

---

# Comparison and verification of tidal feature detection methods

---

## LITERATURE REVIEW



**UNSW**  
SYDNEY

---

Haotian Lyu | z5283856

Supervisor:  
Professor Sarah Brough

---

August 29, 2024

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Structure of Universe</b>	<b>4</b>
2.1	Lambda Cold Dark Matter cosmological model . . . . .	4
2.2	Hierarchical structure formation . . . . .	5
<b>3</b>	<b>Galaxy mergers</b>	<b>6</b>
3.1	Types of galaxy mergers . . . . .	6
3.2	Tidal features . . . . .	7
<b>4</b>	<b>Detection of galaxy mergers</b>	<b>9</b>
4.1	Close pair detection . . . . .	9
4.2	Tidal features detection . . . . .	10
4.2.1	Visual classification . . . . .	10
4.2.2	Machine learning . . . . .	12
4.2.3	Cosmological simulations . . . . .	13
4.3	CAS method . . . . .	15
<b>5</b>	<b>This project</b>	<b>18</b>
<b>6</b>	<b>Data sources and models</b>	<b>19</b>
6.1	Mock images . . . . .	19
6.2	Machine learning model . . . . .	21
6.2.1	Self-Supervised Model . . . . .	21
6.2.2	The Linear Classifier . . . . .	23
<b>7</b>	<b>Thesis contents</b>	<b>24</b>

## 1 Introduction

Galaxies are gravitationally bound systems composed of stars, gas, dust, and dark matter. Galaxy evolution is a very important topic in contemporary astronomy, as it reveals how current galaxies have evolved from early galaxies. The widely accepted formation model of these structures is the ‘hierarchical structure formation model’, which states that large, massive galaxies are formed through the continuous merging of smaller galaxies (Lacey & Cole, 1993). Therefore, studying the mergers of galaxies has become an important element in the study of the evolution of galaxies.

One method to identify merging galaxies is Close pair detection. Pairs of galaxies that are close to each other are likely to interact or merge. By using their projected separation on the sky and velocity difference derived from redshifts from spectroscopic observation (e.g. Robotham et al. 2014), we can determine whether the physical 3D distance between two galaxies is small enough that they are gravitationally bound. This allows us to select galaxies for further study of their likely merger.

Another method used to identify merging galaxies is the detection of tidal features. These are diffuse, non-uniform regions of stars that extend into space from galaxies and are formed by gravitational interactions between the galaxies. During the mergers of the galaxies, the gravitational force pulls out the stellar material into distinctive shapes. Compared to close pair detection, which focuses on ongoing or imminent interactions and mergers, tidal features are the results of current or past mergers and interactions (e.g. Desmons et al. 2023). The morphology, colour, and number of tidal features provide valuable insights into the details of mergers. The number of tidal features indicates the frequency of mergers, while the colour reveals information about the mass of the parent galaxies involved in these interactions (Kado-Fong et al., 2018), and the morphology of the tidal features offers information on the initial angular momentum (e.g. Hendel & Johnston 2015).

In the past, the identification and classification of tidal features usually relied on visual classification (e.g. Tal et al. 2009, Atkinson et al. 2013, Martin et al. 2022). The low surface brightness of tidal features made them difficult to be detected in previous wide-field optical astronomical surveys, resulting in a lack of sufficient images. Consequently, visual classification was deemed sufficient for such classification tasks. However, with the advent of the Vera C. Rubin Observatory’s Legacy Survey of Space and Time (LSST; Ivezić et al. 2019), the images will be sensitive enough to detect tidal features, generating a vast amount of data. LSST is predicted to detect billions of galaxies (Ivezić et al., 2019), which makes it impossible for visual classification to deal with such a large amount of data. To solve this problem, Professor Sarah Brough and her team have begun developing methods that use Self-Supervised Representation Learning to automatically identify and classify galaxies with tidal features (Desmons et al., 2024).

Another method to explore galaxy mergers is the CAS (concentration, asymmetry, smoothness) parameters method. This non-parametric approach measures the shapes of galaxies in images (Conselice et al., 2008), based on the idea that the light distribution of galaxies provides insights into their past and present formation modes (Conselice, 2003). Among the three CAS parameters, asymmetry is considered as the main indicator of galaxy mergers. A higher value of asymmetry means an asymmetric light distribution, which is usually found in spiral galaxies and galaxy mergers (Conselice et al., 2008). The other two parameters also provide insights into galaxy mergers. For example, the comparison between asymmetry and smoothness can also be used as a reference for

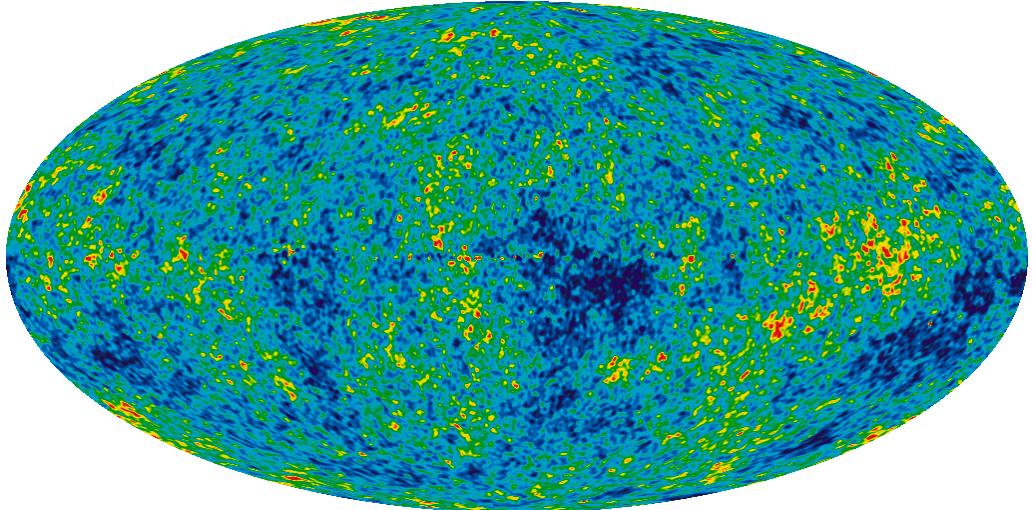
measuring galaxy mergers. Usually, the asymmetry parameter  $A > 0.35$  and a higher value of asymmetry than smoothness can be considered as the criterion of merging galaxies.

In this project, we will evaluate the self-supervised representation learning algorithm developed by [Desmons et al. \(2024\)](#) applied to mock images of galaxies generated by cosmological simulations ([Khalid et al., 2024](#)). By comparing the results from machine learning (ML) with those from visual classification, we aim to understand the advantages and limitations of automatic recognition algorithms in classifying tidal features. Additionally, we will examine how the results from ML and visual classification sit in the CAS parameter space. By analyzing these comparisons, parameters, and the characteristics of the galaxies, we hope to gain a deeper understanding of the key factors in detecting tidal features.

The literature review of this project consists of the following six sections: Section [2](#) describes the most widely accepted model of the universe: the Lambda Cold Dark Matter cosmological model, and also introduces the Hierarchical merger model for the formation of the galaxy structures. Then, Section [3](#) introduces the information about mergers between the galaxies, including the types of the mergers and the tidal features created by those mergers. Section [4](#) introduces the ways to detect galaxy mergers and the previous research on mergers and tidal features. The following Section [5](#) talks about the aims and plans of this project. The data sources that will be used in this project are presented in Section [6](#). Finally, Section [7](#) is an outline for the thesis of the whole project.

## 2 Structure of Universe

In the 1960s, the detection of the cosmic microwave background (CMB) provided significant support for the Big Bang theory (e.g. [Wilson & Penzias 1967](#), [Thorne 1967](#)). Then, the Lambda Cold Dark Matter ( $\Lambda$ CDM) Cosmological Model is developed to explain the formation of structure in the Universe (e.g. [Blumenthal et al. 1984](#), [Riess et al. 1998](#), [Perlmutter et al. 1999](#)). Figure 1 is the image of the cosmic microwave background showing the temperature anisotropies in the Universe.



**Figure 1:** The cosmic microwave background image, reveals temperature fluctuations (shown as colour differences). Original Figure 27 in [Bennett et al. \(2013\)](#).

In this section, we will introduce the most widely accepted cosmological model of the Universe ( $\Lambda$ CDM Model) and the paradigm of cosmic structure formation in this model.

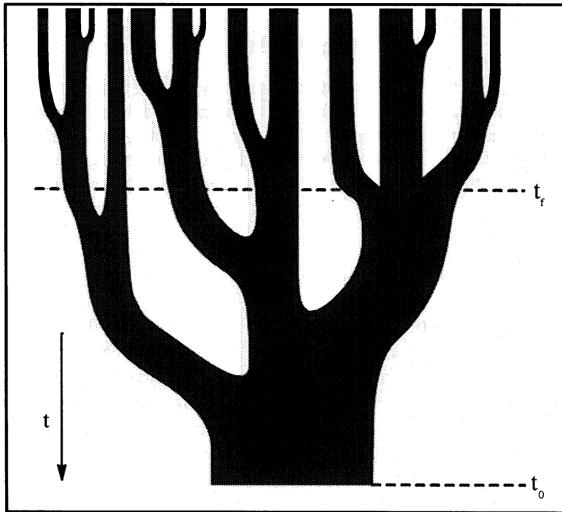
### 2.1 Lambda Cold Dark Matter cosmological model

The Lambda Cold Dark Matter cosmological model is currently the most widely accepted model in cosmology. Lambda ( $\Lambda$ ) is the cosmological constant driven by Dark Energy, and the dark matter is Cold Dark Matter. In this model, approximately 27% of the universe is dark matter and 68% is dark energy, the baryonic matter that makes up the visible galaxies, and stars accounts for only a small fraction of 5% ([Abbott et al., 2019](#)). The introduction of Dark Energy helps the model successfully predict the observed accelerating expansion of the Universe (e.g. [Riess et al. 1998](#), [Perlmutter et al. 1999](#)).

In the  $\Lambda$ CDM model, the early expansion of the Universe creates the initial density perturbations that eventually develop into galaxies and other large-scale structures. Cold dark matter plays an important role since it forms initial density perturbation. Baryons condense in these gravitational potential wells and form stars and eventually galaxy structures. These small perturbations allow matter to gradually collapse and combine in the dark matter halo under the influence of gravity (e.g. [Wechsler & Tinker 2018](#)), and then gradually form the current observed structure of galaxies and clusters of galaxies according to the ‘[hierarchical structure formation paradigm](#)’. The CMB therefore provides a snapshot of the early Universe, showing the initial density perturbations that lead to the formation of current galaxy structures.

## 2.2 Hierarchical structure formation

Hierarchical structure formation is a major prediction from the  $\Lambda$ CDM model (e.g. [Dalal et al. 2008](#)). According to the hierarchical structure formation paradigm, the formation of structures, such as galaxies and galaxy clusters, occurs from the bottom up. This means that larger structures are formed through the continuous merging of smaller structures ([Lacey & Cole, 1993](#)). Figure 2 shows a schematic of hierarchical formation of a massive galaxy structure over time.



**Figure 2:** A schematic of a ‘merger tree’ that describes the formation of a massive galaxy as the result of a series of mergers. The widths of branches represent the masses of the galaxies, and time increases from top to bottom, where  $t_0$  is now and  $t_f$  is the time of formation of the final structure. Originally Figure.6 in [Lacey & Cole \(1993\)](#).

As mentioned in Section 2.1, perturbations in the early Universe caused baryons within dark matter halos to condense under gravity, forming low-mass galaxies. These dark matter halos acted as gravitational scaffolding. As the Universe expanded and gravity continued to act, halos and matter further merged, forming more massive galaxies as well as groups and clusters of galaxies, each other their own dark matter halo ([White & Frenk, 1991](#)).

In the hierarchical structure formation paradigm, galaxy mergers play a significant role in the evolution of galaxies. The mergers of low-mass galaxies over time eventually form massive galaxies. We will describe the mergers of the galaxies in the next section.

### 3 Galaxy mergers

A galaxy merger occurs when two or more galaxies collide and merge due to the gravitational attraction between them. This is an important process in galaxy evolution, and the merger rate is a fundamental measurement of galaxy evolution. In this section, we will introduce the types of galaxy mergers and the tidal features created by those mergers.

#### 3.1 Types of galaxy mergers

The most common classification of galaxy mergers is based on the relative mass of the two merging galaxies. A minor merger is when two merging galaxies have a large mass difference. Usually, if the secondary galaxy is less than 1/4 the mass of the primary galaxy in a merger (e.g. [Robotham et al. 2014](#)), we call this a minor merger. In contrast, major mergers are mergers where the two merging galaxies have comparable masses with mass ratios from 1:1 to 1:4.

Of the two types of mergers, minor mergers occur much more frequently than major mergers. Observations and Simulations indicate that minor mergers are at least three to four times more common than major mergers due to the larger number of lower mass galaxies ([Lin et al. 2004](#), [Lotz et al. 2011](#)). Furthermore, the hierarchical structure formation model indicates that low-mass collisions (minor mergers) play a more significant role in the evolution of galaxies (e.g. [Naab et al. 2009](#)). Simulations by [Martin et al. \(2018\)](#) predict that most of the morphological evolution of galaxies since  $z \sim 1$  has been driven primarily by minor mergers. Due to the large mass difference between the primary and secondary galaxies, a minor merger has a relatively smaller effect on the primary galaxy. However, many minor mergers over time will cause disturbances in the primary galaxy's morphology. Unfortunately, the faintness of low-mass galaxies makes minor mergers difficult to be able to quantify this effect observationally (e.g. [Robotham et al. 2014](#)).

Major mergers are more intense because the mass difference between two merging galaxies is smaller which leads to significant changes in their morphology in just one single merger (e.g. [Santucci et al. 2024](#)). For example, during a major merger of disk galaxies, the tidal forces can cause the originally stable disk to develop a bar-like structure ([Barnes & Hernquist, 1996](#)). Major mergers of spiral galaxies are often considered one of the mechanisms for the formation of elliptical galaxies ([Barnes, 1992](#)), and this has been verified in hydrodynamical simulations ([Springel et al., 2005](#)). Additionally, compared to minor mergers, the remnants or morphological changes resulting from major mergers are more pronounced and easier to detect. Figure 3 is the image of NGC 2623, which is the result of a major merger between two galaxies. We can see that the major merger has caused a significant change in the galaxies morphology, forming two tidal tails that extend 15.3 kpc.

Galaxy mergers not only influence the morphology of galaxies but also star formation and the growth of black holes. Galaxy mergers are predicted to cause the merging of central black holes and cause nuclear inflows of gas which contribute to black hole growth ([Hopkins et al., 2005](#)). Also, major mergers will lead to gas inflow and compression, which can trigger a sharp spike in star formation ([Mihos & Hernquist 1996](#)). We can see this phenomenon well in Figure 3, the active star formation is indicated by speckled patches of bright blue. These patches are clustered at the centre of the merger and along the tidal tails where tidal forces would be at a maximum.

In addition to classifying the merging galaxies based on their mass ratios, we can also classify them



**Figure 3:** Image of NGC 2623, shows a major merger between two galaxies. Original images from ESO.

based on galaxies’ gas richness. A wet merger is a merger between gas-rich galaxies. Typically, the wet merger will lead to a large amount of star formation since a lot of gas is present so can be compressed and flow into the galaxies ([Lin et al., 2008](#)). If the merging galaxies are all gas-poor, this merger is called a dry merger.

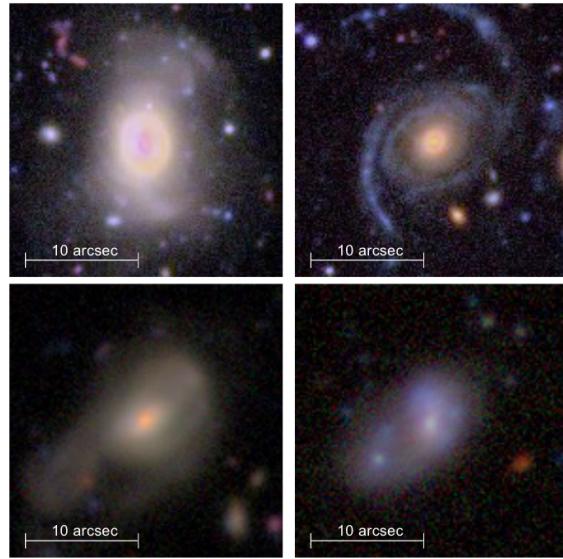
### 3.2 Tidal features

Both major and minor mergers will cause changes in galaxy morphology. Gravitational interactions between merging galaxies will pull stars out and these form diffuse, non-uniform regions of stars that extend into space. These are called tidal features.

Because of the range of angles, speeds, and angular momentum of galaxy mergers, as well as the different properties of the parent galaxies, tidal features also have different morphologies. The examples of the main tidal feature morphologies are shown in Figure 4. These include shells, streams, asymmetric halos, and double nuclei.

The formation of different tidal features has been studied using N-body simulations which provide detailed predictions of the galaxy mergers and interactions that cause different tidal feature morphologies, such as dynamics, mass ratios, and orbits. For example, [Karademir et al. \(2019\)](#)’s simulation found that in minor mergers, tidal streams can be formed from the nearly circular infall of a satellite galaxy with high angular momentum while shells are the result of a nearly radial infall. In contrast, tidal tails are generally considered to be the result of high angular momentum interactions in major mergers. [Pop et al. \(2018\)](#)’s analysis of Illustris simulation ([Nelson et al., 2015](#)) predicted that higher mass galaxies are more likely to have tidal shells, while small satellite galaxies need almost purely radial infall to produce tidal shells.

Therefore, tidal features provide a pathway to explore the details of the mergers and interactions between galaxies. The number of tidal features indicates information about the rate of collisions, and the colour provides information on the mass of the parent galaxies involved in the collisions. For example, [Kado-Fong et al. \(2018\)](#) found the average colours of shells are bluer than their host galaxies but that some shells are red indicating their formation from major mergers. These features



**Figure 4:** Example galaxies with different tidal feature morphologies. Top row from left to right: shells, stream. Bottom row from left to right: asymmetric halo, and double nuclei. Originally Figure.1 in [Desmons et al. \(2024\)](#).

are predicted to exist for a few billion years after a collision, which makes the tidal features good tracers of galaxy mergers.

## 4 Detection of galaxy mergers

There are several ways to observe and study galaxy mergers. One approach is using relative sky position and relative velocity driven by redshifts to detect close pairs of galaxies which may indicate potential galaxy mergers. Another method is finding the tidal features in the galaxy images, that can provide clues of galaxy mergers. The CAS method is another effective way to detect merging galaxies. CAS analyzes the light distribution in galaxies (Conselice et al., 2008). In this Section, we will discuss the basic ideas, advantages, and disadvantages of each method and some existing examples.

### 4.1 Close pair detection

Close pair detection is a method that focuses on finding dynamically close galaxies that by definition are pre-merger states. The close pairs of galaxies can be selected by determining whether they are close both in projected position and relative velocity (Robotham et al., 2014). Then using the frequency of close pairs, the near future merger rates of galaxies can be predicted.

The distance of pairs of galaxies can be described by the physical projected separation  $r_p$  and rest-frame relative velocity along the line of sight  $\Delta v$ :

$$\begin{aligned} r_p &= \theta d_A(z_i) \\ \Delta v &= \frac{c|z_j - z_i|}{1 + z_i} \end{aligned} \quad (1)$$

Where  $z_i$  and  $z_j$  are the redshifts of the two galaxies,  $i$  is the index of the more luminous galaxy, and  $j$  is the secondary galaxy.  $\theta$  is the angular separation between the two galaxies (in arcsec), and  $d_A(z)$  is the angular scale at redshift  $z$  (in kpc/arcsec) (e.g. López-Sanjuan et al. 2011). One definition of close pairs is given by (Patton et al., 2002):

$$\begin{aligned} 5 h^{-1}kpc &\leq r_p \leq 20 h^{-1}kpc \\ \Delta v &\leq 500 km s^{-1} \end{aligned} \quad (2)$$

Galaxies that are this close are considered to be likely to merge in 0.5 Gyr, and it is expected that 50–70 % of them will finally merge (López-Sanjuan et al., 2011). Selecting close pairs requires the redshifts of the galaxies and the measurement of the redshifts is based on spectroscopy. Large-scale high-completeness spectroscopic surveys such as the Galaxy And Mass Assembly (GAMA; Driver et al. 2011) provide complete enough spectroscopic data to robustly select close pairs.

Close pair detection provides per-merger insights, allowing us to study the initial conditions and dynamics of the merging galaxies. By counting the number of close pairs at different redshifts, we can understand how the merger rate evolves over cosmic time.

However, close pair detection also has some limitations. The main one of these is the observational expense of high-completeness spectroscopy that then limits the sample sizes available. Also, due to the potentially low mass of the secondary galaxies, spectroscopy may have difficulty detecting them, which may lead to biases in the identified pairs and then skew the merger rates. Close pairs may also have complex orbital dynamics, so it is difficult to accurately identify interacting galaxies simply based on projected separation and relative velocity, noting that only 50% to 70% were

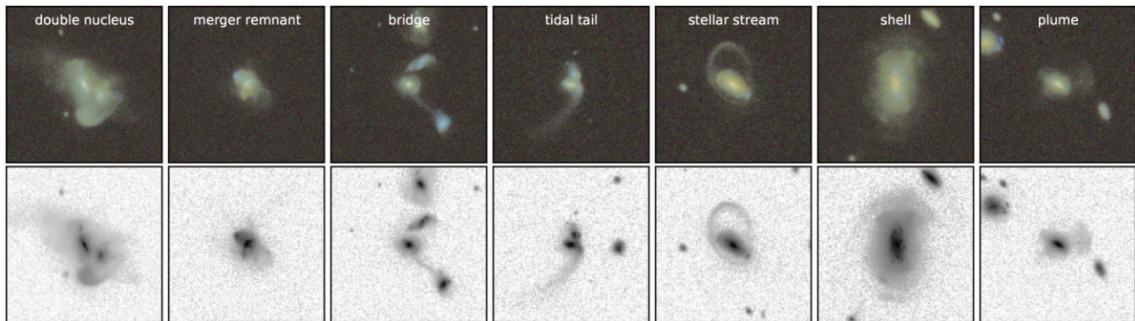
found to eventually merge in simulation by [López-Sanjuan et al. \(2011\)](#). Since close pair detection focuses on the pre-merger states, it also has limitations in studying the properties of galaxies after mergers, such as their morphology.

## 4.2 Tidal features detection

Tidal feature detection differs from close pair detection by focusing on the current or post-merger states of galaxies (e.g. [Desmons et al. 2024](#)). The detection and study of tidal features provide information about different phases of galaxy mergers. It faces challenges due to the low surface brightness of the tidal features, which can easily reach  $\mu_r \geq 27$  mag arcsec $^{-2}$  ([Bilek et al., 2020](#)). Such low surface brightnesses make it difficult for current wide-field optical astronomical surveys to detect, leading to the current tidal feature dataset being smaller than close pair samples constructed by spectroscopic surveys. But such surface brightnesses will be able to be observed by LSST, which reaches depths of  $\mu_r \sim 30.3$  mag arcsec $^{-2}$  ([Martin et al., 2022](#)) and plans to detect billions of galaxies ([Ivezić et al., 2019](#)), significantly enhancing the samples of galaxies with selection of tidal features.

Currently, the identification and classification of tidal features mainly rely on visual classification, but with the development of artificial intelligence, classifications based on machine learning are becoming more common. In this Section, we will discuss these two different methods. Cosmological simulations are also important methods for studying tidal features and these will also be discussed in this Section.

### 4.2.1 Visual classification



**Figure 5:** Example of an extended range of tidal feature morphologies from simulated images of each category. The top row shows the initialized false colour images and the bottom shows the greyscaled surface brightness maps. Originally Figure 6 in [Martin et al. \(2022\)](#).

Over the years, visual classification has been the primary method for tidal feature classification and identification. Many studies have performed visual classification and identification on different datasets, using slightly different classification criteria depending on the research data available and their focus.

In [Tal et al. \(2009\)](#)'s paper, they divided tidal features into four main categories: (1) Shells, (2) Tidal tails, (3) Broad fans of stellar light, (4) Highly disturbed galaxies. Studying a sample of 55 nearby massive elliptical galaxies observed with Cerro Tololo Inter-American Observatory, they found that 73% of galaxies show tidal features.

[Atkinson et al. \(2013\)](#) analyzed a sample of 1781 galaxies from the wide-field component of the Canada–France–Hawaii Telescope Legacy Survey. They classified tidal features into six categories: (1) Streams, (2) Arms, (3) Linear features, (4) Diffuse fans, (5) Shells, and (6) Miscellaneous diffuse structure. The galaxies showing clear tidal features classified with the highest confidence account for 12% of the sample. If also including galaxies with weaker tidal features, the fraction rises to 18%.

[Hood et al. \(2018\)](#)’s study focuses on the tidal features around galaxies in the REsolved Spectroscopy Of a Local VolumE (RESOLVE) survey of very nearly galaxies. They categorized tidal features as “narrow” and “broad”. The narrow category includes the streams, arms, and linear features while the broad category includes shells, fans, and miscellaneous diffuse structures. This led to their finding that  $17 \pm 2\%$  of their sample with 1048 galaxies have tidal features.

[Desmons et al. \(2023\)](#) visually classified and analyzed a sample of 852 galaxies from the Galaxy and Mass Assembly (GAMA) survey, identifying 198 galaxies with tidal features, which represents 23% of the total sample. A comparison with the GAMA spectroscopic close-pair sample revealed that 42 galaxies appeared in both the close-pair and tidal feature samples. They found that while close-pair detection is more effective for identifying early-stage mergers, tidal features are better suited for detecting late-stage mergers.

In preparation for the upcoming LSST, [Martin et al. \(2022\)](#) visually classified a sample of mock images generated by NEWHORIZON cosmological simulation ([Dubois et al., 2021](#)) to study the nature, frequency, and visibility of tidal features and debris in a variety of environments and stellar masses. They further modified and subdivided the [Bílek et al. \(2020\)](#)’s classification, dividing tidal features into the following 7 categories: (1) Stellar streams, (2) Tidal tails, (3) Asymmetric stellar halos (plume), (4) Shells, (5) Tidal bridges, (6) Merger remnants and (7) Double nuclei. Examples of each category in [Martin et al. \(2022\)](#) are shown in Figure 5.

More details of galaxy mergers were revealed through the visual identification of tidal features. The study performed by [Bridge et al. \(2010\)](#) includes  $\approx 27000$  galaxies observed by the Canada–France–Hawaii Telescope Legacy Deep Survey. They found  $\approx 1600$  merging galaxies through visual classification. The merger fraction at different redshifts was  $4.3\% \pm 0.3\%$  at  $z \sim 0.3$  and  $19.0\% \pm 2.5\%$  at  $z \sim 1$ . They also explored the relationship between merging galaxies and star formation rate, finding that interacting galaxies have twice the average star formation rate of non-interacting galaxies.

[Bílek et al. \(2020\)](#) examined how galaxies with tidal features relate to their stellar mass. In 177 nearby massive early-type galaxies from the Mass Assembly of early-Type GaLAXies with their fine Structures survey data source, visual classification revealed that approximately 15 percent of galaxies exhibit shells, streams, and tails, with each category representing a similar fraction. However, the incidence of shells and streams are 1.7 times higher for galaxies with masses over  $10^{11} M_{\odot}$ . [Atkinson et al. \(2013\)](#) also found that the proportion of tidal features is a strong function of galaxy mass, with massive galaxies (stellar masses  $\geq 10^{10.5} M_{\odot}$ ) more likely to have tidal shells and fans. [Kado-Fong et al. \(2018\)](#) selected 1201 galaxies with tidal features from 21,208 galaxies in Hyper Suprime-Cam Subaru Strategic Program (HSC-SSP) images. They found that tidal streams occur in galaxies across the full range of stellar masses in their sample, whereas galaxies with shells are predominantly massive.

The gas content of merging galaxies and the resulting merger type (see Section 3.1) have also been analyzed in relation to the classification of tidal features. In [Tal et al. \(2009\)](#), the survey focused on tidal features in massive elliptical galaxies and found that those in groups and low-density environments are predominantly growing through mergers between gas-poor galaxies, with little star formation. In comparing tidal features in gas-rich and gas-poor galaxies, [Hood et al. \(2018\)](#) found that tidal features are more likely to appear in gas-rich galaxies ( $19\% \pm 2\%$ ), while only  $13\% \pm 3\%$  of gas-poor galaxies have tidal features. However, by looking at the relationship between tidal features and stellar masses, they found that in gas-poor galaxies, tidal features are likely to come from gas-poor mergers, while tidal features in gas-rich galaxies are mainly tails and streams.

In summary, there have been many studies on tidal features based on visual classification, which have increased our understanding of the relationship between galaxy properties and tidal features. However, there are discrepancies between these, e.g. tidal fractions [Tal et al. \(2009\)](#), that are likely due to the depths of the images studied and the sample selection. LSST will greatly improve on this with a large, homogeneous, deep sample

#### 4.2.2 Machine learning

In the area of galaxy surveys, the use of data-driven methods based on artificial intelligence such as deep learning is becoming more and more popular as a solution to long-standing problems ([Huertas-Company & Lanusse, 2023](#)). These methods are also increasingly used in automatically identifying tidal features. They offer a significant advantage over visual identification in that they can process a large number of galaxy images and so are much more efficient than the visual classification method. In this section, we will discuss automated classification based on machine learning.

Nowadays, the most widely used method for image classification problems is supervised learning. The characteristic of supervised learning is using already classified images as training datasets, such as using Convolutional Neural Networks (CNN). CNNs have been used to do galaxy morphological classification (e.g. [Walmsley et al. 2019](#), [Cheng et al. 2020](#)) and to further classify galaxy mergers and tidal features.

[Walmsley et al. \(2018\)](#) used CNNs to classify faint tidal features in the images from the Wide component of the Canada–France–Hawaii Telescope Legacy Survey (CFHTLS-Wide). The performance of the CNN in tidal feature classification was evaluated by comparing its results with those from visual classification. They found that their CNN method had a high completeness of 76% and a low contamination of 20%, which shows the potential of CNNs for tidal feature classification tasks.

A major problem with using supervised learning like CNNs to classify tidal features is the lack of appropriate training data sets. There are two reasons for this. First, due to the limitation of detection depth, past surveys could not form a large tidal feature data set. Second, the labeled data used as training data sets require visual classification, which is then limited by the time that process takes.

To address the first limitation, many supervised learning studies have used simulated images as training data sets. For example, [Bickley et al. \(2021\)](#) used the data set of post-merger galaxies

generated from IllustrisTNG simulation (Nelson et al., 2015) as the training sample and processed this data set with the observational realism code REALSIM (Zhu et al., 2015). Their CNN method achieved a classification accuracy of 88%. Sánchez et al. (2023) also use simulated images as their training set, these approximately 6000 images were generated by the NewHorizon simulations and visually classified by Martin et al. (2022). They used CNNs to achieve good results on simulated images, with an accuracy of 0.84, a precision of 0.72, and a recall of 0.85. However, when these models trained with simulated images are applied to the classification of observed images, the performance drops significantly, which is due to the gap in authenticity between the simulated and observed images. Pearson et al. (2019) compared the performance of deep learning networks trained on simulated images from EAGLE (Crain et al., 2015) and observed images from SDSS (Almeida et al., 2023). They discovered that the networks that are trained and evaluated using the same type of data exhibit optimal performance and the observation-trained network performed better than the simulation-trained network. When using the simulation-trained model for the classification of observed galaxies, they only get an accuracy of 65.2% compared with visual classifications, which shows that only using simulated images to train models will have limitations in real image classification without adding an additional transfer learning fine-tuning stage (Walmsley et al., 2024).

To solve the second problem, many studies have turned to unsupervised learning. Unlike supervised learning, unsupervised learning does not require labeled training samples but classifies images by self-summarizing features. For instance, Hocking et al. (2017) and Martin et al. (2020) used unsupervised learning for galaxy morphological classification, and achieved good consistency with the results of manual classifications.

However, unsupervised learning can only cluster and classify data without matching them to labels. There is also self-supervised learning (SSL), which is a middle point between supervised learning and unsupervised learning. This approach uses a large unlabeled dataset for training and then combines it with a classifier trained on a smaller labeled dataset. This satisfies the task requirements of classifying tidal features well. SSL uses unlabeled images to train the model which converts the images into meaningful low-dimensional representations, then uses a smaller labeled data set to train a linear classifier. SSL can perform classifications by combining the linear classifier with meaningful low-dimensional representations of images. Desmons et al. (2024)'s SSL model is trained on unlabeled HSC-SSP data and was shown to achieve very good results, maintaining low contamination while achieving 96% of completeness. Moreover, it maintained good performance when trained with only 50 labeled galaxies.

#### 4.2.3 Cosmological simulations

In contrast to the previous two methods, cosmological simulations do not detect or classify tidal features, but they are helpful for studying the formation pathways of tidal features and their relationship with host galaxy properties. In this Section, we will discuss several examples of cosmological simulation analyses focused on tidal features or galaxy mergers and how these simulations help support our understanding.

Due to a lack of tidal feature images, many surveys used mock images generated from simulations to train the machine learning model as discussed earlier (e.g. Pearson et al. 2019, Hood et al. 2018, Bottrell et al. 2019). Martin et al. (2018) used the Horizon-AGN hydrodynamical cosmological

simulation to explore the influence of galaxy mergers on galaxy morphology. They focused on massive galaxies over  $10^{10} M_{\odot}$  and found that one-third of the morphological transformations of elliptical galaxies are due to minor mergers, and these become the main driving force after  $z \sim 1$ . The gas content of the merging galaxies determines the morphology of the galaxy merger remnants, gas-rich mergers will lead to the formation of a disk after the merger. ([Lotz et al., 2011](#)) also used cosmological simulations to determine the merger rate of galaxies, and their conclusions were in good agreement with [Martin et al. \(2018\)](#)'s. They compared the merger rates of major and minor mergers and found that at  $z \sim 0.7$ , the merger rate of minor mergers was about three times that of major mergers.

Cosmological simulations are also used to study galaxy mergers in order to study the formation paths resulting in different tidal features. [Karademir et al. \(2019\)](#) concentrated on the influence of infall-orbit configurations on the formation of streams and shells. They used the GADGET-3 simulation ([Springel, 2005](#)) and performed 36 simulations of galaxy mergers between two disc galaxies with different masses. They reached the conclusion that minor mergers can significantly change the morphology of galaxies, while micro mergers (with smaller mass ratios that  $\leq 1 : 10$ ) mainly act to expand the galactic disk. They also found that tidal streams are mainly produced by circular satellite infall, while tidal shells are produced by (nearly) radial satellite infall.

[Pop et al. \(2018\)](#) used the Illustris hydrodynamical cosmological simulation ([Nelson et al., 2015](#)) to study the formation of the shell galaxies. They found that the proportion of galaxies with shells increases with galaxy mass. However, as redshift increases, the number of massive galaxies with shells becomes increasingly rare. They discovered that larger mass ratios are more likely to produce shells after mergers (deviating further from a purely radial infall), and small satellites require almost pure radial infall to form shells, in agreement with [Karademir et al. \(2019\)](#).

[Khalid et al. \(2024\)](#) used four different cosmological hydrodynamical simulations—NEWHORIZON, EAGLE, ILLUSTRISNG, and MAGNETICUM—to generate LSST-like mock images and visually classified the tidal features within them. They found that the fraction of galaxies exhibiting tidal features are similar in different cosmological hydrodynamical simulations:  $40\% \pm 6\%$  in NEWHORIZON,  $37\% \pm 1\%$  in EAGLE,  $32\% \pm 1\%$  in ILLUSTRISNG, and  $32\% \pm 1\%$  in MAGNETICUM. Tidal features are also related to the stellar and halo mass of galaxies. Central galaxies within groups and clusters almost universally display tidal features, whereas only a few satellite galaxies show tidal features.

Simulations can also be combined with real observational data to understand the formation of the tidal features observed in real galaxies and reconstruct the merger history of galaxies. For instance, the goal of [Foster et al. \(2014\)](#) was to explore the cause of the streams and shells of the umbrella galaxy NGC 4651. They obtained image and spectroscopic observations of NGC 4651 and simulated the formation of the stream through N-body simulation. They found that the satellite must have passed through the galactic disk and caused disturbances in the host galaxy. [Martínez-Delgado et al. \(2021\)](#) used numerical simulations combined with images to reconstruct the merger history of the M104 galaxy, which is considered to be the result of a gas-rich merger of two galaxies of comparable mass 3.5 around Gyr ago. They found that the formation of tidal streams on both sides of the galaxy disc was not related to this major gas-rich merger, so they can use the formation of these tidal streams to constrain the time of this major gas-rich merger.

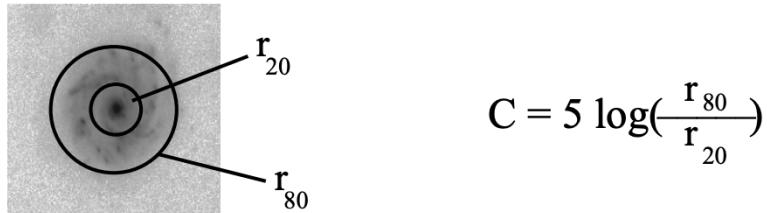
### 4.3 CAS method

Using CAS (concentration, asymmetry, smoothness) parameters is another merger detection method based on galaxy images. This measures the shapes of galaxies by measuring these three parameters ([Conselice et al., 2008](#)) because the light distribution of the galaxies reflects their formation (e.g. [Conselice 2003](#)), CAS can also suggest whether there are traces of mergers in the images. In this Section, we will introduce the specific calculation method of each CAS parameter and the criteria for measuring galaxy mergers. We will also discuss some examples of using this method.

Concentration C is a parameter that measures how light is concentrated in a galaxy. It is typically measured using the ratio of the radii of two circular areas, which respectively contain 20% and 80% of the total light flux of the galaxy. If the central region of a galaxy contains higher ratio of light than another galaxy, then this galaxy corresponds to a higher C value ([Conselice et al., 2008](#)). The calculation of concentration C is:

$$C = 5 \log \left( \frac{r_{80}}{r_{20}} \right). \quad (3)$$

where  $r_{80}$  is the radius of the circular area that contains 80% of the total flux of the galaxy and  $r_{20}$  is the radius of the circular area that contains 20% of the total flux, this index C is also called  $C_{28}$ . The graphical representation of how to calculate concentration parameter C is Figure 6.

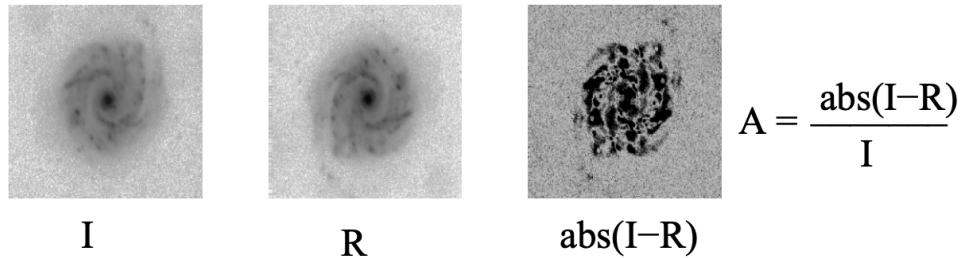


**Figure 6:**  $r_{20}$  is the radius of the circle that contains 20% of the total light flux and  $r_{80}$  is the radius of the circular area that contains 80% of the total flux, C is calculated by the ratio these two radii. Original Figure 3 in [Conselice \(2003\)](#).

Asymmetry A is the parameter that measures the degree of symmetry in a galaxy. It is calculated by subtracting the rotated image of the galaxy from the original image. The rotated image is generated by rotating the galaxy image 180 degrees about the galaxy's central point. The central point of the galaxy is identified by finding the point where the asymmetry is minimized. The residual value between the original and rotated images will be normalized based on the intensity of the original image. The asymmetry A formula ([Conselice et al., 2008](#)) is given by:

$$A = \min \left( \frac{\sum |I_0 - I_{180}|}{\sum |I_0|} \right) - \min \left( \frac{\sum |B_0 - B_{180}|}{\sum |I_0|} \right) \quad (4)$$

where  $I_0$  is the pixels of the original image and the  $I_{180}$  is that of the image rotated by 180 degrees. The second term is added for the correction of the background noise, where  $B_0$  is the light from a blank sky area and the  $B_{180}$  is the rotated version of  $B_0$ . The correction term will be normalized by the original image to keep it the same as the first term. Also, this term is minimized in the same way as the original galaxy image term. The graphical representation of how asymmetry parameters are calculated is given in Figure 7.

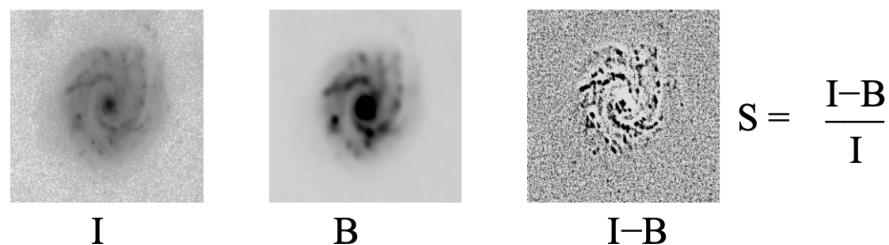


**Figure 7:** The graphical representation of the calculation of asymmetry parameter A. The left image is the original image of a galaxy. The middle shows the image of the left panel rotated by 180 degrees. The right panel is the absolute value of the residuals generated from subtracting the rotated image from the original one. Then, the asymmetry parameter is calculated based on the absolute value of the residuals. Original Figure 3 in [Conselice \(2003\)](#).

The smoothness (also called clumpiness) parameter S measures the proportion of light in a galaxy that is contained in clumpy light concentrations. Smooth systems, like elliptical galaxies, contain light at low spatial frequencies, while clumpy galaxies have a comparatively large amount of light at high spatial frequencies. Therefore, in order to obtain the smoothness parameter, it is necessary to only keep the high spatial frequency part of the galaxy through image processing. [Conselice \(2003\)](#) proposed that a residual map can be obtained by subtracting the blurred image from the original galaxy image, and these residual values contain only the high-frequency part of the galaxy. The residual values are normalized by the original image as same as the calculation of the asymmetry parameter. The most common formula for calculating the smoothness parameter is described by [Conselice et al. \(2008\)](#) as follows:

$$S = 10 \left\{ \left( \frac{\sum(I_{x,y} - I_{x,y}^\sigma)}{\sum I_{x,y}} \right) - \left( \frac{\sum(B_{x,y} - B_{x,y}^\sigma)}{\sum I_{x,y}} \right) \right\} \quad (5)$$

where  $I_{x,y}$  is the original image's light intensity at position  $(x, y)$ , and  $I_{x,y}^\sigma$  is the intensity of blurred image with smoothing filter of size  $\sigma$  at position  $(x, y)$ . Similar to the calculation of the asymmetry parameter, S have the correction for background noise, using the same smoothing kernel  $\sigma$  to blur the background intensity  $B_{x,y}$  to get  $B_{x,y}^\sigma$ . The correction term will also be normalized by the original image. The smoothing kernel  $\sigma$  is determined by the radius of the galaxies  $r$ , as  $\sigma = 0.2 \times 1.5r$  ([Conselice, 2003](#)). Figure 8 shows the graphical representation of the smoothness calculation.



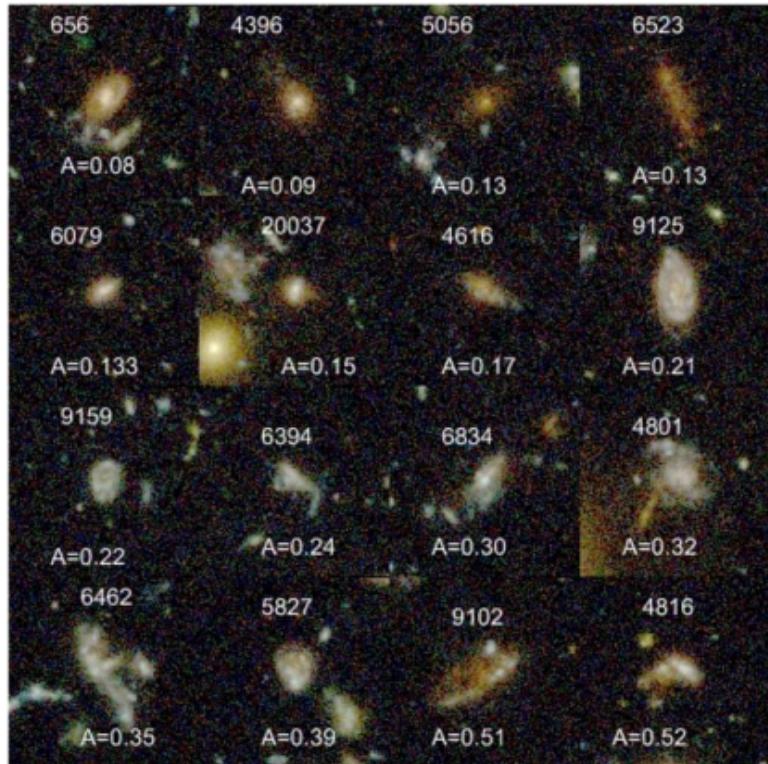
**Figure 8:** The graphical representation of the calculation of smoothness parameter S. The left image is the original image of a galaxy. The middle panel is the image blurred by a smoothing filter of size  $\sigma$ . The right panel shows the absolute value of the residuals generated from subtracting the blurred image from the original one. Then, the smoothness parameter S is calculated based on the absolute value of the residuals. Original Figure 3 in [Conselice \(2003\)](#).

These three different parameters provide us with different information about the galaxy. Typically, the concentration parameter C ranges from 2 to 3 for disc galaxies, exceeds 3.5 for massive elliptical galaxies and spans the entire range for peculiar galaxies. The asymmetry parameter A can provide information about galaxy morphology and mergers. High A values represent asymmetric light distributions, usually found in spiral galaxies or galaxy mergers. Typically,  $A = 0 - 0.05$  represents elliptical galaxies, while  $A = 0.1 - 0.3$  represents disk galaxies or irregular galaxies.  $A > 0.3$  usually indicates galaxies involved in mergers (Conselice et al., 2008). Conselice et al. (2000) found that  $A > 0.35$  locates nearly all major galaxy mergers in the nearby universe. Past studies (e.g. Bershady et al. 2000, Courteau et al. 2007) have shown that concentration is highly correlated with galaxy mass and morphology. The smoothness parameter S can be used to measure star formation activity in galaxies. Typically, galaxies with  $S < 0.1$  are non-star-forming, such as elliptical galaxies, while star-forming galaxies, such as discs and irregulars, have  $S = 0.1 - 1$ .

For the measurement of galaxy mergers that this project is concerned with, Conselice (2003) gave the following criteria:

$$A > 0.35 \quad \text{and} \quad A > S \quad (6)$$

The first condition is based on the properties of the asymmetry parameter as mentioned before: the asymmetry parameter should be higher than the limit found in local mergers. The second condition  $A > S$  ensures that the asymmetric light is not dominated by clumpy star-forming regions of the galaxies. Figure 9 shows some examples of galaxies with their corresponding value of asymmetry parameter.



**Figure 9:** Example galaxy images with their corresponding asymmetry parameters. Plotted on the top of each image is the ID number from Coe et al. (2006), and the bottom number is the computed value of the asymmetry parameter. Original Figure 11 in Conselice et al. (2008).

## 5 This project

Galaxy mergers play an important role in a galaxy's evolution. By understanding galaxy mergers, we can further investigate the evolution of galaxies and the formation of universe structures. The detection and study of tidal features are crucial for understanding galaxy mergers, but visual classification and detection methods for tidal features are relatively inefficient, particularly for the next generation of imaging surveys in the pipeline. With the massive influx of galaxy images about to come from LSST, traditional visual classification is no longer feasible.

To handle a large number of galaxy images from the upcoming LSST and automatically identify low surface brightness tidal features, [Desmons et al. \(2024\)](#) have developed an automated classification model based on self-supervised learning. This model is designed to be applied to LSST-like images and has been tested on existing datasets such as HSC-SSP and SDSS, showing promising results. However, it needs to be further tested to verify its performance, which is the key focus of this project.

The primary aim of this project is to evaluate the performance of automatic methods for identifying tidal features and to compare and validate these against existing methods. Thereby investigating the factors influencing the classification of tidal features. [Khalid et al. \(2024\)](#) have used cosmological simulations to generate mock LSST-like images of galaxies. These mock images are the main image source of this project. These have been adjusted to match the input format of [Desmons et al. \(2024\)](#)'s model. The details of mock images will be introduced in Section 6.

In this project, I will evaluate [Desmons et al. \(2024\)](#)'s self-supervised representation learning (SSL) algorithm. I first need to visually classify the mock images into different confidence levels of different tidal features. To test the performance of the SSL model, I will apply [Desmons et al. \(2024\)](#)'s model on the mock images to get its classification result. I will then compare the machine learning (ML) results with those from visual classification. By analyzing the similarities and differences between the visual and ML classifications for the range of galaxy properties, We will gain a deeper understanding of the algorithm's parameters that affect the classification of tidal features, as well as any limitations regarding the galaxy samples to which it can be applied.

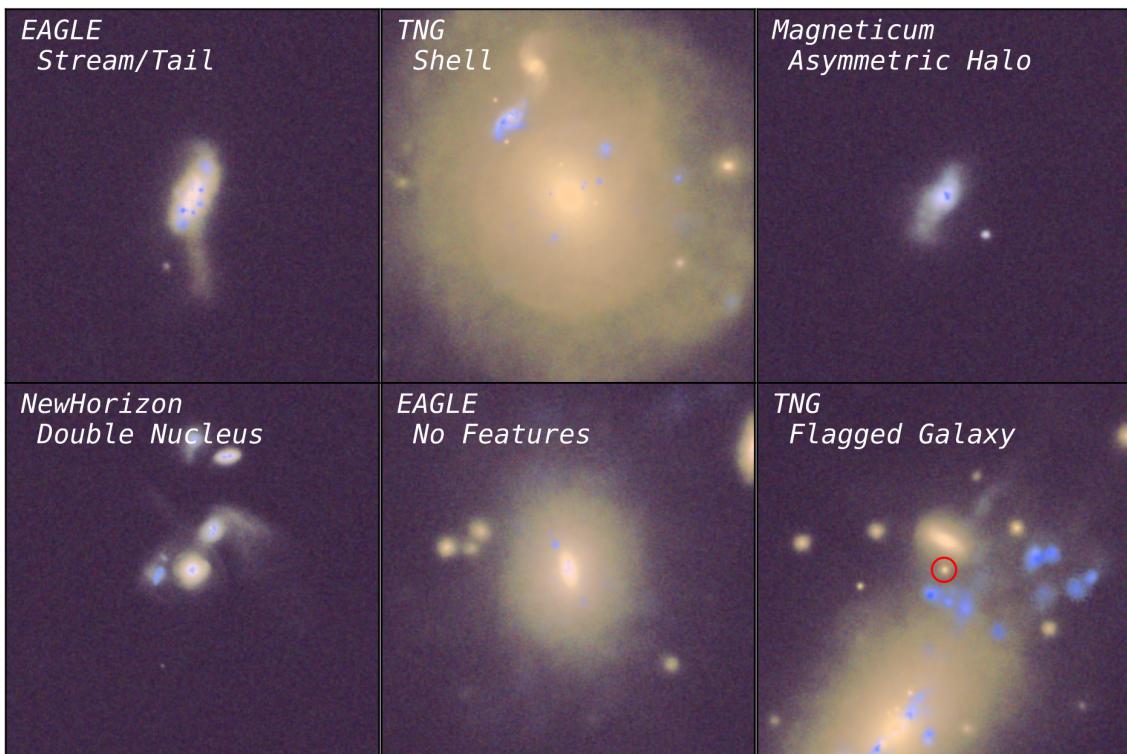
In addition, we will explore where the galaxies in these mock images sit in the CAS (concentration, asymmetry, smoothness) parameter space ([Bershady et al., 2000](#)). We will compare the results from using CAS parameters with those from both self-supervised learning and visual classification to assess the relationship of CAS parameters with these other two identifications of tidal features. This will allow us to determine whether comparisons can be made between results obtained using these different approaches to finding galaxy mergers.

## 6 Data sources and models

The data for this Honours project will be the mock images generated by Khalid et al. (2024) from Illustris hydrodynamical cosmological simulations (Nelson et al., 2015). The tidal-feature detecting model that will be used in this project is the self-supervised machine learning model developed by Desmons et al. (2024). We will develop our own code to measure the CAS parameters of the galaxies in the mock images.

### 6.1 Mock images

We will use mock images generated from cosmological simulations to test the classification of the machine-learning algorithm from Desmons et al. (2024). The mock images of 1826 galaxies have already been generated from the “IllustrisTNG” (Springel et al., 2017) cosmological simulation by Khalid et al. (2024). Example mock images are shown in Figure 10. These galaxies have stellar masses  $M_*$  ranging from  $3.16 \times 10^9 M_\odot$  to  $6.50 \times 10^{11} M_\odot$ , with an average mass of  $2.67 \times 10^{10} M_\odot$ .



**Figure 10:** Examples of Mock images from 4 different hierarchical cosmological simulations including IllustrisTNG showing different tidal features. Original Figure 1 in Khalid et al. (2024). In this project, we will only use the TNG mock images due to the limited time available.

It is important to note that the mock images visually classified in Khalid et al. (2024) are different from the mock images that will be inputted into the self-supervised learning algorithm in this project. In Khalid et al. (2024), the mock images were placed close by, at a distance of  $z \sim 0.025$  (average  $\sim 105\text{Mpc}$ ) and were sized 2400 pixels  $\times$  2400 pixels, as shown in Figure 10. As these images are produced to be LSST-like, the arcsecond-to-pixel scale is  $0.2''/\text{pixel}$ . However, the Desmons et al. (2024) model was developed for galaxies identified by GAMA (Galaxy And Mass Assembly; Driver et al. 2011). This sample has a median distance of  $z \sim 0.2$ , and the input image size of Desmons et al. (2024)’s model is  $128 \times 128$  pixels, so  $25.6 \times 25.6''$ , designed to be consistent

with the types of images that LSST will provide. So, [Khalid et al. \(2024\)](#) re-made the mock images applying a distance of  $z \sim 0.2$ , and image size of  $128 \times 128$  pixels.

To construct LSST-like mock images, [Khalid et al. \(2024\)](#) produced mock images to match the  $0.2''/\text{pixel}$  spatial resolution and expected 10-year surface brightness limits of the LSST:  $\mu_g \sim 30.3 \text{ mag arcsec}^{-2}$ ,  $\mu_r \sim 30.3 \text{ mag arcsec}^{-2}$ ,  $\mu_i \sim 29.7 \text{ mag arcsec}^{-2}$ ,  $\mu_z \sim 28.4 \text{ mag arcsec}^{-2}$ , and  $\mu_y \sim 28.1 \text{ mag arcsec}^{-2}$  ([Yoachim, 2024](#)). The mock images used in [Khalid et al. \(2024\)](#) only include *gri*-bands but to match the input image format of [Desmons et al. \(2024\)](#)'s model, he re-made mock images include the 5 HSC bands (*grizy*). The re-made *grizy* images are made to  $128 \times 128$  pixels which is roughly equivalent to  $87 \text{ kpc} \times 87 \text{ kpc}$  at  $z \sim 0.2$  (using cosmological parameters of Plank2015 ([Ade et al., 2016](#))). [Khalid et al. \(2024\)](#) used Friends-Of-Friends to extract all relevant star particles in a 1Mpc cube centered on the simulated galaxy and calculated the spectral energy distribution of the stars. They then smoothed those pixels containing only a few star particles and generated mock images by collapsing the cube along one of its axes. The mock images are then rescaled to match the LSST pixel size and Gaussian noise is added to apply the surface brightness limit.

Figure 11 shows an example of the re-made mock images including different tidal features.

Since the mock images have been re-made and the distance and resolution both affect the identification of tidal feature details (e.g. [Martin et al. 2022](#)). In order to better compare the machine learning classification results with visual classifications, I have visually re-classified the mock images. I used the same classification criteria as [Khalid et al. \(2024\)](#), which follows a simplified version of the [Bilek et al. \(2020\)](#) scheme. Tidal features are classified into four categories:

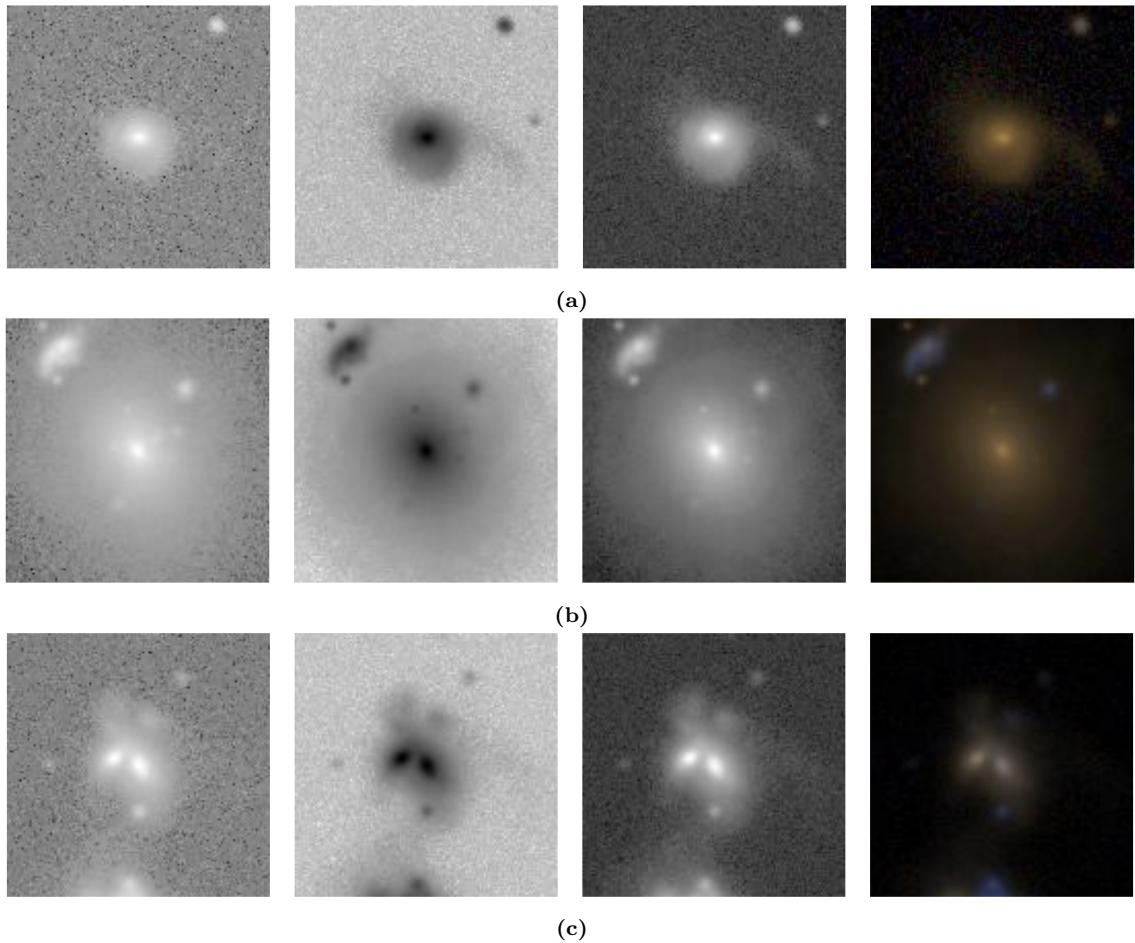
- **Streams/Tails:** Prominent, elongated structures orbiting or expelled from the host galaxy. These usually have similar colours to the host galaxy.
- **Shells:** Concentric radial arcs or ring-like structures around a galaxy.
- **Asymmetric Stellar Halos:** Diffuse features in the outskirts of the host galaxy, lacking well-defined structures like stellar streams or tails.
- **Double nuclei:** Galaxies that are visibly merging but where both objects are still clearly separated.

Table 1 shows the classification Confidence levels together with their corresponding descriptions.

Confidence level	Description
0	No tidal feature detected.
1	Hint of tidal feature detected, classification difficult.
2	Even chance of correct classification of tidal feature presence and/or morphology.
3	High likelihood of the tidal feature being present and morphology being obvious.

**Table 1:** Descriptions of the classification confidence levels ([Khalid et al. 2024](#)).

The IllustrisTNG data, produced by hierarchical cosmological simulations, is stored in the UNSW supercomputer Katana. Katana is a shared computational cluster located on the UNSW campus, designed to provide access to computational resources for groups working with non-sensitive data



**Figure 11:** Examples of tidal features in the mock images re-made at a distance of  $z \sim 0.025$  and  $128 \times 128$  pixels. The four images from left to right in each row are generated according to *g*-band grayscale, 5-band grayscale, 5-band inverse grayscale, and *gri*-colour respectively. Images a, b, and c show different galaxies, 11a is classified as having a Tail (confidence level 3) and Asymmetric Halo (confidence 3). 11b is classified as having a Shell (confidence 3), 11c is classified as having a Double Nucleus (confidence 3) and Asymmetrical Halo (confidence 3). The description of the confidence levels is given in Table 1.

([Katana, 2024](#)).

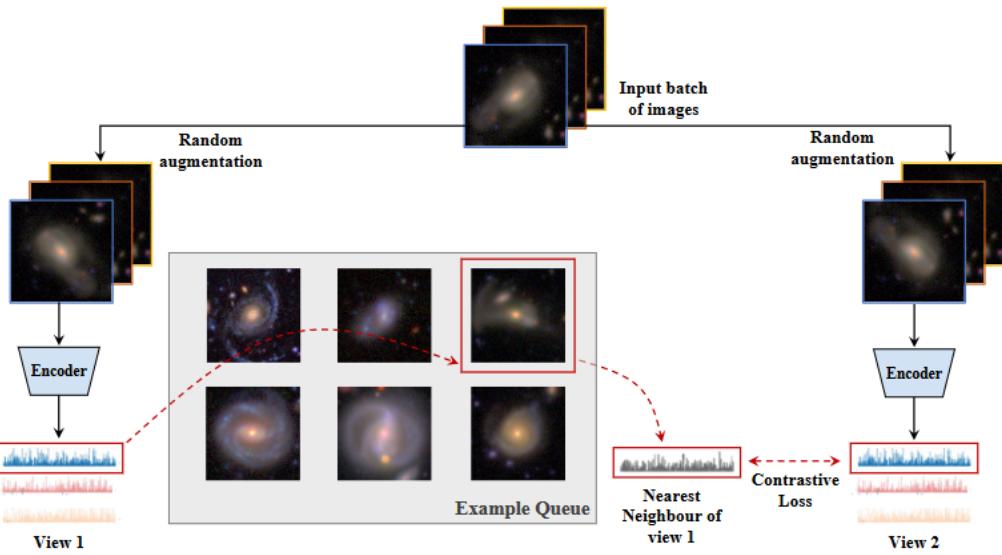
## 6.2 Machine learning model

The machine learning model from [Desmons et al. \(2024\)](#) consists of two components: a self-supervised model used for pre-training, and a linear classifier used for classification.

### 6.2.1 Self-Supervised Model

The self-supervised model uses a type of self-supervised learning called Nearest Neighbour Contrastive Learning of visual Representations (NNCLR; [Dwibedi et al. 2021](#)), Figure 12 shows the structure of the NNCLR used here.

Self-supervised learning models create various augmented versions of an image, then the model pairs images and categorizes them as positive or negative pairs. For instance, if  $x_i$  is the original image, for any pair  $(x_i, x_j)$ , if  $x_j$  is an augmented version of  $x_i$ , then the pair is labeled as a positive pair. Conversely, if  $x_j$  is not an augmented version of  $x_i$ , then it is a negative pair.



**Figure 12:** The schematic of Nearest Neighbour Contrastive Learning of Visual Representations. Original Figure 2 in Desmons et al. (2024)

For every image in the sample, the encoder networks in the model generate a 128-dimensional representation. This encoder is trained to produce similar representations for positive pairs and dissimilar representations for negative pairs, allowing for the clustering of similar samples while distancing dissimilar ones.

The encoder uses a contrastive loss function to create the representations for each image (Eqn. 7), but this contrastive learning has problems since it is based on the different augmented versions of the same image to create positive pairs, the objects with large differences that still belong to the same class (such as galaxy with different tidal features) will not be linked.

$$L_i = -\log \left( \frac{\exp(\text{sim}(z_i, z_i^+))}{\exp(\text{sim}(z_i, z_i^+)) + \sum_{z^-} \exp(\text{sim}(z_i, z^-))} \right) \quad (7)$$

To address this problem, NNCLR does not simply create positive and negative pairs based on whether they are augmented versions of the original image. Instead, it creates a queue of samples and uses whether they are the nearest neighbors in the queue to define the positive and negative pairs for a given image. The contrastive loss function of NNCLR is:

$$L_i^{NNCLR} = -\log \left( \frac{\exp(\text{NN}(z_i, Q) \cdot z_i^+ / \tau)}{\sum_k \exp(\text{NN}(z_i, Q) \cdot z_k^+ / \tau)} \right) - \log \left( \frac{\exp(\text{NN}(z_k, Q) \cdot z_i^+ / \tau)}{\sum_k \exp(\text{NN}(z_k, Q) \cdot z_i^+ / \tau)} \right) \quad (8)$$

where  $Q$  is the queue of the sample and  $\text{NN}$  is the nearest neighbour operator that:

$$\text{NN}(z, Q) = \arg \min_{i \in Q} \|z - i\|_2 \quad (9)$$

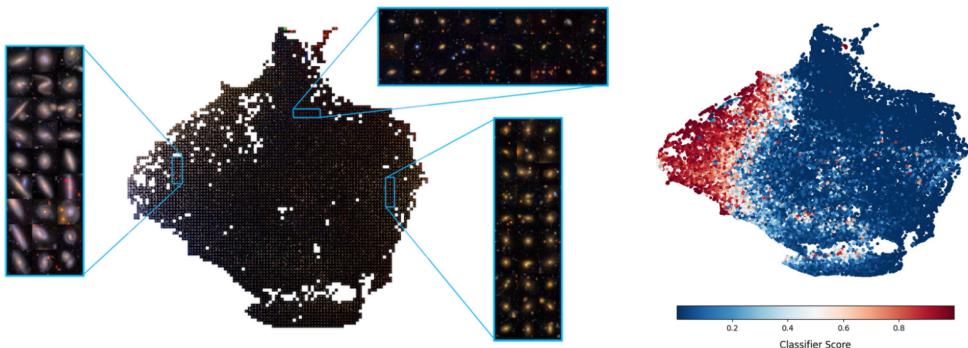
where  $\|x\|_{l_2}$  is  $l_2$ -normalisation of  $x$ , which is:

$$\|x\|_2 = \sqrt{\sum_{k=1}^n |x_k|^2} \quad (10)$$

This model uses ResNet-50 as the encoder followed by a global pooling layer. After that, there are two 128-sized fully connected layers with L2 kernel regularisation and a penalty of 0.0005. A batch-normalisation layer follows each fully connected layer and ReLu activation follows the batch-normalisation layer ([Desmons et al., 2024](#)).

### 6.2.2 The Linear Classifier

The training of the self-supervised model does not require labels for data, the encoder simply transforms images into meaningful low-dimensional representations based on their augmentations. However, since we need to classify galaxies based on tidal features, [Desmons et al. \(2024\)](#) has trained a linear classifier using a small visually-classified dataset ([Desmons et al., 2024](#)) and the encoder to classify the images. The encoded representations from the Self Supervised Model are sent to a fully connected layer with a sigmoid activation, which outputs a single number between 0 and 1 to measure the likelihood of a single image having tidal features. In this project, we will compare the output of the linear classifier with the visual classification.



**Figure 13:** Left figure is the 2D Uniform Manifold Approximation and Projection (UMAP) projection of the representations from the self-supervised model. Made by binning the space into  $100 \times 100$  cells and randomly selecting a sample from that cell to plot in the corresponding cell location. The right figure is the same 2D UMAP projection without binning, coloured according to the scores assigned to each galaxy by the linear classifier. Original Figure 8 in [Desmons et al. \(2024\)](#)

## 7 Thesis contents

Abstract

Contents

1. Introduction
2. Background Theory (Section 2-5 of this literature review)
  - (a) Structure of Universe
  - (b) Galaxy mergers
  - (c) Detection of galaxy mergers
3. Data (Section 6.1 of this literature review)
4. Methods
  - (a) Visual Classification
  - (b) Machine learning
  - (c) CAS parameters
5. Results
6. Discussion
7. Conclusion

Acknowledgments

## References

- Abbott T. M. C., et al., 2019, *The Astrophysical Journal*, 872, L30
- Ade P. A. R., et al., 2016, *Astronomy & Astrophysics*, 594, A13
- Almeida A., et al., 2023, *Astrophysical Journal Supplement Series*, 267, 44
- Atkinson A. M., Abraham R. G., Ferguson A. M. N., 2013, *The Astrophysical Journal*, 765, 28
- Barnes J. E., 1992, *The Astrophysical Journal*, 393, 484
- Barnes J. E., Hernquist L., 1996, *The Astrophysical Journal*, 471, 115
- Bennett C. L., et al., 2013, *The Astrophysical Journal Supplement Series*, 208, 20
- Bershady M. A., Jangren A., Conselice C. J., 2000, *The Astronomical Journal*, 119, 2645
- Bickley R. W., et al., 2021, *Monthly Notices of the Royal Astronomical Society*, 504, 372
- Bílek M., et al., 2020, *Monthly Notices of the Royal Astronomical Society*, 498, 2138
- Blumenthal G. R., Faber S. M., Primack J. R., Rees M. J., 1984, *Nature*, 311, 517
- Bottrell C., et al., 2019, *Monthly Notices of the Royal Astronomical Society*, 490, 5390
- Bridge C. R., Carlberg R. G., Sullivan M., 2010, *The Astrophysical Journal*, 709, 1067
- Cheng T.-Y., et al., 2020, *Monthly Notices of the Royal Astronomical Society*, 493, 4209
- Coe D., Benítez N., Sánchez S. F., Jee M., Bouwens R., Ford H., 2006, *The Astronomical Journal*, 132, 926
- Conselice C. J., 2003, *The Astrophysical Journal Supplement Series*, 147, 1
- Conselice C. J., Bershady M. A., Gallagher J. S., 2000, *Astronomy and Astrophysics*, 354
- Conselice C. J., Rajgor S., Myers R., 2008, *Monthly Notices of the Royal Astronomical Society*, 386, 909
- Courteau S., McDonald M., Widrow L. M., Holtzman J., 2007, *The Astrophysical Journal*, 655, L21
- Crain R. A., et al., 2015, *Monthly Notices of the Royal Astronomical Society*, 450, 1937
- Dalal N., White M., Bond J. R., Shirokov A., 2008, *The Astrophysical Journal*, 687, 12
- Desmons A., Brough S., Martínez-Lombilla C., De Propris R., Holwerda B., López-Sánchez A. R., 2023, *Monthly Notices of the Royal Astronomical Society*, 523, 4381
- Desmons A., Brough S., Lanusse F., 2024, *Monthly Notices of the Royal Astronomical Society*, 531
- Driver S. P., et al., 2011, *Monthly Notices of the Royal Astronomical Society*, 413, 971
- Dubois Y., et al., 2021, *Astronomy and Astrophysics*, 651, A109

- Dwibedi D., Aytar Y., Tompson J., Sermanet P., Zisserman A., 2021, With a Little Help from My Friends: Nearest-Neighbor Contrastive Learning of Visual Representations, doi:10.48550/arXiv.2104.14548, <https://arxiv.org/abs/2104.14548>
- Foster C., et al., 2014, *Monthly Notices of the Royal Astronomical Society*, 442, 3544
- Hendel D., Johnston K. V., 2015, *Monthly Notices of the Royal Astronomical Society*, 454, 2472
- Hocking A., Geach J. E., Sun Y., Davey N., 2017, *Monthly Notices of the Royal Astronomical Society*, 473, 1108
- Hood C. E., Kannappan S. J., Stark D. V., Dell'Antonio I. P., Moffett A. J., Eckert K. D., Norris M. A., Hendel D., 2018, *The Astrophysical Journal*, 857, 144
- Hopkins P. F., Hernquist L., Cox T. J., Di Matteo T., Martini P., Robertson B., Springel V., 2005, *The Astrophysical Journal*, 630, 705
- Huertas-Company M., Lanusse F., 2023, *Publications of the Astronomical Society of Australia*, 40, e001
- Ivezić Ž., et al., 2019, *The Astrophysical Journal*, 873, 111
- Kado-Fong E., et al., 2018, *The Astrophysical Journal*, 866, 103
- Karademir G. S., Remus R.-S., Burkert A., Dolag K., Hoffmann T. L., Moster B. P., Steinwandel U. P., Zhang J., 2019, *Monthly Notices of the Royal Astronomical Society*, 487, 318
- Katana 2024, Katana — UNSW Research, <https://research.unsw.edu.au/katana>
- Khalid A., Brough S., Martin G., Kimmig L. C., Lagos C. D. P., Remus R. S., Martinez-Lombilla C., 2024, *Monthly Notices of the Royal Astronomical Society*, 530, 4422
- Lacey C. G., Cole S., 1993, *Monthly Notices of the Royal Astronomical Society*, 262, 627
- Lin L., et al., 2004, *The Astrophysical Journal*, 617, L9
- Lin L., et al., 2008, *The Astrophysical Journal*, 681, 232
- López-Sanjuan C., et al., 2011, *Astronomy and Astrophysics*, 530, A20
- Lotz J. M., Jonsson P., Cox T. J., Croton D., Primack J. R., Somerville R. S., Stewart K., 2011, *The Astrophysical Journal*, 742, 103
- Martin G., Kaviraj S., Devriendt J. E. G., Dubois Y., Pichon C., 2018, *Monthly Notices of the Royal Astronomical Society*, 480, 2266
- Martin G., Kaviraj S., Hocking A., Read S. C., Geach J. E., 2020, *Monthly Notices of the Royal Astronomical Society*, 491, 1408
- Martin G., et al., 2022, *Monthly Notices of the Royal Astronomical Society*, 513, 1459
- Martínez-Delgado D., et al., 2021, *Monthly Notices of the Royal Astronomical Society*, 506, 5030
- Mihos J. C., Hernquist L., 1996, *The Astrophysical Journal*, 464, 641
- Naab T., Johansson P. H., Ostriker J. P., 2009, *The Astrophysical Journal*, 699, L178

- Nelson D., et al., 2015, [Astronomy and Computing](#), 13, 12–37
- Patton D. R., et al., 2002, [The Astrophysical Journal](#), 565, 208
- Pearson W. J., Wang L., Trayford J. W., Petrillo C. E., van der Tak F. F. S., 2019, [Astronomy & Astrophysics](#), 626, A49
- Perlmutter S., et al., 1999, [The Astrophysical Journal](#), 517, 565
- Pop A.-R., Pillepich A., Amorisco N. C., Hernquist L., 2018, [Monthly Notices of the Royal Astronomical Society](#), 480, 1715
- Riess A. G., et al., 1998, [The Astronomical Journal](#), 116, 1009
- Robotham A. S. G., et al., 2014, [Monthly Notices of the Royal Astronomical Society](#), 444, 3986
- Sánchez H. D., et al., 2023, [Monthly Notices of the Royal Astronomical Society](#), 521, 3861
- Santucci G., et al., 2024, [Monthly Notices of the Royal Astronomical Society](#), 528
- Springel V., 2005, [Monthly Notices of the Royal Astronomical Society](#), 364, 1105
- Springel V., Di Matteo T., Hernquist L., 2005, [The Astrophysical Journal](#), 620, L79
- Springel V., et al., 2017, [Monthly Notices of the Royal Astronomical Society](#), 475, 676
- Tal T., Dokkum v., Nelan J. E., Bezanson R., 2009, [The Astronomical Journal](#), 138, 1417
- Thorne K. S., 1967, [The Astrophysical Journal](#), 148, 51
- Walmsley M., A., Mann R., Lintott C., 2018, [Monthly Notices of the Royal Astronomical Society](#), 483, 2968
- Walmsley M., et al., 2019, [Monthly Notices of the Royal Astronomical Society](#), 491, 1554
- Walmsley M., et al., 2024, Scaling Laws for Galaxy Images, <https://arxiv.org/abs/2404.02973>
- Wechsler R. H., Tinker J. L., 2018, [Annual Review of Astronomy and Astrophysics](#), 56, 435
- White S. D. M., Frenk C. S., 1991, [The Astrophysical Journal](#), 379, 52
- Wilson R. W., Penzias A. A., 1967, [Science](#), 156, 1100
- Yoachim P., 2024, Surface Brightness Limit Derivations, <https://smtn-016.lsst.io/>
- Zhu J.-Y., Krähenbühl P., Shechtman E., Efros A. A., 2015, [arXiv \(Cornell University\)](#)