# Generative Adversarial Networks for Galaxy Classification Using Convolutional Neural Networks

Polina Petrova
ppetrova@umass.edu

Aryan Singh
arysingh@umass.edu

## Abstract

*The classification of galaxies poses significant challenges due to inherent data imbalances among various morphological classes. This research addresses the pressing issue of accurately classifying galaxies using Convolutional Neural Networks (CNNs), where traditional approaches often struggle with underrepresented categories. To mitigate this problem, we propose a novel combined model utilising a Conditional Generative Adversarial Network (cGAN) to generate synthetic images that augment the Galaxy Zoo 2 dataset, which contains approximately 250,000 labelled galaxy images. Our methodology begins with enhancing image resolution using Super Resolution GANs (SRGAN), followed by training the cGAN to produce additional samples for each class, thereby addressing class imbalance without resorting to deeper networks. We draw on foundational works, including studies on deep generative models for galaxy image simulations and prior CNN applications in galaxy classification, to inform our approach. Our model's performance will be rigorously evaluated using metrics such as accuracy, precision, recall, and F-1 score, comparing results from training on both real and synthetic data against traditional CNNs trained on the original dataset. Through this work, we aim to contribute to the ongoing efforts to improve automated galaxy classification and provide a scalable solution to the prevalent issue of data imbalance in astronomical research.*

## 1. Introduction

With rapid development in technology comes a phenomenon that is referred to as the *data deluge*, – a consequence of an overwhelming influx of data and insufficient resources to process it. Astronomers, in particular, contend with this challenge; estimates suggest the observable universe holds between 100 billion and 200 billion galaxies [2]. A single photograph of a small sky section can capture up to 25,000 galaxies, and the resulting daily data volume overwhelms the limited pool of experts available to classify them [7]. To address this challenge, astrophysicists launched the *Galaxy Zoo* project in 2007, inviting citizen scientists to help classify over 900,000 galaxies, marking a transformative moment in data processing through public participation [6].

Following the success of the *Galaxy Zoo*, advancements in machine learning spurred efforts to automate galaxy classification. Modern approaches employ Convolutional Neural Networks (CNNs) to recognise patterns with minimal human input. However, a persistent challenge is the quality and distribution of available data. While the *Galaxy Zoo* project provided a substantial dataset, imbalances in class representation can skew training, causing models to favour more frequent classes. Our research specifically addresses this class imbalance in the Galaxy Zoo 2 dataset – a collection of categorised images taken from the Sloan Digital Sky Survey (SDSS) – where certain galaxy types are underrepresented. Building deeper networks is a common workaround to address this issue, but we argue that this approach only sidesteps the core problem.

Ideally, a large, balanced dataset would improve model accuracy, but limitations in space imaging and classification complexity make this difficult. To tackle this, we propose a novel solution using Generative Adversarial Networks (GANs) to generate synthetic images for underrepresented galaxy classes in the Galaxy Zoo 2 dataset, which we then use to train a CNN. We expect that by augmenting our data with synthetic images, we can improve classification accuracy across all classes. Our evaluation will compare the GAN-CNN model's performance against traditional CNNs trained on imbalanced data, particularly examining accuracy gains in classifying underrepresented classes. With this combined GAN-CNN model, we aim to address class imbalance directly, eliminating the need for deeper networks as a compensatory measure.

## 2. Related work

A literature survey of past work on this topic. Introduce the baselines you will compare to and the weakness you plan to address. *This section should be nearly complete.*

Lahav et al. (1996) offers one of the first discussions of neural networks in the galaxy classification problem.

This study clarifies the role of Artificial Neural Networks (ANNs) in galaxy classification by demonstrating their ability to replicate human classification using ESO-LV galaxy data. ANNs achieve comparable accuracy to human experts, operating within 2 T-type units. The authors argue that ANNs provide a robust statistical framework, improving on linear methods through their capacity for non-linear modelling. While the paper does not cover all classification methods, it emphasises the potential of unsupervised algorithms to discover new features in galaxy data without external guidance. It also highlights the importance of integrating dynamic properties and multiwavelength data to enhance the classification process, as we now have in the Galaxy Zoo 2 dataset.

Sharma et al. (2019) proposed a comparison study between classic machine learning algorithms, such as ANNs and Random Forests (RFs) against a deep CNN. Their implementation of CNN resulted in better acccuracry than traditional methods – at 89 per cent for stellar spectral class prediction. An issue cited was the need for a large training set to implement an effective deep network, which they accounted for by using an auto-encoder and more convolution layers.

Lanusse et al. presented a framework combining deep learning and physical modelling to enhance galaxy morphology analysis, decoupling galaxy structure from observational noise and PSF effects. This hybrid approach included a conditional generative model that produced accurate galaxy features like size and ellipticity, outperforming traditional models in capturing complex details. Beyond image simulations, such generative models also offer a data-driven solution for building signal priors, valuable for inverse imaging problems like denoising, deconvolution, and deblending.

Walmsley et al. (2020) address the limitations of previous deep learning approaches to galaxy morphology prediction, which often ignore uncertainty or rely on confidently labelled data. The authors propose a Bayesian CNN model that utilises a generative model of Galaxy Zoo volunteer responses, allowing for detailed morphological analysis with sparse labels (about 10 responses per galaxy). Their method incorporates Monte Carlo Dropout to provide well-calibrated predictions of expected volunteer responses, enhancing insights into the connections between morphology and astrophysical phenomena. Additionally, the authors introduce an active learning strategy, iteratively selecting the most informative galaxies for human labelling using a custom acquisition function based on BALD. This approach significantly improves performance compared to random selection.

## 3. Method

Describe the methods you intend to apply to solve the given problem.

## 4. Results

State and evaluate your results upto the milestone.

## 5. Conclusion

State your conclusions upto the milestone. This section might be empty for now but your final report should contain the conclusions of your project.

## References

[1] M. Walmsley et al. Galaxy zoo: probabilistic morphology through bayesian cnns and active learning. *Monthly Notices of the Royal Astronomical Society*, 491(2):1554–1574, 2019.

[2] E. Howell. How many galaxies are there?, 2018. Accessed: Oct. 30, 2024. 1

[3] O. Lahav, A. Nairn, L. Sodré, and M. C. Storrie-Lombardi. Neural computation as a tool for galaxy classification: methods and examples. *Monthly Notices of the Royal Astronomical Society*, 283(1):207–221, 1996.

[4] F. Lanusse, R. Mandelbaum, S. Ravanbakhsh, C.-L. Li, P. Freeman, and B. Póczos. Deep generative models for galaxy image simulations. *Monthly Notices of the Royal Astronomical Society*, 504(4):5543–5555, 2021.

[5] V. Lukic and M. Brüggen. Galaxy classifications with deep learning. *Proceedings of the International Astronomical Union*, 12(S325):217–220, 2016.

[6] C. McGourty. Scientists seek galaxy hunt help, 2007. Accessed: Oct. 30, 2024. 1

[7] E. Sauers. Webb telescope reveals more galaxies in a snapshot than hubble's deepest survey, 2023. Accessed: Oct. 30, 2024. 1

[8] K. Sharma, A. Kembhavi, T. Sivarani, S. Abraham, and K. Vaghmare. Application of convolutional neural networks for stellar spectral classification. *Monthly Notices of the Royal Astronomical Society*, 491(2):2280–2300, 2019.