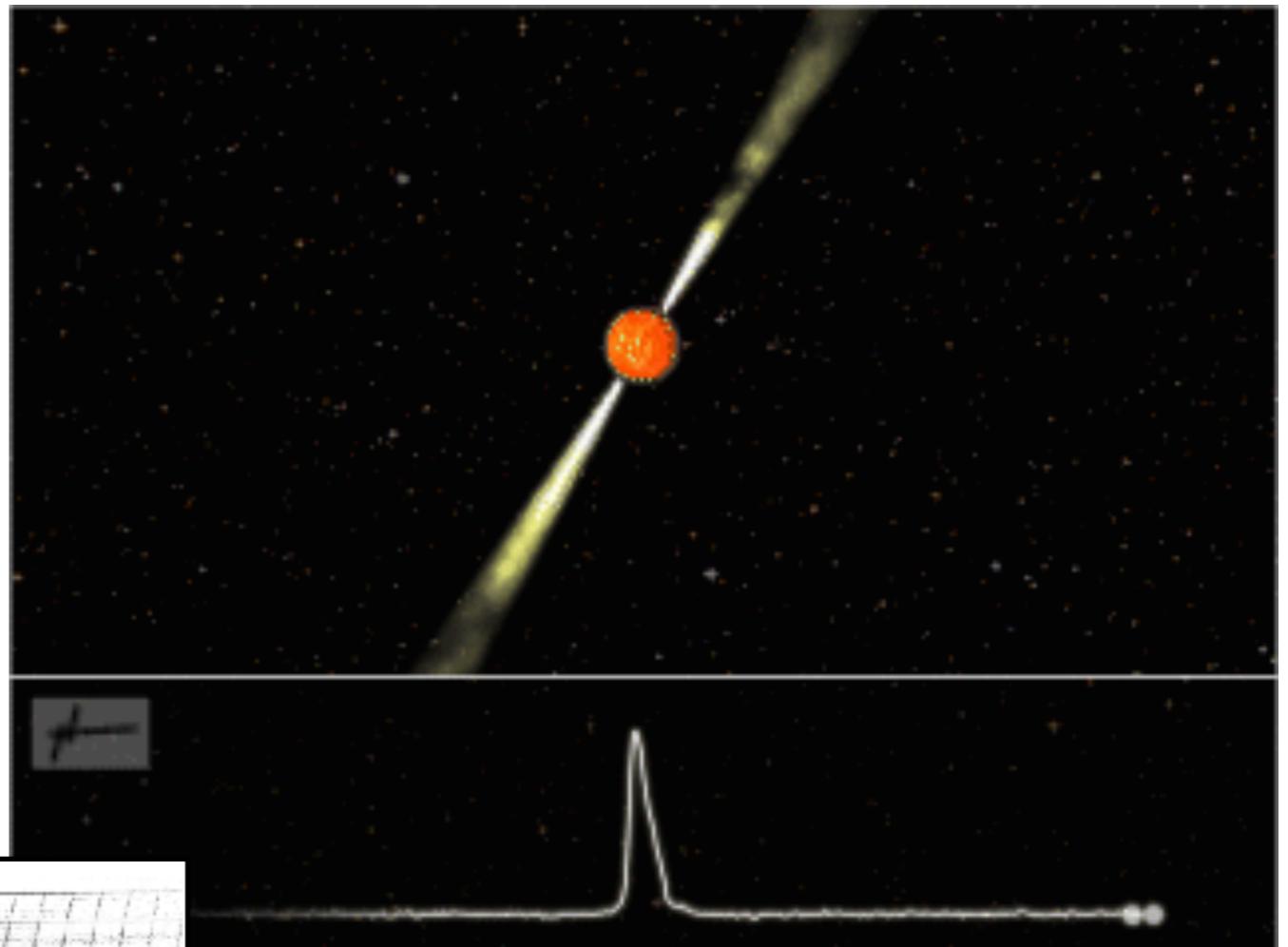
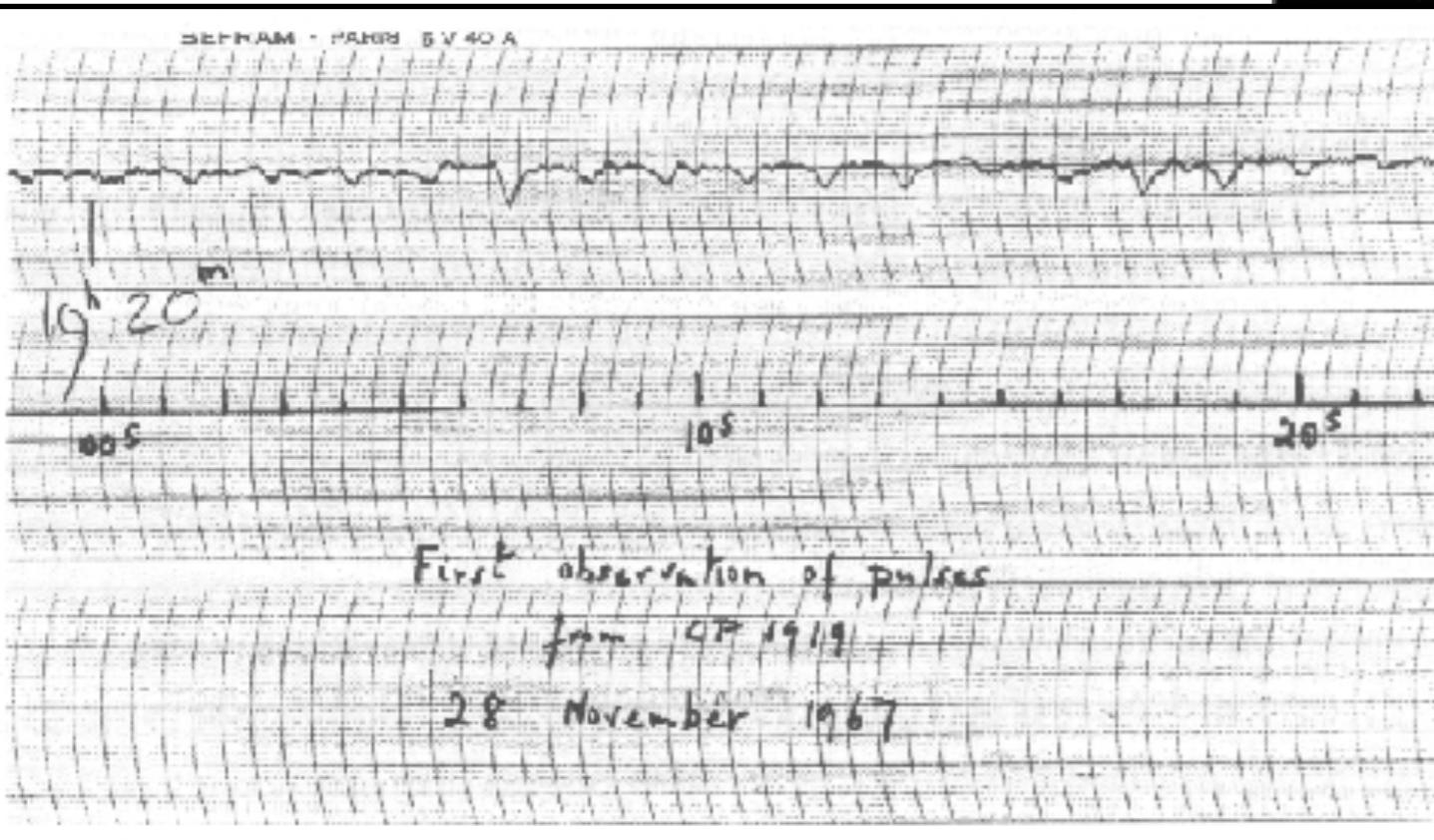


An artist's impression of a pulsar. Image credit: [Joeri van Leeuwen](#), License: [CC-BY-AS](#)

WHAT IS A PULSAR?

PSR B0329+54: This pulsar is a typical, normal pulsar, rotating with a period of 0.714519 seconds, i.e. close to 1.40 rotations/sec.

<http://www.jb.man.ac.uk/distance/frontiers/pulsars/section1.html>

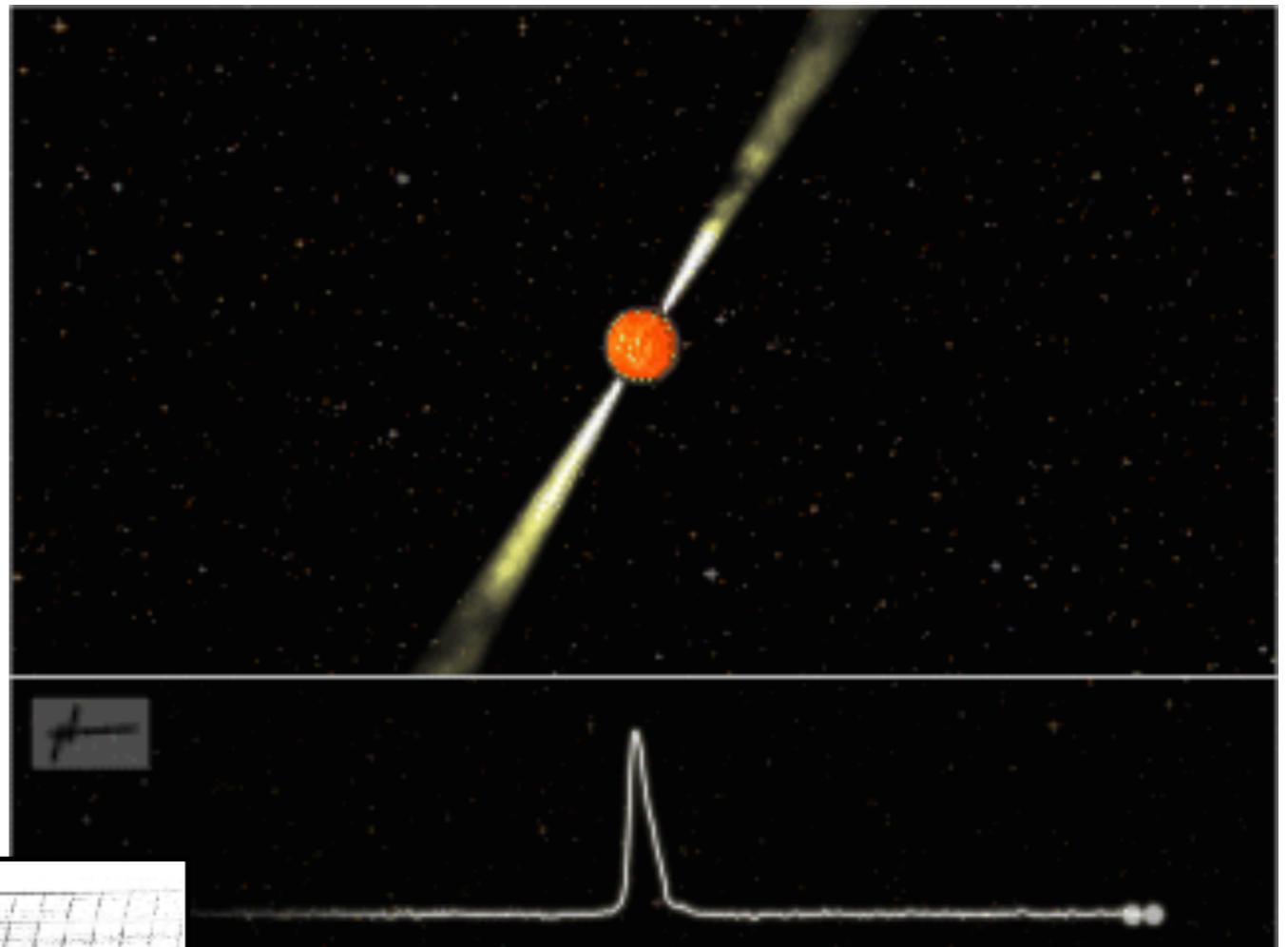
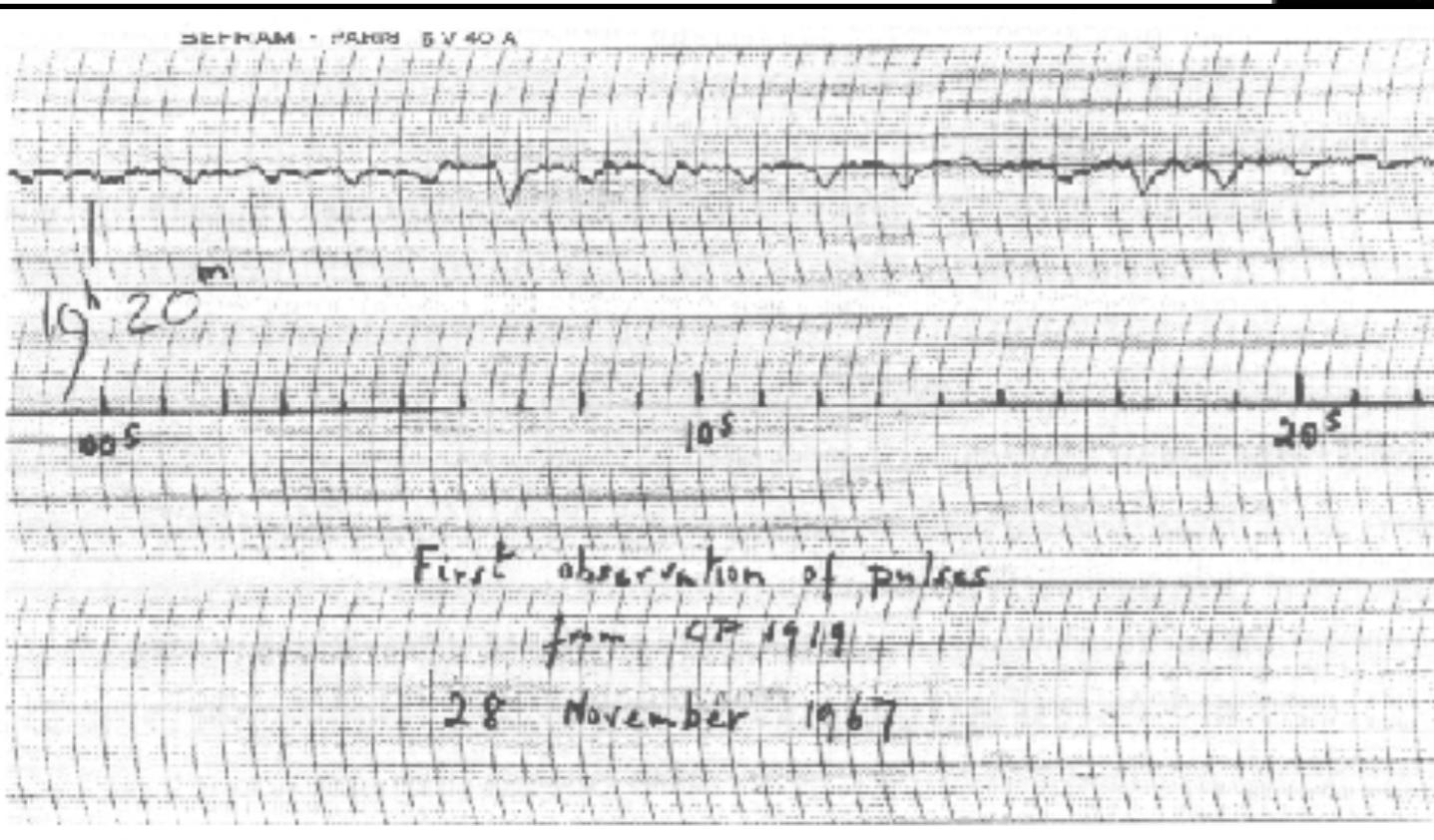


An artist's impression of a pulsar. Image credit: [Joeri van Leeuwen](#), License: [CC-BY-AS](#)

WHAT IS A PULSAR?

PSR B0329+54: This pulsar is a typical, normal pulsar, rotating with a period of 0.714519 seconds, i.e. close to 1.40 rotations/sec.

<http://www.jb.man.ac.uk/distance/frontiers/pulsars/section1.html>

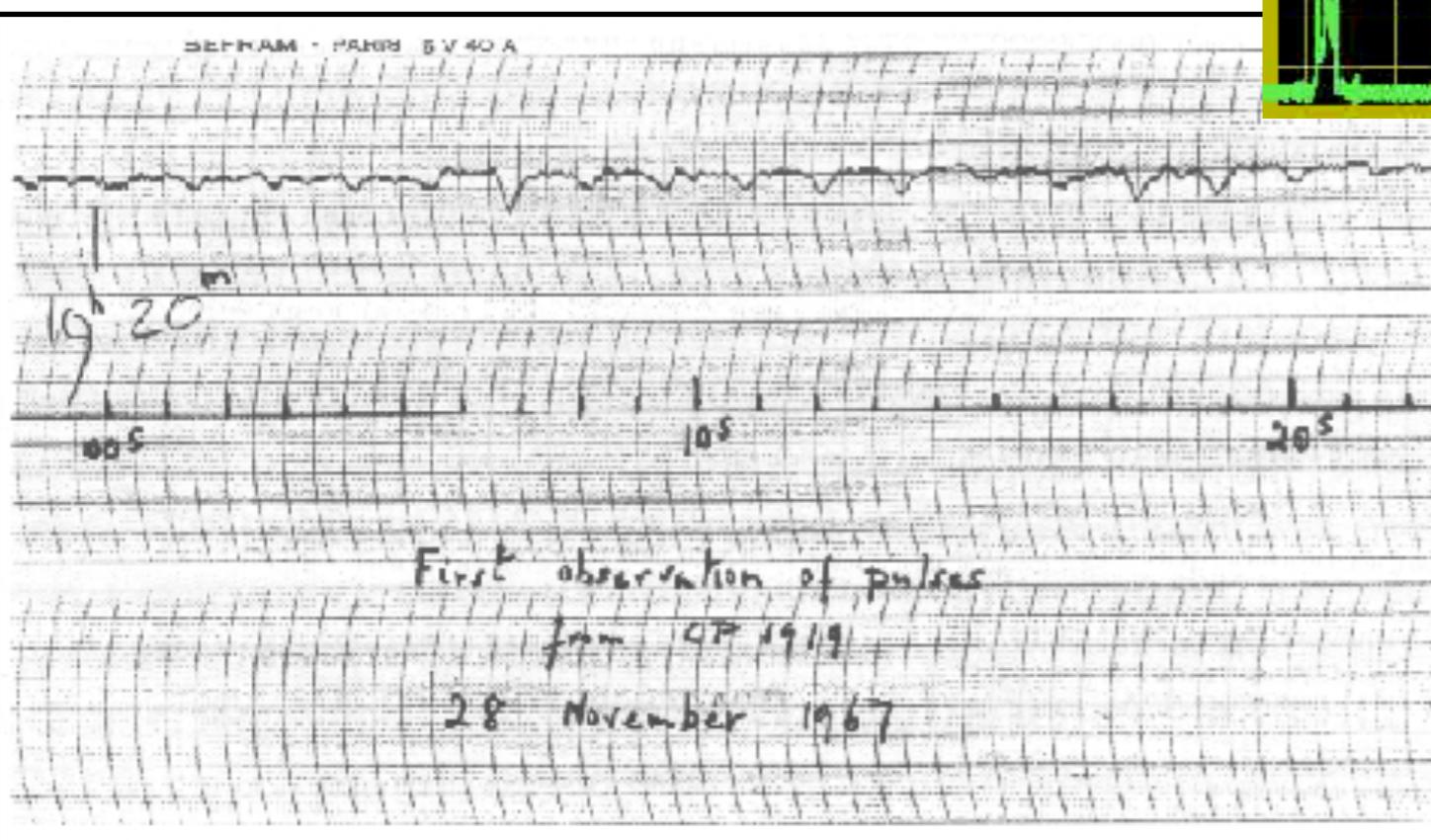
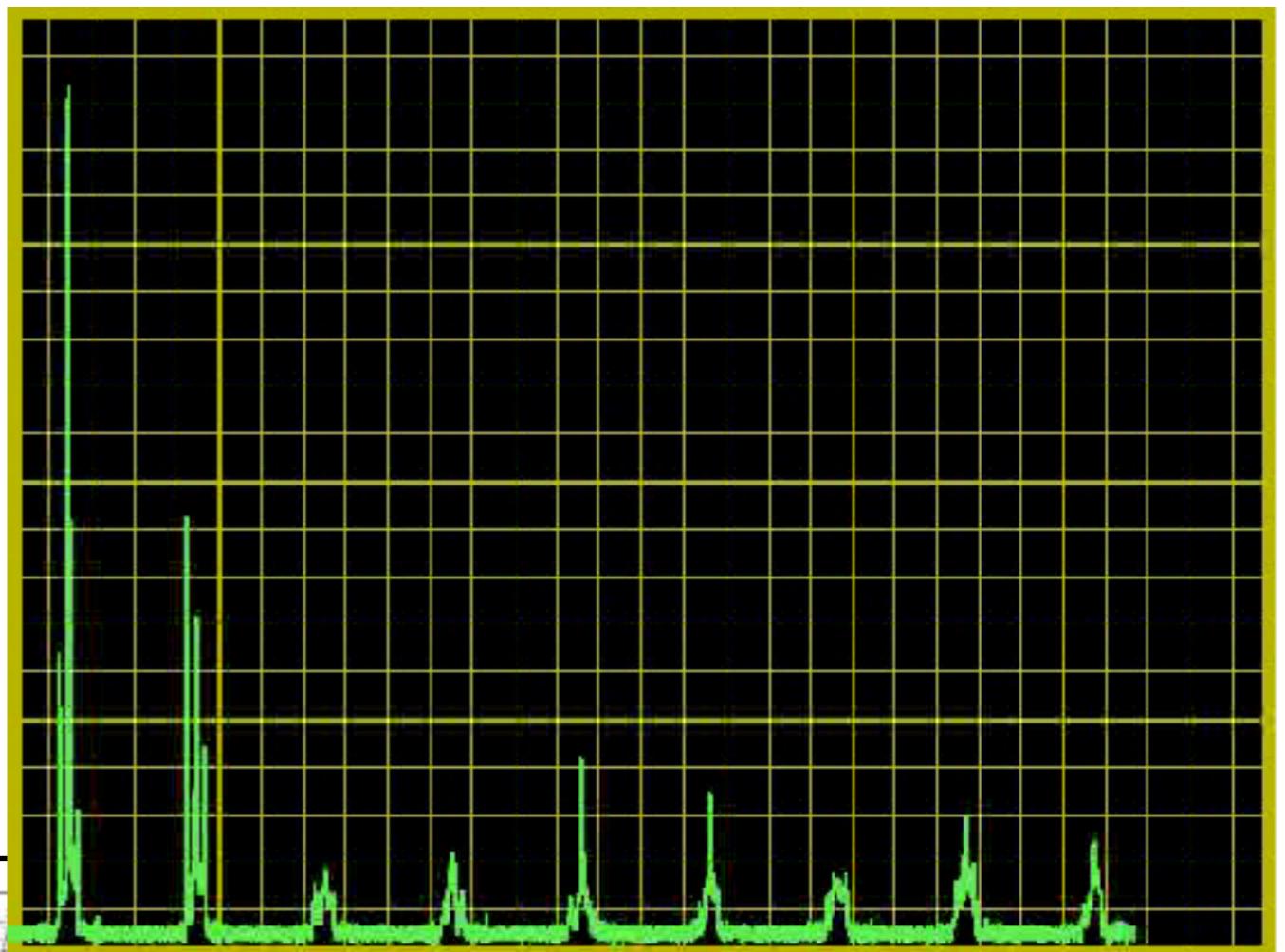


An artist's impression of a pulsar. Image credit: [Joeri van Leeuwen](#), License: [CC-BY-AS](#)

WHAT IS A PULSAR?

PSR B0329+54: This pulsar is a typical, normal pulsar, rotating with a period of 0.714519 seconds, i.e. close to 1.40 rotations/sec.

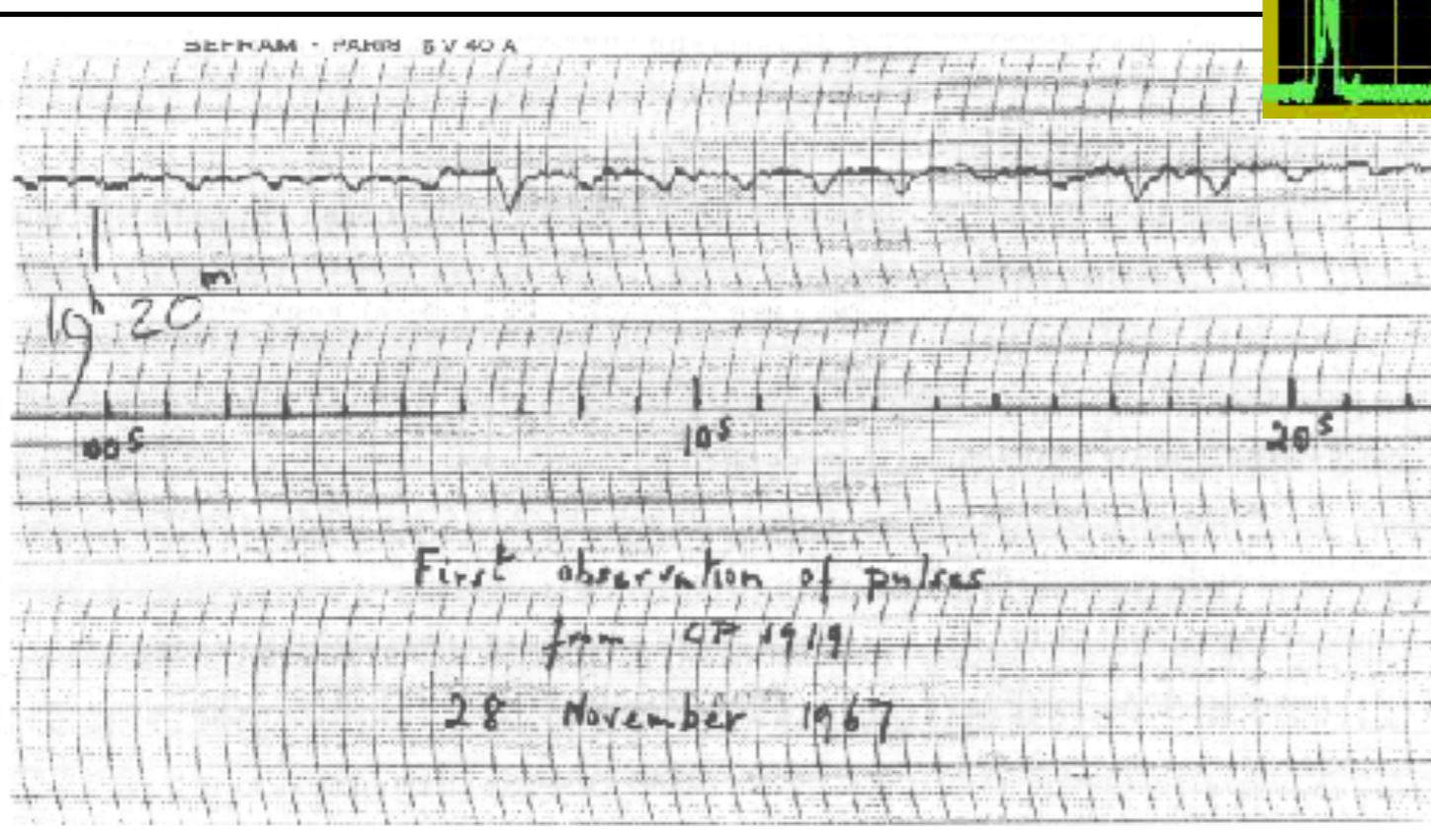
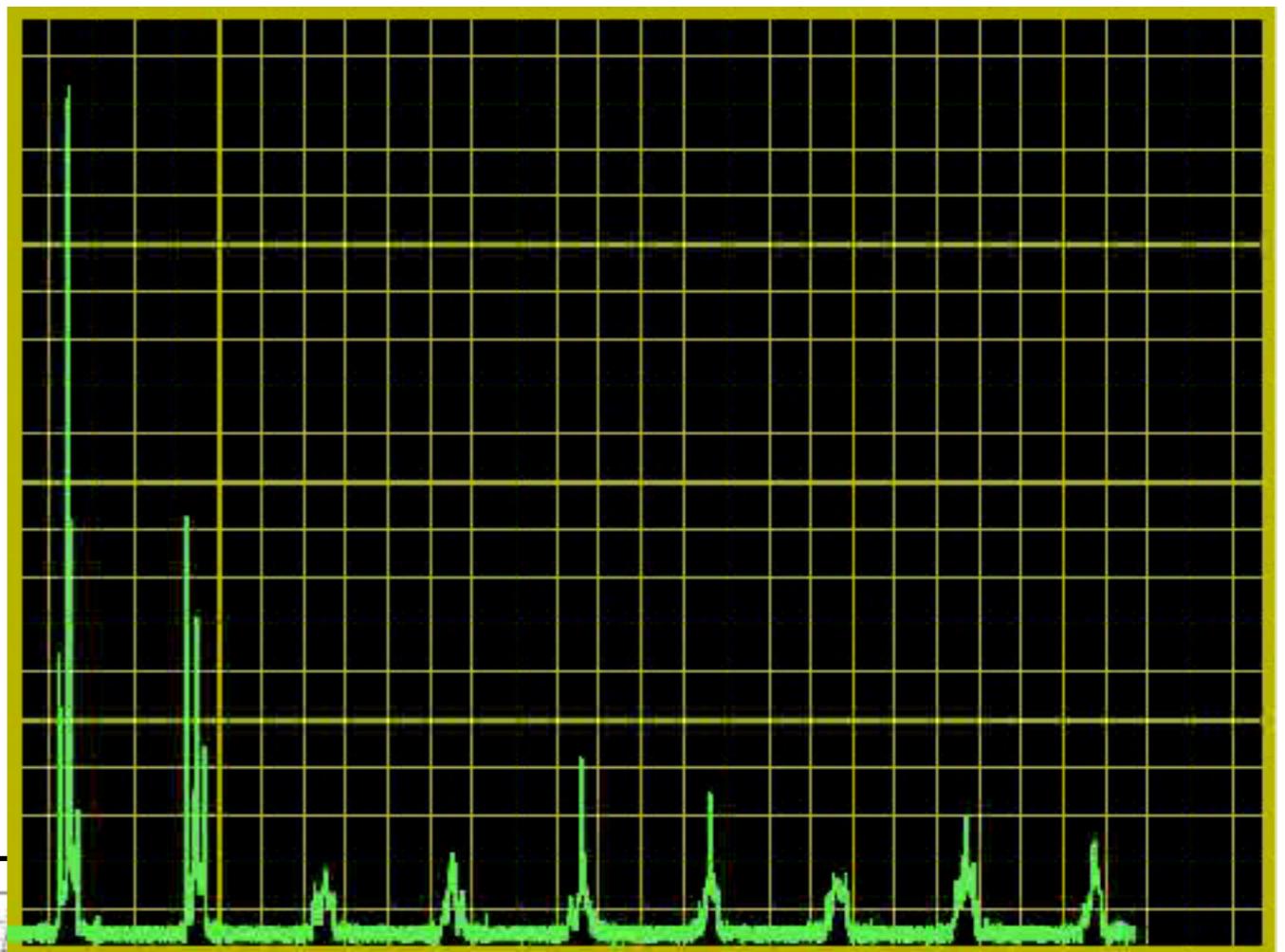
<http://www.jb.man.ac.uk/distance/frontiers/pulsars/section1.html>



WHAT IS A PULSAR?

PSR B0329+54: This pulsar is a typical, normal pulsar, rotating with a period of 0.714519 seconds, i.e. close to 1.40 rotations/sec.

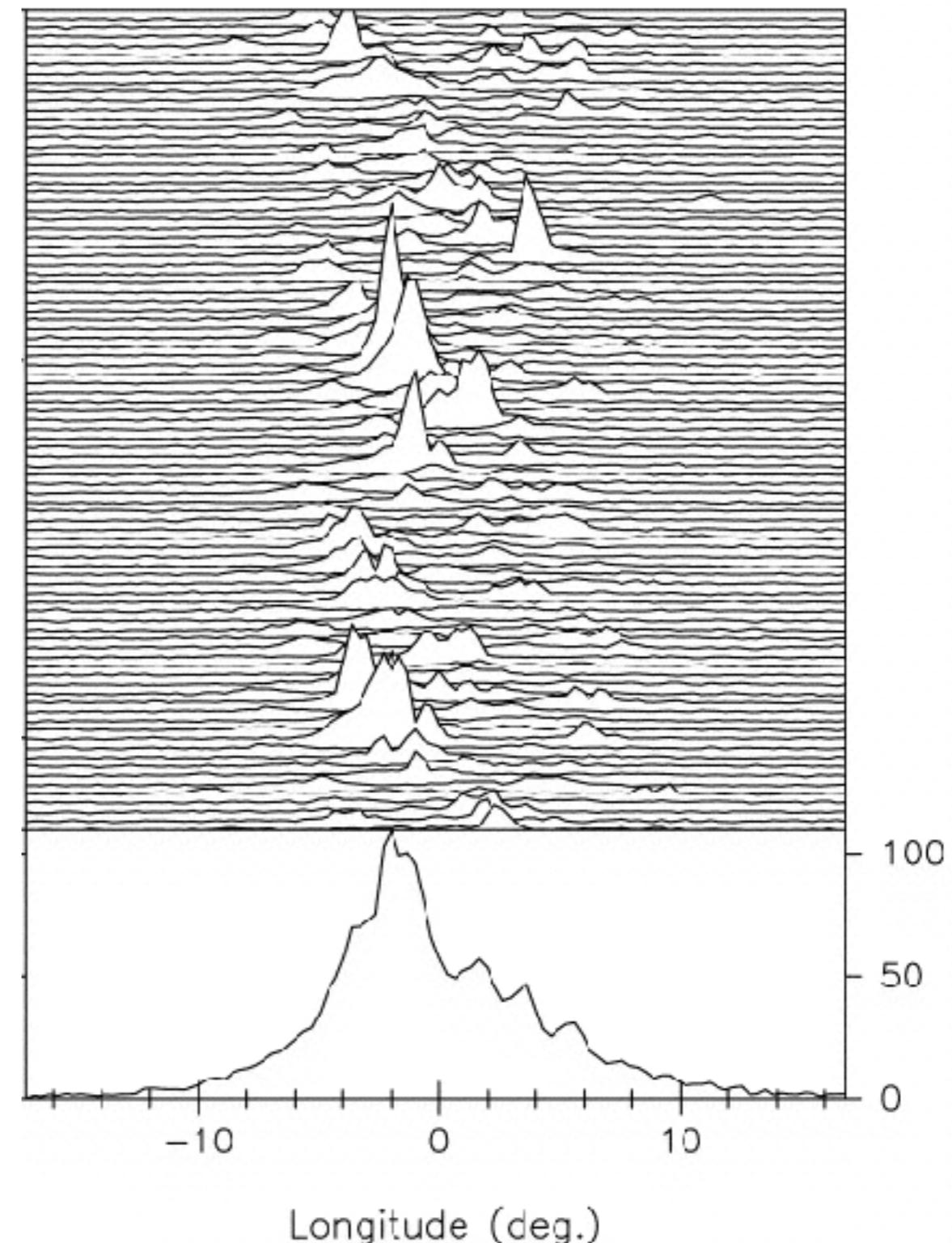
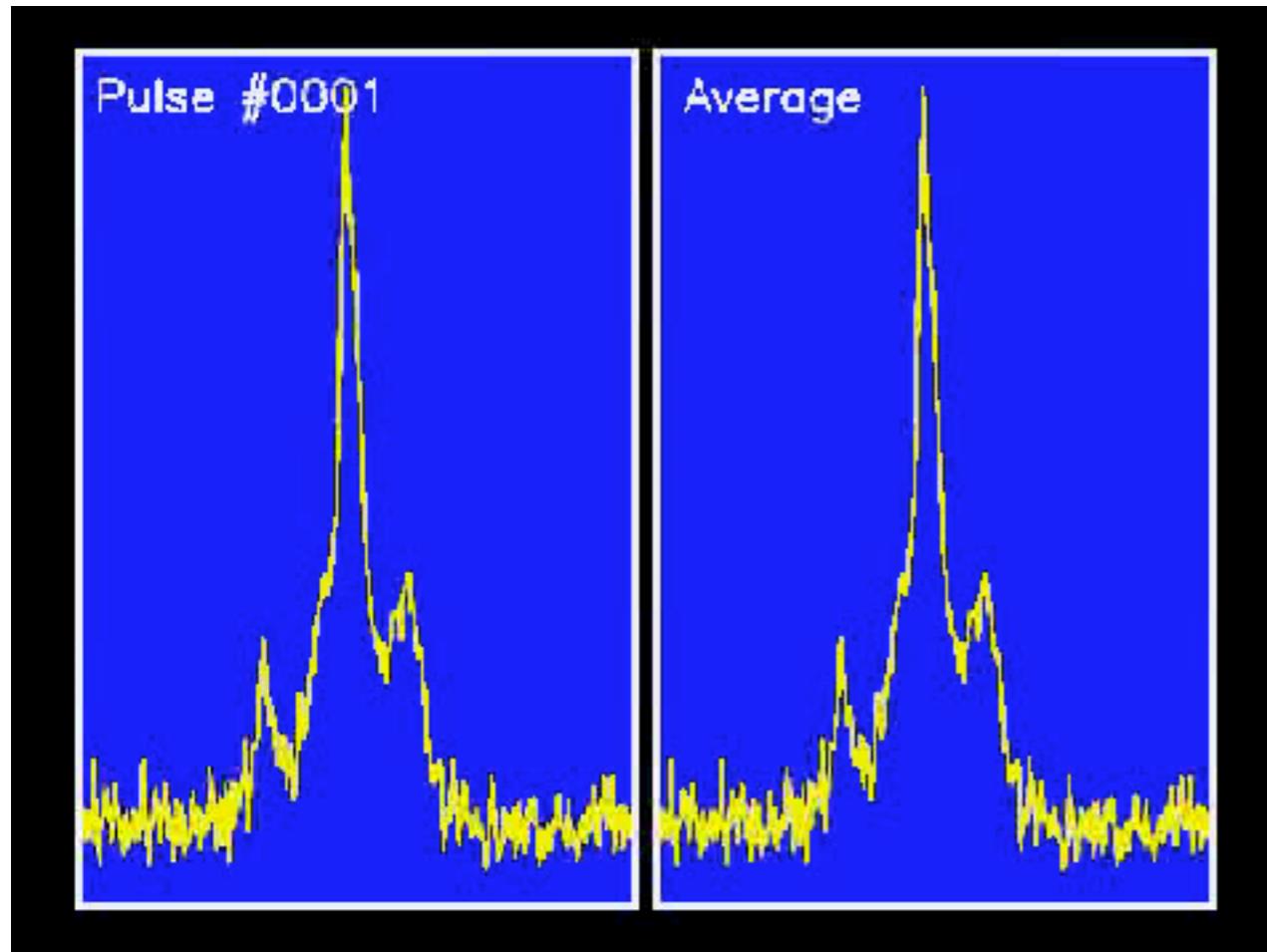
<http://www.jb.man.ac.uk/distance/frontiers/pulsars/section1.html>



INTEGRATED PROFILE

PSR B0329+54: This pulsar is a typical, normal pulsar, rotating with a period of 0.714519 seconds, i.e. close to 1.40 rotations/sec.

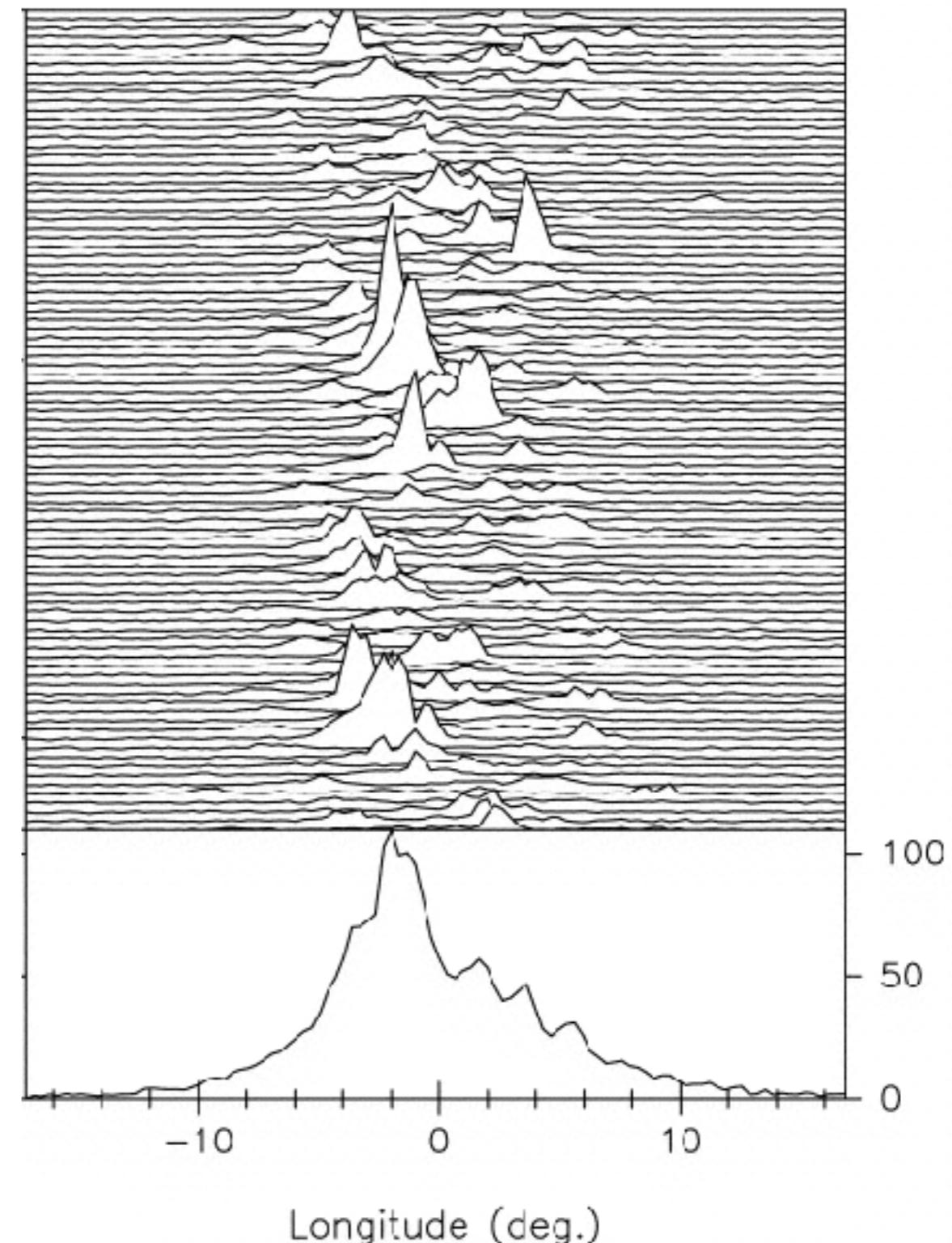
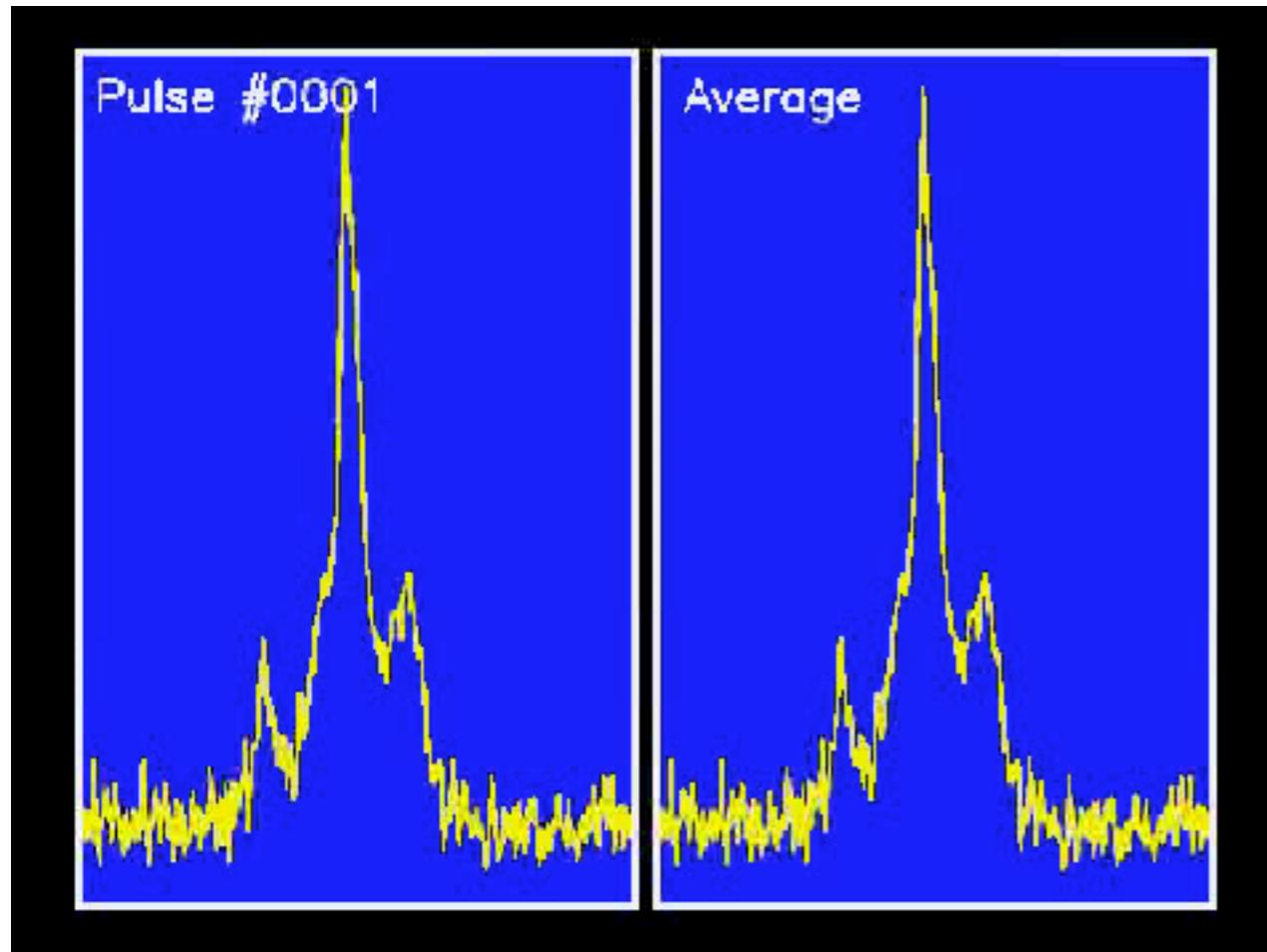
<http://www.jb.man.ac.uk/distance/frontiers/pulsars/section1.html>



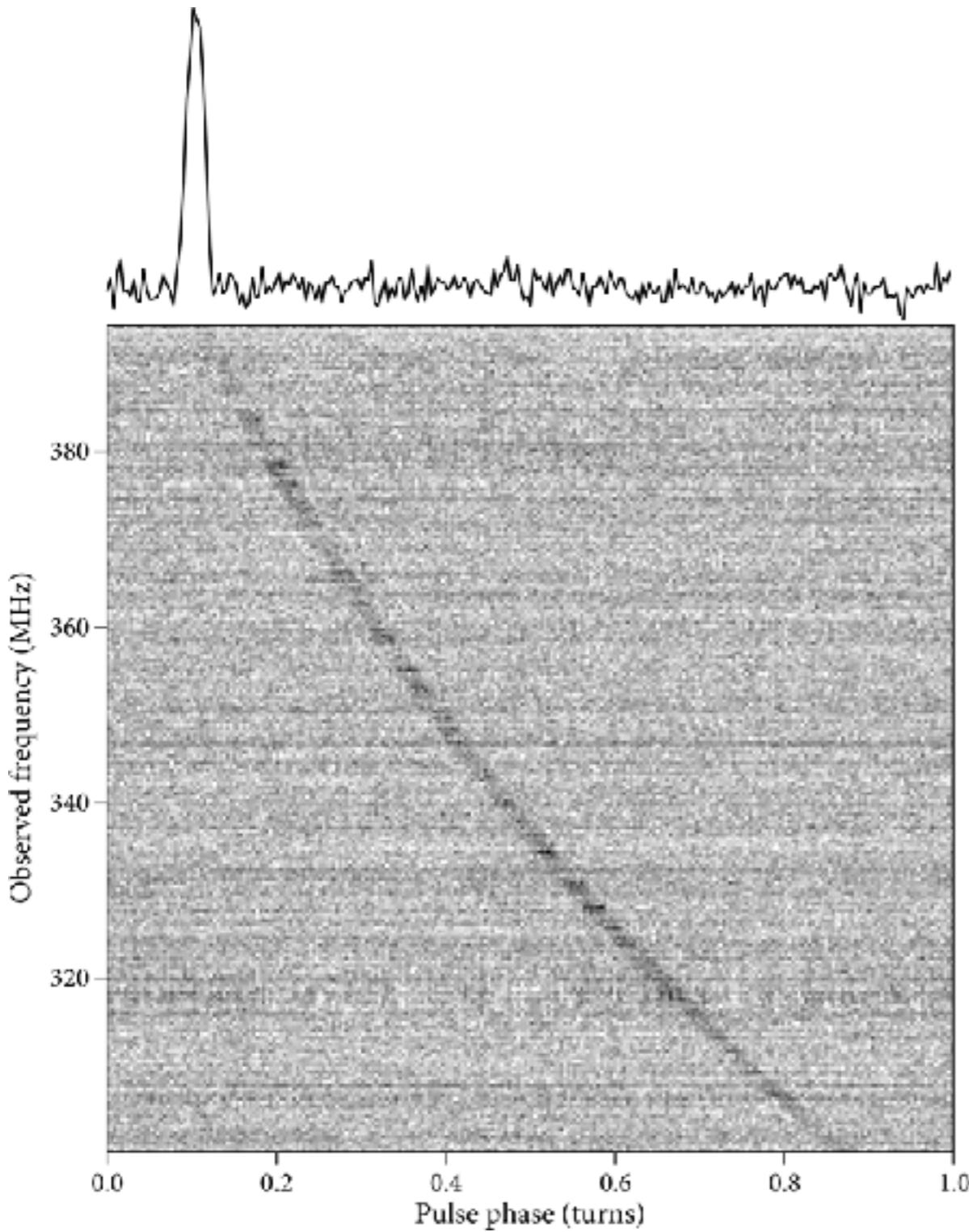
INTEGRATED PROFILE

PSR B0329+54: This pulsar is a typical, normal pulsar, rotating with a period of 0.714519 seconds, i.e. close to 1.40 rotations/sec.

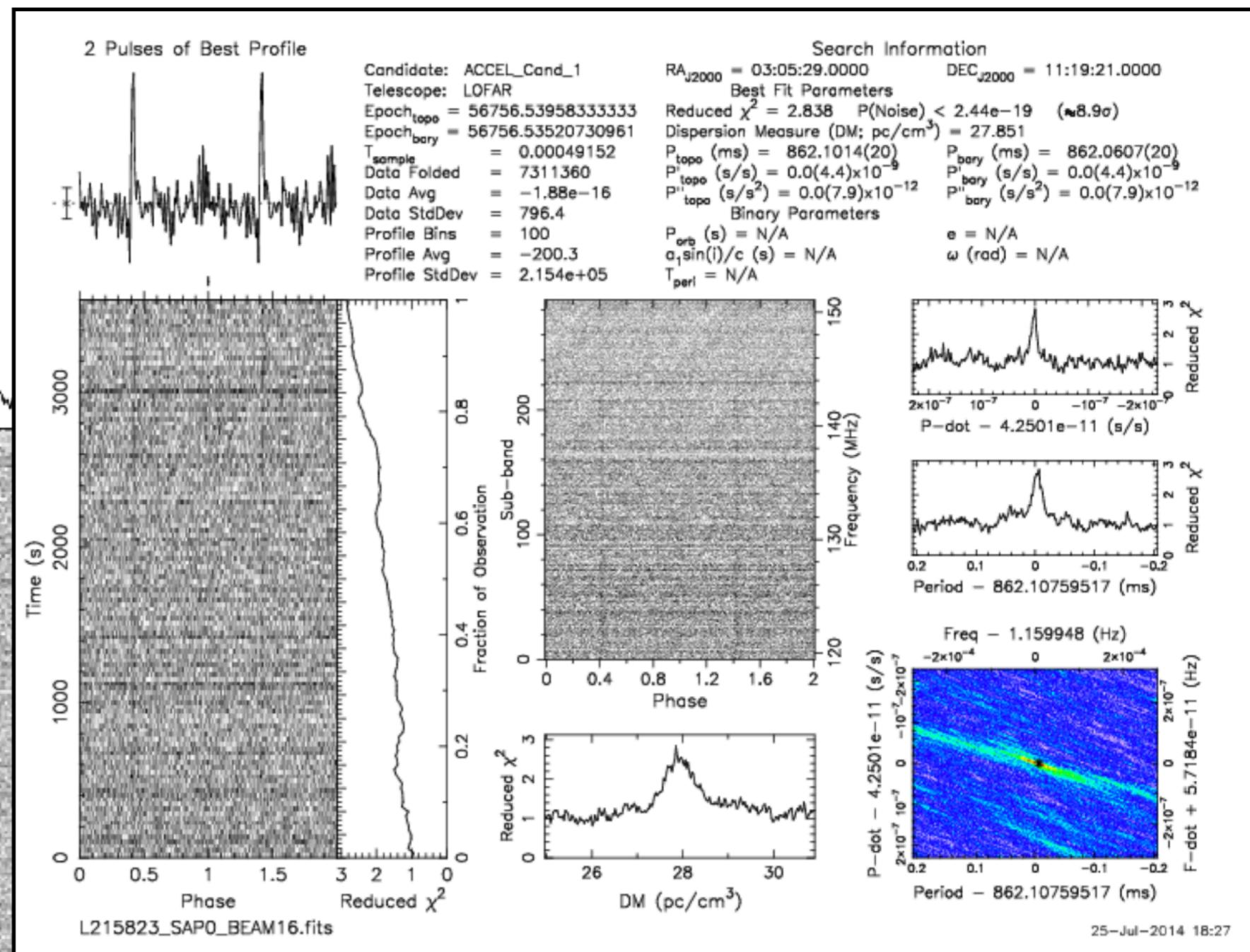
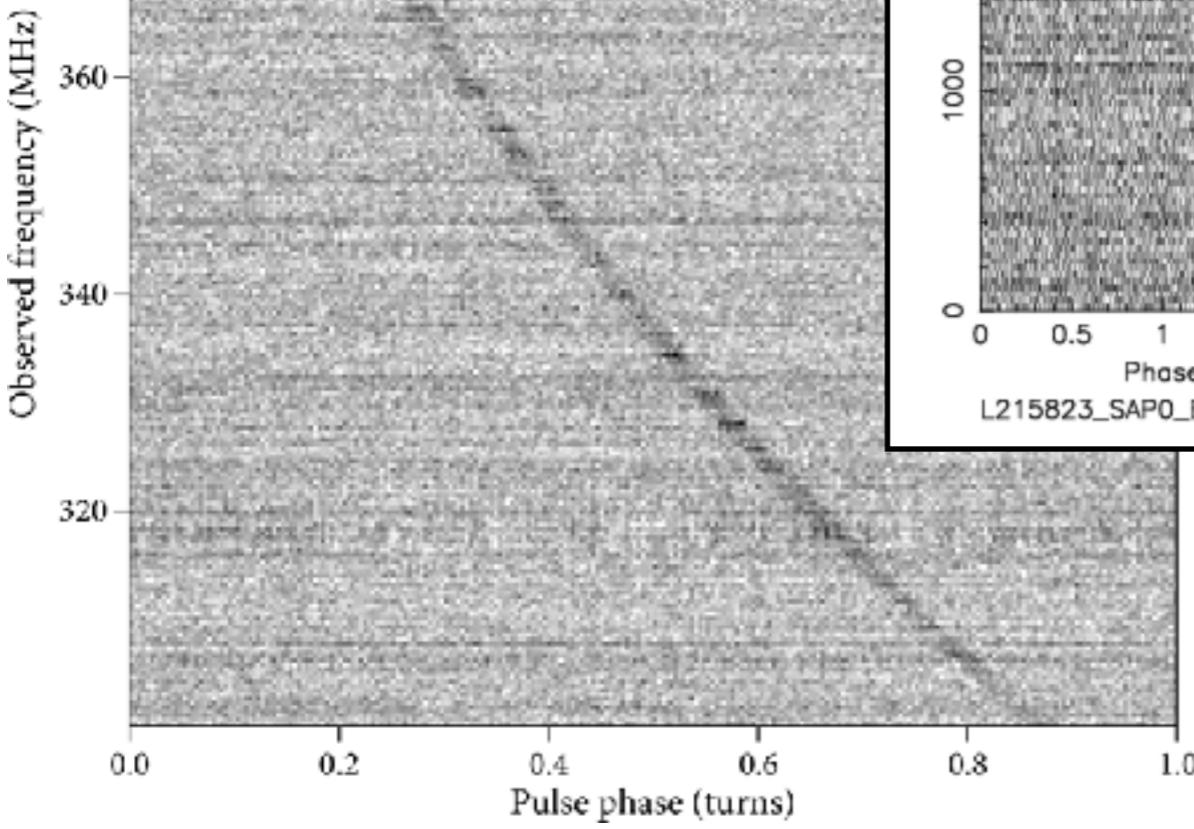
<http://www.jb.man.ac.uk/distance/frontiers/pulsars/section1.html>



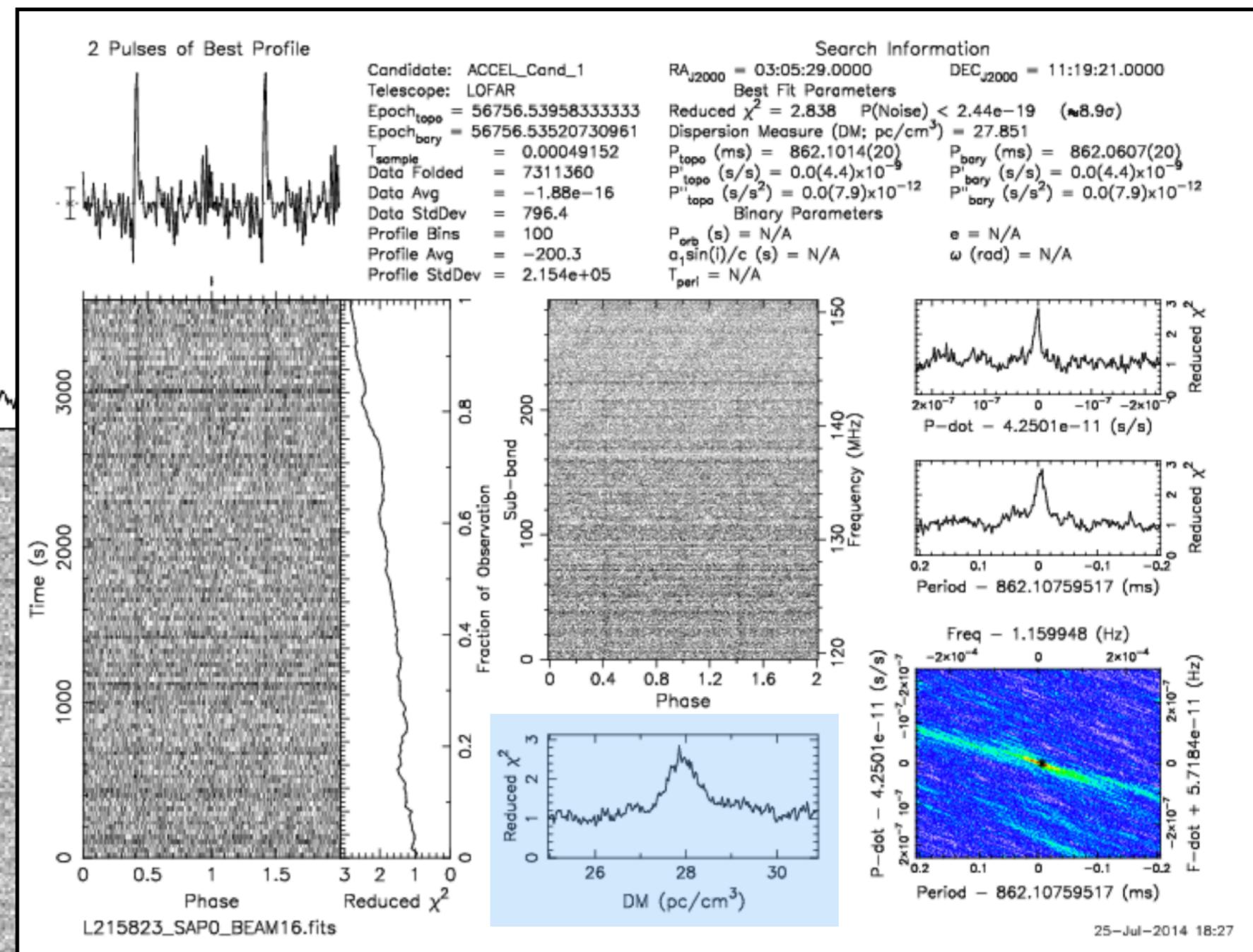
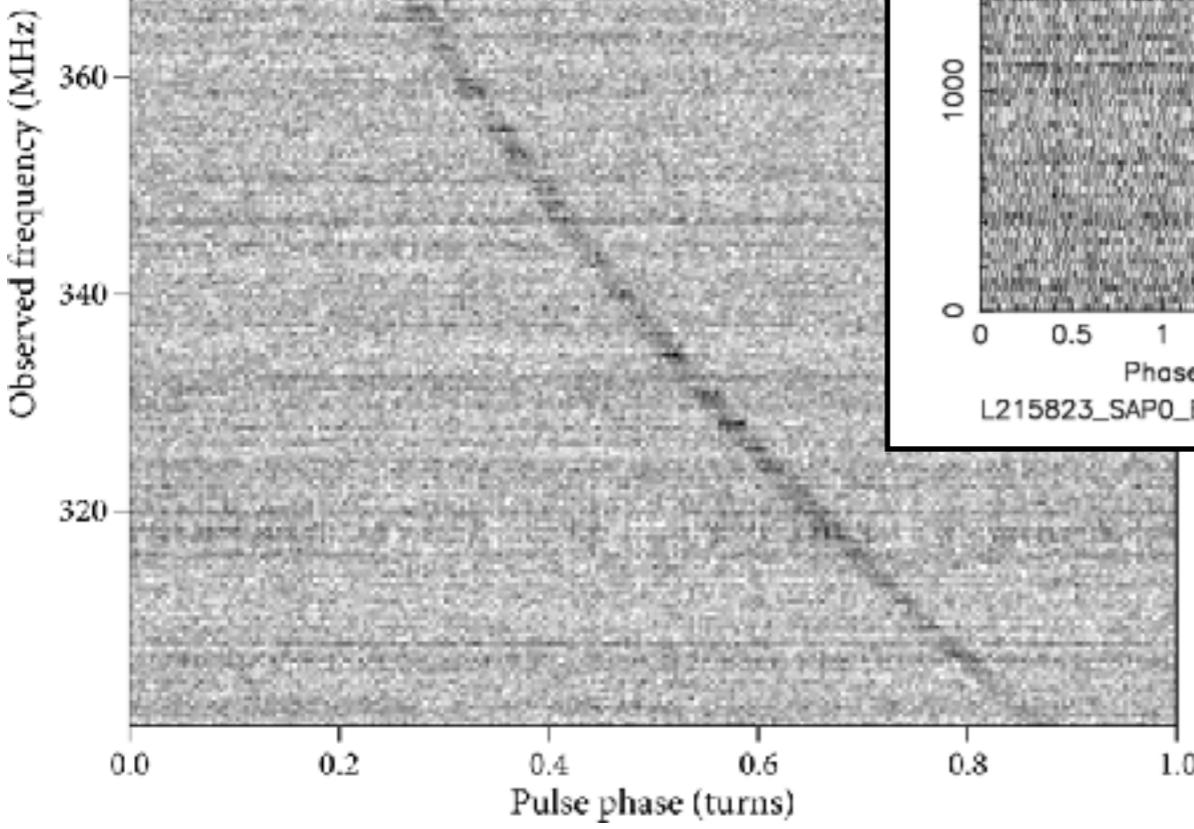
DM-SNR CURVE



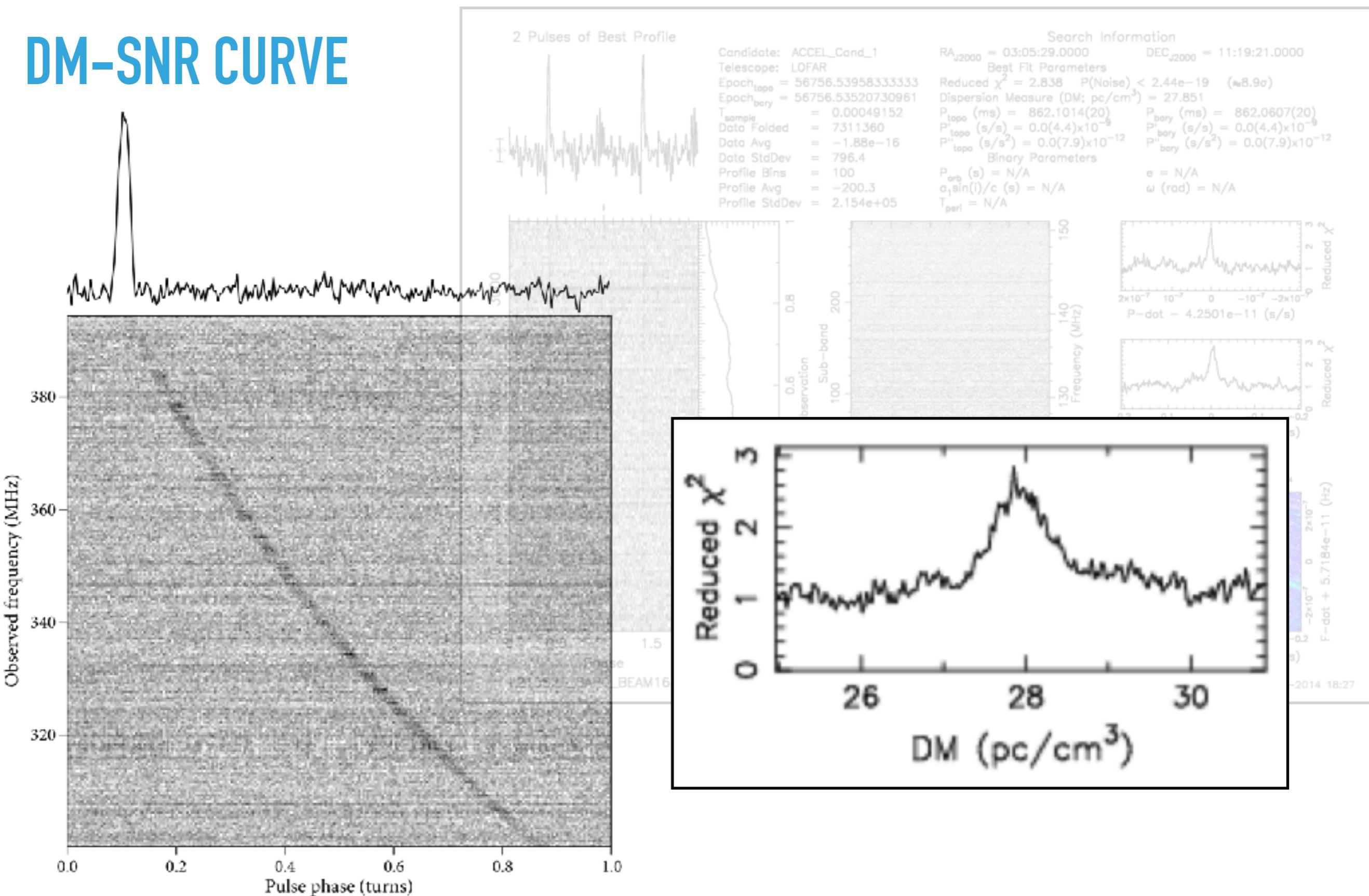
DM-SNR CURVE



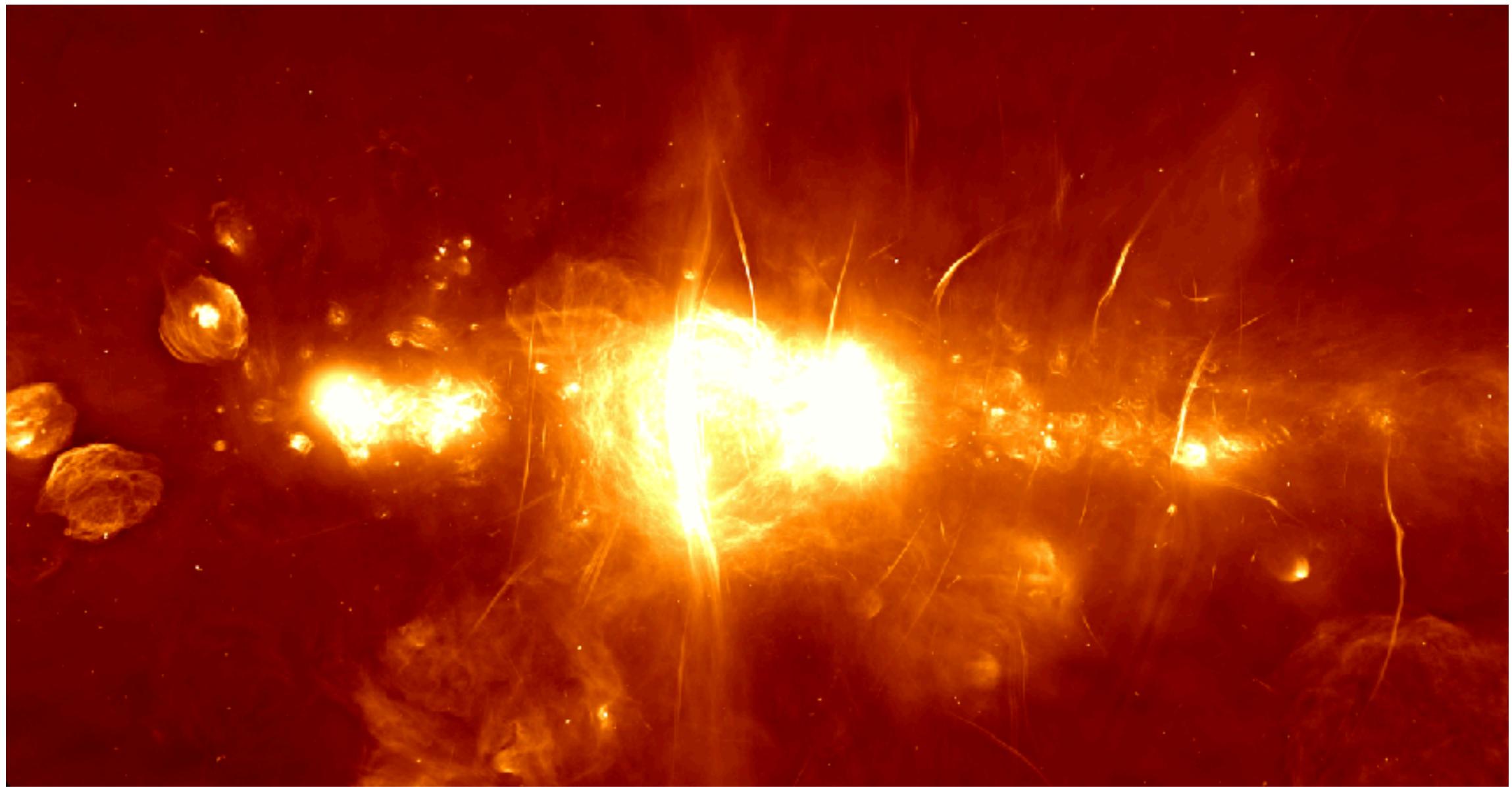
DM-SNR CURVE



DM-SNR CURVE



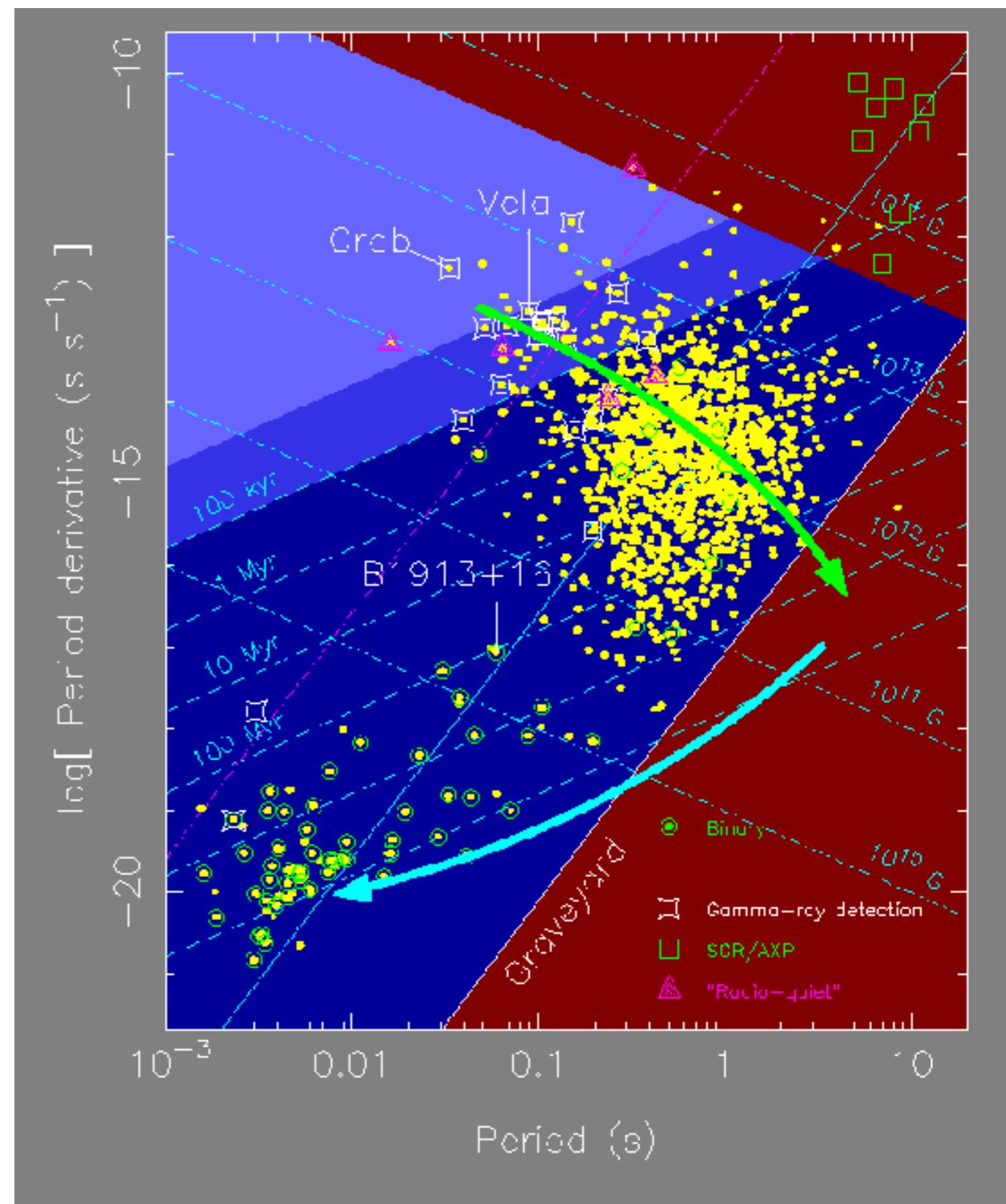
WHAT IS A PULSAR?

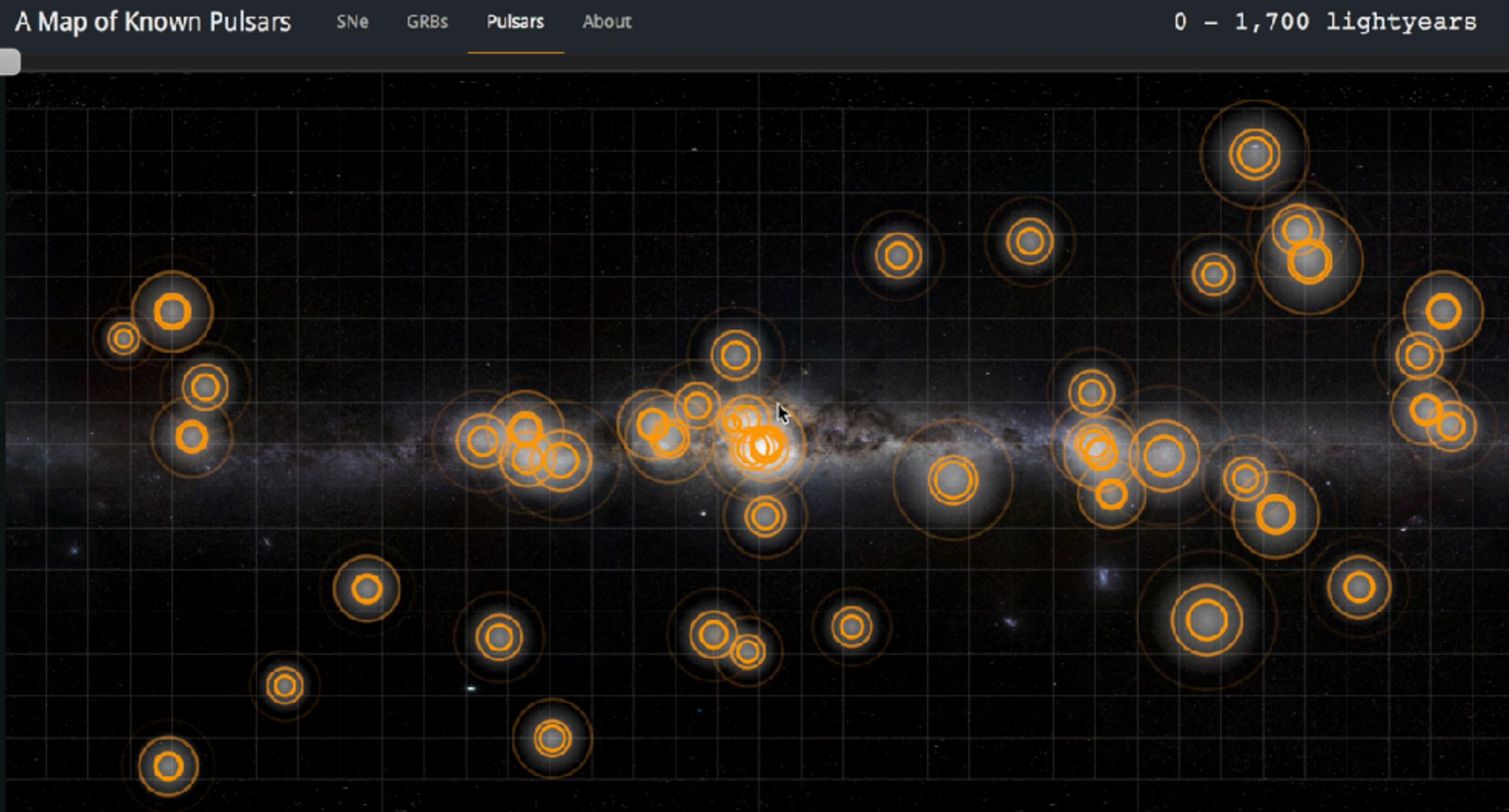


CRAB NEBULA, A.D. 1054

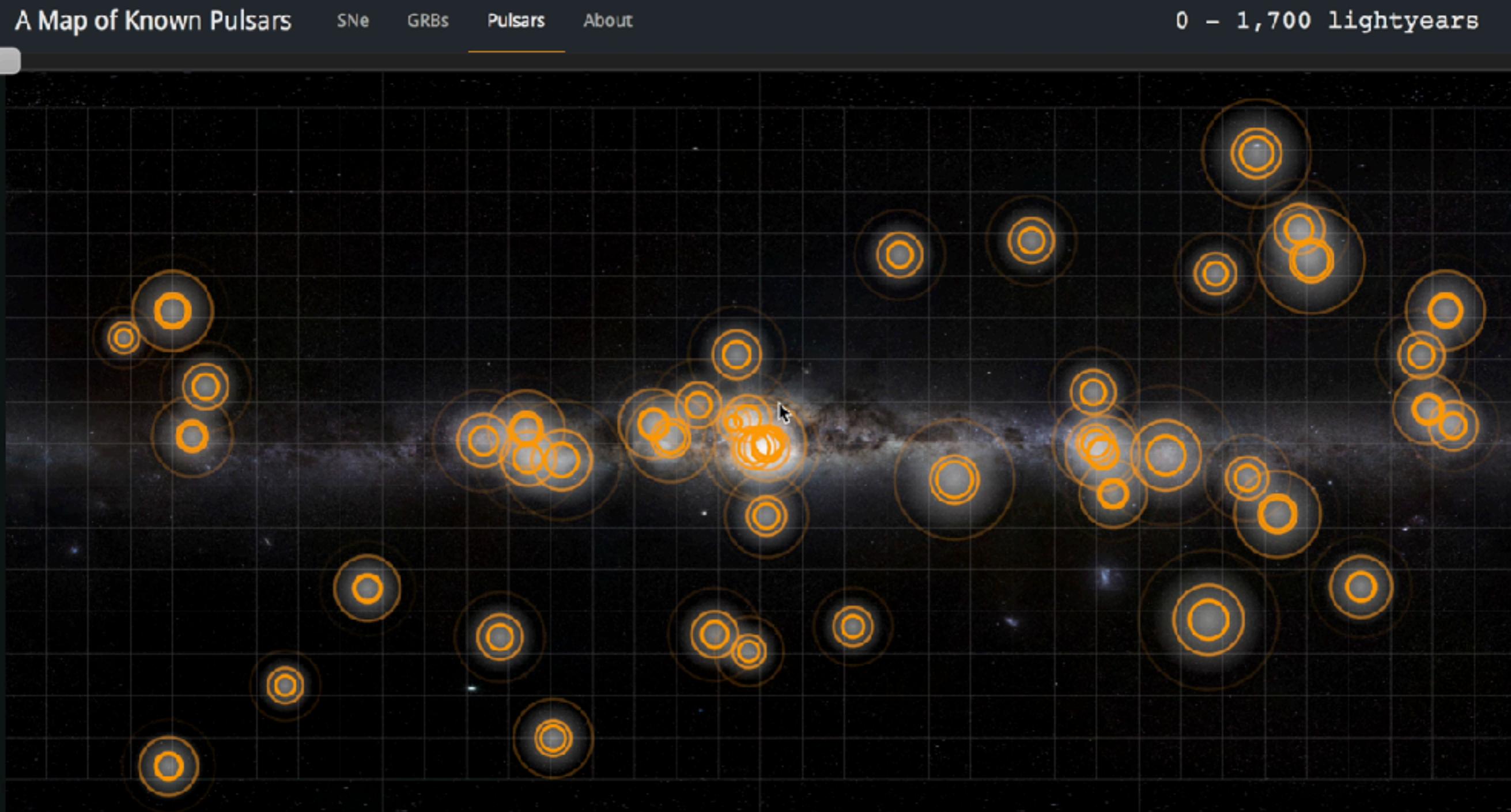


WHAT IS A PULSAR?



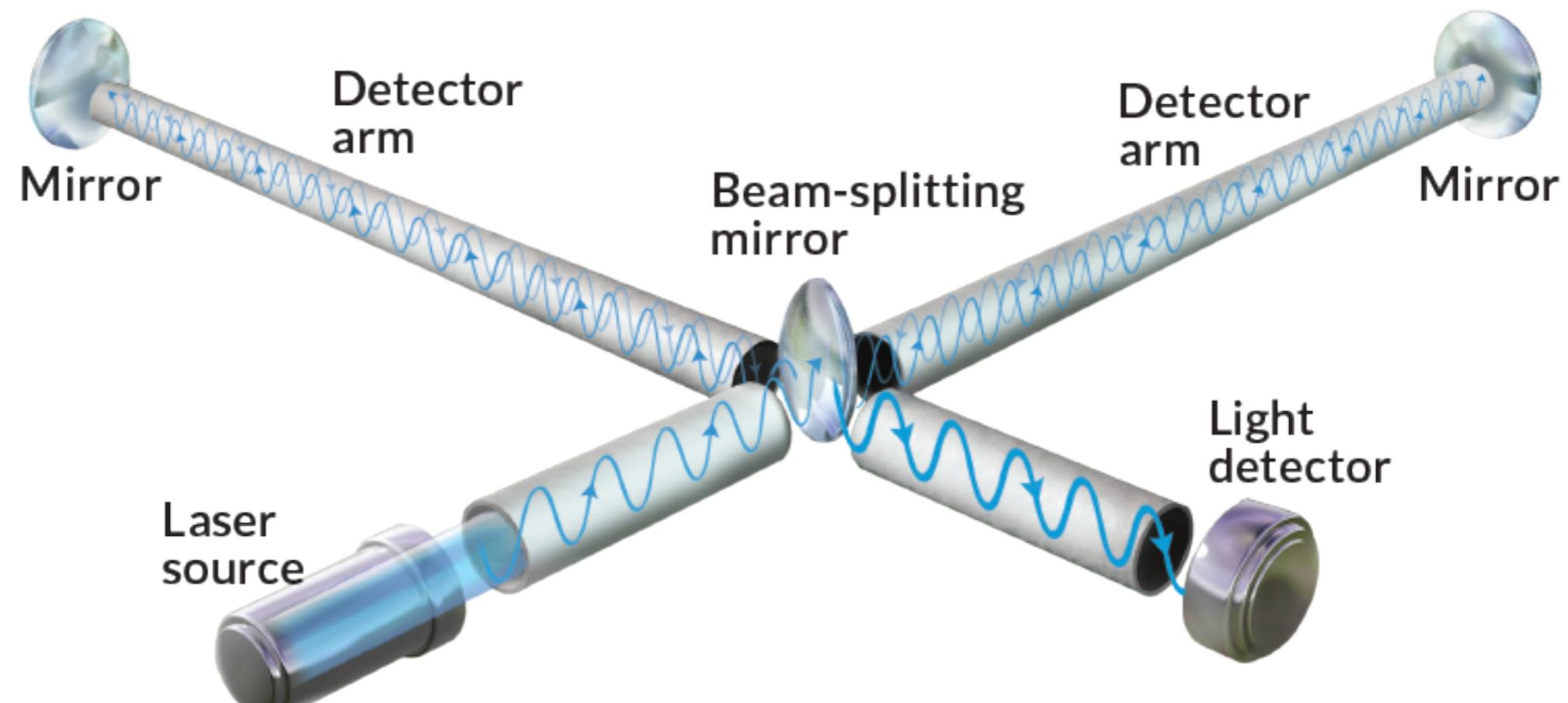


<http://www.ishivvers.com/maps/pulsars.html>



<http://www.ishivvers.com/maps/pulsars.html>

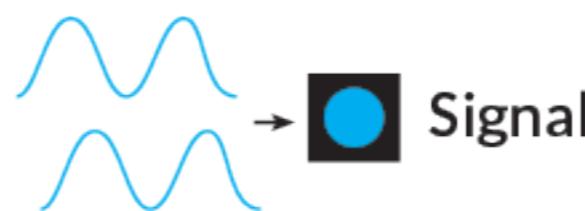
MEASURING GRAVITATIONAL WAVES



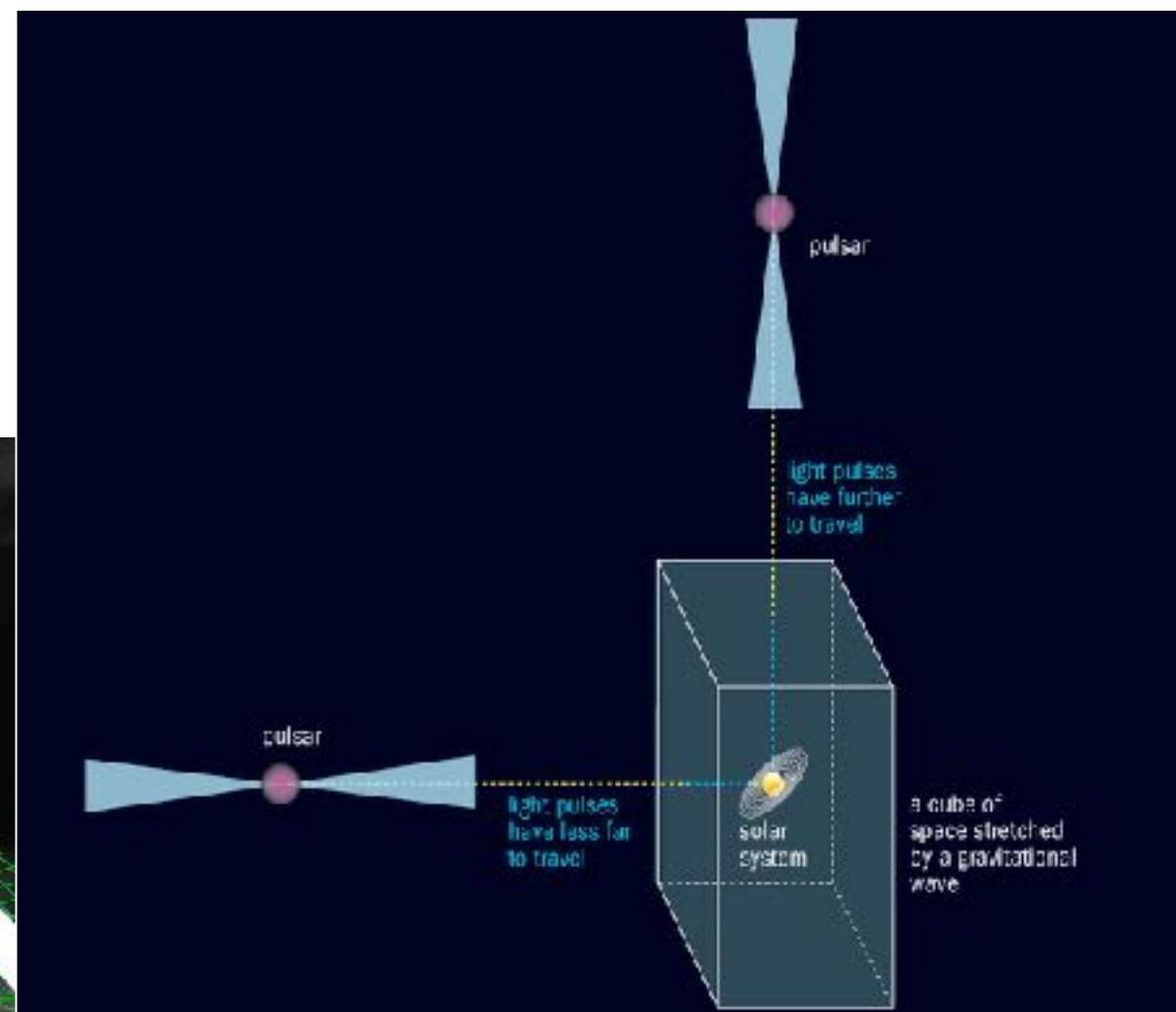
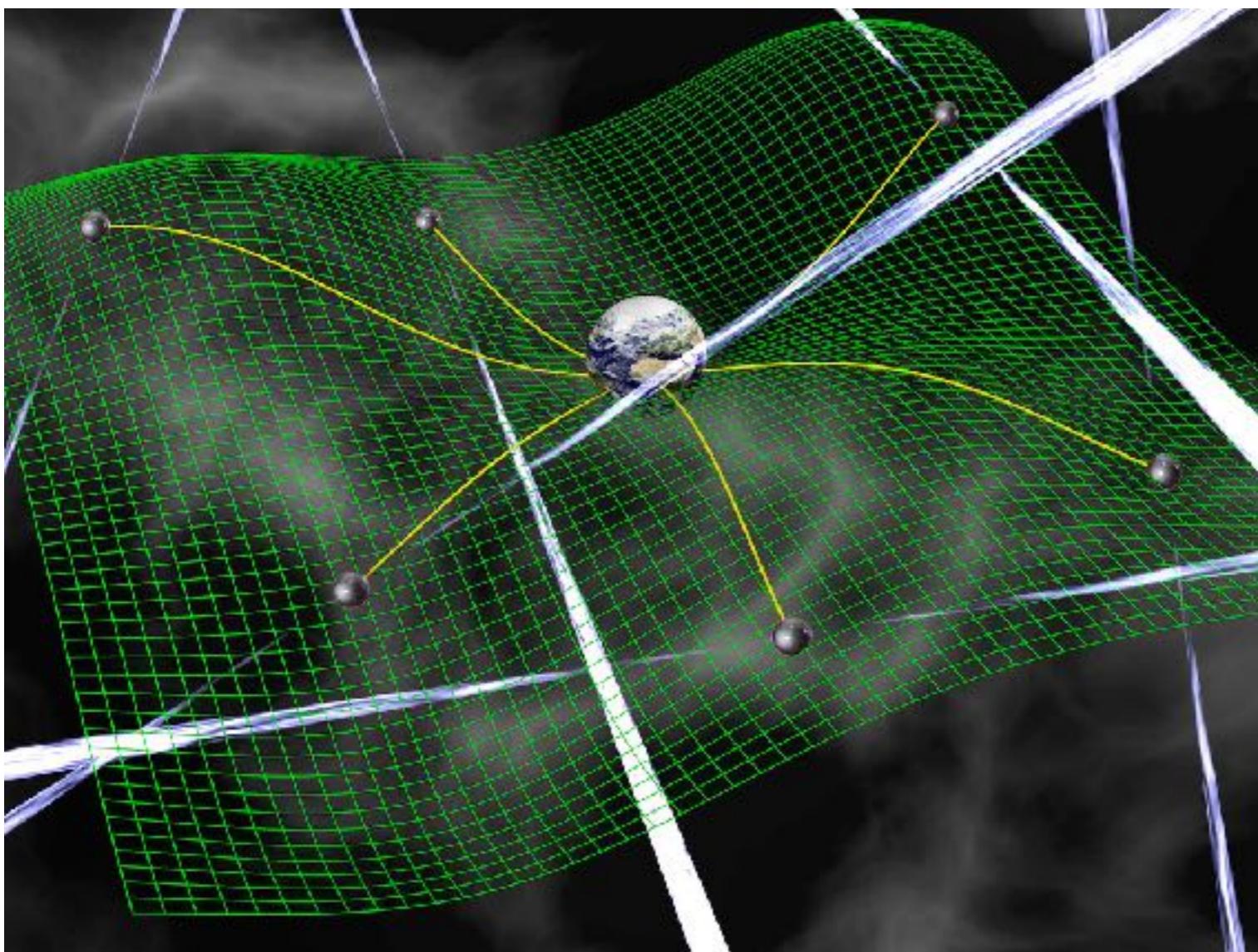
Normal situation



Gravitational wave detection

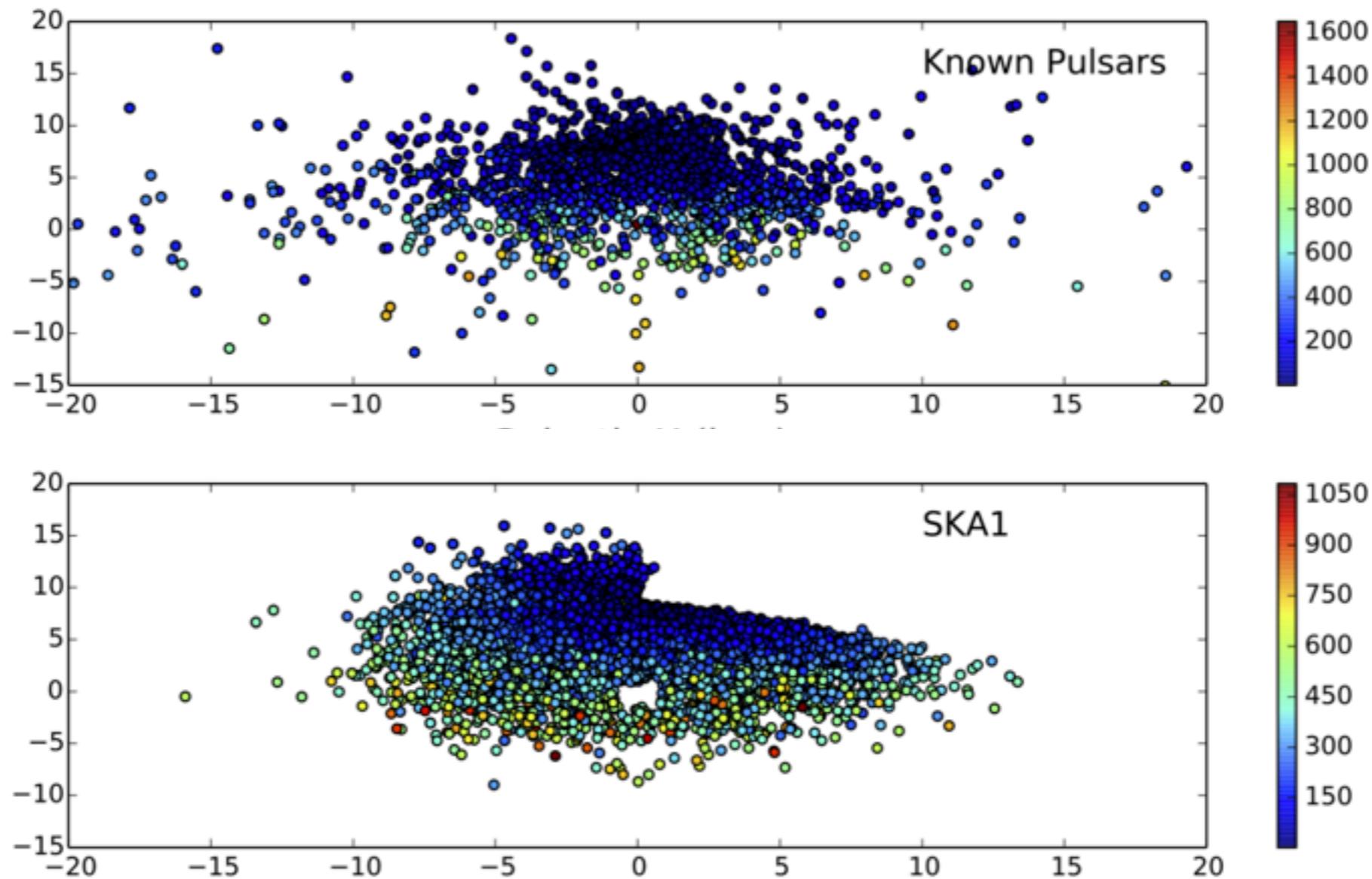


MEASURING GRAVITATIONAL WAVES



[http://live.iop-pp01.agh.sleek.net/2014/09/25/hunting- gravitational-waves-using-pulsars/](http://live.iop-pp01.agh.sleek.net/2014/09/25/hunting-gravitational-waves-using-pulsars/)

PULSARS WITH THE SKA



SKA1-MID is expected to find **10,000 new pulsars**

Survey	Year	Candidates	Per Sq. Degree
2nd Molonglo Survey(Manchester et al. 1978)	1977	2,500	~0.1
Phase II survey (Stokes et al. 1986)	1983	5,405	~1
Parkes 20 cm survey (Johnston et al. 1992)	1988	~ 150,000	~188
Parkes Southern Pulsar Survey (Manchester et al. 1996)	1991	40,000	~2
Parkes Multibeam Pulsar Survey (Manchester et al. 2001)	1997	8,000,000	~5,161
Swinburne Int. Lat. Survey (Edwards et al. 2001)	1998	> 200,000	~168*
Arecibo P-Alfa all configurations (Cordes et al. 2006; Lazarus 2012; P-Alfa Consortium 2015)	2004	> 5,000,000	~16,361*
6.5 GHz Multibeam Survey (Bates et al. 2011a; Bates 2011)	2006	3,500,000	~77,778 †
GBNCC survey (Stovall et al. 2014)	2009	> 1,200,000	~89*
Southern HTRU (Keith et al. 2010)	2010	55,434,300	~1,705
Northern HTRU (Barr et al. 2013; Ng 2012)	2010	> 80,000,000	~2,890*
LOTAAS (Cooper, private communication, 2015)	2013	39,000,000	~2,000

Table 1. Reported folded candidate numbers. Note * indicates a lower bound on the number of candidates per square degree, calculated from incomplete candidate numbers. † indicates very long integration times, with further details supplied in Tables [2](#) & [3](#).

- We need to deal with the **rising candidate numbers**.
- **Data rate is also rising**... we can no longer store all the data and repeatedly re-process it - we need to be right the first time around (or as often as possible).

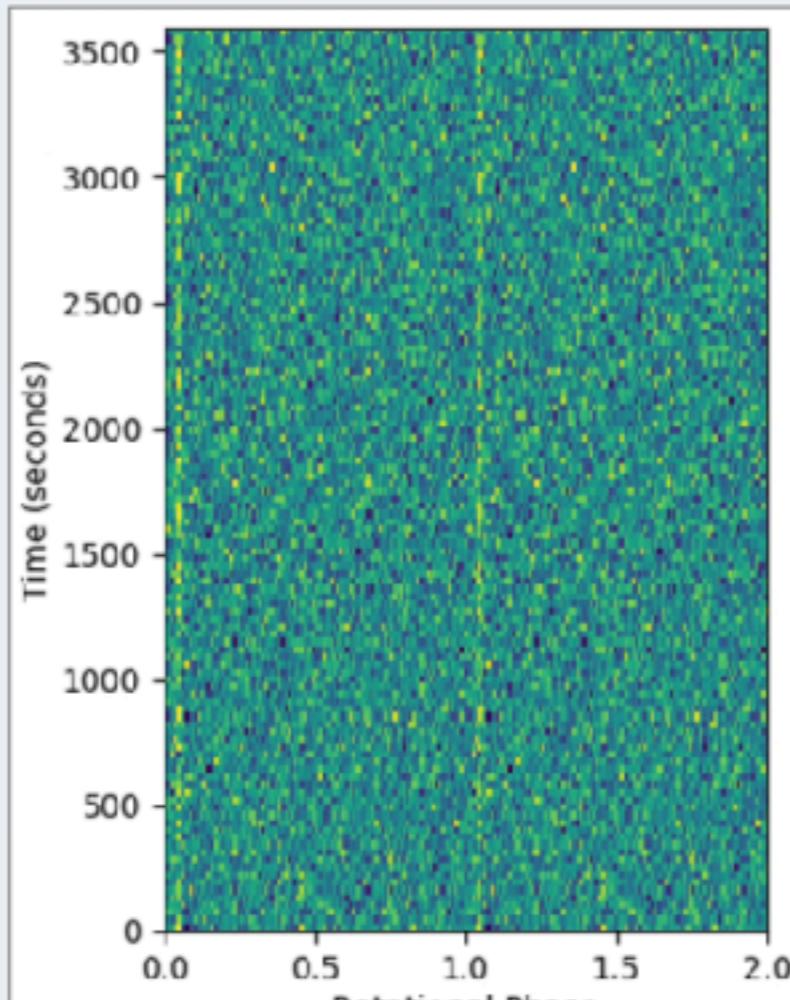
Survey	Year	Candidates	Per Sq. Degree
2nd Molonglo Survey(Manchester et al. 1978)	1977	2,500	~0.1
Phase II survey (Stokes et al. 1986)	1983	5,405	~1
Parkes 20 cm survey (Johnston et al. 1992)	1988	~ 150,000	~188
Parkes Southern Pulsar Survey (Manchester et al. 1996)	1991	40,000	~2
Parkes Multibeam Pulsar Survey (Manchester et al. 2001)	1997	8,000,000	~5,161
Swinburne Int. Lat. Survey (Edwards et al. 2001)	1998	> 200,000	~168*
Arecibo P-Alfa all configurations (Cordes et al. 2006; Lazarus 2012; P-Alfa Consortium 2015)	2004	> 5,000,000	~16,361*
6.5 GHz Multibeam Survey (Bates et al. 2011a; Bates 2011)	2006	3,500,000	~77,778 †
GBNCC survey (Stovall et al. 2014)	2009	> 1,200,000	~89*
Southern HTRU (Keith et al. 2010)	2010	55,434,300	~1,705
Northern HTRU (Barr et al. 2013; Ng 2012)	2010	> 80,000,000	~2,890*
LOTAAS (Cooper, private communication, 2015)	2013	39,000,000	~2,000

Table 1. Reported folded candidate numbers. Note * indicates a lower bound on the number of candidates per square degree, calculated from incomplete candidate numbers. † indicates very long integration times, with further details supplied in Tables [2](#) & [3](#).

- There are far more non-pulsars than true pulsars... this is a machine learning **class imbalance problem**.
- Typically there are ~10,000 non-pulsars for every true pulsar.



Pulsar Hunters

[ABOUT](#)[**CLASSIFY**](#)[TALK](#)[COLLECT](#)

You should sign in!

[TASK](#)[**TUTORIAL**](#)

Which characteristic best describes the vertical feature?

Continuous

Patchy

Halving/interrupted

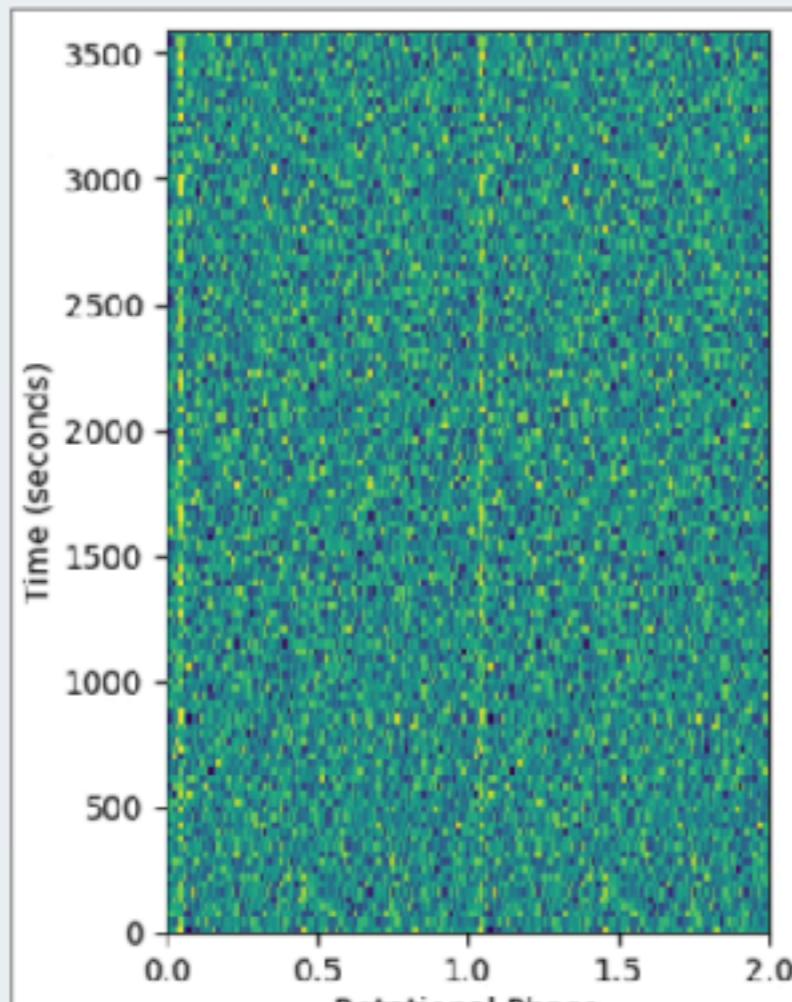
Curved

No feature

Done



Pulsar Hunters

[ABOUT](#)[**CLASSIFY**](#)[TALK](#)[COLLECT](#)

You should sign in!

[TASK](#)[**TUTORIAL**](#)

Which characteristic best describes the vertical feature?

Continuous

Patchy

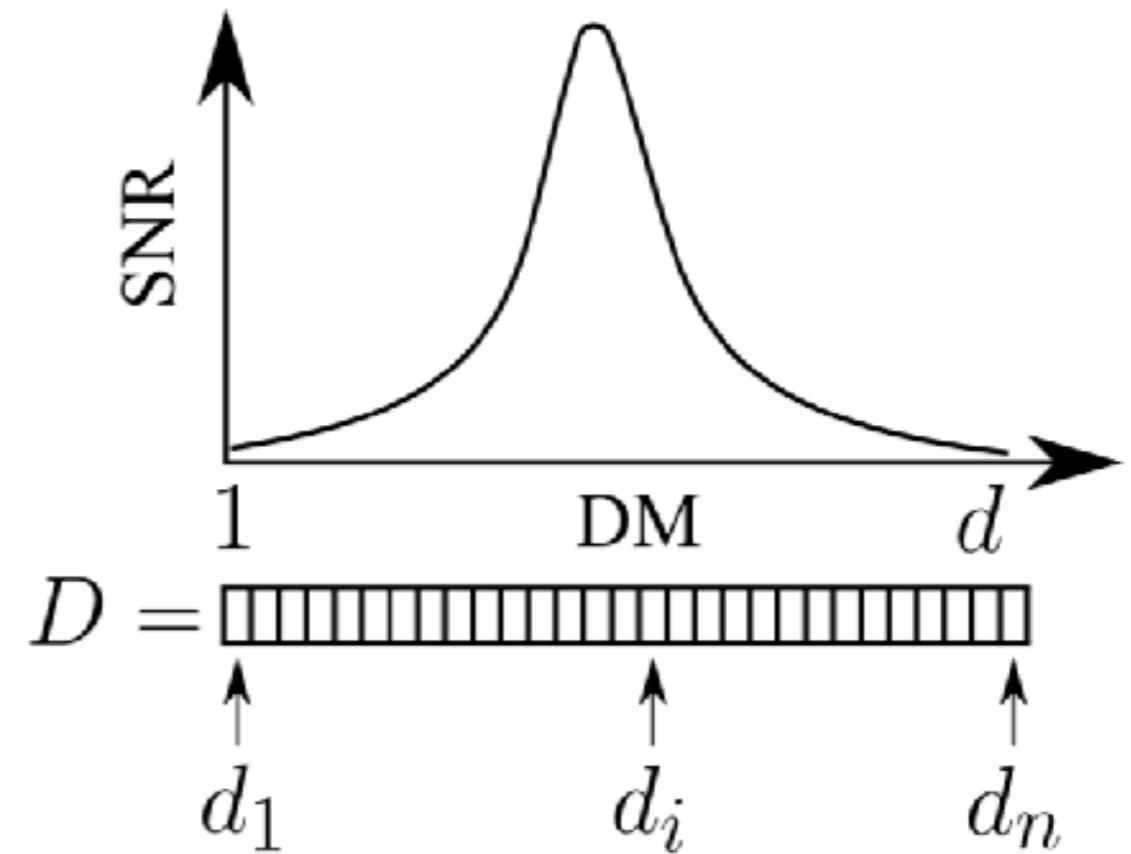
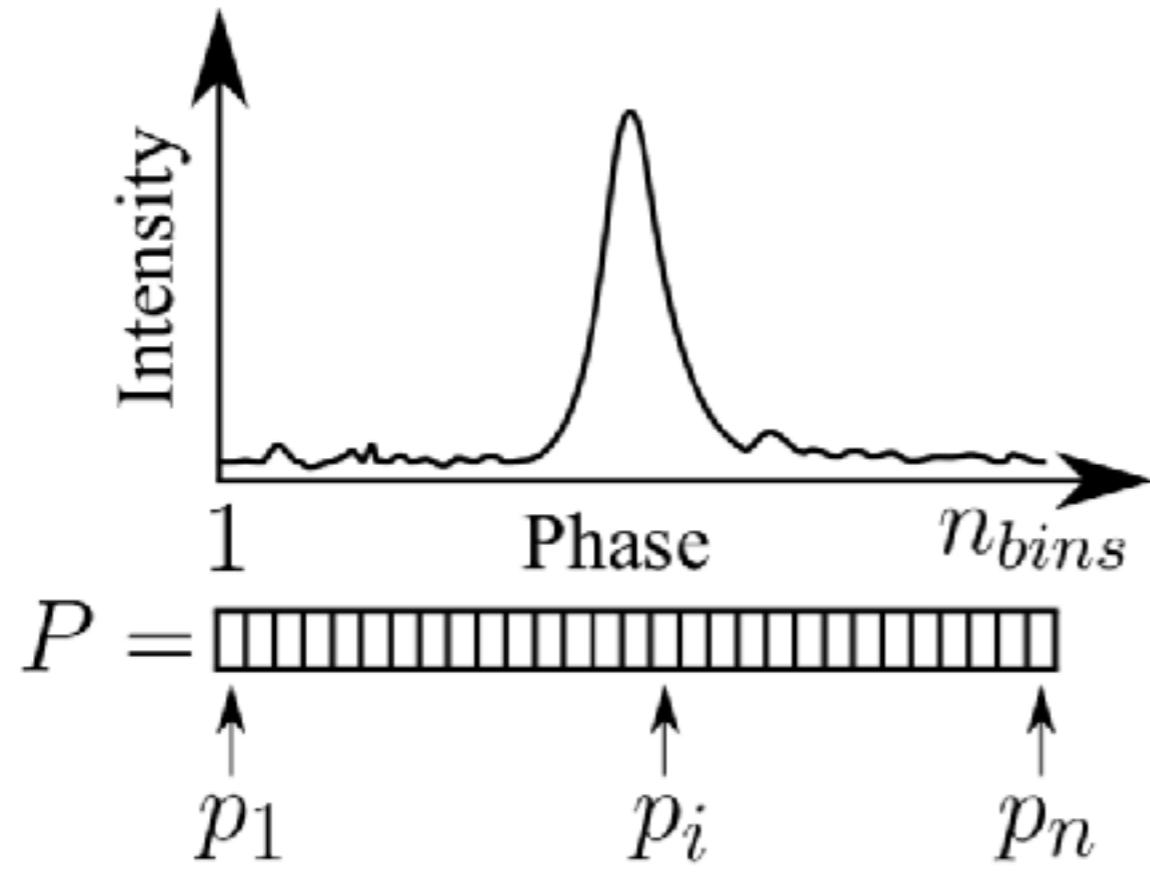
Halving/interrupted

Curved

No feature

Done

For each pulsar candidate we have **8** pieces of data - these are the **features**.

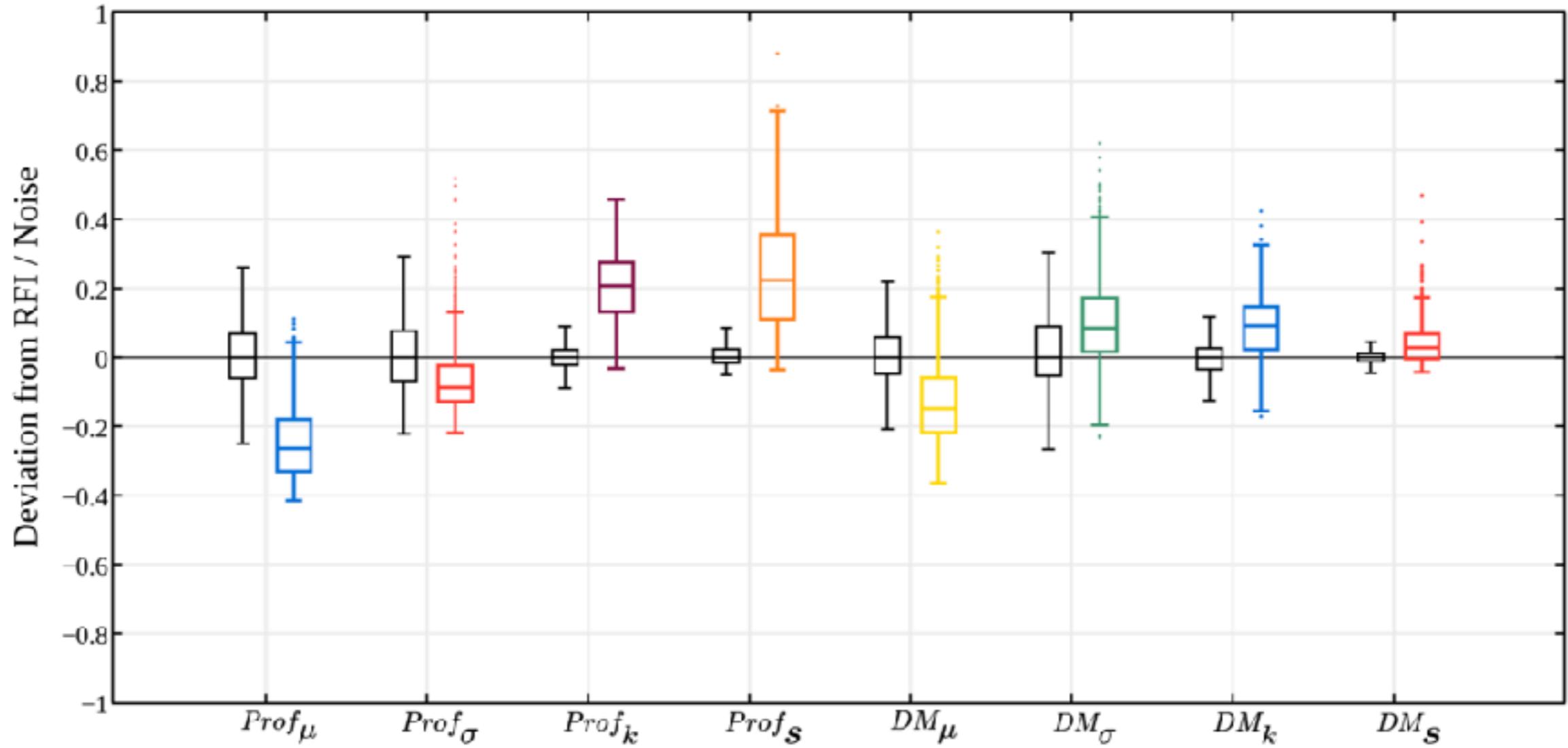


Integrated Profile

- Mean
- Standard deviation
- Kurtosis
- Skewness

DM-SNR Curve

- Mean
- Standard deviation
- Kurtosis
- Skewness

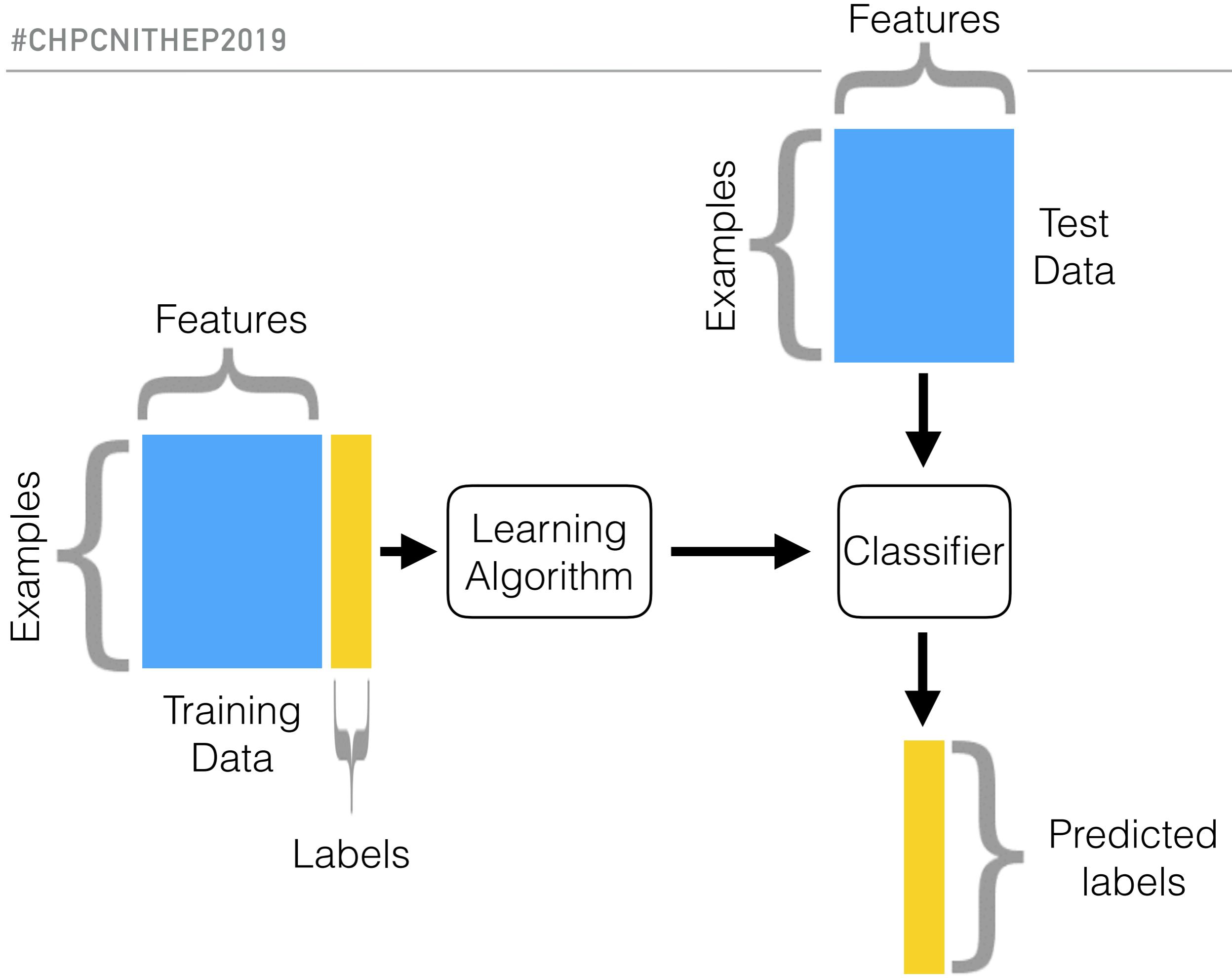


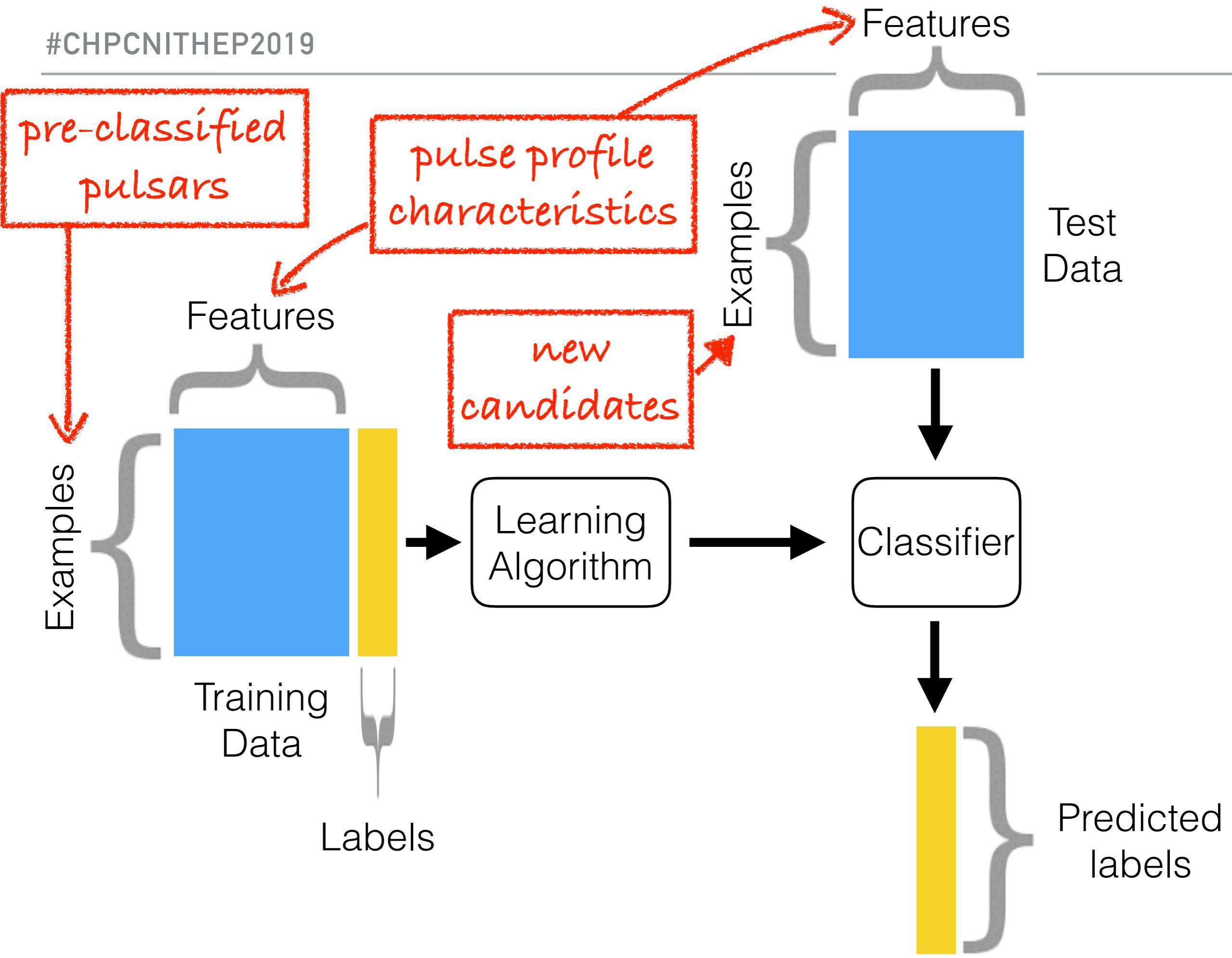
Integrated Profile

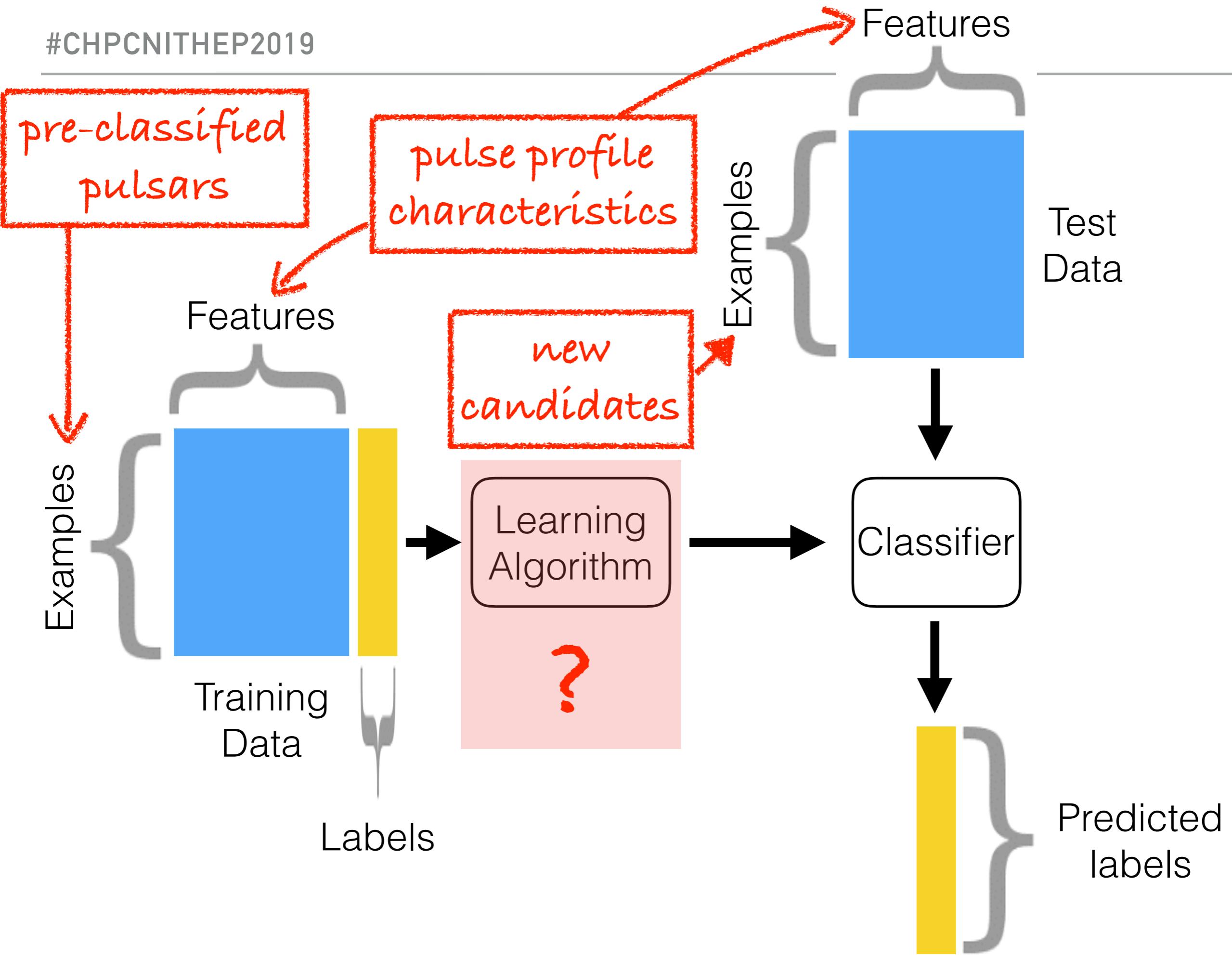
- Mean
- Standard deviation
- Kurtosis
- Skewness

DM-SNR Curve

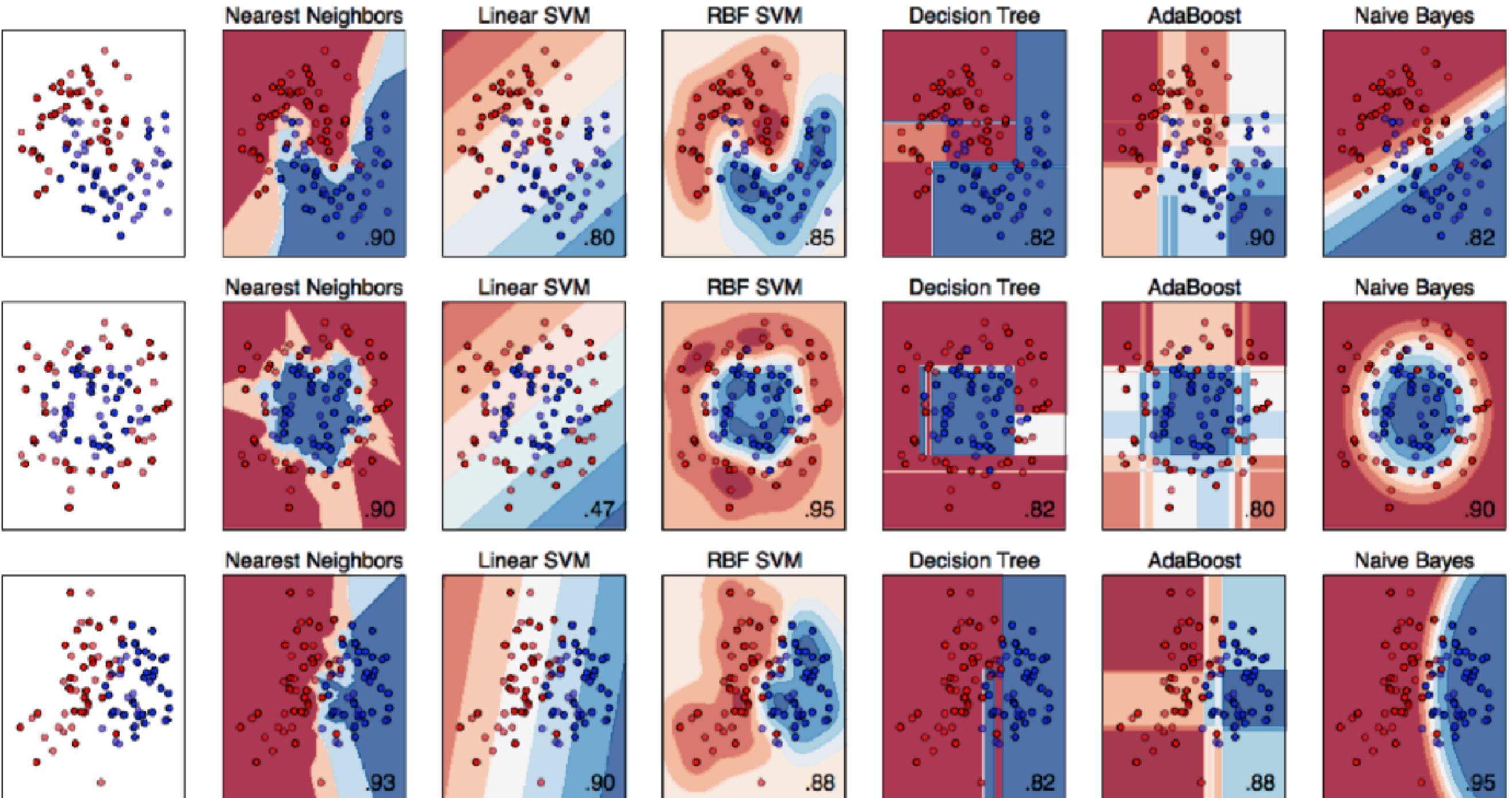
- Mean
- Standard deviation
- Kurtosis
- Skewness





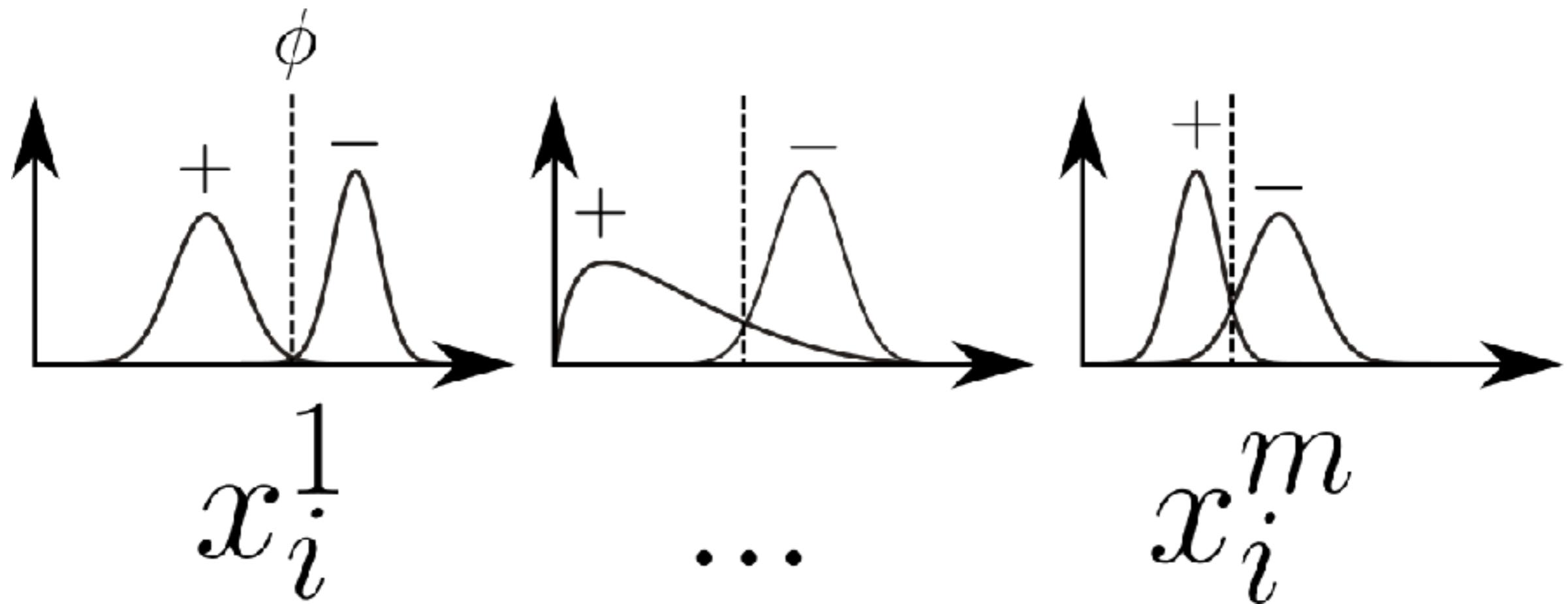


WHAT KIND OF MACHINE LEARNING?

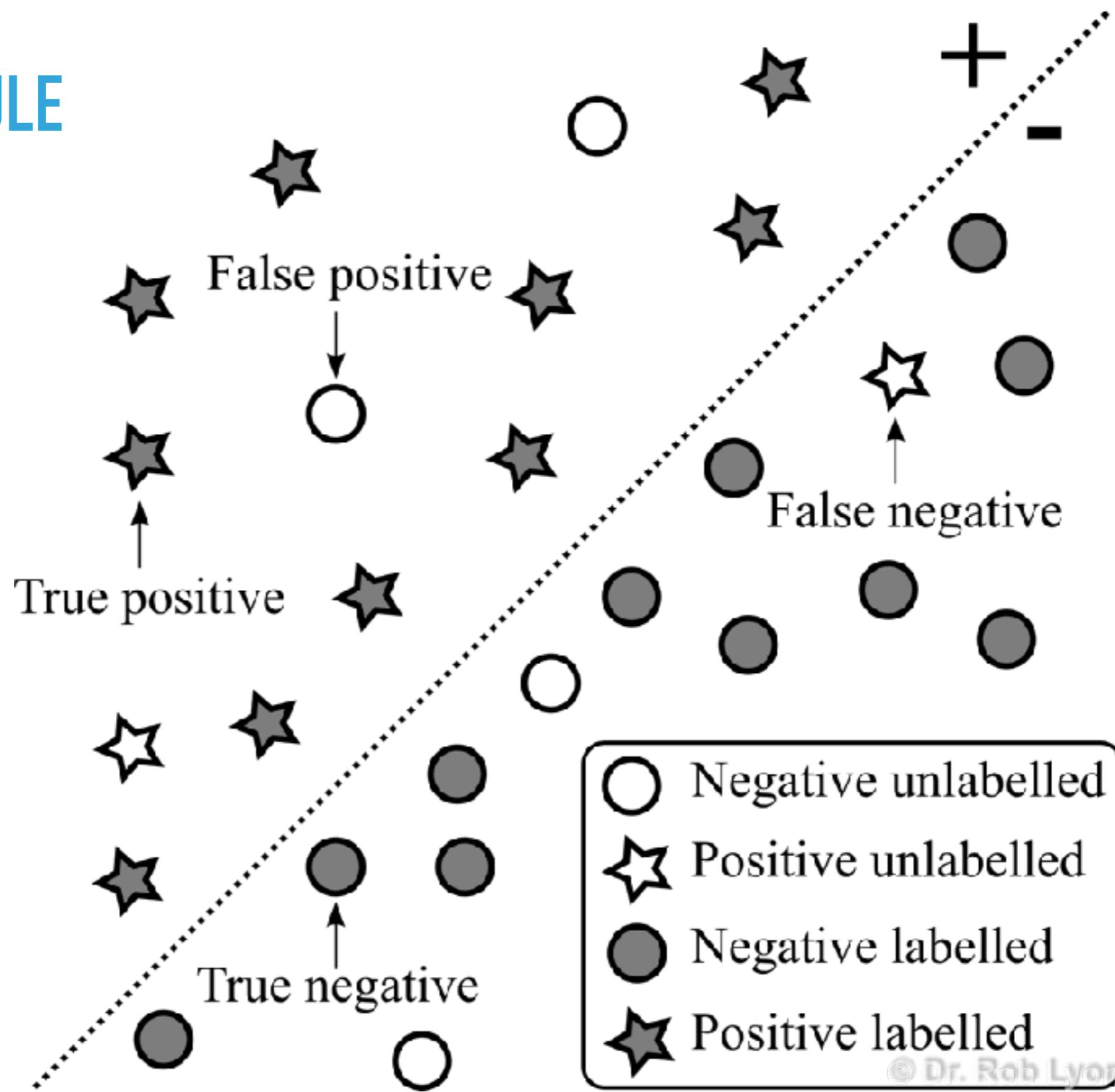


DECISION RULE

if $s(x) > \phi$ THEN predict $y = +ve$
if $s(x) < \phi$ THEN predict $y = -ve$



DECISION RULE



		Truth
Prediction	Positive	Negative
Positive	True Positive (TP)	False Positive (FP) Type I Error
Negative	False Negative (FN) Type II Error	True Negative (TN)

$$Precision = \frac{TP}{TP + FP}$$

% of **+ve** predictions that are truly **+ve**

$$Recall = \frac{TP}{TP + FN}$$

% of true **+ves** that are predicted **+ve**

$$Specificity = \frac{TN}{TN + FP}$$

% of true **-ves** that are predicted **-ve**

EVALUATING CLASSIFIER PERFORMANCE

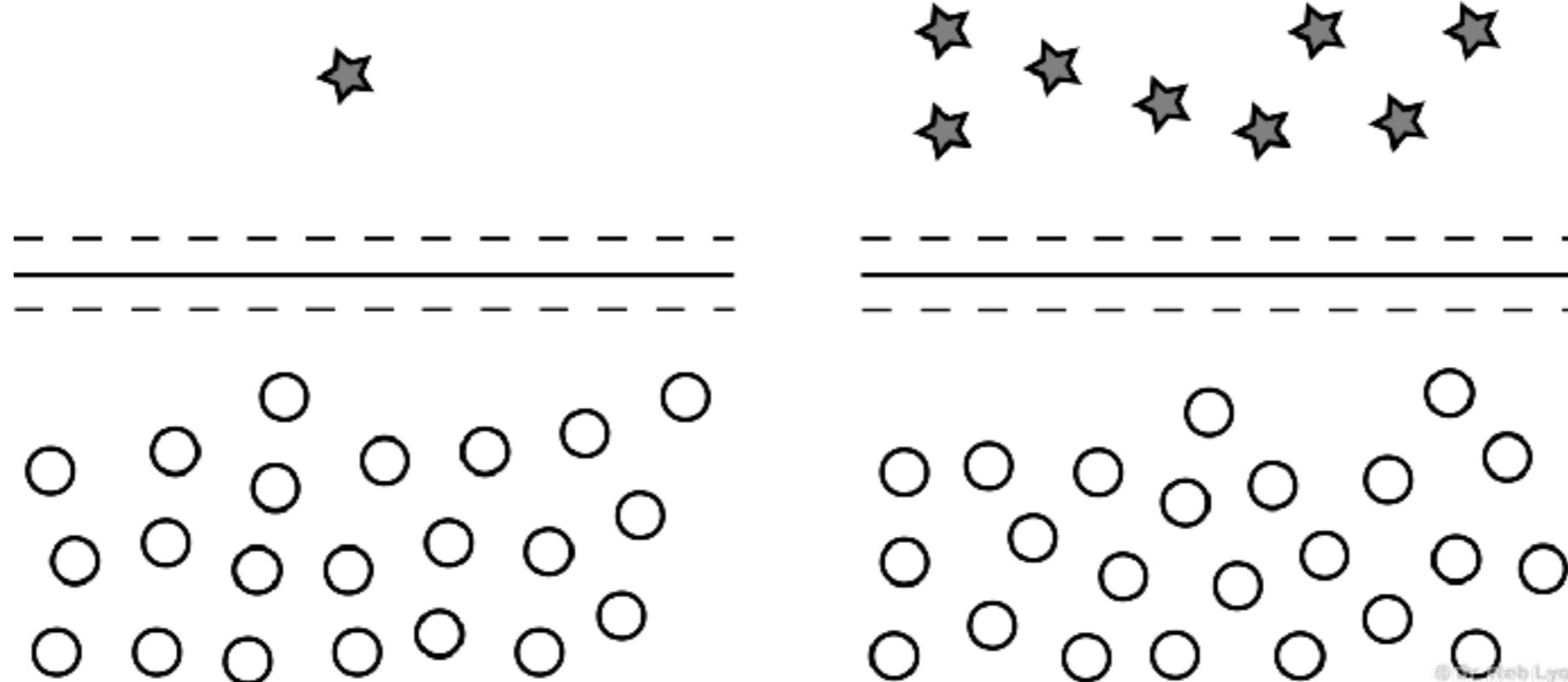
Statistic	Description	Definition
Accuracy	Measure of overall classification accuracy.	$\frac{(TP+TN)}{(TP+FP+FN+TN)}$
False positive rate (FPR)	Fraction of negative instances incorrectly labelled positive.	$\frac{FP}{(FP+TN)}$
G-Mean	Imbalanced data metric describing the ratio between positive and negative accuracy.	$\sqrt{\frac{TP}{TP+FN} \times \frac{TN}{TN+FP}}$
Precision	Fraction of retrieved instances that are positive.	$\frac{TP}{(TP+FP)}$
Recall	Fraction of positive instances that are retrieved.	$\frac{TP}{(TP+FN)}$
F-Score	Measure of accuracy that considers both precision and recall.	$2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$
Specificity	Fraction of negatives correctly identified as such.	$\frac{TN}{(FP+TN)}$

Lyon et al. 2016

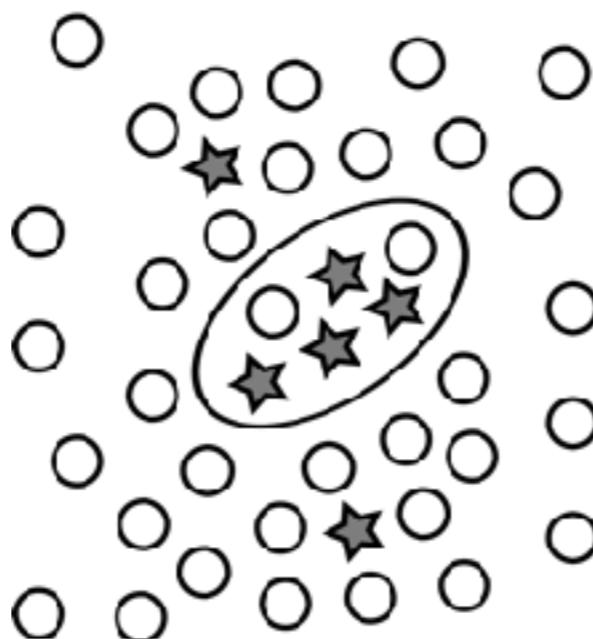
For imbalanced classification problems, we can also treat **rarity** as **importance** and try to minimise the ***expected cost*** of mis-classification:

$$R = (1/p_{-ve}) \times FP + (1/p_{+ve}) \times FN$$

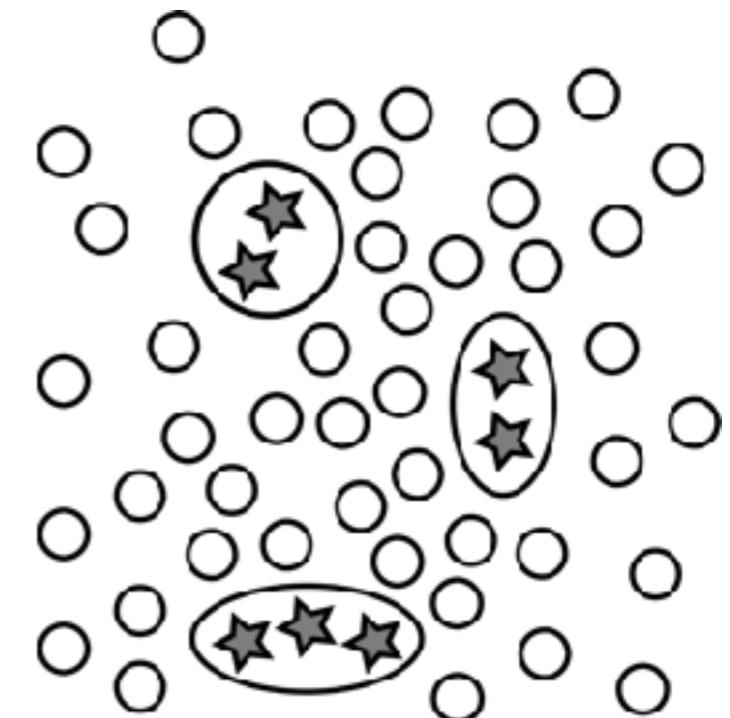
IMBALANCE IS NOT ALWAYS A PROBLEM



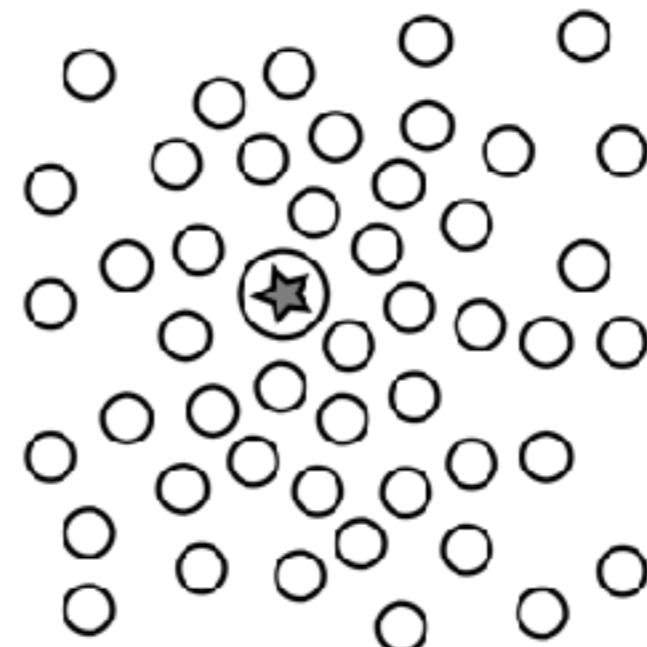
...BUT IT OFTEN IS.



a)

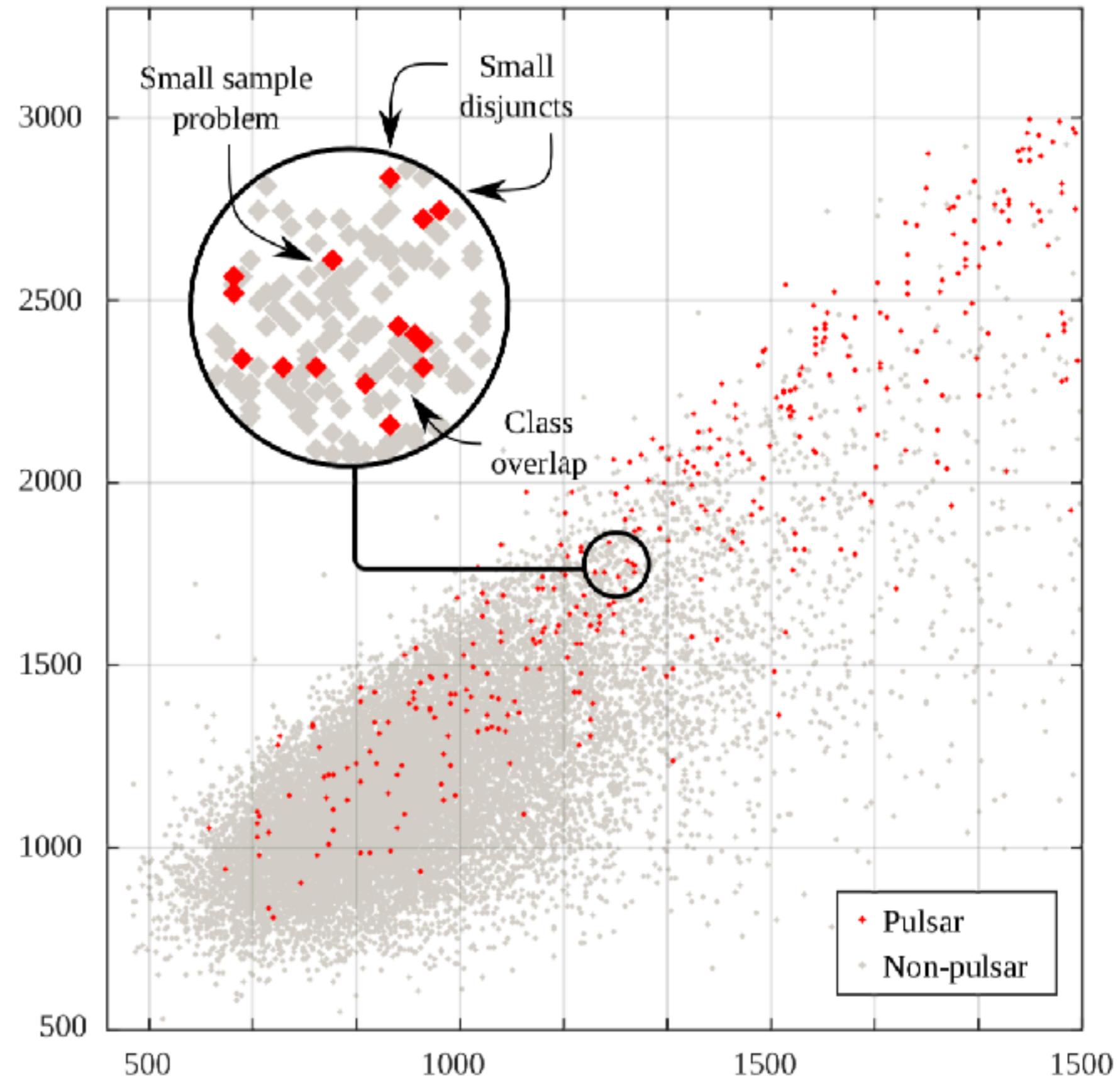


b)



c)

...BUT IT OFTEN IS.



POSSIBLE APPROACHES

- Do nothing.
- Force balance in the dataset.
- Modify algorithm.
- Write a new algorithm.
- Treat it as an anomaly detection problem.
- Get more data.

POSSIBLE APPROACHES

- Oversample minority class
- Undersample majority class
- Introduce synthetic examples

- Do nothing.
- Force balance in the dataset.
- Modify algorithm.
- Write a new algorithm.
- Treat it as an anomaly detection problem.
- Get more data.

POSSIBLE APPROACHES

- Weight input data
- Modify the loss function
- Calibrate probability estimates

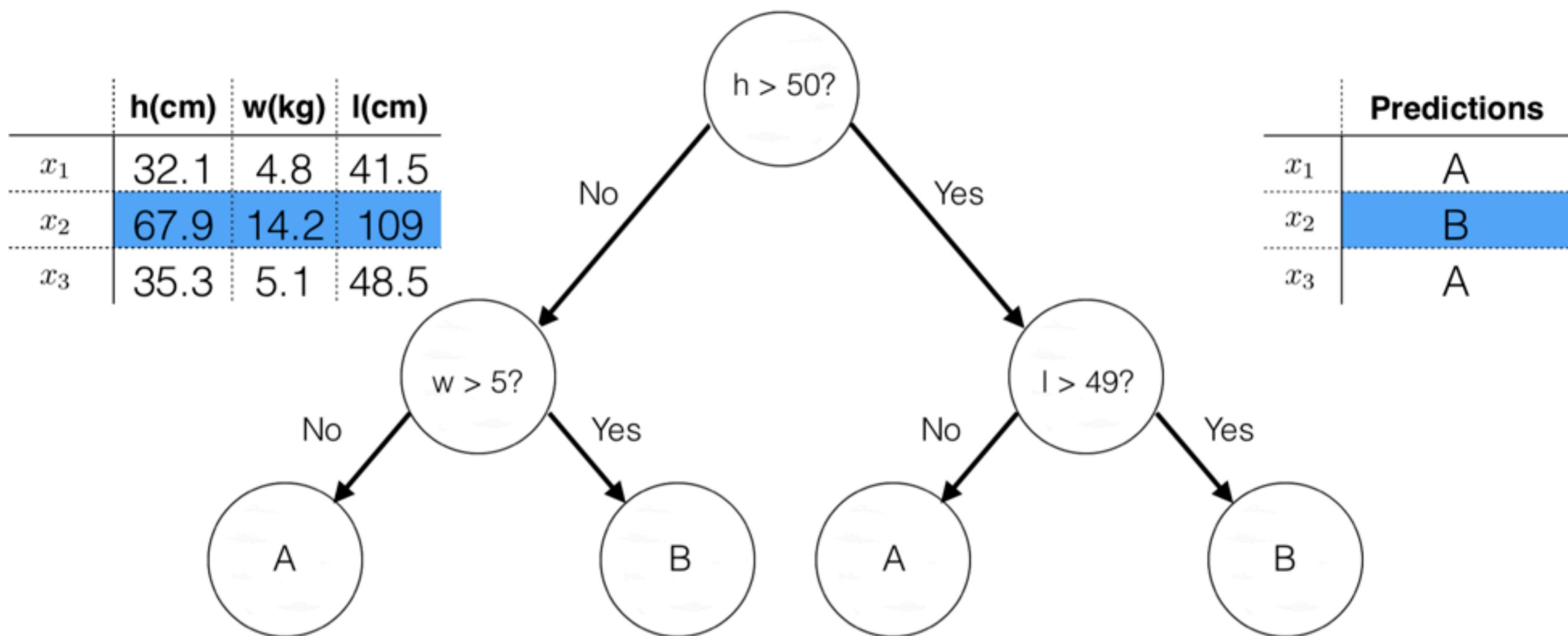
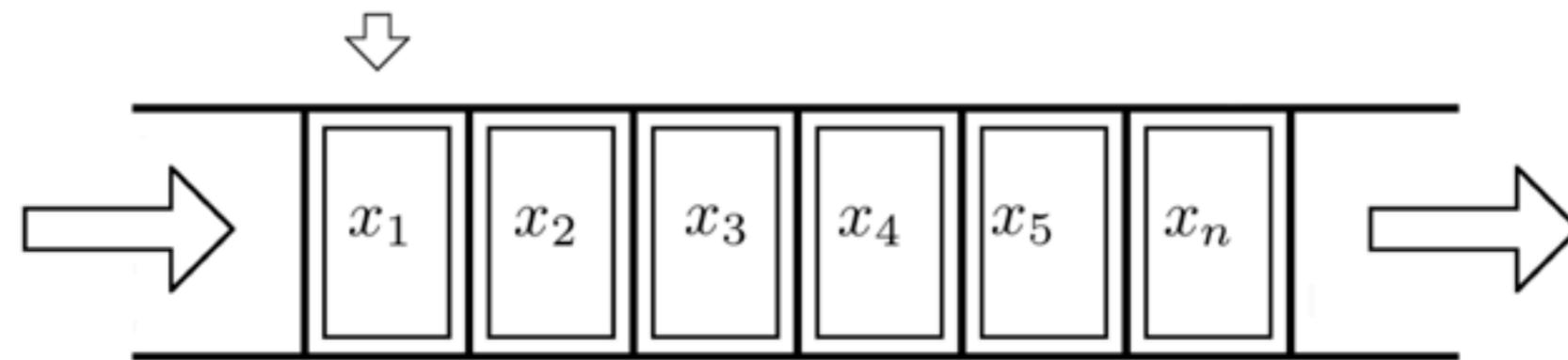
- Do nothing.
- Force balance in the dataset.
- Modify algorithm.
- Write a new algorithm.
- Treat it as an anomaly detection problem.
- Get more data.

POSSIBLE APPROACHES

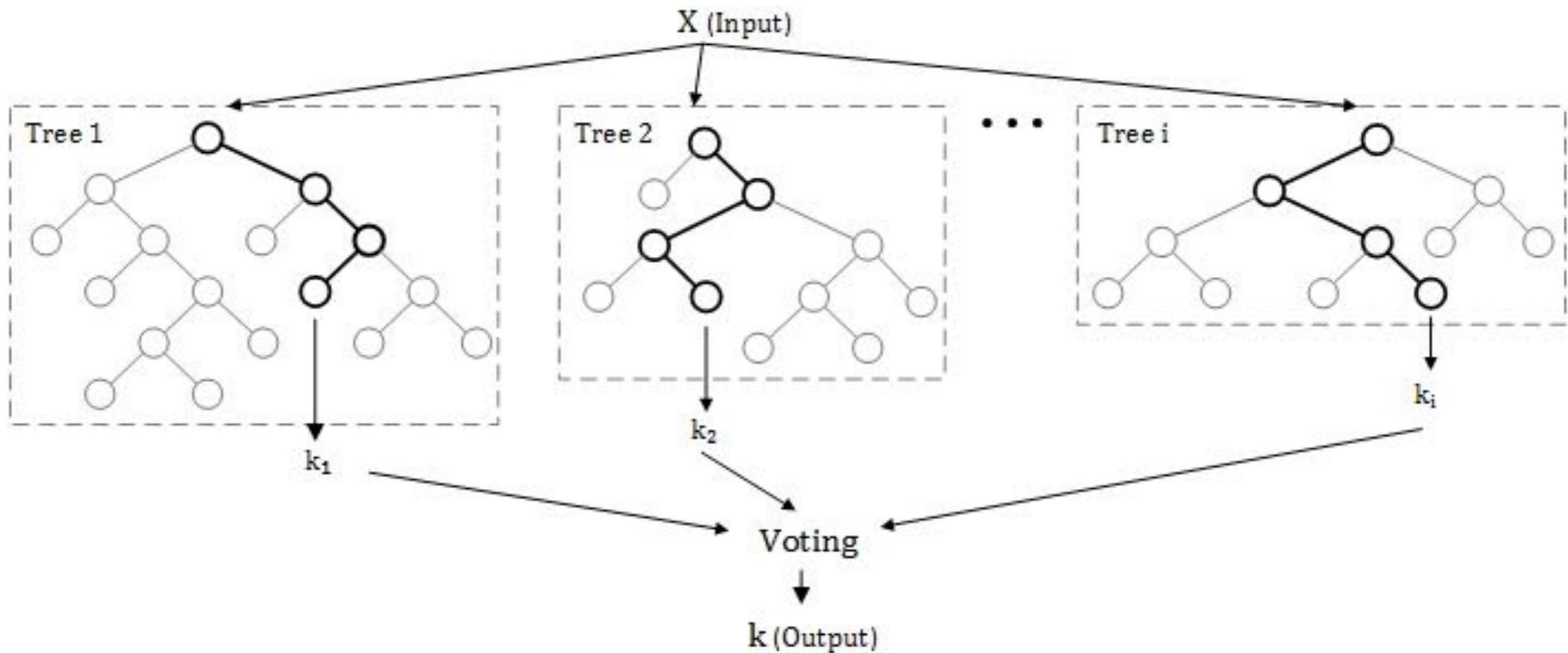
- Do nothing.
- Force balance in the dataset.
- Modify algorithm.
- Write a new algorithm.
- Treat it as an anomaly detection problem.
- Get more data.

- Not always possible
 - Probably part of the reason you have an imbalance in the first place...

DECISION TREES

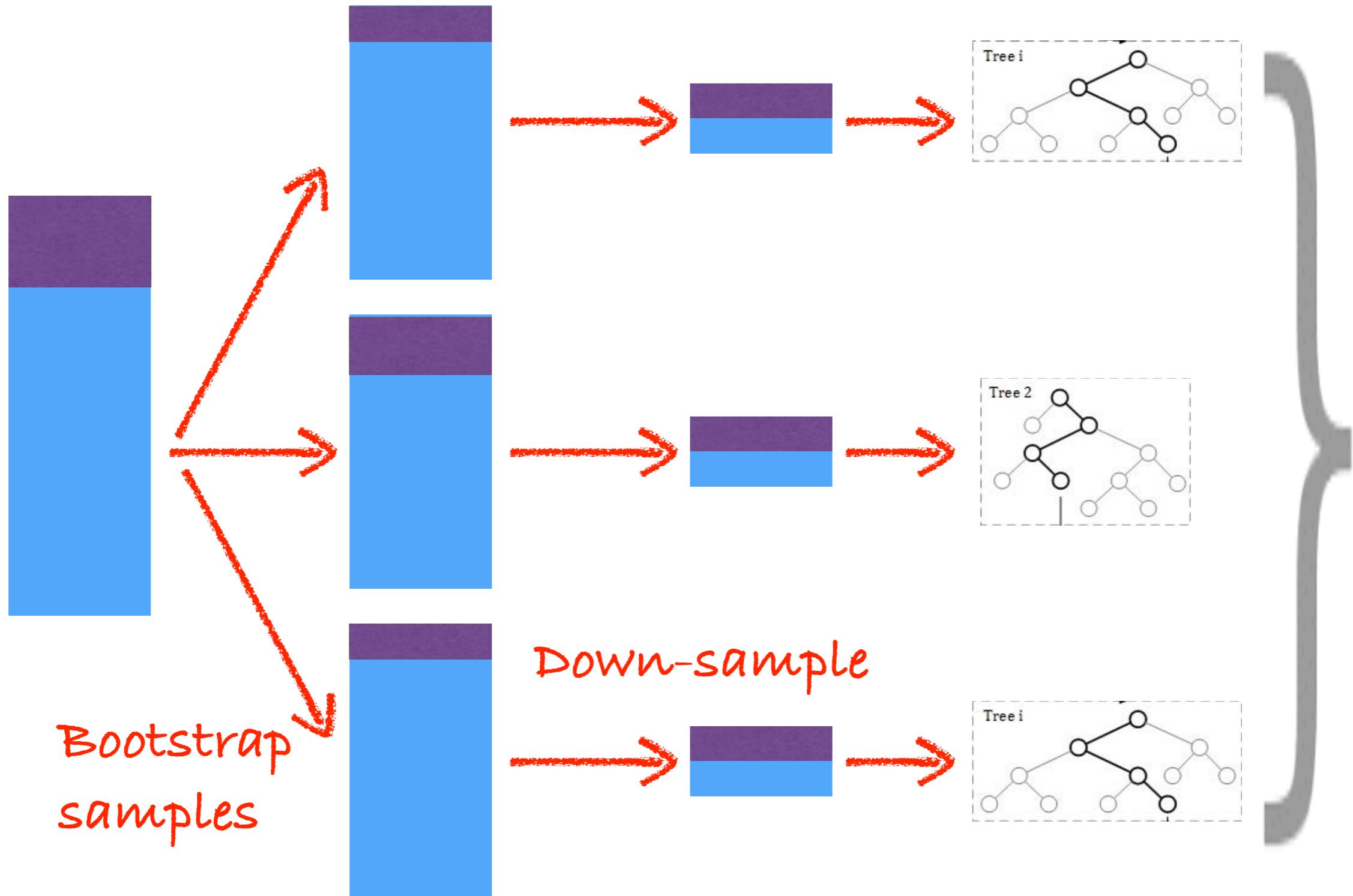


RANDOM FOREST CLASSIFIERS



Each tree uses **a subset of the data** and **a subset of the features**
(this is known as *bagging*)

BAGGED UNDERSAMPLING

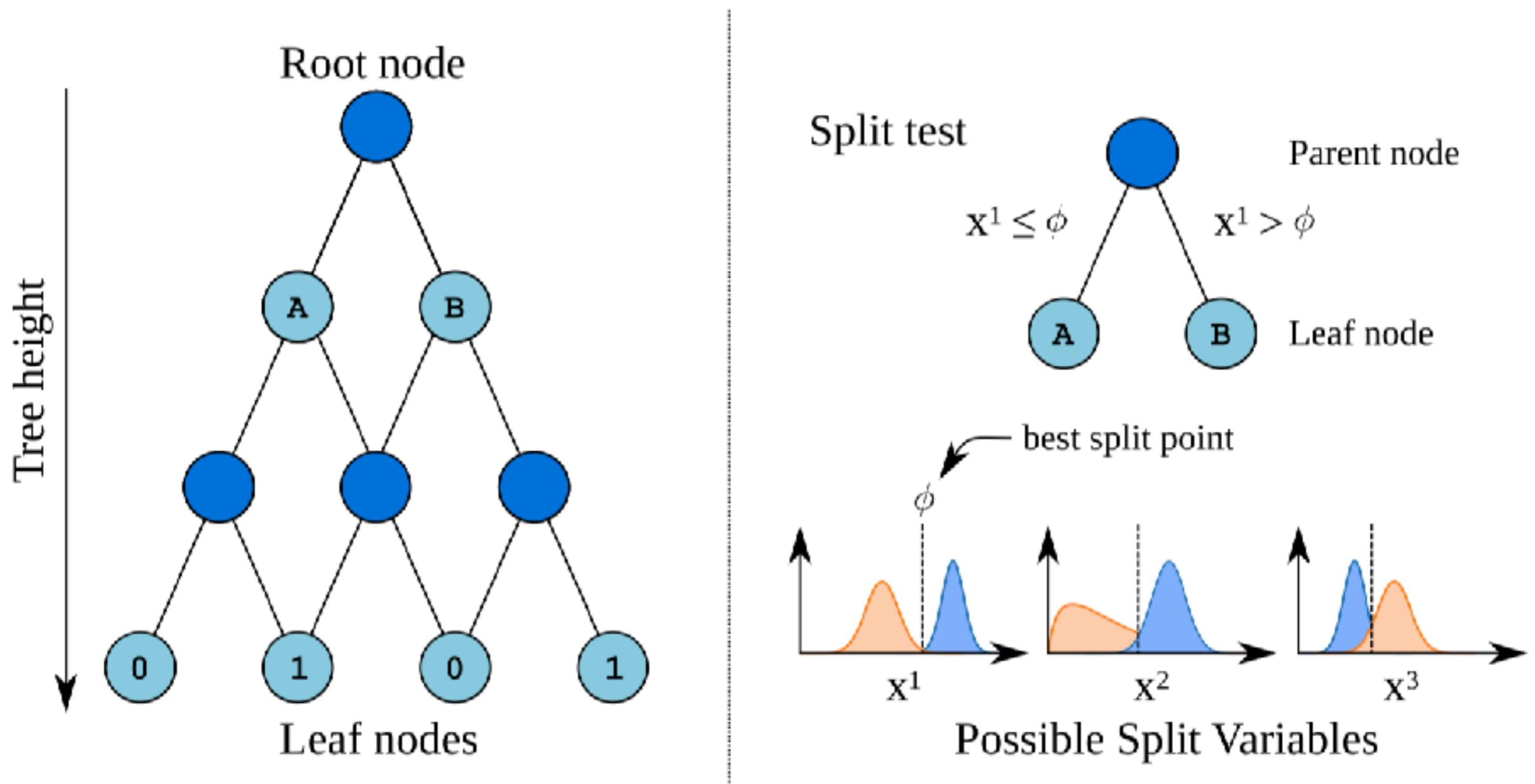


SKA PULSAR CLASSIFICATION

1,000 candidates ***per second***

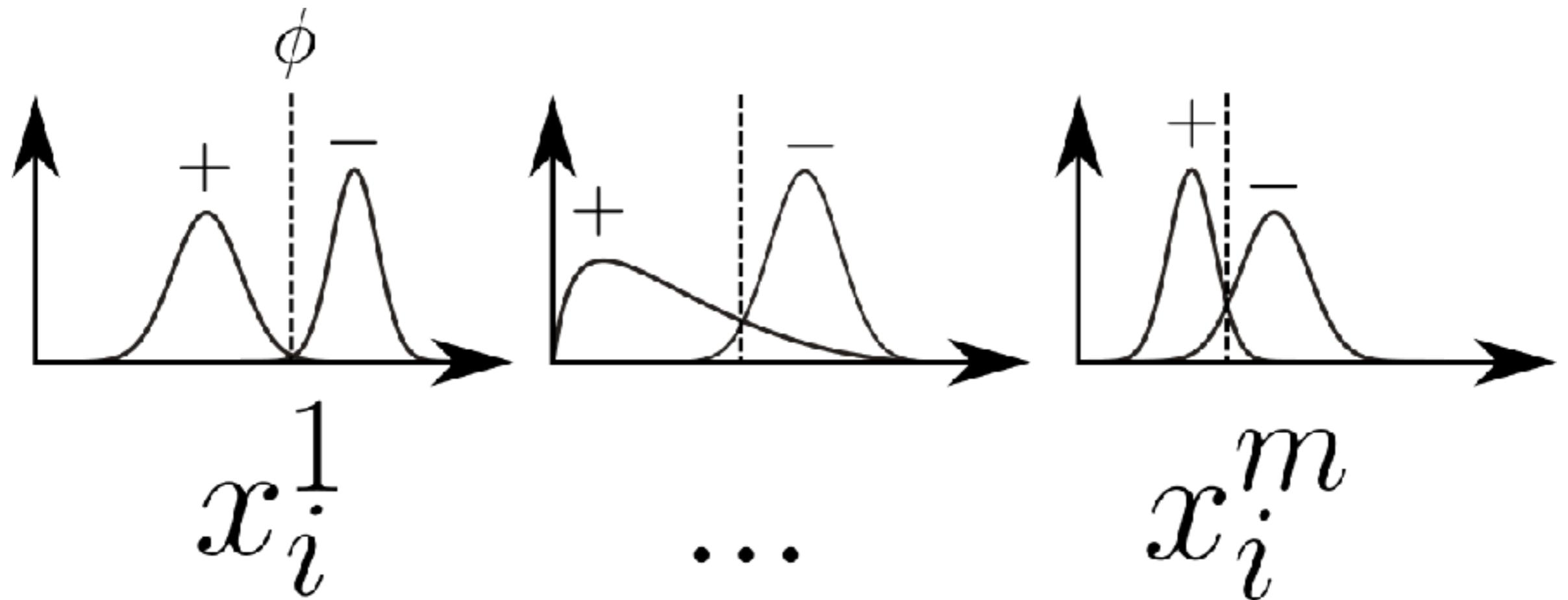
> 80,000,000 candidates ***per day***

SKA PULSAR CLASSIFICATION: GH-VFDT STREAM CLASSIFIER



Gaussian Hellinger Very Fast Decision Tree

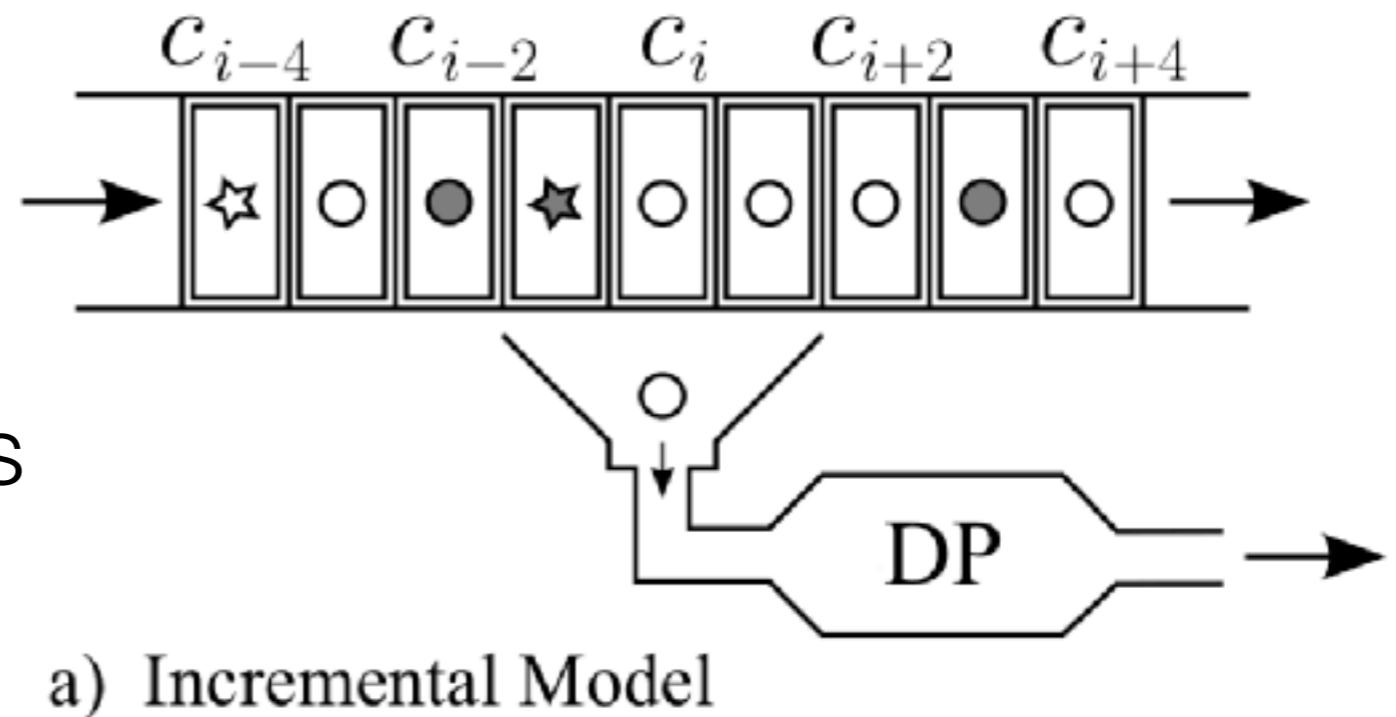
SKA PULSAR CLASSIFICATION: GH-VFDT STREAM CLASSIFIER



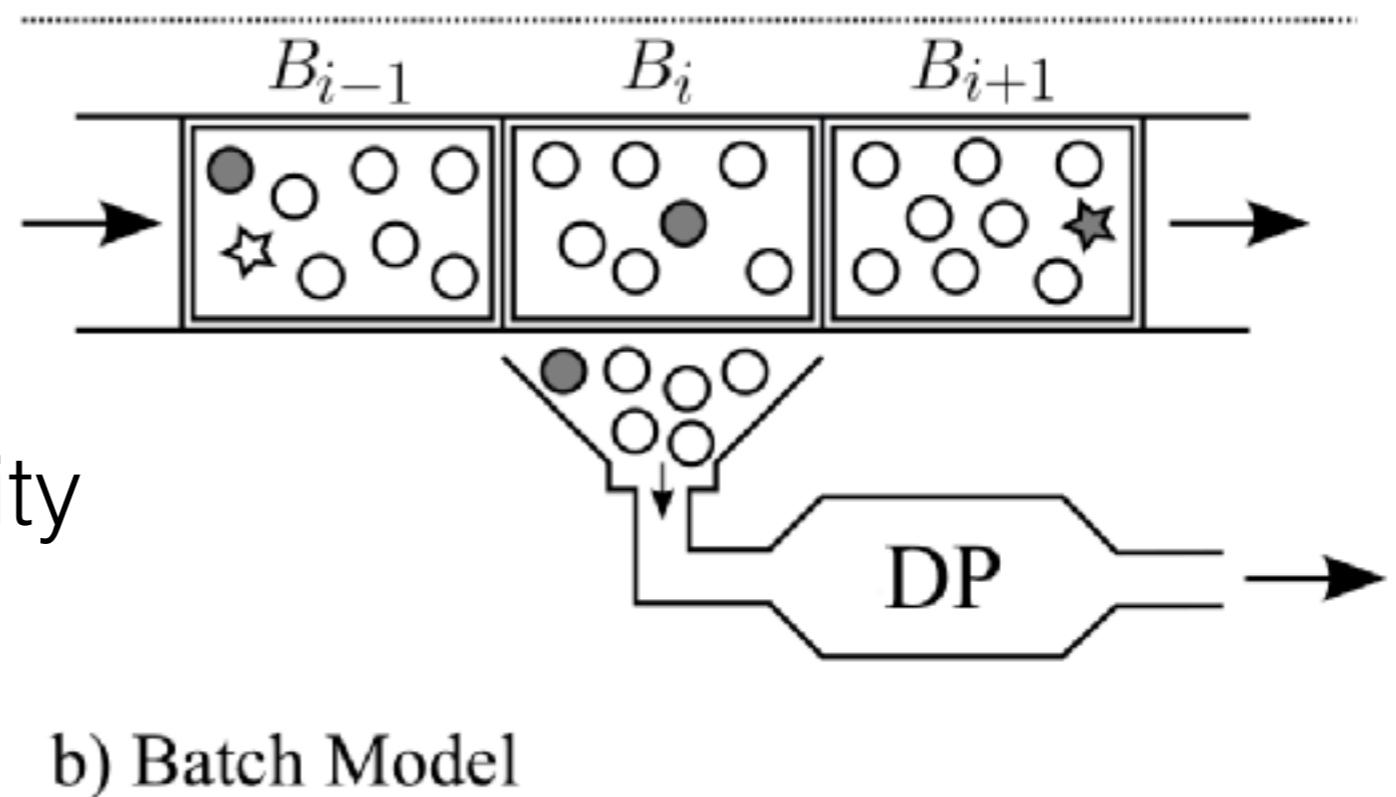
Hellinger Distance = separation between two distributions

REAL-TIME PROCESSING

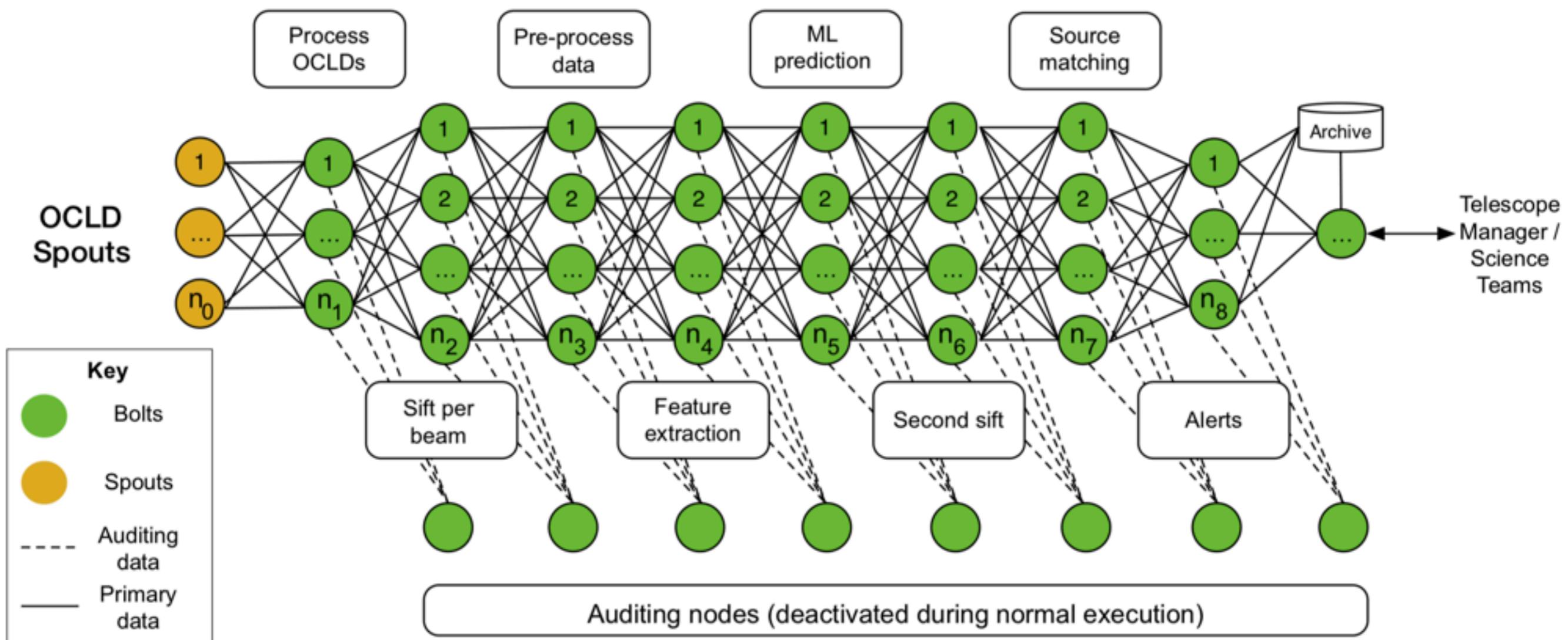
good for
isolated outliers



good for
localised similarity



SKA REAL-TIME PROCESSING



<https://github.com/scienceguyrob/GHVFDT>

TUTORIAL



<https://github.com/as595/NITheP>