# Detection of Galaxy Tidal Features using Self-Supervised Machine Learning

Alice Desmons, Sarah Brough, Francois Lanusse

## Gap in Knowledge

Tidal features are diffuse regions of stars that extend in the outskirts of galaxies. They are the result of stellar material being pulled out of merging galaxies due to immense gravitational forces. Studying these tidal features can help us better understand the process of galaxy evolution, but to do so we need a large sample of galaxies with tidal features. We turn to machine learning to achieve this.



*Figure 1:* Example of galaxies with tidal features taken from the HSC-SSP Survey
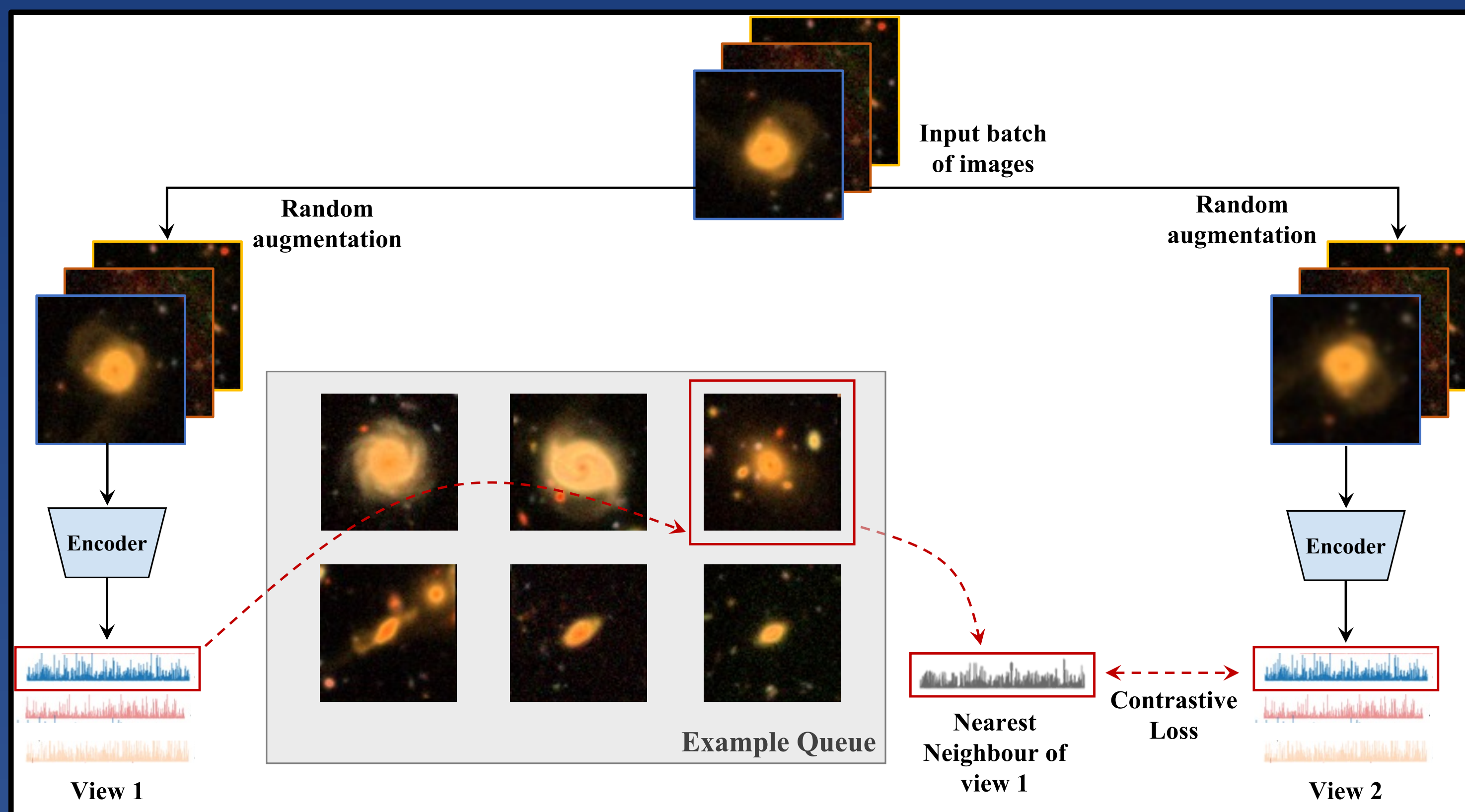


*Figure 2:* Illustration of the self-supervised model architecture

## Results: Self-Supervised vs. Supervised Performance

We compare the performance of our self-supervised model to that of a fully supervised model designed by Pearson et al. (2019). We do this by comparing the Area Under the Curve (**AUC**) of a Receiver Operating Characteristic (**ROC**) Curve for the two models as a function of the amount of unique training examples (**Figure 3**). Both models are trained on the same data and our model not only has a consistently higher ROC AUC but also maintains its performance regardless of whether we use 600 or 50 examples to train the classifier.

## Results: Similarity Search

One advantage of self-supervised models over supervised models is the ability to use just one labelled example to find examples of similar galaxies from the full unlabelled dataset. As shown in **Figure 4**, given just one "query image" the model can look through the dataset of ~45,000 galaxies and find examples of similar images. This is a useful tool for finding more of one type of galaxy.

## The Model

Our model consists of two sub-models: a self-supervised encoder, and a supervised classifier.

### The Self-Supervised Encoder

This model does no require labelled data, and instead uses **data augmentations** to learn to encode galaxy images into lower-dimensional representations. The encoder uses a **contrastive loss** function which forces the loss to be minimised when two representations are based on differently augmented versions of the same image. **Figure 2** illustrates the model architecture. In this case, the loss between view 2 and the nearest neighbour of view 1 will be minimised. The augmentations we use include flipping the images, adding random noise, and randomly cropping the images.

### The Supervised Classifier

This is a simple linear model which takes in the encoded representations output by the encoder; and returns a number between 0 and 1. Outputs close to 1 indicate a galaxy that likely possesses tidal features. This part of the model requires a small labelled training set.

## The Dataset

To train the self-supervised encoder we use an unlabelled dataset of ~45,000 galaxies. These galaxies were selected from the Ultradeep region of the Hyper Suprime-Cam Subaru Strategic Program (**HSC-SSP**) survey and have magnitudes $15 < i < 20$ mag.

To train the supervised classifier we use the Ultradeep HSC-SSP dataset visually classified in Desmons et al. (2023). This datasets consists of **300** galaxies with tidal features and **300** galaxies without tidal features.
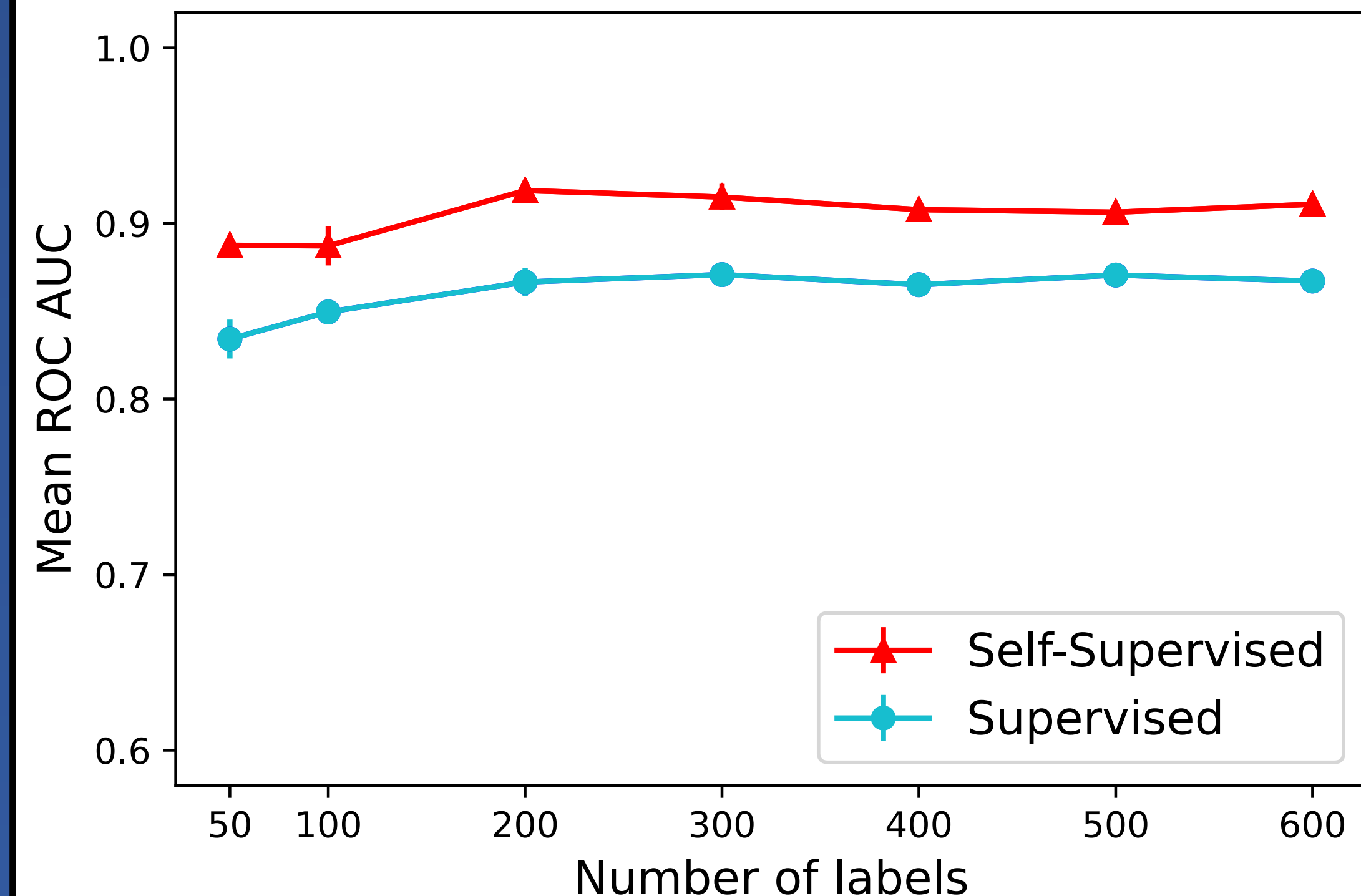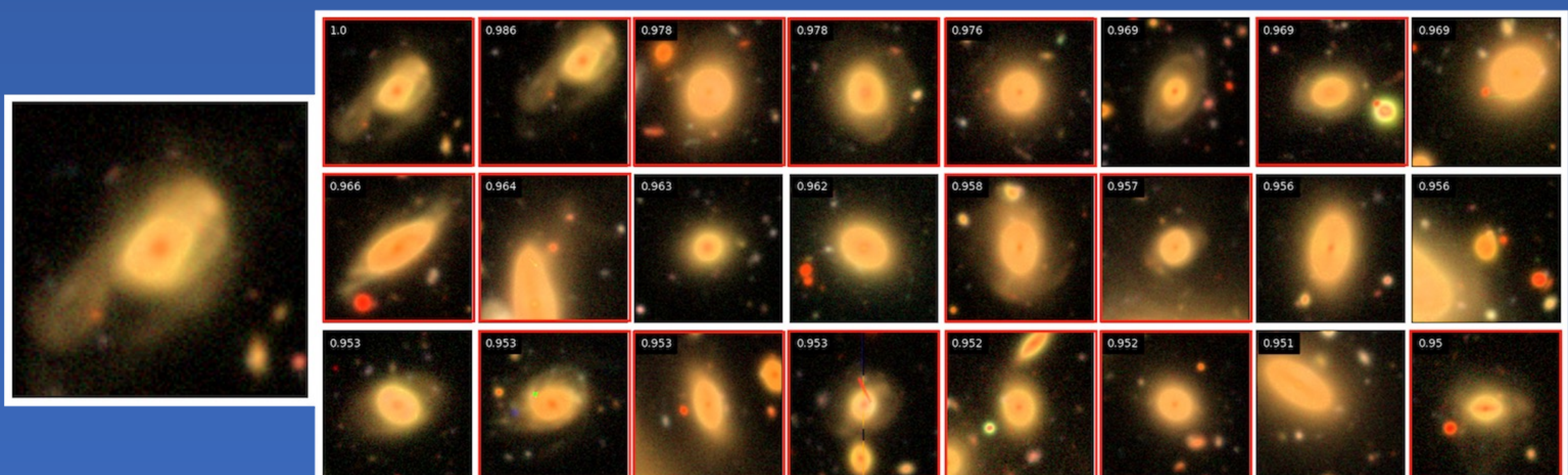


*Figure 3: Average ROC AUC as a function of labels used for training for a supervised and self-supervised model. Each point is an average of ten runs.*

Check out the code and documentation for the model to:
- Load and use the pre-trained models
- Train one or both models using your own data



*Figure 4:* Results from a similarity search using a query image, displayed on the left, alongside the top 24 galaxies with the highest similarity scores. The red outlines indicate images which would be visually classified as hosting tidal features.

a.desmons@unsw.edu.au