



UNIVERSITY
OF AMSTERDAM



MSc PHYSICS AND ASTRONOMY

TRACK: GRAVITATION AND ASTROPARTICLE PHYSICS (GRAPPA)

MASTER'S THESIS

Neural Subhalo Density Estimation

Probabilistic Image Segmentation for Strong Gravitational Lensing

by

ELIAS DUBBELDAM

13418505

July 2022

60 ECTS

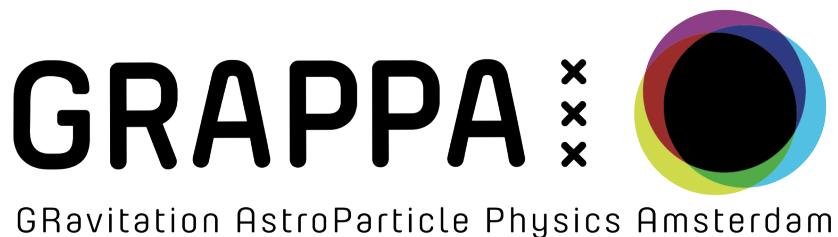
September 2021 - July 2022

Supervisor:

Dr. Christoph Weniger

Second Examiner:

Prof. Dr. Gianfranco Bertone



Abstract

Strong gravitational lensing images provide a pristine probe of the nature of dark matter. To date, only a few heavy individual subhalos have been detected by analyzing their imprint on the image. Current analyses are likelihood-based, and both time-consuming and computationally challenging. Therefore, compromises are required, such as ignoring other subhalos and adopting a particular form of the noise and source realization. Other methods that are including multiple subhalos do not exploit all the information coming from the observation but require compressing it into a summary statistic. In this work, we aim at inferring the parameters of multiple subhalos directly from the full image. To achieve this, we estimate the subhalo density in space and mass instead of the individual subhalo parameters. We combine marginal neural ratio estimation (MNRE), a neural simulation-based inference technique, with probabilistic image segmentation. MNRE enables us to directly estimate marginal posteriors of an arbitrary complex model by training a neural network on data simulated from the model. With probabilistic image segmentation, we estimate the probability of a subhalo being in a particular pixel. With this approach, we determine the subhalo density in mock observations of gravitational lensing systems. We verify the accuracy of the method on single subhalo measurements, and discuss how the network resolves high-mass subhalos but does not pick up the signal from low-mass subhalos when multiple subhalos are modeled. Furthermore, we show that with calibration we are able to correctly interpret the density predictions as probabilities. These results suggest that our method is a powerful approach to measurements of multiple subhalos in a single strong lensing system.

Acknowledgments

First and foremost I would like to express my gratitude to Dr. Christoph Weniger for his guidance, advice, and encouragement. Enthusiastically, he has come up with an overflow of ideas to work on. I would also like to thank Noemi Anau Montel for her tireless supervision, insights, and instant replies that helped me greatly; Benjamin Kurt Miller for his answers to questions about `swyft`; Dr. Adam Coogan and Kosio Karchev for their discussions about lensing, the other members of the unDark group and GRAPPA for their inspiring research environment; and finally my fellow students for making the time spent during long sessions in the library, and outside of it, a lot more fun.

This thesis was carried out on the Lisa Compute Cluster at SURFsara. We acknowledge the use of `astropy` (The Astropy Collaboration et al., 2013), `jupyter` (Kluyver et al., 2016), `matplotlib` (Hunter, 2007), `numpy` (Harris et al., 2020), `plotly` (Plotly Technologies Inc., 2015), `pytorch` (Paszke et al., 2019), `scikit-learn` (Pedregosa et al., 2018) and `tqdm` (da Costa-Luis et al., 2022).

Contents

1	Introduction	2
2	Modeling subhalos in strong lensing systems	9
2.1	Main lens	10
2.2	Subhalos	11
2.3	Source	13
2.4	Instrumentation	13
3	Statistical concepts	15
3.1	Likelihood-free inference	15
3.1.1	Simulation-based inference	16
3.1.2	Marginal Neural Ratio Estimation	16
3.2	Image Segmentation	18
3.2.1	U-Net	18
4	Subhalo Density Estimation	21
4.1	Pixel posterior probability	22
4.2	Training pipeline	24
4.2.1	Simulating	24
4.2.2	Pixelating	24
4.2.3	Mapping	24
4.3	Evaluation	25
4.3.1	Normalization	26
4.3.2	Validation & calibration	26
5	Results	28
5.1	Mock HST observation	28
5.2	Single subhalo inference	29
5.3	Multiple subhalo inference	32
6	Discussion & conclusion	37
6.1	Conclusion	37
6.2	Discussion and future work	38
Appendices		42
A.1	Additional subhalo density predictions	42
A.2	Normalizing during training	47

Chapter 1

Introduction

Dark matter

Our Universe is intriguing, with its vast amount of exotic processes we observe. Cosmologists gave themselves the daunting task to describe it with one single model. The lambda cold dark matter (Λ CDM) model, the standard model of cosmology, is successfully explaining many cosmological observations (Dodelson & Schmidt, 2021). According to it, our Universe consists of three key components: baryonic matter, dark matter, and dark energy. The model describes how in the early Universe, during inflation, quantum fluctuations seeded the structure formation. These fluctuations were amplified by gravity, and matter was pulled from under-dense regions towards over-dense regions, making the over-dense regions denser and the under-dense regions less dense. These over-dense regions eventually collapsed into gravitationally bound structures. All these structures combined are what we call the cosmic web. Overall, dark matter dominates the cosmic web, while baryonic matter collapsed down to smaller volumes, triggering star and galaxy formation.

Although dark matter is a well-established concept within modern physics (Bertone & Hooper, 2018), its nature is still hypothetical. It comprises currently $\sim 85\%$ of the Universe's mass. Pioneering work has been performed by Zwicky (1933), who found that the kinematics of galaxy clusters hinted at the existence of dark matter. Due to the study on rotation curves by Rubin & Ford (1970); Roberts (1966), it became evidently clear that dark matter was needed to explain astrophysical observations. More recent evidence was found through studies on the light element abundance (Cyburt, 2004), bullet cluster (Markevitch et al., 2004), weak gravitational lensing (Heymans et al., 2012), galaxy clustering (Anderson et al., 2014), and cosmic microwave background anisotropies (Planck Collaboration et al., 2016). These observations confirmed the cold dark matter (CDM) description of the Λ CDM model. According to the model, CDM behaves as a massive, neutral particle, that dominantly, if not solely, interacts with baryonic matter through gravity. It has a non-relativistic speed, hence the ‘cold’ prefix.

There have been several CDM candidates proposed throughout the years. Each of the candidates explains the abundance of dark matter in the Universe with their own production mechanism. The most widely considered candidate is the weakly interactive massive particle (WIMP) (Roszkowski et al., 2018). WIMPs were in thermal equilibrium in the early Universe but stopped interacting (significantly) with any other particles as the Universe expanded. This is the *thermal freeze-out* of WIMPs. The mass range where the particle is considered is 1 GeV - 100 TeV, making the particle decouple at very early times without erasing the structure on small scales. Another well-motivated CDM candidate with lower masses is the axion (Duffy & van Bibber, 2009) because it also solves the strong CP problem (Peccei & Quinn, 1977a,b). To act as CDM, it should have a mass of order $\mu\text{eV}/c^2$.

CDM predictions have been successful at explaining the observed distribution of matter on large

scales (\gtrsim Mpc) across all epochs. However, the agreement with observations on smaller, i.e. galactic and sub-galactic, scales has been less clear. The most well-known tension of CDM is the missing satellite problem (Kravtsov, 2009). Numerical cosmological simulations are used to predict the evolution of the distribution of matter in the Universe. In these simulations, dark matter clusters hierarchically in so-called *halos*. Due to gravity, these halos host a galaxy, like our Milky Way. These halos are orbited by a large number of smaller virialized *subhalos*, or *substructures* (Wechsler & Tinker, 2018).^{1,2} However, in these simulations, there seems to be a mismatch between the observed number of dwarf galaxies and the predicted number of subhalos that should host these dwarf galaxies (Klypin et al., 1999; Moore et al., 1999).

There are two possible solutions to the small-scale tensions. Star formation should be suppressed within low-mass halos and subhalos such that the low number of observed dwarf galaxies can be explained. Alternatively, the abundance of low-mass halos and subhalos has to be suppressed such that there is no mismatch with the observations. Numerous baryonic processes could result in the first solution. Feedback from massive stars, supernovae, or black holes and photoevaporation during reionization could suppress the star formation (Kravtsov, 2009; Bullock, 2010). This solution would have the consequence that subhalos are (even) more difficult to detect because they would host fewer bright galaxies.

The second solution requires modifying the microphysics of dark matter. The abundance of small-scale substructures should be suppressed, while the large-scale predictions should be unaffected. There are multiple proposed modifications to CDM. As an alternative to the non-relativistic CDM, the ultra-relativistic hot dark matter (HDM) has been proposed, possibly in the form of the Standard Model neutrino. But due to its light mass, this was quickly ruled out (White et al., 1983). As a middle ground between CDM and HDM, warm dark matter (WDM) (Bode et al., 2001; Lovell et al., 2014) was proposed, with sterile neutrinos (Boyarsky et al., 2019) or gravitinos (Bond et al., 1982) as its main candidates. With its mass around the keV scale, the particles are relativistic when they decouple but are non-relativistic during the radiation-dominated era. This gives rise to the *free-streaming* of dark matter particles out of density perturbations due to their thermal velocity. This suppresses the structure formation at small scales. There are two more recent alternatives to CDM. The first is self-interacting dark matter (SIDM) (Tulin & Yu, 2018), which has non-gravitational interactions with itself. Because of the scattering, there are fewer dark matter particles at small scales. Finally, the large de Broglie wavelength of ultra-light bosonic, or ‘fuzzy’ dark matter (Hu et al., 2000; Hui et al., 2017), can also suppress the small-scale structures.

The various dark matter models can be differentiated by their free-streaming length, which tells how much the small-scale structure is suppressed (Schneider et al., 2012). It only depends on the particle mass, because it determines the time of decoupling and therefore the speed of the particles during structure formation. A dark matter particle with a lower mass would lead to a higher speed and decoupling at later times. This results in the suppression of larger scales. Similarly, higher masses would lead to earlier decoupling and the suppression would only be limited to smaller scales. The typical suppression scale is parametrized by the half-mode mass M_{hm} . The half-mode mass describes when the halo mass function (HMF) is suppressed by a factor of two with respect to the CDM HMF.

Techniques to detect and characterize dark matter can be divided into three main categories: direct detection, indirect detection, and astrophysical methods. Direct detection methods aim to ob-

¹Substructure generally refers to bound dark matter overdensities on sub-galactic scale within larger structures, while the definition of (sub)halos can be more exact (Despali et al., 2016). We will not go into that discussion and use the terms subhalo and substructure interchangeably.

²Because dark matter clusters hierarchically, the subhalos could have sub-subhalos, and so on. See Giocoli et al. (2010, Figure 1) for a schematic overview of this.

serve the result of interactions between dark matter candidates and particles in an experimental setup. There has been no well-established dark matter detection, which leads to strong constraints on its mass ([XENON Collaboration et al., 2019](#)). Indirect detection relies on observing the product of decay or self-annihilation of dark matter particles. In high-density dark matter regions, i.e. (sub)halos, dark matter particles could produce gamma rays or particle-antiparticle pairs which (constituents) could be detected on Earth. For instance, the GeV excess at the Galactic center could be explained by dark matter ([Fermi LAT Collaboration et al., 2017](#)). There are two drawbacks to direct and indirect detection methods. The first one is that for such a detection, one has to assume a specific dark matter model while we do not know the nature of dark matter particle physics. The other drawback is that a perfect understanding of all other physical processes, such as any other interactions in the experiment or Galactic center, is needed before one can imply anything about the presence of dark matter.

When relying on astrophysical methods, these problems can be rephrased to a simpler one. With these methods, one exploits the gravitational imprint of dark matter on other astrophysical objects. These interactions are relatively simple and better understood. However, one would still have to simulate baryonic physics correctly. There is a variety of astrophysical objects and processes, such as stellar streams ([Banik et al., 2021](#)), Lyman alpha forest ([Baur et al., 2016](#)), ‘dressed’ black holes ([Kavanagh et al., 2020](#)), weak gravitational lensing, ([Mondino et al., 2020](#)) and strong gravitational lensing that are being used. In this work, we will use the latter as a probe for probe dark matter.

Gravitational lensing

During the Newtonian era, it has already been considered that light could be deflected by gravity. Around 1784, unpublished calculations of Cavendish showed how massive bodies could deflect, or ‘bend’, light around themselves ([Will, 1998](#)). [Soldner \(1801\)](#) was the first who published these results, although not fully correct, because the local curvature of space-time was neglected. We had to wait until the general theory of relativity (GR) of [Einstein \(1916\)](#) made the correct prediction. GR describes how light paths are bent due to a curved spacetime. Large masses act therefore as a lens, and light gets gravitationally lensed by this mass. The modern theory of lensing was developed in the 60s (e.g. [Liebes, 1964](#); [Refsdal & Bondi, 1964](#); [Refsdal, 1964](#)), resulting in the first observation of a lensed quasar by [Walsh et al. \(1979\)](#).

Gravitational lensing is divided into the regimes of strong and weak lensing. Weak lensing, reviewed in [Schneider \(2006\)](#), occurs when the lensing mass is located far away from the *line of sight*, the straight line between the observer and the light’s source. Weak lensing causes small modifications to the light’s direction. It is extremely common, every single light ray is weakly lensed up to some point. Objects from all scales, from galaxies and galaxy clusters to cosmological lensing, act as weak lenses. No individual weak-lensed sources can be identified, due to the small modifications. Therefore it can only be studied in a statistical sense.

Strongly-lensed sources, on the contrary, can readily be identified. It occurs when the lensing mass, typically a galaxy or galaxy cluster, is located close to the line of sight. Light from all sides around the lens gets deflected to the line of sight, as shown by the straight black lines in Figure 1.1. It produces multiple images of the background object at the observer’s location. What we observe depends on the size of the source. Small, point-like sources, such as quasars or other AGNs, are lensed into multiple point-like components (e.g. [Fassnacht et al., 1999](#)). Larger sources, such as extended galaxies, are observed as extended arcs or Einstein rings (e.g. [Lin et al., 2009](#)).

Gravitational lensing has various applications in cosmology and astrophysics. For example, it is possible to study the most distant objects because the lens magnifies the background source light ([Deane](#)

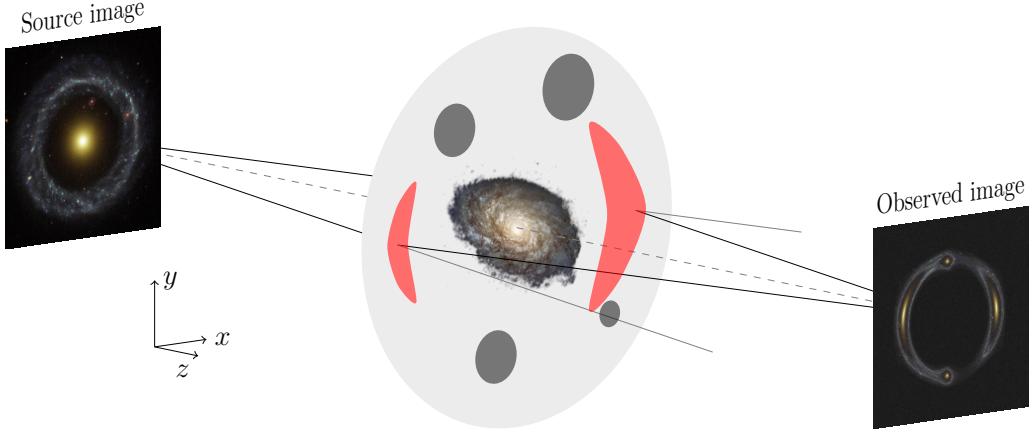


Figure 1.1: Not to scale sketch of strong galaxy-galaxy lensing. Light, coming from the source, travels around the lensing galaxy to the observed image. In this sketch, light from the lensing galaxy is not included in the observation. By analyzing how the deflected light rays are distorted by subhalos, we can learn more about dark matter. The main halo (light gray circles) hosts the central lensing galaxy and the subhalos (dark gray circles). The black straight lines indicate the deflected light curves. The grey solid line shows the non-deflection light curves if there was no lens present. The gray dashed line indicates the line of sight.

et al., 2013); small objects such as planets can be detected with the effect of microlensing (Tsapras, 2018); the mass of supermassive black holes can be measured (Winn et al., 2004); the radial mass density profile of galaxies can be studied (Fong et al., 2018); with observed time-delays the Hubble constant can be measured (Wong et al., 2020). Furthermore, strong gravitational lensing provides a direct way of probing the dark matter distribution in the lens galaxy. This is because gravitational lensing is only sensitive to the mass distribution, whether it is baryonic or dark matter. The source flux, strongly lensed by the main halo’s mass distribution, has small distortions due to the subhalos. By analyzing the gravitational effect of the subhalos, we can eventually extract the subhalo mass function (SHMF) and test the various dark matter models.

Besides subhalos, line-of-sight (LOS) halos leave also a small imprint on the lensing. These halos are not part of the main halo and could be anywhere between the source and observer. They are not affected by tidal effects of the main lens or the presence of baryonic matter of the lensing galaxy and could therefore be a more direct probe of small-scale structure suppression. Furthermore, they are (far) more abundant than subhalos. This work, however, will focus on the development of new methods to measure dark matter subhalos. To simplify the problem, we will not model LOS halos.

The lensed images are clearly separate from each other with lensed quasar sources. Analyzing subhalos with quasar lensing was first done by Mao & Schneider (1998), where the flux-ratio from the projected images were compared with each other. Constraints on the abundance of substructure were found by Dalal & Kochanek (2002), by analyzing a small sample of strongly lensed quasars. Recent constraints using quasars are for example from Zelko et al. (2022); Gilman et al. (2020).

If the background source is a galaxy, the arc or Einstein ring has small variations due to the subhalos. The analysis of galaxy-galaxy lensing has been developed by Koopmans (2006); Vegetti & Koopmans (2009a,b). A smooth lens and source galaxy are modeled and fit to an observation, and subhalos are added to the model to improve the fit. With this gravitational imaging technique, Vegetti et al. (2010a,b, 2012); Hezaveh et al. (2016) have detected several individual $\sim 10^9 M_\odot$ subhalos. However, the large

uncertainty in the detections does not yet allow for a discriminating test for dark matter models.

The technique used to detect subhalos is a *likelihood-based* method. That means that one needs to marginalize over the nuisance parameters to infer the properties of the subhalos. Nuisance parameters are part of the model, but not of immediate interest. To detect subhalos in a strong lensed system, the lens and source parameters are the nuisance parameters. The marginalization over all the lens and source parameters is traditionally done with Markov-chain Monte Carlo (MCMC) or nested-sampling methods (Skilling, 2004, 2006). This becomes very slow and time-consuming for high-dimensional problems. Therefore one needs to compromise between an intractable high number of source and lens parameters or less complex modeling.

The complexity of the lens is reduced by modeling only one, or at most two, subhalos. These detection methods probe typically the most massive subhalo(s), which is less interesting than a low-mass subhalo if the goal is to constrain small-scale structures. Furthermore, it is not taken into account that the collective effect of close-by multiple low-mass subhalos could also be interpreted as a massive subhalo. If N_{sub} subhalos would be modeled, the marginalization would become much more complex due to the *label switching problem*. Because the subhalos are interchangeable, the ordering of the subhalos' parameters is arbitrary. This makes the joint posterior highly multimodal.

The methods to overcome these challenges can be divided into two classes. The first is transdimensional Bayesian inference (Brewer et al., 2015; Daylan et al., 2018). This method takes models with a different number of subhalos into account and maps them to the parameter space of multiple low-mass subhalos. The output, probabilities for subhalo catalogs, can be used to determine the population properties of the subhalos. This method is ‘transdimensional’ because it goes through different subhalo models. However, this technique still requires a large computational cost.

The other approach is measuring the collective effect of multiple subhalos through a summary statistic, the power spectrum, as proposed by Hezaveh et al. (2016). The power spectrum of the residuals between the observation and fitted reconstruction without subhalos can be related to the properties of the subhalo population (Chatterjee & Koopmans, 2018; Bayer et al., 2018; Cyr-Racine et al., 2019). The advantage of describing the population instead of individual subhalos is that the collective effect of otherwise non-detectable subhalos can be measured. Low-mass subhalos can have a too low imprint on the gravitational lensing to be detected individually, but their collective impact could be significant. Another advantage is that we directly measure population properties, which are more interesting for testing dark matter theories. However, all information gets compressed into a summary statistic, which might leave valuable statistics of the subhalo information out.

Another challenge of MCMC or nested sampling is the reproducibility. The marginalization is done for a single observation, and specific marginalizations can not be redone without rerunning the full procedure again. Therefore, statistics, such as biases and systematics, can not directly be explored. This makes the subhalo analysis slow and computationally costly. Furthermore, likelihood-based methods take a particular configuration of the noise and lens parameters, so that they are analytically marginalized. This makes it hard to implement more flexible sources.

There are currently on the order of a hundred observations available that are suitable for subhalo detection in galaxy-galaxy lenses. Most of these observations are from the SLACS (Shu et al., 2017) and BELLS (Cornachione et al., 2018) optical surveys, taken with the Hubble Space Telescope (HST). However, near-future telescopes such as LSST from the Vera C. Rubin Observatory (LSST Science Collaboration et al., 2009), Euclid (Refregier et al., 2010), James Webb Space Telescope (Gardner et al., 2006), and the Extremely Large Telescope (Simon et al., 2019) are expected to detect hundreds of thousands of observations of higher quality (Collett, 2015). This calls for the development of methods that can efficiently analyze large samples of observations with fast automated pipelines because the speed of existing methods to analyze these observations will not be sufficient.

Machine learning techniques have already proved itself to speed up and improve the quality of the analysis of strongly lensed images. Several works have applied it to the reconstruction of the source image and lens parameters. [Hezaveh et al. \(2017\)](#) and [Levasseur et al. \(2017\)](#) used convolutional neural networks (CNNs) to infer the parameters of the main lens. CNNs were combined with recurrent neural networks to infer the lens parameters by [Morningstar et al. \(2018\)](#) and the whole system (both lens parameters and source image) by [Morningstar et al. \(2019\)](#). [Chianese et al. \(2020\)](#) proposed using variational autoencoders to infer the whole system. [Karchev et al. \(2022\)](#) used variational inference to describe the source light with Gaussian processes. Finally, both [Wagner-Carena et al. \(2021\)](#); [Pearson et al. \(2021\)](#) used Bayesian neural networks to model the lens parameters.

CNNs have also been used to infer subhalos from strongly lensed images. [Rivero & Dvorkin \(2020\)](#) trained a binary classifier with a CNN to determine if a strongly lensed image contains substructure or not. This method has a lower mass limit of $\sim 10^9 M_{\odot}$, similar to the subhalos detected with traditional methods. Similar to the approach with the power spectrum, [Varma et al. \(2020\)](#) used the effect of multiple subhalos to measure the low-mass cutoff of the SHMF. [Alexander et al. \(2020\)](#) trained CNNs to distinguish between different types of dark matter while [Alexander et al. \(2021\)](#) used unsupervised techniques to infer model-independent subhalos. On top of that, the use of machine learning techniques has made it possible to measure multiple subhalos individually at once. [Lin et al. \(2020\)](#) used CNNs to produce a probability map, indicating whether there is a subhalo at a specific location in the lens. The network was also able to reject subhalos if there were no subhalos. [Ostdiek et al. \(2022a,b\)](#) elevated this approach by using an architecture designed to produce such ‘maps’. They use image segmentation to predict the probability of a subhalo being in a pixel, while also discriminating between subhalo mass classes.

Simulation-based inference

Recent developments in machine learning have opened up a new branch of machine learning techniques to infer dark matter properties from strongly-lensed images. With the method of simulation-based inference (SBI), there is no need to compute the likelihood ([Cranmer et al., 2020](#)). This *likelihood-free* method makes it therefore possible to relax the requirements that are necessary with likelihood-based methods. SBI samples observations from a stochastic simulator (the lensing model). These samples function as the likelihood and are being passed to the neural network as training data. The networks are trained to estimate either the posterior ([Lueckmann et al., 2017](#)), the likelihood ([Papamakarios et al., 2019](#)), or a likelihood ratio ([Hermans et al., 2020](#)). These methods have all their own specific strengths and weaknesses.

The latter method, known as neural ratio estimation (NRE) ([Hermans et al., 2020](#); [Miller et al., 2020](#)), has the advantage that the inference task is rephrased to a binary classification task. Binary classification is the most ‘simple’ task that one can solve, which generally results in less complex network architecture and fewer training iterations. NRE is also suitable to directly calculate *marginals*. With this procedure, the network is trained on a specific subset of the model parameters, the parameters of interest, while marginalizing over the nuisance parameters. This method can further be extended by doing *truncations*. The training consists of multiple ‘rounds’ where in the next round, only data is considered in the most relevant parameter space for a particular observation. Since the data is ‘focused’ in future rounds, there is less training data needed to obtain the same inference quality compared without truncation. This full procedure is called truncated marginal neural ratio estimation (TMNRE) ([Miller et al., 2021](#)).

An early precursor is Approximate Bayesian Computation (ABC). This likelihood-free method was

applied on strongly lensed images by [Birrer et al. \(2017\)](#); [Gilman et al. \(2018\)](#); [He et al. \(2022\)](#). With ABC, a hand-crafted summary statistic is chosen to be minimized, similar to the likelihood-based method with power spectra. A drawback of these methods is that it is unclear how much information is discarded when focusing on a summary statistic. It is hard to define what the most optimal summary statistic would be. Another drawback of ABC is scalability, because each observation requires an independent chain of hundreds of thousands of tailored simulations.³

The usage of machine learning and SBI has already found its way into measuring substructure with lensing. [Brehmer et al. \(2019\)](#) used a likelihood ratio to estimate the slope and normalization of the SHMF. [Wagner-Carena et al. \(2022\)](#) used neural posterior estimation to measure the normalization of the SHMF. [Anau Montel et al. \(2022\)](#) showed how TMNRE can be used to determine the cutoff of the SHMF in the WDM model. Besides these works that measured subhalo population statistics, [Coogan et al. \(2022\)](#) demonstrated how TMNRE is able to produce posteriors of the location and mass of a single subhalo while modeling other subhalos as well. SBI has also found its use in other astrophysical contexts, such as in cosmology ([Alsing et al., 2018, 2019](#); [Cole et al., 2021](#)), gravitational waves ([Dax et al., 2021](#); [Delaunoy et al., 2020](#)), and the Galactic Center gamma-ray excess ([Mishra-Sharma & Cranmer, 2022](#)).

This thesis focuses on measuring the parameters of multiple subhalos simultaneously with SBI. Two main challenges arise when doing this. Firstly, the exact number of subhalos is unknown. So one could not specify how many parameters should be estimated. The second problem is the label switching problem. Previous works, as explained earlier, approached the analysis of multiple subhalos by estimating a single or a few population statistics. Instead, we approach the problem, similar to [Ostdiek et al. \(2022a,b\)](#), by using image segmentation. Such ‘traditional’ image segmentation predicts a ‘binary’ outcome: a pixel belongs to a certain class, or not. Additionally, we apply MNRE, similar to [Anau Montel et al. \(2022\)](#); [Coogan et al. \(2022\)](#). By combining image segmentation with MNRE, we are able to directly estimate the probability that there is a subhalo in a particular pixel. The probability of all pixels combined can be interpreted as a subhalo density estimation.

The structure of this thesis is as follows. In Chapter 2, we lay down the physical framework of strong gravitational lensing. We discuss how the lens, source, and subhalos are modeled in the simulator. We introduce the necessary statistical concepts in Chapter 3, where we explain separately how image segmentation and MNRE work. We will combine these two concepts in Chapter 4, where we explain how the subhalo density estimation is trained and evaluated. The result of these methods on HST mock images is in Chapter 5. Finally, the nuances of the method and results, and future development of the framework are discussed in Chapter 6.

³One could raise this also as an argument against TMNRE. However it is not the same: TMNRE uses tailored simulations, but narrows the prior, and therefore increasing the (potential) inference quality per simulation. ABC requires tailored simulations, but it is not able to narrow the prior.

Chapter 2

Modeling subhalos in strong lensing systems

In this chapter, we review how strongly lensed images are modeled. We closely follow the setup from [Meneghetti \(2021\)](#), other detailed works on gravitational lensing are e.g. [Schneider et al. \(1992\)](#); [Treu \(2010\)](#).

We assume that the lens is small compared to the overall dimensions of the optical system.⁴ From this assumption, we do two usual approximations. The first is that mass densities are low enough such that gravitational effects can be ignored. This implies that the metric is fully described by a scalar gravitational potential ψ , which is the Newtonian approximation of general relativity. The second approximation is that of a ‘thin lens’. This assumes that all the lens mass lies in a single plane, the *image plane* or *lens plane*. Additionally, the source light is emitted from a single plane as well, the *source plane*. Because of the above assumptions, the lensing deflections are small and the image plane covers only a small part of the sky. Therefore, the coordinate system can be treated as Cartesian. The two-dimensional coordinates $\xi = (\xi_x, \xi_y)$ and $x = (x, y)$ are used to describe the source and lens plane, respectively. Within this model, the lensing happens ‘instantly’ in the lensing model.

The matter density of the lens $\rho(\xi, z)$ is projected onto the lens plane, and can therefore be described as the surface density

$$\Sigma(\xi) = \int dz \rho(\xi, z), \quad (2.1)$$

where z is the direction perpendicular to the lens plane. The lens plane coordinates ξ of light rays are related to the source-plane coordinates x through the *lens* or *raytrace equation*

$$x = \xi - \alpha(\xi). \quad (2.2)$$

Here α is the *deflection* or *displacement field*, the deflection of all mass elements integrated over the lens plane,

$$\alpha(\xi) = \frac{4G}{c^2} \frac{D_{LS}}{D_L D_S} \int d^2(D_L \xi') \frac{\xi - \xi'}{|\xi - \xi'|^2} \Sigma(\xi'). \quad (2.3)$$

We have introduced the distances D_{LS} (from the lens to source), D_L (from observer to lens) and total distance D_S (from observer to source).⁵ G and c are the gravitational constant and speed of light. The expression for the displacement field, Equation (2.3), can be simplified by introducing the convergece κ and critical surface density Σ_{cr}

⁴Similar to almost all astrophysical systems due to the large cosmological distances that light travels.

⁵Due to the expansion of the Universe, one has to relate the comoving distances to angular diameter distances. We use the flat cosmology ([Planck Collaboration et al., 2016](#)) and relate the comoving distances depening on the redshift of the lens and source to angular diameter distances with the `astropy` package ([The Astropy Collaboration et al., 2013, 2018](#)).

$$\kappa(\xi) \equiv \frac{\Sigma(\xi)}{\Sigma_{\text{cr}}}, \quad \Sigma_{\text{cr}} = \frac{c^2}{4\pi G} \frac{D_S}{D_L D_{\text{LS}}}. \quad (2.4, 2.5)$$

Where $\kappa(\xi)$ is defined so that it the divergence of $\alpha(\xi)$. The lens equation can be written down as a coordinate transformation $\frac{dx}{d\xi}$. It can be shown that the Jacobian of the lensing transformation is fully determined by the convergence κ and two external shear components (which will be introduced shortly). The ratio between these parameters differentiates strong from weak lensing. See [Meneghetti \(2021, Section 2.4\)](#) for a discussion on this.

Because lensing creates multiple images of the same source, it might look like photons are being created or destroyed. This is not the case, as lensing only alters the trajectory of the photons. Light emitted from the source is described by the surface brightness $\beta_{\text{source}}(x)$. The total surface brightness at the source $\beta_{\text{source}}(x)$ is the same as the total surface brightness at the lens $\beta_{\text{lens}}(\xi)$:

$$\beta_{\text{lens}}(\xi) = \beta_{\text{source}}(x(\xi)). \quad (2.6)$$

So lensing conserves surface brightness and energy. From this expression we can see that lens and source are degenerate, there are multiple configurations of the lens and source that could produce the same observation.

The altering of trajectory of the surface brightness $\beta(x)$ is determined by the deflection field α through the lens equation. The total deflection field is commonly split into the components of the main lens, external shear, and substructure. The total deflection field is a superposition of the individual components, and can be computed by summing the respective displacement fields:

$$\alpha = \alpha_{\text{lens}} + \alpha_{\text{ext}} + \sum_{i=1}^{N_{\text{sub}}} \alpha_{\text{sub},i}. \quad (2.7)$$

In the sections below we will describe how these displacement fields and the surface brightness at the source are modeled.

2.1 Main lens

From cosmological N -body simulations, [Navarro et al. \(1996, 1997\)](#) found that the dark matter distribution is well described by Navarro-Frenk-White (NFW) profile at galactic scales. This profile has extensively been used to describe dark matter. However, this profile does not generalize to strong lensing, as the region where strong lensing probes is much smaller than the typical virial radius of the main halo. At large scales, the mass budget is mainly dominated by dark matter. At smaller scales, like the scales of strong lensing, the baryonic bulge dominate the mass budget. The singular isothermal ellipsoid (SIE) model, or the simplified version, the singular isothermal sphere (SIS) model, take this into account. These models are often used to describe the main lens in automated analysis of strong lensing, such as in [Brehmer et al. \(2019\)](#); [Hezaveh et al. \(2017\)](#); [Ostdiek et al. \(2022a,b\)](#); [Legin et al. \(2021\)](#). However, observations indicated that a more flexible model is required ([Bolton et al., 2008](#)).

For that reason, we use the singular power law sphere (SPLE), an extension SIE/SIS, to describe the main lens. This model takes a more representative combination of dark and baryonic matter into account ([Suyu et al., 2009](#)). The SPLE displacement field, given by

$$\alpha^{\text{SPLE}}(R, \phi) = \theta_E \frac{2\sqrt{q}}{1+q} \left(\frac{\theta_E}{R(\xi)} \right)^{\gamma-2} e^{i\phi} {}_2F_1 \left(1, \frac{\gamma-1}{2}; \frac{5-\gamma}{2}; -\frac{1-q}{1+q} e^{2i\phi} \right), \quad (2.8)$$

has a closed form in the complex notation $\alpha = \alpha_x + i\alpha_y$ (O’Riordan et al., 2020; Tessore & Metcalf, 2015). The elliptical coordinates (R, ϕ) are related to the Cartesian coordinates ξ through a transformation. This transformation is parametrized by ellipticity or axes ratio q , the ellipse orientation φ and the lens’ position ξ_0 (Karchev et al., 2022; Coogan et al., 2022):

$$\begin{pmatrix} R_x \\ R_y \end{pmatrix} = \begin{pmatrix} q^{1/2} & 0 \\ 0 & q^{-1/2} \end{pmatrix} \begin{pmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{pmatrix} \begin{pmatrix} \xi_x - \xi_{0,x} \\ \xi_y - \xi_{0,y} \end{pmatrix}. \quad (2.9)$$

They map from the Cartesian to elliptical coordinates through the usual transformation

$$R = \sqrt{R_x^2 + R_y^2}, \quad \phi = \arctan(R_y/R_x). \quad (2.10)$$

The hypergeometric function ${}_2F_1$ is calculated through an interpolation of a pretabulate, as described in Chianese et al. (2020, Appendix A).⁶ The final parameters of the displacement field are the Einstein radius θ_E and the slope γ . Note that the displacement field reduces to the simple expression $\alpha = \theta_E e^{i\phi}$ if the lens is circular ($q = 1$) and isothermal ($\gamma = 2$).

Additional to the lens parameters, we include an external shear component to account for weak lensing across the light path. It describes all lenses that are not modeled or properly taken into account. The displacement field

$$\alpha^{\text{shear}}(\xi) = \begin{pmatrix} \gamma_1 & \gamma_2 \\ \gamma_2 & -\gamma_1 \end{pmatrix} \xi \quad (2.11)$$

is modeled by introducing the external shear components γ_1, γ_2 .

Taking both the SPLE and external shear into account, we end up with eight parameters that describe the main lens: $\theta_{\text{lens}} = \{\xi_{0,x}, \xi_{0,y}, \varphi_{\text{lens}}, q_{\text{lens}}, \theta_E, \gamma, \gamma_1, \gamma_2\}$, where we have added the subscripts \cdot_{lens} to the parameters φ and q .

2.2 Subhalos

Dark matter subhalos leave a small imprint on the gravitationally lensed image. They are distributed within the virial radius of the main lens halo and lie therefore on the lens plane according to the thin-lens approximation. They are described by three components. Firstly, the internal density profile, which describes how matter is distributed inside the subhalo. Secondly, the SHMF, which tells how many halos of a given mass there are. And finally, its spatial distribution, which tells how the halos are distributed inside the main halo. We will discuss the distributions below.

The mass within the subhalos is modeled with the truncated NFW profile (Baltz et al., 2009)

$$\rho_{\text{tNFW}}(r) = \frac{\rho_s}{r/r_s} \frac{1}{(1+r/r_s)^2} \frac{1}{1+(r/r_t)^2}. \quad (2.12)$$

The first term, the NFW profile, is described by the radial distance from the center of the halo $r = \sqrt{\xi_x^2 + \xi_y^2 + z^2}$, the density normalization ρ_s and the scale radius r_s . The truncation⁷ is described by the second term and it depends, specified by the truncation radius $r_t \equiv \tau r_s$, on the history of the subhalo. Typical truncation scale τ values follow a skewed distribution with typical values between 4 – 10 (Gilman et al., 2020). We fix it on the mean $\tau = 6$ for simplicity. The standard NFW profile has

⁶The hypergeometric function is a solution for many second-order linear ordinary differential equation.

⁷Not to confuse with truncation in TMNRE.

an infinite total mass. Therefore, we have chosen the truncated NFW profile, which makes the profile decay as r^{-5} for large radii. This results in the finite total mass

$$m_\tau = 4\pi\rho_s r_s^3 \frac{\tau^2}{(\tau^2 + 1)^2} \left[(\tau^2 - 1) \ln \tau + \tau\pi - (\tau^2 + 1) \right]. \quad (2.13)$$

With a fixed truncation scale, the truncated NFW is described by the same parameters as the standard NFW profile, the scale radius r_s and density normalization ρ_s . However, studies on the subhalo mass typically measure the subhalo's mass m_{200} and the concentration c_{200} .⁸ These variables determine the NFW free parameters with

$$r_s = \frac{1}{c_{200}} \left[\frac{3m_{200}}{4\pi 200 \rho_{\text{cr}}(z_{\text{lens}})} \right]^{1/3}, \quad (2.14)$$

$$\rho_s = \rho_{\text{cr}}(z_{\text{lens}}) \frac{1}{3} \frac{c_{200}^3}{\log(1 + c_{200}) - c_{200}/(1 + c_{200})}. \quad (2.15)$$

We fix the concentration $c_{200} = 15$, following Richings et al. (2021). We note that sampling the concentration, by taking the scatter in the mass-concentration into account, would improve our inference results, as higher concentrations lead to significantly stronger lensing signals (Amorisco et al., 2022).

Masses of subhalos are determined by the SHMF, which could be expressed by (Giocoli et al., 2010)

$$\frac{1}{M} \frac{dN_{\text{sub}}(m_{200}, z_{\text{lens}})}{d \log m_{200}} = (1 + z_{\text{lens}})^{1/2} A_M m_{200}^\alpha \exp \left[-\beta \left(\frac{m_{200}}{M} \right)^3 \right], \quad (2.16)$$

with M the mass of the main lens. The free parameters, the normalization A_M , slope α and exponential cutoff β , of the SHMF can be fit to the values of the cosmological simulation EAGLE (Despali & Vegetti, 2017). The number of subhalos in a lensing system would be determined by sampling from a Poisson distribution with parameter N_{sub} . N_{sub} would be calculated by integrating the SHMF over the appropriate range of the subhalo mass ($m_{\text{low}}, m_{\text{high}}$). However, if we are modeling multiple subhalos, we simplify the sampling of number of subhalos and their masses. We will stick to the a uniform distribution for both N_{sub} and m_{sub} .

The spatial distribution of the subhalos is assumed to be spherically symmetric around the center of the main lens. It has been shown that the subhalos are distributed following the Einasto profile (Springel et al., 2008)

$$\rho_{\text{Einasto}}(r) \propto r^2 \exp \left\{ -\frac{2}{\alpha} \left[\left(\frac{r}{R_S} \right)^\alpha - 1 \right] \right\}, \quad (2.17)$$

with index $\alpha = 1.1$ and the scale radius R_S as a function of the main lens mass M . However, with the Einasto profile, the virial radius of a typical main halo is much larger than the image plane, as can be seen in Despali & Vegetti (2017, Figure 3). Therefore, we can approximate the distribution of subhalos to be uniform in the lens plane.

The modeling of LOS halos will not be discussed, as they are not considered in our model. However, we would like to stress that these can straightforwardly be taken into account with our simulator. These would be modeled as low-mass halos without any subhalos. Practically, they are modeled by projecting them on the lens plane. They would have a similar imprint on the gravitational lensing as subhalos.

Given the truncated NFW, mass and position, the displacement field of the subhalos can be computed. We refer to (Baltz et al., 2009, Appendix A) for the complete calculation. With most subhalo

⁸The subscript 200 refers to the definition of the a halo: the mass enclosed in a sphere with an average density of 200 times the critical density at redshift $z = 200$.

parameters fixed to values from cosmological parameters, we only vary the subhalos' local parameters of mass and position: $\boldsymbol{\theta}_{\text{subs}} = \{\boldsymbol{\theta}_{\text{sub},1}, \boldsymbol{\theta}_{\text{sub},2}, \dots, \boldsymbol{\theta}_{\text{sub},N_{\text{sub}}}\}$, with $\boldsymbol{\theta}_{\text{sub},i} = (\xi_{x,i}, \xi_{y,i}, m_{200,i}) = (x_{\text{sub},i}, y_{\text{sub},i}, m_{\text{sub},i})$. In the last equality we have renamed the conventional subhalo parameters definition to names that need less context.

2.3 Source

The source model describes the emission from the background galaxy. It produces a surface brightness at the source plane. Modeling the surface brightness of a realistic-looking galaxy is not easy, therefore a more simplified model is often chosen. One has to decide between two types of source modeling: non-parametric and flexible sources or parametric and simplified.

Non-parametric sources are not described by a specified set of parameters. Instead, the surface brightness of each pixel is modeled by a complex model. [Vegetti & Koopmans \(2009a\)](#) mapped an adaptive grid onto the source plane, [Birrer et al. \(2017\)](#) used shapelets, [Chianese et al. \(2020\)](#) variational autoencoders and [Karchev et al. \(2022\)](#) Gaussian processes to model the source. Non-parametric modeling is more flexible because it is not constrained to a specific set of parameters. It is able to better capture the complex morphology of the surface brightness of real galaxies without assumptions on the kind of galaxy. It is therefore an active field of research.

Parametric sources model the surface brightness of a galaxy as a profile, described by a function dependent on a specified set of parameters. The profile is usually not flexible enough to describe the surface brightness of real galaxies. The big variety of galaxies is then considered as uncertainty of the model. However, they are an efficient approximation of galaxy sources in general and describe the general shape and size of the source. Due to their simplicity, they are often used in data-driven analysis of dark matter in strongly lensed images (e.g. [Anau Montel et al., 2022](#); [Coogan et al., 2022](#); [Ostdiek et al., 2022b,a](#); [Brehmer et al., 2019](#)).

The typical used parametric source model is the Sérsic profile ([Sérsic, 1963](#))

$$\beta(\mathbf{x}) = I_e \exp \left\{ -k_n \left[\left(\frac{R(\mathbf{x})}{R_e} \right)^{1/n} - 1 \right] \right\}, \quad (2.18)$$

with $R(\mathbf{x})$ the elliptical radial coordinate, similar to Equations (2.9) and (2.10), but now with \mathbf{x} and \mathbf{x}_0 instead of ξ and ξ_0 . The normalization factor I_e is the surface brightness at the half-light radius R_e . The other normalization k_n depends on the index n and is related to an implicit transcendental⁹ equation with gamma functions¹⁰ $2\gamma(2n, k_n) = \Gamma(2n)$. We use the series expansion

$$k_n \approx 2n - \frac{1}{3} + \frac{4}{405} n^{-1} + \frac{46}{25515} n^{-2} + O(n^{-3}) \quad (2.19)$$

from [Ciotti & Bertin \(1999\)](#) to estimate the normalization factor.

The source, modeled with a Sérsic profile, has in total seven parameters: $\boldsymbol{\theta}_{\text{source}} = \{\mathbf{x}_{0,x}, \mathbf{x}_{0,y}, \varphi_{\text{source}}, q_{\text{source}}, n, R_e, I_e\}$.

2.4 Instrumentation

As our modeling is intended to be applied to observations from telescopes, we have to model the instrumental uncertainties. Light reaching a telescope is transformed into digital images by the software of

⁹A [transcendental function](#) is an analytic function that does not have a polynomial structure

¹⁰To be precise, they are the [lower incomplete](#) gamma function γ and [complete](#) (or ordinary) gamma function Γ .

the telescope. The telescope consists of light-receiving pixels, small photon-counting detectors, that are exposed over an extended period to a restricted area of the sky. However, the light is spread out over the detectors due to atmospheric distortions and defects in the optical system. This should be taken into account by convoluting a point spread functions (PSF) with the received surface brightness for each pixel. Modeling the PSF is nevertheless not so straightforward. Its effect varies for each pixel and is different for each filter from the telescopes. Moreover, it can vary with time and temperature for the same instrument. To solve this, we use a resolution that is slightly lower than the expected resolution of the telescope. This allows us to neglect the PSF of the instrument, as this should be within the pixel size. Nevertheless, to account for pixelation issues and simulate the integration of light across the pixel areas, we initially generate mock data with a higher resolution and subsequently downsample it to the adopted resolution by local averaging.

The instrumental effects are modeled by adding noise to each pixel in the form of uncorrelated Gaussian noise. The mean of the Gaussian is the simulated received surface brightness μ , and the width of the Gaussian is the user-defined noise level σ . So we model the instrumental parameters for now with just a single parameter σ . Finally, we can write down our simulator as

$$p(\mathbf{x}_{\text{obs}} | \boldsymbol{\theta}_{\text{lens}}, \boldsymbol{\theta}_{\text{source}}, \boldsymbol{\theta}_{\text{subs}}, \sigma) = \mathcal{N}(x | \mu(\boldsymbol{\theta}_{\text{lens}}, \boldsymbol{\theta}_{\text{source}}, \boldsymbol{\theta}_{\text{subs}}), \sigma^2) \quad (2.20)$$

The exact values that are chosen for the parameters are discussed when the results are presented, Chapter 5.

There are various aspects of strongly lensed observations not taken into account. We have not discussed the light from the foreground galaxy, as it is in general subtracted in real data analysis. Furthermore, telescopes can often differentiate light of different wavelengths into multiple energy bands or channels. The different light bands are also not modeled. Further discussion about these topics is in Section 6.2.

Chapter 3

Statistical concepts

In this chapter, we discuss in detail two separate concepts which are built upon in the next chapters. These concepts do not have a direct connection to dark matter or gravitational lensing. Therefore, the sections can be read independently of the rest of this thesis. In the first section, we will discuss the technicalities of MNRE. In the second section, we will go through the U-Net, the neural network that we will use.

3.1 Likelihood-free inference

To learn about the properties of dark matter, we want to measure the subhalos in strong gravitational lensing galaxies. Practically, this means that we are interested in the posterior probability $p(\boldsymbol{\theta}|\mathbf{x})$ of the parameters describing the subhalos. That is, the probability of a set of parameters $\boldsymbol{\theta}$ given some observation \mathbf{x} . With Bayes theorem, the posterior is defined as

$$p(\boldsymbol{\theta}|\mathbf{x}) = \frac{p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{x})}. \quad (3.1)$$

Bayes theorem consists of three terms:

- the **likelihood** $p(\mathbf{x}|\boldsymbol{\theta})$ tells the probability that some parameters $\boldsymbol{\theta}$ have produced an observation \mathbf{x} .
- the **prior** $p(\boldsymbol{\theta})$ reflects our existing knowledge. It spans the entire parameter space where we think the parameters $\boldsymbol{\theta}$ are. A uniform distribution is often used.
- the **evidence** $p(\mathbf{x}) = \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$ normalizes the numerator by the probability of the observation \mathbf{x} .

In physics and astronomy, we are often only interested in a subset of the parameters $\boldsymbol{\theta}$. Therefore we split the parameters into physical parameters $\boldsymbol{\vartheta}$ that we are interested in, and nuisance parameters (consisting of physical parameters and latent variables) $\boldsymbol{\eta}$, i.e. $\boldsymbol{\theta} \equiv (\boldsymbol{\vartheta}, \boldsymbol{\eta})$. We are only interested in the marginal posterior $p(\boldsymbol{\vartheta}|\mathbf{x})$, which means we have to integrate the joint posterior $p(\boldsymbol{\vartheta}, \boldsymbol{\eta}|\mathbf{x})$ over the nuisance parameters $\boldsymbol{\eta}$:

$$p(\boldsymbol{\vartheta}|\mathbf{x}) = \int p(\boldsymbol{\vartheta}, \boldsymbol{\eta}|\mathbf{x}) d\boldsymbol{\eta} = p(\boldsymbol{\vartheta}) \frac{\int p(\mathbf{x}|\boldsymbol{\vartheta}, \boldsymbol{\eta})p(\boldsymbol{\eta}) d\boldsymbol{\eta}}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\boldsymbol{\vartheta})}{p(\mathbf{x})} p(\boldsymbol{\vartheta}). \quad (3.2)$$

If there are a lot of nuisance parameters, this problem gets the curse of dimensionality. The marginalization becomes intractable because one would have to integrate over a high number of dimensions.

The intractability of the likelihood makes it infeasible to use traditional approaches such as MCMC and nested sampling. These methods compute the joint posterior, so one would still need to marginalize all the nuisance parameters to obtain the marginal posterior.

3.1.1 Simulation-based inference

With SBI, one can overcome these challenges. Instead of calculating the likelihood, one sample from a simulator, the physical model. This is done by drawing parameters θ from the prior, and passing them as input to the simulator to generate simulations:

$$\mathbf{x}_{\text{sim}} \sim \text{Simulator}(\boldsymbol{\theta}). \quad (3.3)$$

By drawing from the simulator, we obtain a probabilistic model of what parameters θ produce an observation \mathbf{x} . This is the same as a likelihood $p(\mathbf{x}|\theta)$ is telling us. The simulator ‘mimics’ therefore the role of the likelihood without actually calculating the likelihood. SBI is, therefore, a likelihood-free method.

ABC ([Sisson et al., 2018](#)), an early SBI method, is such a likelihood-free method. It is a rejection algorithm where the distribution of the parameters of interest ϑ is recorded only for the simulations \mathbf{x}_{sim} that are sufficiently similar to the analyzed observation \mathbf{x}_{obs} . The simulation and observation are ‘similar’ when the condition $\rho(\mathbf{x}_{\text{sim}}, \mathbf{x}_{\text{obs}}) < \epsilon$ is met, where ρ is a measure for the distance and ϵ is a user-defined tolerance.¹¹ This method becomes exact in the limit $\epsilon \rightarrow 0$. However, too small values of ϵ require too many simulations, as the acceptance probability will become too low. On the other hand, large ϵ will decrease the inference quality. The approximated posterior can be inaccurate or biased. For a further discussion about the shortcomings of ABC, we refer to [Cranmer et al. \(2020\)](#).

A revolution in machine learning (ML), in particular deep neural networks, has made it possible to address the above issues ([Cranmer et al., 2020](#)). With neural SBI, the simulator generates a set of N sample-parameters pairs $\{(\mathbf{x}_{\text{sim}}^{(1)}, \boldsymbol{\theta}^{(1)}), (\mathbf{x}_{\text{sim}}^{(2)}, \boldsymbol{\theta}^{(2)}), \dots, (\mathbf{x}_{\text{sim}}^{(N)}, \boldsymbol{\theta}^{(N)})\}$. A neural network can use these pairs as training data to estimate the posterior, the likelihood or the likelihood-to-evidence ratio. Instead of sampling from the posterior, neural SBI directly learns an estimator by training neural networks on the full input data. This makes the need for handcrafted summary statistics unnecessary. Furthermore, it has the advantage over ABC that it is *amortized*. Amortized methods learn to estimate all posteriors that are supported by the prior, while non-amortized methods are designed to approximate only a single posterior. One of the methods that enjoy these benefits is the NRE ([Hermans et al., 2020; Miller et al., 2020](#)) algorithm, where the likelihood-to-evidence ratio is estimated. As discussed in the introduction, NRE can be extended with marginalization and truncation ([Miller et al., 2021](#)). In this work, we will only extent with marginalization. Below we will discuss the technicalities of this technique, MNRE.

3.1.2 Marginal Neural Ratio Estimation

The goal with NRE is to directly learn the likelihood-to-evidence ratio

$$r(\mathbf{x}, \boldsymbol{\vartheta}) \equiv \frac{p(\mathbf{x}|\boldsymbol{\vartheta})}{p(\mathbf{x})} = \frac{p(\mathbf{x}, \boldsymbol{\vartheta})}{p(\mathbf{x})p(\boldsymbol{\vartheta})} = \frac{p(\boldsymbol{\vartheta}|\mathbf{x})}{p(\boldsymbol{\vartheta})}. \quad (3.4)$$

Given that one knows the prior $p(\boldsymbol{\vartheta})$, one can directly find the posterior from the ratio without calculating the likelihood through Equation (3.2).

¹¹A summary statistic first chosen to describe to data. These summary statistics are compared against each other by the distance measure.

Following [Hermans et al. \(2020\)](#), this ratio is obtained by comparing samples jointly drawn $(\mathbf{x}, \boldsymbol{\vartheta}) \sim p(\mathbf{x}, \boldsymbol{\vartheta})$ (labeled with class $y = 1$) and marginally drawn $(\mathbf{x}, \boldsymbol{\vartheta}) \sim p(\mathbf{x})p(\boldsymbol{\vartheta})$ (labeled with class $y = 0$) from the simulator. The jointly drawn parameters have a matching pair of sample and parameter. This is not the case for marginally drawn pairs, where the sample is not produced by the parameter passed to the simulator. The optimal binary classifier $d^*(\mathbf{x}, \boldsymbol{\vartheta})$ describes the probability of class 1:

$$\begin{aligned} d^*(\mathbf{x}, \boldsymbol{\vartheta}) &= p(y = 1 | \mathbf{x}, \boldsymbol{\vartheta}) \\ &= \frac{p(\mathbf{x}, \boldsymbol{\vartheta} | y = 1)p(y = 1)}{p(\mathbf{x}, \boldsymbol{\vartheta} | y = 0)p(y = 0) + p(\mathbf{x}, \boldsymbol{\vartheta} | y = 1)p(y = 1)} \\ &= \frac{p(\mathbf{x}, \boldsymbol{\vartheta})}{p(\mathbf{x})p(\boldsymbol{\vartheta}) + p(\mathbf{x}, \boldsymbol{\vartheta})}. \end{aligned} \quad (3.5)$$

By rewriting this we obtain the likelihood ratio as

$$r(\mathbf{x}, \boldsymbol{\vartheta}) \stackrel{\text{Eq. (3.4)}}{=} \frac{p(\mathbf{x}, \boldsymbol{\vartheta})}{p(\mathbf{x})p(\boldsymbol{\vartheta})} \stackrel{\text{Eq. (3.5)}}{=} \frac{d^*(\mathbf{x}, \boldsymbol{\vartheta})}{1 - d^*(\mathbf{x}, \boldsymbol{\vartheta})}. \quad (3.6)$$

This is known in the literature as the likelihood ratio trick ([Cranmer et al., 2016](#)).

Now we have rephrased the goal to estimate the optimal binary classifier instead of the likelihood-to-evidence ratio. To estimate the optimal classifier $d^*(\mathbf{x}, \boldsymbol{\vartheta})$ we train a binary classifier $d_\phi(\mathbf{x}, \boldsymbol{\vartheta}) \in [0, 1]$, with ϕ the network parameters, using stochastic gradient descent to minimize the binary-cross entropy function

$$L = - \int [p(\mathbf{x}, \boldsymbol{\vartheta}) \log(d_\phi(\mathbf{x}, \boldsymbol{\vartheta})) + p(\mathbf{x})p(\boldsymbol{\vartheta}) \log(1 - d_\phi(\mathbf{x}, \boldsymbol{\vartheta}))] d\boldsymbol{\vartheta} d\mathbf{x} \quad (3.7)$$

This binary-cross entropy loss function is typically used in classification problems. We compare a jointly drawn pair with a marginally drawn pair. So the network learns to classify whether a parameter matches with an observation or not. Minimizing the loss with respect to the network parameters gives

$$\frac{\partial}{\partial \phi} L[d_\phi(\mathbf{x}, \boldsymbol{\vartheta})] = -\frac{\partial}{\partial \phi} \int [p(\mathbf{x}, \boldsymbol{\vartheta}) \log(d_\phi(\mathbf{x}, \boldsymbol{\vartheta})) + p(\mathbf{x})p(\boldsymbol{\vartheta}) \log(1 - d_\phi(\mathbf{x}, \boldsymbol{\vartheta}))] d\boldsymbol{\vartheta} d\mathbf{x} \quad (3.8)$$

$$= - \int \left[\frac{p(\mathbf{x}, \boldsymbol{\vartheta})}{d_\phi(\mathbf{x}, \boldsymbol{\vartheta})} - \frac{p(\mathbf{x})p(\boldsymbol{\vartheta})}{1 - d_\phi(\mathbf{x}, \boldsymbol{\vartheta})} \right] \frac{\partial d_\phi(\mathbf{x}, \boldsymbol{\vartheta})}{\partial \phi} d\boldsymbol{\vartheta} d\mathbf{x} \quad (3.9)$$

Demanding $\frac{\partial}{\partial \phi} L[d_\phi(\mathbf{x}, \boldsymbol{\vartheta})] = 0$, we can set the expression within the square brackets to zero, resulting in the desired expression, Equation (3.5),

$$d_\phi(\mathbf{x}, \boldsymbol{\vartheta}) \approx \frac{p(\mathbf{x}, \boldsymbol{\vartheta})}{p(\mathbf{x})p(\boldsymbol{\vartheta}) + p(\mathbf{x}, \boldsymbol{\vartheta})}. \quad (3.10)$$

In practice, we parametrize the classifier $d_\phi(\mathbf{x}, \boldsymbol{\vartheta}) = \sigma(f_\phi(\mathbf{x}, \boldsymbol{\vartheta}))$. With $f_\phi(\mathbf{x}, \boldsymbol{\vartheta})$ the output of the neural network and σ the logistic sigmoid $\sigma(x) = \frac{e^x}{1+e^x}$. Combining this with Equation (3.6), we have implied that

$$f_\phi(\mathbf{x}, \boldsymbol{\vartheta}) = \log r(\mathbf{x}, \boldsymbol{\vartheta}). \quad (3.11)$$

Using the logarithmic ratio improves numerical stability during training.

The above procedure can be extended to the estimation of marginals. Nuisance parameters, the parameters that need to be marginalized over, are sampled but not presented to the classifier. So they are being passed to the simulator, but they are not labeled with either class $y = 0, 1$. Because the

parameters are randomly drawn from the prior, the nuisance parameters are marginalized by random sampling. With the model $p(\mathbf{x}, \boldsymbol{\theta}) = p(\mathbf{x}, \boldsymbol{\eta}, \boldsymbol{\vartheta}) = p(\mathbf{x}|\boldsymbol{\eta}, \boldsymbol{\vartheta})p(\boldsymbol{\eta}, \boldsymbol{\vartheta})$, the loss functions becomes

$$L = - \int [p(\mathbf{x}, \boldsymbol{\eta}, \boldsymbol{\vartheta}) \log (d_\phi(\mathbf{x}, \boldsymbol{\vartheta})) + p(\mathbf{x})p(\boldsymbol{\eta}, \boldsymbol{\vartheta}) \log (1 - d_\phi(\mathbf{x}, \boldsymbol{\vartheta}))] d\boldsymbol{\vartheta} d\mathbf{x} d\boldsymbol{\eta} \quad (3.12)$$

$$= - \int [p(\mathbf{x}, \boldsymbol{\vartheta}) \log (d_\phi(\mathbf{x}, \boldsymbol{\vartheta})) + p(\mathbf{x})p(\boldsymbol{\vartheta}) \log (1 - d_\phi(\mathbf{x}, \boldsymbol{\vartheta}))] d\boldsymbol{\vartheta} d\mathbf{x}. \quad (3.13)$$

Resulting in the same loss function if there were no nuisance parameters, Equation (3.7).

Finally, one can directly obtain the posteriors by multiplying the results of the network with the posterior:

$$p(\boldsymbol{\vartheta}|\mathbf{x}) = e^{f_\phi(\mathbf{x}, \boldsymbol{\vartheta})} p(\boldsymbol{\vartheta}). \quad (3.14)$$

3.2 Image Segmentation

Image segmentation is a technique used for object detection by classifying every pixel in an image (Minaee et al., 2020). It divides the image into different parts, where each part belongs to a certain class. The classes are user-defined and are dependent on the task. Typical tasks of image segmentation are detecting traffic (Cordts et al., 2016) (to improve for example software of self-driving cars) or remote sensing (Yuan et al., 2021), the classification of aerial photos. A typical segmenting model has some characteristic features. Firstly, there are no fully connected layers, all layers are fully convolutional. Secondly, there is an encoder-decoder structure. By both taking the information from the encoder and decoder into account, the model is able to extract features across different scales.

3.2.1 U-Net

A well recognized neural network architecture for image segmentation is the U-Net, originally designed for biomedical image segmentation by Ronneberger et al. (2015). Adoptions of it are widely used, for example in Imagen (Saharia et al., 2022), the latest text-to-image diffusion model of Google. A U-Net has the typical encoder-decoder structure, also known as contracting-expanding or downsampling-upsampling structure. This architecture makes it possible to be sensitive to small objects in images. Below, we will discuss the most important features of the network, as shown in Figure 3.1.

The contracting path is similar to a ‘traditional’ convolutional neural network. It takes as input the image and can be divided into blocks, consisting of two sequential convolutions (the blue arrow in Figure 3.1) and the downsampling (the red arrow).

The blue arrow indicates three operations. The first operation is a 2D convolution. It applies a 3×3 kernel with stride, padding, and dilatation set to 1. The kernel consists of weights that are being learned. The second operation is batch normalization (Ioffe & Szegedy, 2015). This was not in the original U-Net by Ronneberger et al. (2015), as batch normalizations were not developed yet. It regularizes the network by normalizing each training mini-batch, allowing faster training and higher learning rates. Finally, the rectified linear unit (ReLU) (Nair & Hinton, 2010) is used as the activation function. The downsampling is performed after the convolutions, depicted by the red arrow. It halves the number of pixels with a 2×2 max pooling operation with a stride of 2.

Note that sequential convolutional steps are set such that the number of layers is doubled after the first block, and that the number of pixels is preserved. This is indicated by the numbers on top of each layer in black, and the height and width at the beginning of each row in gray, respectively.

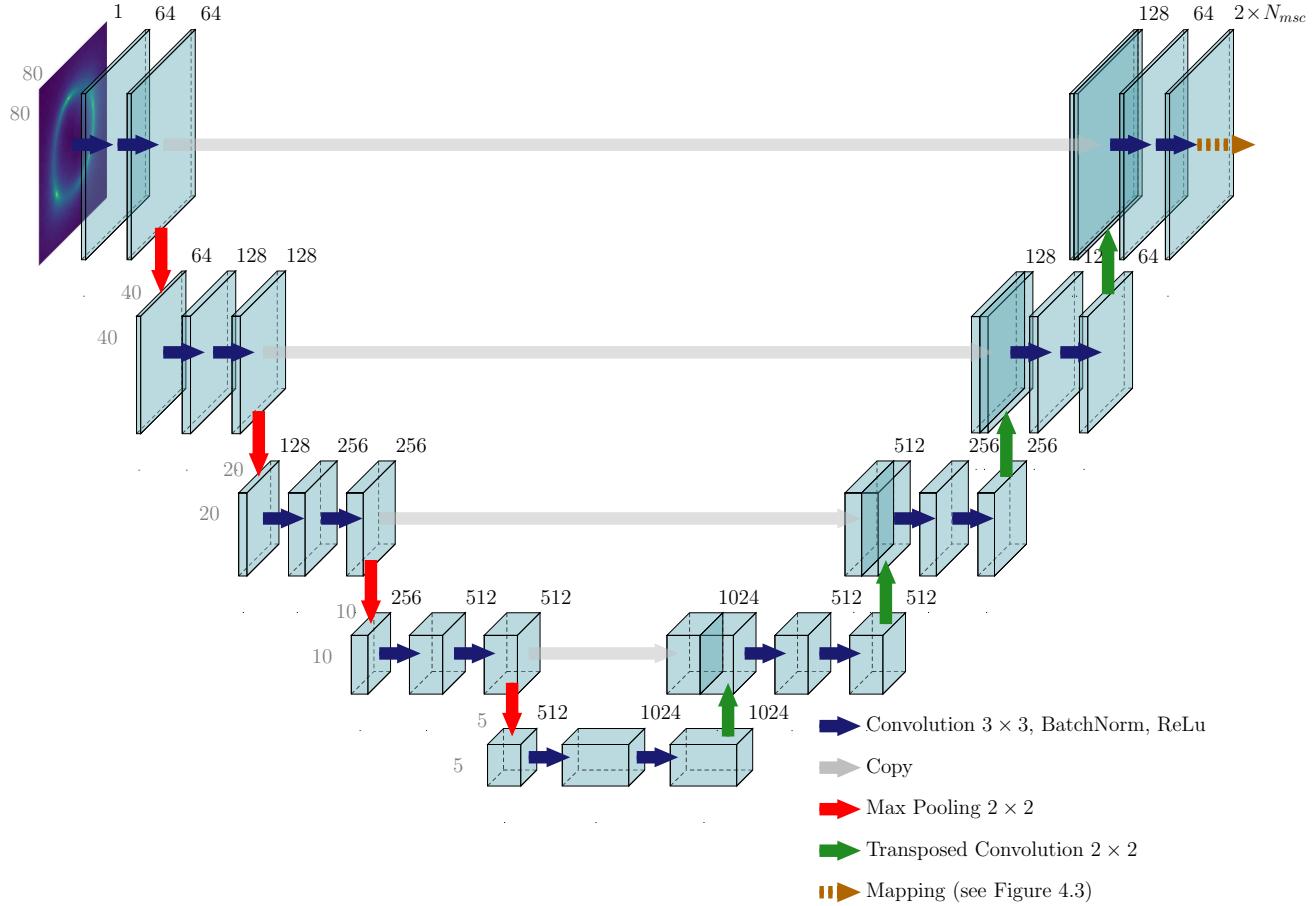


Figure 3.1: The U-Net architecture. The image takes a single normalized image, and returns output layers, with the same size as the input image, for each mass channel. The network consists of two parts. The first part is going down, the contracting path. It has sequential ordinary convolutional blocks (blue arrows) and max pooling layers (red). The second part is going up, the expanding path. The max pooling layers are now replaced by transposed convolutions (green). These convolutions upsample the observations again to the same size as the initial input. The U-Net is characterized by the concatenation (gray) of the contracting and expanding outputs. The numbers in black above the layers indicate the number of filters for that layer. The values in gray denote the size of those layers.

The expanding path starts after four downsampling blocks. Each expanding block consists of an upsampling step (green arrow), a concatenation (grey arrow), and again two sequential convolutions.

The upsampling step (in green) is a transposed convolution. It has, similar to the max pooling operation, a 2×2 kernel size and stride of 2. Therefore, it halves the number of layers and doubles the number of pixels. The output of each second convolution during the contracting path is stored and now concatenated to the output of the corresponding layer of the upsampling. This is depicted by the gray array. Finally, the expanding block ends with two sequential convolutions, similar to the contracting path.

The combination of the contracting and expanding path makes the U-Net architecture behave well for image segmentation tasks. The first convolutional layers learn about the small-scale information in the image. The information in the image is ‘boiled down’ to its large-scale features when repeating the convolutional steps and downsampling. This information is passed on to the final output by the upsampling path. The concatenation of the small and large-scale features is the main U-Net strength. During the upsampling path, the network is able to localize new features with this combined information.

The number of input and output layers must be specified by the user. The number of channels of a normal photo is three (the RGB layers). That is of the same order as a typical space telescope,¹² although its channels range from optical to infrared wavelengths. In most cases, the characteristics of certain objects are more pronounced in specific channels. Hence, having more input layers helps the network to distinguish between objects. The number of output layers is determined by the number of classes the user wants to differentiate between. For example, if a network has to differentiate if a pixel has a pedestrian or not, two classes are needed. We are interested if a pixel contains a subhalo of a certain mass, where we differentiate the masses between N_{msc} number of mass channels. Therefore the number of output channels for us is $2 \times N_{\text{msc}}$. This will be discussed in more detail in the next section.

¹²See for example the instrumentation of [HST](#) (three channels), the [Euclid VIS](#) (one channel) and the [James Webb Space Telescope](#) (three channels)

Chapter 4

Subhalo Density Estimation

This chapter uses the concepts of SBI and the U-Net, introduced in the previous chapter, to develop a definition of the pixel posterior probability. With the pixel posterior probability, one can estimate a subhalo density in the lens.

Our goal is to approximate the posterior of multiple subhalos

$$p(\boldsymbol{\vartheta}_{\text{subs}} | \mathbf{x}_{\text{obs}}) = p(\boldsymbol{\vartheta}_{\text{sub},1}, \boldsymbol{\vartheta}_{\text{sub},2}, \dots, \boldsymbol{\vartheta}_{\text{sub},N_{\text{sub}}} | \mathbf{x}_{\text{obs}}), \quad (4.1)$$

with the parameters of interest $\boldsymbol{\vartheta}_{\text{sub},i} = (x_{\text{sub},i}, y_{\text{sub},i}, m_{\text{sub},i})$, i.e. the parameters of the i^{th} subhalo. $m_{\text{sub},i}$ is the mass of the subhalo, $(x_{\text{sub},i}, y_{\text{sub},i})$ is the position of the subhalo in the lens plane. When determining the posterior of multiple subhalos $p(\boldsymbol{\vartheta}_{\text{subs}} | \mathbf{x}_{\text{obs}})$, two problems arise:

- (i) **The unknown number of subhalos** makes it impossible to exactly determine how many posteriors should be calculated. Current subhalo models contain a stochastic number of subhalos. Therefore one can only tell in advance the expected number of subhalos $\langle N_{\text{sub}} \rangle$, but not the actual number of subhalos N_{sub} .
- (ii) **The interchangeability of the subhalos** causes the label switching problem. This problem arises when the ordering of parameters is arbitrary. There is no physical ordering in the subhalos, as all subhalos are described by the same parameters $\boldsymbol{\vartheta}_{\text{sub}} = (x_{\text{sub}}, y_{\text{sub}}, m_{\text{sub}})$. The target, the true subhalos' position and mass, is therefore invariant under a permutation. In the case that there are N_{sub} subhalos, and the posterior of a parameter has a peak at a particular point in the parameter space, it will necessarily have peaks at all $N_{\text{sub}}!$ points related by symmetry.

Let us illustrate the above statement with an example. Consider that we inferring only the mass of three subhalos (m_1, m_2, m_3) with true masses (m_A, m_B, m_C) , respectively. Ideally, the joint posterior $p(m_1, m_2, m_3 | \mathbf{x})$ would have the highest posterior density region around $(m_1, m_2, m_3) = (m_A, m_B, m_C)$. However, as there is no difference for the network between the $6!$ permutations of the ordering of (m_A, m_B, m_C) , the estimated joint posterior would similarly have high posterior density regions around $(m_1, m_2, m_3) = (m_A, m_C, m_B)$ and all the other permutations. Therefore, the joint posterior would have $6!$ modes.

Originally, solutions of this problem have been studied in the context of mixture models and MCMC (Celeux, 1998; Stephens, 2000; Jasra et al., 2005) and more recently in the context of gravitational wave astronomy (Buscicchio et al., 2019). Solutions to this problem include imposing an artificial identifiability constraint (Diebolt & Robert, 1994), where an artificial constraint such as $m_1 < m_2 < m_3$ is placed, and relabelling algorithms (Stephens, 1997), where the labels are correctly re-labeled with k -means clustering.

A simplified visualization of these problems is shown in Fig. 4.1. The solutions that are used to solve the label switching problem can not directly be applied to our problem of subhalo measurements. Typical solutions solve the label switching problem in the case where the number of interchangeable objects is fixed and relatively low.¹³ However, the number of subhalos is not fixed in a lensing system. Furthermore, the number of subhalos can be significantly higher than the number of parameters which the label switching problems were designed for. Therefore, we introduce another approach to solve the label switching problem. This will be discussed in the remainder of this chapter.

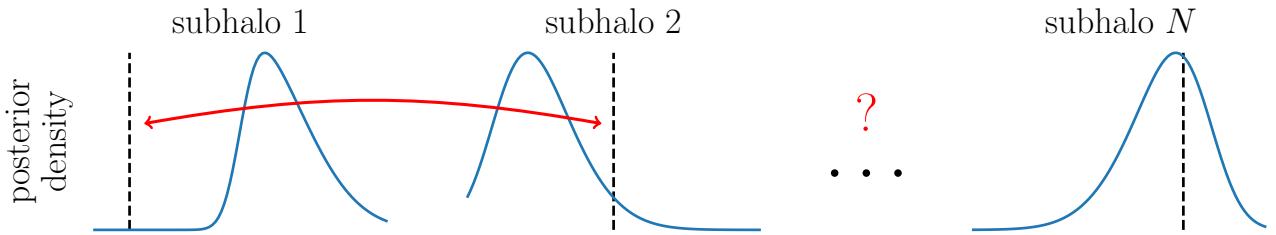


Figure 4.1: When determining posteriors of multiple subhalos one runs into problems because of: (i) the unknown number of subhalos (and therefore posteriors), one can only determine $\langle N_{\text{sub}} \rangle$ (visualized by the three dots and question mark); (ii) the label switching problem causes the prediction to not match the true value (a possible wrong match visualized by the red arrow).

4.1 Pixel posterior probability

Instead of trying to solve the problems when calculating an unknown number of interchangeable subhalos, we rephrase the goal. We divide the parameter space of the prior in a hypergrid of pixels. The hypergrid, visualized in Figure 4.2, has a dimensionality of $\mathbb{Z}^{N_{\text{par}}}$, with N_{par} the number of parameters of interest. The dimensions of the hypergrid are the parameters of interest ($x_{\text{sub}}, y_{\text{sub}}, m_{\text{sub}}$), hence $N_{\text{par}} = 3$.

The number of pixels in each direction of the hypergrid are $(N_{\text{pix},x}, N_{\text{pix},y}, N_{\text{msc}})$. Where N_{pix} is the number of pixels we have along that observation dimension, and N_{msc} is the number of mass channels we want to differentiate the subhalos between. The number of pixels along the (x, y) direction is determined by the number of pixels of the observation. The number of mass channels is set by the user itself. With $N_{\text{msc}} = 1$, one would consider all subhalos without differentiating their masses. With $N_{\text{msc}} > 1$, one divides the prior mass range into multiple mass channels. The goal is then to determine whether a subhalo belongs to a mass channel, i.e. if its mass is within the mass ranges of that mass channel.

For each pixel on the grid, we will determine the probability that *there is* a subhalo in that pixel. That probability of that pixel is what we will call the *pixel posterior probability*. Because it describes the probability that there is a subhalo in a pixel, it can be interpreted as the zeroth-order marginal posterior distribution, where it is marginalized over the values between the pixel boundaries. All the pixels together describe the subhalo number density on the hypergrid. The pixel posterior probability is a newly defined quantity, and therefore needs a careful introduction. We describe the pixel posterior probability mathematically as

$$p(\tilde{\vartheta}_{ijk}^1 | \mathbf{x}_{\text{obs}}) \quad (4.2)$$

¹³This makes sense because MCMC problems are by construction seldom high-dimensional. High-dimensional MCMC problems are often computationally intractable.

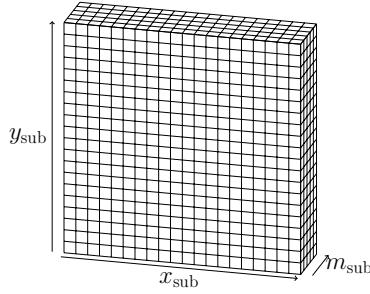


Figure 4.2: Visualization of the hypergrid where the predicted pixel posterior probability would live on. The number of pixels in the (x, y) -directions, $(N_{\text{pix},x}, N_{\text{pix},y})$, is determined by the number of pixels of the observation. The number of pixels in the m -direction, N_{msc} , is determined by the user. It specifies the mass channels where the network should differentiate between.

Let us discuss the sub- and superscripts around ϑ in more detail.:

- $\tilde{\cdot}$ represents a variable in pixel space, i.e. it is describing a pixel on the hypergrid.
- \cdot^1 means that we are talking about the probability that *there is* a subhalo in that pixel. This is contrary to \cdot^0 , describing the probability that *there is not* a subhalo in that pixel. The usage of this will become clear in Section 4.2.3.
- \cdot_{ijk} are the indices of the pixel on the hypergrid. It ranges from the first pixel on the hypergrid to the final pixel, i.e. $(0, 0, 0) < (i, j, k) \leq (N_{\text{pix},x}, N_{\text{pix},y}, N_{\text{msc}})$

To obtain the pixel posterior probability we implement the U-Net into the MNRE scheme. This work combines these techniques. In order for that to work, we introduce four new concepts: (i) the pixelating of the observation by the U-Net; (ii) the mapping to select the ‘correct’ pixels to train on; (iii) the normalization to interpret the posteriors as probabilities; (iv) the calibration to ensure that the probabilities reflect correct values. We will discuss all these concepts in the sections below. These concepts are, together with the simulating step, connected by the flowchart shown in Figure 4.3.

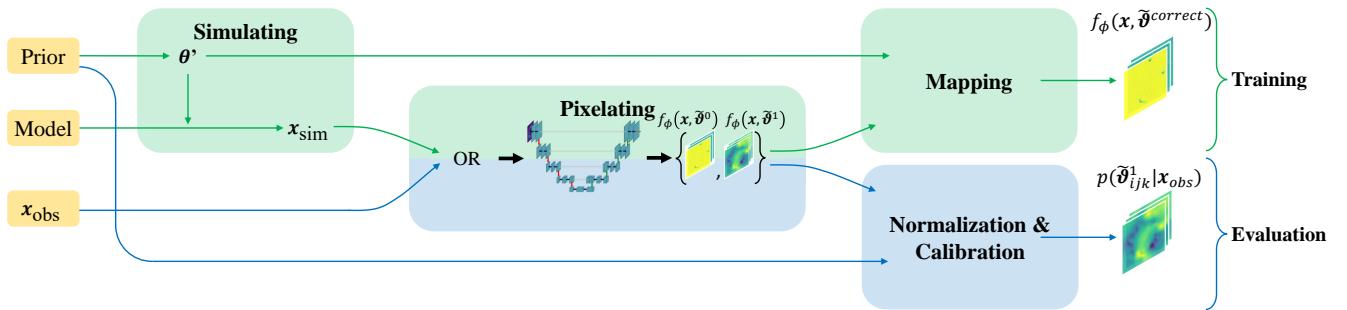


Figure 4.3: Flowchart of this work, visualizing the training (in green) and evaluation (in blue) steps. Simulated observation x_{sim} are passed to the U-Net, which results in two hypergrids $\{f_{\phi}(x, \tilde{\vartheta}^0), f_{\phi}(x, \tilde{\vartheta}^1)\}$. The mapping procedure uses those hypergrids to select the ‘correct’ pixels $f_{\phi}(x, \tilde{\vartheta}^{\text{correct}})$, which are passed to the loss function. After training the evaluation starts. An actual (mock) observation x_{obs} is passed to the U-Net. Instead of mapping, the two hypergrids are now the ingredients for the normalization and calibration procedures, resulting in the pixel posterior probability $p(\tilde{\vartheta}^1_{ijk} | x_{\text{obs}})$.

4.2 Training pipeline

4.2.1 Simulating

The training data for the network consists of simulated data. This data is produced by the model, Equation (2.20), which functions as the simulator. The simulator takes as input all the model parameters θ , which are sampled from their priors:

$$\theta' \sim p(\theta), \quad (4.3)$$

$$\mathbf{x}_{\text{sim}} = \text{Simulator}(\theta'). \quad (4.4)$$

With θ' the sampled parameters. The exact distribution and ranges chosen for the priors are listed in Section 5.1, where the results are discussed.

4.2.2 Pixelating

The simulated data is the training data for the network. It is normalized before it is passed through the network. The standard Z-score normalization

$$z = \frac{x - \mu}{\sigma} \quad (4.5)$$

is used. Here we are using z for the normalized data, x for the input data and finally μ and σ for the mean and standard variance of all the data.¹⁴ After the normalization, the image is passed through the U-Net. See Section 3.2.1 for a further discussion about the network architecture. The network's output is chosen such that it represents two hypergrids. The first hypergrid relates to the probabilities that there *is* a subhalo in a pixel of the hypergrid, while the second relates to *no* subhalo in a pixel. This means that there are $2 \times N_{\text{msc}}$ output channels, where each channel is an image with the same size as the input image. The outputted hypergrids are the (logarithmic) ratios

$$\left\{ f_\phi(\mathbf{x}, \tilde{\vartheta}^0), f_\phi(\mathbf{x}, \tilde{\vartheta}^1) \right\} = \text{Network}(\mathbf{x}). \quad (4.6)$$

Note that we are now using the parameters of interest ϑ and not all parameters θ . Although the simulation or observation is determined by all parameters θ , the network is trained to only distinguish between the parameters of interest ϑ .

4.2.3 Mapping

The two hypergrid are now used to select the ‘correct’ pixel to train on. ‘Correct’ in the way that: if there is *no* subhalo in a pixel, we select $f_\phi(\mathbf{x}, \tilde{\vartheta}^0)$ for that pixel. If there *is* a subhalo, we select $f_\phi(\mathbf{x}, \tilde{\vartheta}^1)$ for that pixel. We do this by defining a binary target map, which maps the true parameters of the subhalo on the binary target map:

$$\vartheta_{\text{subs}} \rightarrow \tilde{\mathbf{z}}. \quad (4.7)$$

The binary target map indicates the pixels that contain a subhalo. We do not take multiple subhalos in the same pixel into consideration. Taking multiple subhalos in the same pixel into account would make the problem unnecessarily more complicated, as it is very unlikely that this would happen. A pixel has

¹⁴Practically, we are using a [parallel algorithm](#) to calculate the variance.

a value of one when the subhalo's parameters fall within the boundaries of that pixel, otherwise, it has a value of zero:

$$\tilde{z}_{ijk} = \begin{cases} 0 & \text{if no subhalo in pixel } (i, j, k), \\ 1 & \text{if subhalo in pixel } (i, j, k). \end{cases} \quad (4.8)$$

The correct pixels are now selected following

$$f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}_{ijk}^{\text{correct}}) = \begin{cases} f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}_{ijk}^0) & \text{if } \tilde{z}_{ijk} = 0, \\ f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}_{ijk}^1) & \text{if } \tilde{z}_{ijk} = 1. \end{cases} \quad (4.9)$$

At this point MNRE and image segmentation come together. Each pixel is now correctly segmented, it is assigned with the correct ratio estimator. These pixels are now passed to the binary-cross entropy loss, Equation (3.7). What practically happens is that the correct pixel is compared against a randomly drawn other pixel. This is the binary classification problem in NRE. The loss updates the weights of the network and a new batch of simulations is put through the pixelating and mapping procedure.

4.3 Evaluation

If the loss has converged after a sufficient number of training iterations,¹⁵ the evaluation stage starts, shown in blue in Figure 4.3. Instead of a simulated image, an actual observation \mathbf{x}_{obs} can be put through the network.¹⁶ The network weights are not updated anymore. We can not even do the mapping procedure because we do not know what the ‘correct’ pixels would be. Instead, we can directly obtain the pixel posteriors by multiplying the network output with the corresponding pixel prior, following Equation (3.14),

$$p(\tilde{\boldsymbol{\vartheta}}_{ijk}^c | \mathbf{x}_{\text{obs}}) = e^{f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}_{ijk}^c)} p(\tilde{\boldsymbol{\vartheta}}_{ijk}^c), \quad c = \{0, 1\}. \quad (4.10)$$

The pixel prior $p(\tilde{\boldsymbol{\vartheta}}_{ijk}^1) = 1 - p(\tilde{\boldsymbol{\vartheta}}_{ijk}^0)$ depends on the subhalo modeling in the simulator. Since the subhalo position is uniformly sampled, this is straightforward. In the case that the subhalo mass is also uniformly sampled, the prior is just

$$p(\tilde{\boldsymbol{\vartheta}}_{ijk}^1) = \frac{\langle N_{\text{sub}} \rangle}{N_{\text{pix},x} \times N_{\text{pix},y} \times N_{\text{msc}}}. \quad (4.11)$$

If the subhalo mass is sampled from the SHMF, each mass channel k is weighted by the probability $M_{\text{frac},k}$ that a subhalo has a mass that falls within the mass boundaries of channel k :

$$p(\tilde{\boldsymbol{\vartheta}}_{ijk}^1) = \frac{\langle N_{\text{sub}} \rangle}{N_{\text{pix},x} \times N_{\text{pix},y} \times M_{\text{frac},k}}. \quad (4.12)$$

This channel weight is calculated my integrating the SHMF dn/dM from the lower to upper mass boundary ($m_{\text{low},k}, m_{\text{high},k}$) of the channel:

$$M_{\text{frac},k} = \frac{M_k}{\sum_k M_k} = \frac{\int_{m_{\text{low},k}}^{m_{\text{high},k}} \frac{dn(M')}{dM} dM'}{\int_{m_{\text{low}}}^{m_{\text{high}}} \frac{dn(M')}{dM} dM'}. \quad (4.13)$$

¹⁵The requirements for stopping training are discussed when the results are presented in Section 5.1.

¹⁶In this work however, we will not apply the pipeline on real observations. Instead, we apply it on a simulated mock observation.

Where $(m_{\text{low}}, m_{\text{high}})$ are the lower and upper mass region of the subhalo mass prior. The numerator is the mass that is contained in mass channel k according to the SHMF, and is normalized by the total mass that the SHMF contains if integrated from the lower to upper subhalo mass that we consider.

Now that we are able to obtain the pixel posterior $p(\tilde{\boldsymbol{\vartheta}}_{ijk}^c | \mathbf{x}_{\text{obs}})$ directly with Equation (4.10) we verify that it can be interpreted as a probability. To do that, we normalize and calibrate.

4.3.1 Normalization

During training, the network receives no information that the output should represent probabilities. Therefore, we normalize the probability that there is a subhalo in that pixel with

$$p_{\text{norm}}(\tilde{\boldsymbol{\vartheta}}_{ijk}^1 | \mathbf{x}_{\text{obs}}) = \frac{p(\tilde{\boldsymbol{\vartheta}}_{ijk}^1 | \mathbf{x}_{\text{obs}})}{p(\tilde{\boldsymbol{\vartheta}}_{ijk}^0 | \mathbf{x}_{\text{obs}}) + p(\tilde{\boldsymbol{\vartheta}}_{ijk}^1 | \mathbf{x}_{\text{obs}})}, \quad (4.14)$$

\leftrightarrow

$$p_{\text{norm}}(\tilde{\boldsymbol{\vartheta}}_{ijk}^0 | \mathbf{x}_{\text{obs}}) + p_{\text{norm}}(\tilde{\boldsymbol{\vartheta}}_{ijk}^1 | \mathbf{x}_{\text{obs}}) = 1. \quad (4.15)$$

This normalization forces the pixel posterior probability to be $0 < p_{\text{norm}}(\tilde{\boldsymbol{\vartheta}}_{ijk}^1 | \mathbf{x}_{\text{obs}}) < 1$. To avoid notational clutter, we drop the subscript \cdot_{norm} . One can now confuse the normalized pixel posterior probability and non-normalized pixel posterior, but one can assume that we always refer to the normalized case if not specified otherwise.

4.3.2 Validation & calibration

To verify that the pixel posterior probability actually represents correct probabilities we perform a ‘binned coverage test’. With the results of that test, we perform afterwards calibration if necessary. The test is a binned version of the original ‘continuous’ coverage test. This test calculates the expected empirical coverage for each confidence level. We refer to Section 2.3 of Cole et al. (2021) for an illustrative explanation how these expected coverage tests are calculated, and Hermans et al. (2021) for its implications on SBI.

A pixel with an estimated probability of p should contain an actual subhalo in $p \times 100\%$ of the cases. This is what the binned coverage test shows in the top right panel of Figure 4.4. Let us discuss the top row of Figure 4.4 to explain how we obtain the coverage plot. First, we produce N_{calib} predictions from simulated observations.¹⁷ The value of all these $N_{\text{calib}} \times N_{\text{pix},x} \times N_{\text{pix},y} \times N_{\text{msc}}$ predicted pixels are put into a histogram, shown in green. Secondly, we count how many times these pixels actually contain a subhalo. These counts are also put into a histogram, shown in orange. Finally, we compute the empirical pixel posterior probability by dividing the histogram with all the pixels by the histogram with pixels only containing a subhalo.

The empirical pixel posterior probability, shown in blue, should lie on the diagonal straight line indicated by the dashed black line, i.e. $ppp_{\text{emp}} = ppp_{\text{pred}}$, where ppp stands for pixel posterior probability and subscripts \cdot_{emp} and \cdot_{pred} for empirical and predicted, respectively. This means that all pixels with a value of ppp_{pred} should contain a subhalo in ppp_{emp} of the times. The predictions are overconfident if the empirical pixel posterior probability lies *under* the diagonal. This is because the pixels having a specific prediction contain a subhalo less often than they should. If the empirical pixel posterior probability lies *above* the diagonal, the predictions are conservative. The predictions contain a subhalo more often

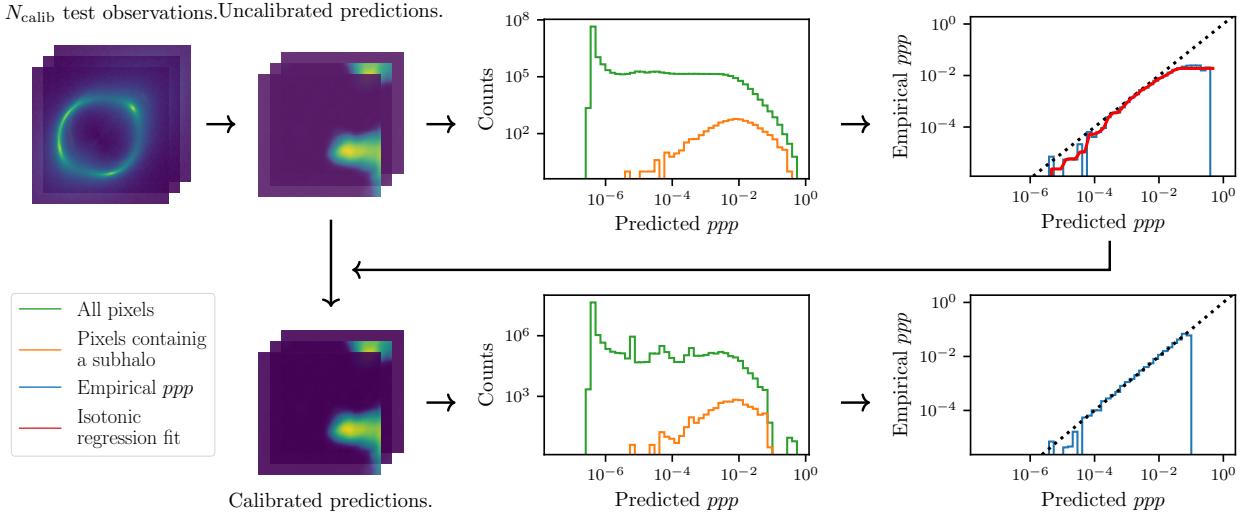


Figure 4.4: Example of the binned coverage test. All pixels of N_{calib} predictions are put into a histogram (blue), the subset of those pixels containing a subhalo selected (orange). Dividing these two histograms gives the empirical pixel posterior probability (green). Pixel posterior probability is abbreviated as ppp .

compared to how often they actually contain a subhalo.

In the particular mock example of Figure 4.4 we see that the empirical pixel posterior probability in the top right panel does not lie on the diagonal for high values. Therefore, we perform calibration. The calibration ‘updates’ the values of the predicted pixels. We do this by fitting the uncalibrated empirical pixel posterior probability, and using this fitted function as the transformation of uncalibrated pixel values to calibrated values.

The fitting technique that we use is *isotonic regression*. This technique, originally developed by [Robertson et al. \(1988\)](#) and reviewed by [Nuesch \(1991\)](#), has been generalised by [Zadrozny & Elkan \(2001, 2002\)](#) to calibrate predictions from various supervised machine learning algorithms. It is a general technique in the sense that the only restriction is that the mapping function is isotonic, i.e. monotonically increasing ([Niculescu-Mizil & Caruana, 2005](#)). That means that we find a non-decreasing function $f(\text{ppp}_{\text{pred}})$ where the isotonic regression task is to define the function such that

$$\min_f \sum_i (f(\text{ppp}_{\text{pred},i}) - \text{ppp}_{\text{emp},i})^2, \quad (4.16)$$

where i is being summed over the bins of the histogram. This function is subject to $f(\text{ppp}_{\text{pred},i}) \leq f(\text{ppp}_{\text{pred},j})$ for all $i < j$, such that the function can only increase for higher ppp_{pred} . We implement this technique with `scikit-learn`.¹⁸ The result of the isotonic regression is shown by the red line in the top right panel of Figure 4.4.

The fit is used to calibrate the predicted pixel posterior probabilities. The calibrated values are obtained through $\text{ppp}_{\text{calibrated}} = f(\text{ppp}_{\text{uncalibrated}})$, with f the fitted function. If we now repeat the validation step, as shown on the bottom row of Figure 4.4, we see that the empirical pixel posterior probability lies now on the diagonal.

¹⁷The number of simulations N_{calib} should be sufficiently high to get enough statistics. However, too high values result in a longer duration of the calibration stage.

¹⁸See the [documentation](#) of `scikit-learn` for examples of isotonic regression.

Chapter 5

Results

5.1 Mock HST observation

The simulations, that functions as the training data for our network, should be representative of the data that telescopes would produce. We chose to adapt our pipeline to apply on HST data. In this work, however, we still apply the pipeline on mock observations. HST data has an expected resolution of $0.04''$. However, we use a slightly lower resolution of $0.05''$ to neglect the PSF, as discussed in Section 2.4. We chose an image size of 100×100 pixels such that the image covers an area of $5'' \times 5''$ on the sky. The Gaussian noise is set such that the signal-to-noise ratio is ~ 30 , similar to HST data.

The parameters that describe the lensing system, i.e. the lens and source parameters $\theta_{\text{lens}}, \theta_{\text{source}}$ are listed in Table 5.1. Since we are using the same simulator as Coogan et al. (2022), we chose the same values for the lensing parameters so that we can do a comparison in specific cases. An example of an image that is produced with these parameters (together with subhalo modeling) is shown in top left panel of Figure 5.1.

The lensing and subhalo modeling is built with pytorch (Paszke et al., 2019) so that we can use GPUs to simulate a large number of observations. pytorch is auto-differentiable, making it possible in future work to use analysis methods that require autodifferentiation. The analysis pipeline uses swyft¹⁹ (Miller et al., 2020, 2021) to implement MNRE. The open-source software swyft is similarly built on pytorch and pytorch-lightning (Falcon & Cho, 2020)

The training, implemented with swyft, is initialized as follows. The Adam optimizer (Kingma & Ba, 2014) is used to minimize the loss with an initial learning rate of 6×10^{-3} . The batch size is, rather small, set to 16. This is to prevent memory overflow of the GPUs when predicting a high number of pixels. An additional advantage is that a low batch size might increase inference quality (Keskar et al., 2017). When the validation loss has not improved 5 sequential epochs, the learning rate is reduced by a factor of 0.1. The training is stopped after 30 epochs, or when the validation loss has not been improved for 10 sequential epochs.

¹⁹<https://github.com/undark-lab/swyft>

Table 5.1: Source and lens parameters that are chosen to model the gravitational lensing system. The value for the subhalo parameters are specified in the main text.

Parameter	True value	Description
Main lens – Section 2.1		SPLE
$\xi_{0,x} [\text{''}]$	-0.05	lens center x -axis
$\xi_{0,y} [\text{''}]$	-0.1	lens center y -axis
$\varphi_{\text{lens}} [\text{°}]$	1.0	rotation angle
q_{lens}	0.75	axis ratio
γ	2.1	SPLE slope
$\theta_{\text{Ein}} [\text{''}]$	1.5	Einstein radius
γ_1	0.005	1 st external shear component
γ_2	-0.01	2 nd external shear component
Source – Section 2.3		Sérsic
$x_{0,x} [\text{''}]$	0.0	source center x -axis
$x_{0,y} [\text{''}]$	0.0	source center y -axis
$\varphi_{\text{source}} [\text{°}]$	0.75	position angle
q_{source}	0.5	axis ratio
n	2.3	index
$R_e [\text{''}]$	2.0	effective radius
I_e	0.6	surface intensity

5.2 Single subhalo inference

We start with the simplified case where we only consider the position of a single subhalo as the free parameters, i.e. $\boldsymbol{\vartheta}_{\text{sub}} = (x_{\text{sub}}, y_{\text{sub}})$. We do this as a sanity check so that we are able to compare our results with other existing methods to measure subhalos with gravitational lensing. We are using the same modeling pipeline as [Coogan et al. \(2022\)](#). By choosing the same values for the model parameters, we can compare our results with the case where also they infer only the position of a single subhalo. We train the network with a set of 40 000 images, this is the same as is being done by [Coogan et al. \(2022\)](#) although we are not truncating. Because only the subhalo position is varied, the images are very similar, with only sub-percentage differences.

After applying the training pipeline, we infer the pixel posterior probability of a mock observation. The red dot in the mock observation, displayed in the left panels of Figure 5.1 shows the true position of a $10^9 M_\odot$ subhalo. With its position on the Einstein ring, the subhalo should be measured better compared to a subhalo that is not on the Einstein ring. The pixel posterior probability, depicted in the middle panel, is predicted for each of the $N_{\text{msc}} \times N_{\text{pix}} \times N_{\text{pix}} = 1 \times 100 \times 100 = 10^4$ pixels. The colorbar, shown on a logarithmic scale, contains an orange bar which indicates the value of the prior. Since we are uniformly sampling the position and a fixed mass, the value of the prior,

$$p(\tilde{\boldsymbol{\vartheta}}_{ijk}^1) = \frac{\langle N_{\text{sub}} \rangle}{N_{\text{pix},x} \times N_{\text{pix},y} \times N_{\text{msc}}} = \frac{1}{100 \times 100 \times 1} = 10^{-4}, \quad (5.1)$$

is the same for each pixel. From the predicted pixel posterior probability, we can accurately determine the position of the subhalo. One would not see this directly from the middle panel, but keep in mind

that this plot has a logarithmic scale. The network is less confident about the pixel posterior probability on the bottom left side of the true subhalo position, outside of the Einstein ring. The values in this outer region are, nevertheless, still lower than the prior value of the pixel posterior probability. Besides the true subhalo position, the network gives it a small chance that there is a subhalo around $(x_{\text{sub}}, y_{\text{sub}}) \approx (1.1, 0.1)''$, but this value is still well below the prior. Besides that, the network is very confident that there is no subhalo anywhere else.

Although the pixel posterior probability in the middle panel is technically no two-dimensional posterior, it can be interpreted as the marginal posterior of the subhalo coordinates $(x_{\text{sub}}, y_{\text{sub}})$. It would be a binned two-dimensional posterior, from where we are also able to compute the binned one-dimensional posteriors of $(x_{\text{sub}}, y_{\text{sub}})$ by marginalizing over the other coordinate. These marginal posteriors are shown in the right panels of Figure 5.1, but now on a normal scale instead of a logarithmic scale. We see what we also have seen from the middle panel. The width of the peak spans a few tenths of arcseconds. The marginalized pixel posterior probability (in blue) is centered around the true value (in red), and much higher than the prior (in orange) around the true value.

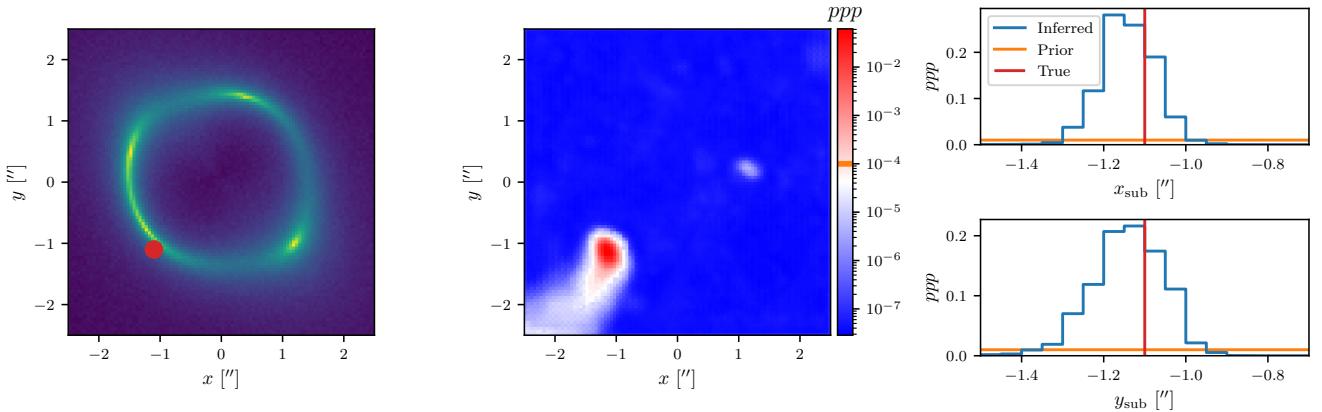


Figure 5.1: Inferring the pixel posterior probability (PPP) of a $10^9 M_\odot$ subhalo while varying its position. The highest density posterior regions centers around the true subhalo's position. **(Left)** Mock observation with the subhalo's true position marked by the red dot. **(Middle)** Predicted pixel posterior probability of all the 100×100 pixels. The orange bar indicates the prior value for each pixel. **(Right)** Marginalized pixel posterior probability.

The values of the pixel posterior probability in Figure 5.1 are calibrated. This calibration is done through the use of the empirical pixel posterior probability. Before going into the discussion of the empirical pixel posterior probability, let remind us quickly how it is calculated. We put all the pixel posterior probabilities of $N_{\text{calib}} = 5000$ predictions in a histogram, the green graph in the left panels of Figure 5.2. From these predictions, we now only consider the pixels that are actually a subhalo, the ‘true’ pixels. This histogram is shown in blue in the left panels. Dividing the histogram containing subhalos (blue) by the histogram of all pixels (green) we obtain the empirical pixel posterior probability shown in the middle panels. The calibration is done by fitting the empirical pixel posterior probability with isotonic regression. We give the predicted pixel posterior probabilities the value of the empirical pixel posterior probability through this fit, and we validate that this works by recalculating the empirical pixel posterior probability after the calibration. See Section 4.3.2 for an in-depth discussion about its technicalities, and Section 6.2 for a discussion about its shortcomings.

In the left panels of Figure 5.2 we can see that the histogram of the predicted pixel posterior probabilities has an asymmetric shape, centered around the prior, which is shown in orange. There are many

more low-value predictions compared to high-value predictions. This makes sense, as there is only one pixel out of the 10^4 pixels in each observation that contains a subhalo. The blue histogram, the pixels that actually contain a subhalo, confirms this statement. For high-value predicted pixel posterior probabilities, the two histograms share the same shape, while for lower values the histograms are different. The blue histogram dies out around the value of the prior. Most of the pixels that contain a subhalo have a pixel posterior probability well above the prior. There are around $\mathcal{O}(10^1)$ pixels out of the 5×10^7 predicted pixels that contain a subhalo but have a pixel posterior probability that is below the prior.

The middle panels Figure 5.2 shows the empirical pixel posterior probability. We can see from the top panel, that the predictions are overconfident only for high values. There is good coverage for values around the prior. We have seen from the blue histogram in the left panels that low-value predicted pixels do often not contain a subhalo. This has the consequence that we can not really assess how good the coverage is for lower values since we can only compute the empirical pixel posterior probability for pixels that contain a subhalo. By calibrating and recalculating the empirical pixel posterior probability, we see that calibration works as it should. The calibrated empirical pixel posterior probability is shown in the bottom panel. The high-valued predictions are now on the diagonal and therefore not overconfident. The gap between the highest bar and the other bars is because there are apparently no predicted pixel posterior probability with that value, as can be seen in the left bottom panel.

The right panels show the sum of all the pixel posterior probabilities of a prediction, for each of the 5000 test observations. The true and prior value of the sum is the same, $N_{\text{sub}} = \langle N_{\text{sub}} \rangle = 1$, and is shown in orange. The sums of the uncalibrated predictions peak at a value slightly higher than it should, including a tail of values $\gtrsim 2$. From the bottom right panel, we see that this is improved after the calibration. The sums peak around the true value, and the width is equal on both sides. This is necessary to interpret all the pixel posterior probabilities together as a subhalo density estimation.

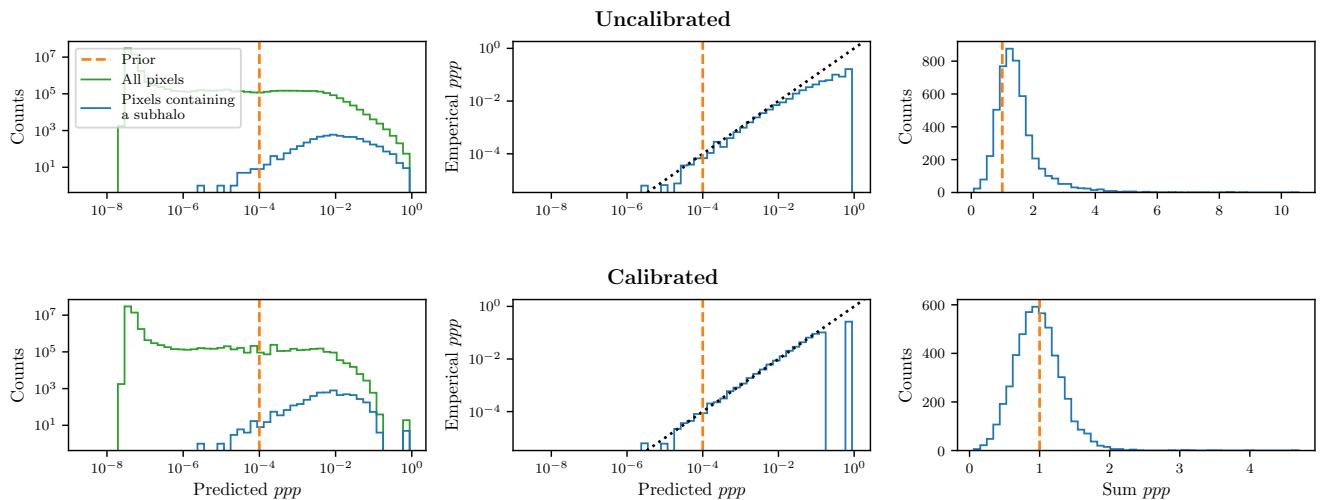


Figure 5.2: Validation and calibration of the single subhalo's position inference, by analyzing the predictions of $N_{\text{calib}} = 5000$ test observations. **(Left)** Histogram of all the predicted pixel posterior probability (ppp) (green) and predicted pixels that actually contain a subhalo (blue). Pixels with a high predicted value have a larger ratio of true subhalos. There are a lot of pixels with a low prediction, but they do not contain a subhalo. **(Middle)** Empirical pixel posterior probability, which is calculated by dividing the histograms in the left panels. Only high predicted values are overconfident. **(Right)** Sum of the predicted pixel posterior probabilities for each test observation. The calibration improves the inference by moving the center of the sum to the true values.

5.3 Multiple subhalo inference

In the previous section we have discussed the results of the position inference of a single subhalo. However, our method is designed to do more complicated inference tasks: determine the position and mass of multiple subhalos. In this section we discuss the inference results when varying the number of subhalos and their masses. The parameters of interest are $\boldsymbol{\vartheta}_{\text{subs}} = \{\boldsymbol{\vartheta}_{\text{sub},1}, \boldsymbol{\vartheta}_{\text{sub},2}, \dots \boldsymbol{\vartheta}_{\text{sub},N_{\text{sub}}}\}$, with $\boldsymbol{\vartheta}_{\text{sub},i} = (x_{\text{sub},i}, y_{\text{sub},i}, m_{\text{sub},i})$ and a varying number of N_{sub} subhalos. Similarly as [Ostdiek et al. \(2022b\)](#), we obtain a three-dimensional map of the subhalo parameters. We will discriminate between $N_{\text{msc}} = 12$ mass classes and consider the limits of subhalos masses $(m_{\text{low}}, m_{\text{high}}) = (10^7, 10^{10}) M_{\odot}$, such that the edges of our mass bins are $\{7.00, 7.25, 7.50, \dots 9.50, 9.75, 10.00\} \log_{10} M_{\odot}$. We train the network on 240 000 simulations.

As discussed in Section 2.2, we should use the SHMF to properly model the number of subhalos and their masses. However, as this work is still applied to mock observations, we will stick to the more simplified case where we model the subhalos uniformly over space, mass, and number density. The number of subhalos is sampled between 0 and 25, such that we have

$$N_{\text{sub}} \sim \mathcal{U}(0, 25), \quad m_{\text{sub}} \sim \mathcal{U}(10^7, 10^{10}) M_{\odot}. \quad (5.2, 5.3)$$

We stress that these sampling distributions are not an attempt to properly model subhalos, but have been simplified to increase the inference results. This will be discussed in more detail in Section 6.2. However, the sampled values are similar with other works. The number sampling is a bit higher than the values of 3 and 4 from [Coogan et al. \(2022\)](#) and [Anau Montel et al. \(2022\)](#) but the same as the number sampling of [Ostdiek et al. \(2022b\)](#). The upper limit of the subhalo mass is chosen such that only ‘dark’ halos are modeled, as higher halos would be visible because of their baryons. Furthermore, the chance of a heavier subhalo in this lensing system would be really low. The lower subhalo mass limit not well substantiated. Lower-mass subhalos would have a minimal imprint on the gravitational lensing and would be below our current sensitivity of subhalos, as we will discuss below. The described sampling results in a value of the prior for each pixel of

$$p(\tilde{\boldsymbol{\vartheta}}_{ijk}^1) = \frac{\langle N_{\text{sub}} \rangle}{N_{\text{pix},x} \times N_{\text{pix},y} \times N_{\text{msc}}} = \frac{(25 - 0)/2}{100 \times 100 \times 12} \approx 1.04 \cdot 10^{-4}. \quad (5.4)$$

A mock observation and its predicted pixel posterior probabilities are shown in Figure 5.3. The top left panel shows the mock observation, the red dots indicate the positions of the five subhalos. The size of the dots indicates their masses. The size of the white circles indicate masses of $\{10^8, 10^9, 10^{10}\} M_{\odot}$. The pixel posterior probability of all the mass channels is separately shown in the 3×4 ‘grid’ bottom panel, but also together in the three-dimensional ‘lavalamp’ in the top right panel. Both the lavalamp and grid figure show the same pixel posterior probability, but due to their different visualization, we can highlight different characteristics.

The gridded bottom panel shows the pixel posterior probability for each mass channel independently. To keep the characteristics of the independent mass channels visible, we are changing the colorbar scale for each row. Because of the logarithmic scaling and difference between colorbars, some characteristics might look large at first sight, while they are not if one keeps the different colorbars in mind.²⁰ We note the following characteristics of the predictions.

²⁰We have made the comprise for one colorbar for each row. One colorbar for all subfigures would result in less pronounced characteristics, and some subfigures would be mainly one single color. A colorbar for each subfigure separately would make the comparison between the subfigures too hard.

- In regions of high-mass channels where there is no subhalo, most pixels have a prediction well below the prior. These are the dark blue regions. The network is very confident that there is no subhalo there.
- High-probability regions around a subhalo spread both in the space and mass direction. For example, there is a high-probability region at $(x, y) \approx (-1.1, -1.1)''$ in all the higher mass channels, while there is only one subhalo with a mass of $m_{\text{sub}} = 10^{9.54} M_{\odot}$. The spread across mass channels is what we will refer to as ‘mass leaking’.
- The inverse of mass leaking also occurs. There is a high-probability region around the subhalo located at $(x_{\text{sub}}, y_{\text{sub}}, m_{\text{sub}}) = (1.5'', 0.5'', 10^{9.16} M_{\odot})$. But in the highest mass channel, there is a very low-probability region at those spatial coordinates. This indicates that the network learns that there is a subhalo, and although it is not sure about its exact position it can clearly exclude the highest mass channels.
- If we go to lower mass channels, the very low- and very high-probability regions get less pronounced. The range of the colorbar is smaller for the rows describing those mass channels. The network is less confident if there is or there is not a subhalo.
- Although the regions get less pronounced for lower mass channels, the effect is still there. The imprint of the heavier subhalos is still possible, even in the lowest mass channel.
- The lowest mass channels have predictions more or less around the prior (the orange line in the colorbar). The signal of any subhalo does not get picked up by the network.

The effect of mass leaking and other general shapes of the pixel posterior probability are better visualized in three dimensions than the two dimensions in the grid of Figure 5.3. With the lavalamp figure in the top right panel, one can get a quick insight into the general pixel posterior probability behavior. An interactive version of the lavalamp figure is available.²¹ When making these plots, one has to do a few compromises. The first compromise is the visibility between pixels on the inside and outside. Therefore, the transparency increases with lower-value predictions all the way to zero so that the lowest predictions are not even shown. Secondly, instead of plotting all the $N_{\text{msc}} \times N_{\text{pix}} \times N_{\text{pix}} = 12 \times 100 \times 100 = 1.2 \times 10^5$ pixels, we draw surfaces of equal value. From the interactive lavalamp plot, we can note a few characteristics of the predictions that we could not see with the two-dimensional grid figure. We can see from the normal scale figure that high-mass subhalos get a much higher probability than low-mass subhalos. From the logarithmic scale, we can see that although low-mass subhalos generally do not get picked up by the network, surfaces of higher values start to form around them.

Additional predictions of mock observations such as Figure 5.3 can be found in Appendix A.1. The inference setup is the same, but the subhalo configuration (number, position, and mass) is different than the discussed example. We do not note any different behavior as already discussed. However, the additional figures can give a better qualitative insight into how the inference behaves.

²¹<https://dm-lensing.github.io/>

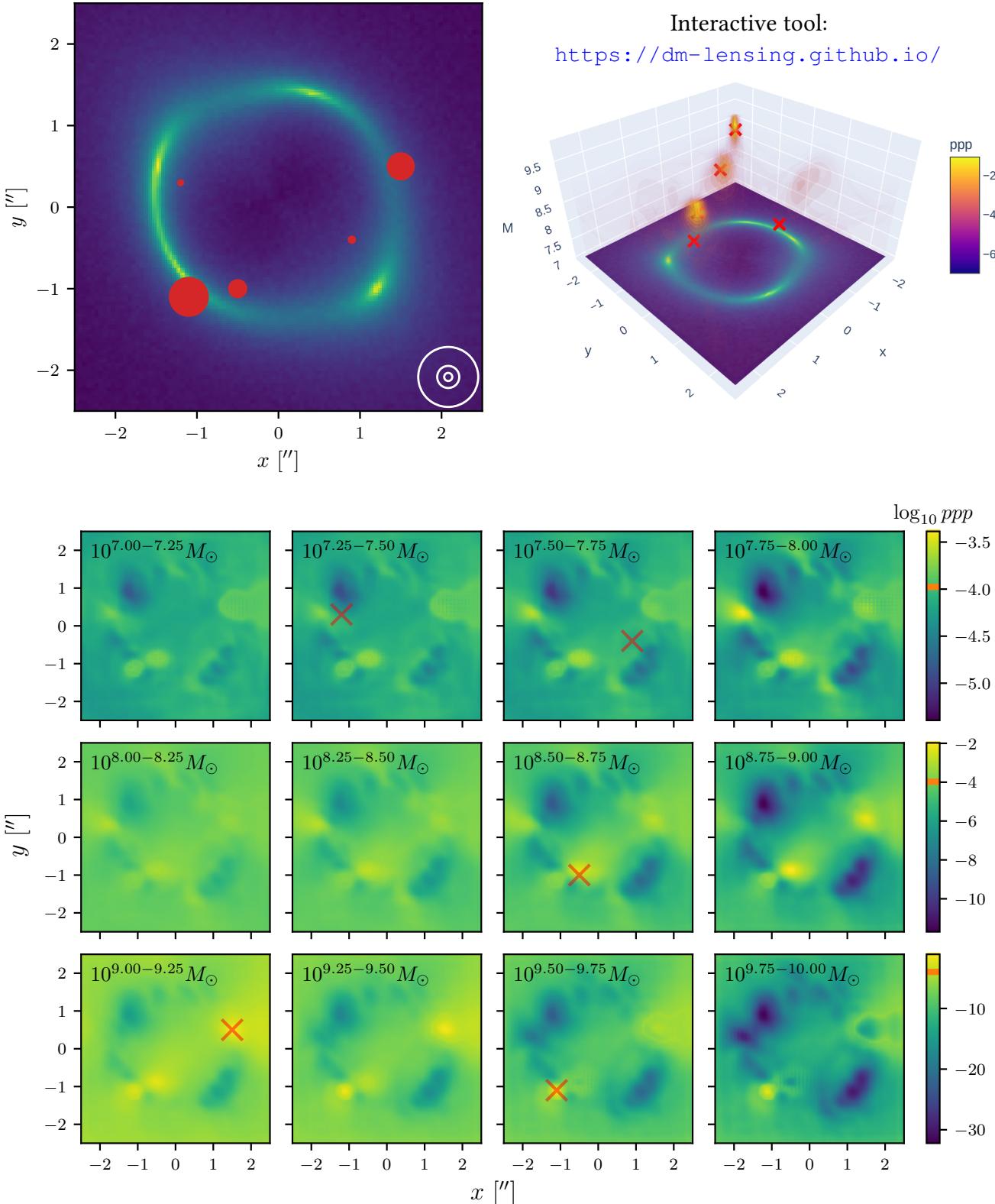


Figure 5.3: Inference of the pixel posterior probability (ppp) of multiple subhalos while varying their position and mass. High-probability regions form around the high-mass subhalos, and pixels from lower-mass channels get predicted values around the prior. **(Top left)** Mock observation with multiple subhalo modeled. The size of the red dots indicates the mass of the subhalo, with the white circles as a reference for $m_{\text{sub}} = \{10^8, 10^9, 10^{10}\} M_\odot$. **(Top right)** The ‘lavalamp’ figure shows the pixel posterior probabilities in three dimensions. The true subhalo positions and masses are indicated by the red crosses. **(Bottom)** The pixel posterior probabilities split out per mass-channel. Note the different colorbar ranges per row. The orange bar indicates the value of the prior.

The validation and calibration results are shown in Figure 5.4. The panels show the same as Figure 5.2, but now for the multiple subhalo case. There are a few similarities but also some subtle differences with the single subhalo position case from the previous section. The left panels are similar to the single subhalo case for high pixel posterior probability values, where the histograms both decrease. This makes sense, as both cases have way more pixels that do not have a subhalo than pixels that do have a subhalo. However, the histograms peak now around the value of the prior. This means that there are a lot of pixels where the network can not tell whether there is a subhalo or not. The middle panel shows how the calibration is less needed for the multiple subhalo case than for the single subhalo case. This is because the uncalibrated empirical pixel posterior probability is already quite close to the diagonal. However, calibrating can only increase our inference quality, so we do it anyway.

The right panels show again the sum of the pixel posterior probability of 5000 predictions of mock observations. The shape of the prior is because of the uniform sampling of the number of subhalos. The interpretation of the combined pixel posterior probability as a subhalo density breaks for the case if there are no subhalos simulated. The lowest sum of pixels, and therefore the lowest bin in the histogram, is $\sum_{ijk} p(\tilde{\vartheta}_{ijk}^1 | \mathbf{x}_{\text{obs}}) = 0.3804$. This bin has also the highest value.

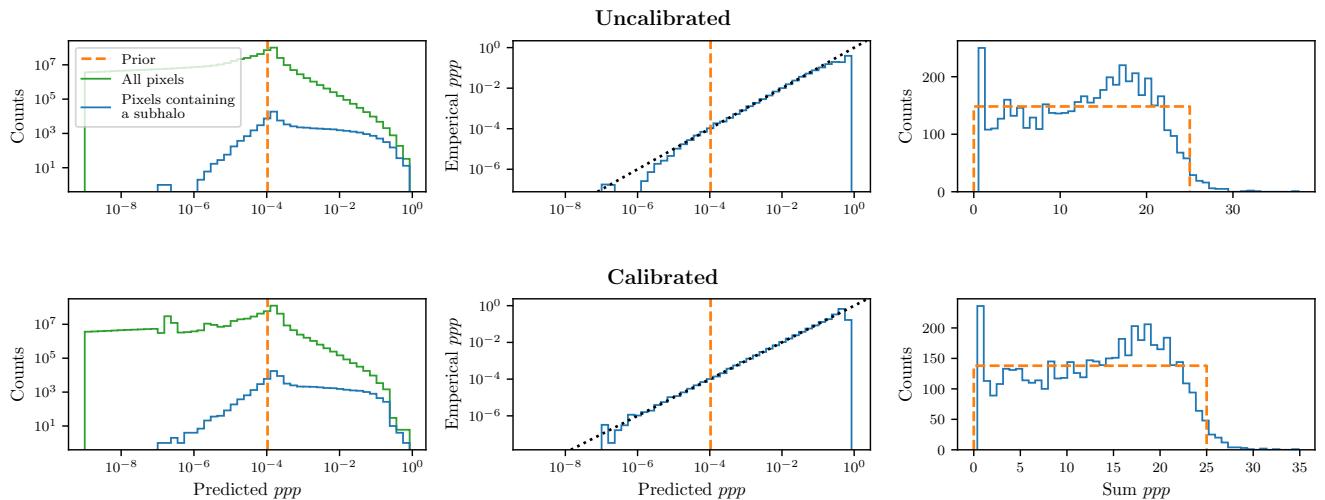


Figure 5.4: Validation and calibration for inference of the mass and position of multiple subhalos, by analyzing the predictions of $N_{\text{calib}} = 5000$ test observations. (**Left**) Histogram of all the predicted pixel posterior probability (ppp) (green) and predicted pixels that actually contain a subhalo (blue). Histograms peak around the prior, indicating that there are a lot of pixels where the network can not give a confident prediction. (**Middle**) Empirical pixel posterior probability, which is calculated by dividing the histograms in the left panels. Calibrating does not improve much, as the uncalibrated prediction has already a good coverage. (**Right**) Sum of the predicted pixel posterior probabilities for each test observation. The prior gets its shape because the number of subhalos is sampled from $\mathcal{U}(0, 25)$.

We can investigate the sensitivity of the network in more detail by splitting out the bottom left panel of Figure 5.4 for each mass channel. This is shown in Figure 5.5. We can validate if our results from the examination of a single prediction (Figure 5.3) generalizes to the predictions in general, as the histograms contain information of $N_{\text{calib}} = 5000$ predictions of mock observations. For lower mass channels, the shape of the histogram of all pixels (green) is mainly the same as the histogram of the pixels containing a subhalo (blue). The histograms center around the prior. This confirms that the network can not tell whether there is a subhalo or not. The difference in the shape of the histograms in the higher mass channels is significant. There is a larger ratio of high pixel posterior probabilities that

contain a subhalo. This is what we would expect, as these subhalos should have a larger effect on the gravitational lensing.

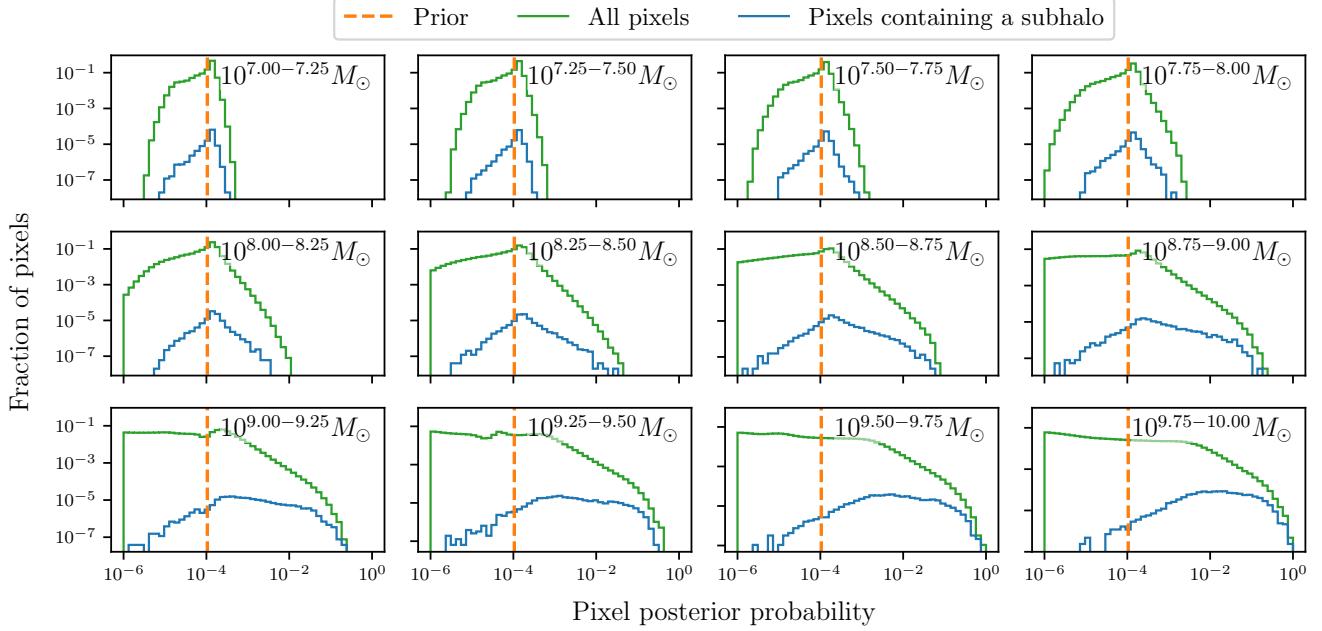


Figure 5.5: Histograms indicating the sensitivity across mass channels from $N_{\text{calib}} = 5000$ test observations. The blue histogram is a subset of the green histogram and shows the distribution of predicted pixel posterior probabilities. Predictions of lower mass channels are typically the same as the prior value, indicating that the network is not sensitive to subhalos of these masses. Predictions of higher mass channels are typically higher, indicating that subhalos are easier to measure when the subhalo's mass is high.

Chapter 6

Discussion & conclusion

6.1 Conclusion

A better understanding of the fundamental nature of dark matter can be obtained through measurements of dark matter subhalos on sub-galactic scales. Analyzing gravitational strongly lensed images is a probe to these subhalos. Near-future telescopes will make a high amount of high-quality observations available, creating the need for a fast, flexible and automated pipeline to analyze these images efficiently. Traditional techniques are not suited for these tasks, as they require calculating the likelihood. This makes these analysis techniques time-consuming and/or not flexible because the number of parameters is practically limited. We are using, as a solution for this, neural SBI. By using this likelihood-free method, we can have an arbitrary complex model because we do not calculate the likelihood explicitly. Instead, we sample training data for a neural network from a simulator. The network learns to estimate a likelihood ratio which can directly be translated to posteriors.

However, one has to feed into the network for which parameters the posteriors should be estimated. This is not possible when one wants to estimate the posterior of multiple subhalos. The number of subhalos is stochastic, so one could not tell the network how many posteriors to estimate. Furthermore, the interchangeability of the subhalos causes the label switching problem. Therefore, we have presented a method than combines MNRE with image segmentation to calculate the pixel posterior probability. This is the probability that there is a subhalo in a user-defined $(x_{\text{sub}}, y_{\text{sub}}, Mm_{\text{sub}})$ pixel. Through this approach, we are able to train on, and estimate a varying number of subhalos.

The presented method has been applied to mock observations, with two different inference tasks. We summarize the key findings of these tests as follows:

- With our method we are able to accurately localize the position of a single heavy ($10^9 M_{\odot}$) subhalo, as presented in Figure 5.1. The marginalization of the pixel posterior probability can be interpreted as the one-dimensional marginal posteriors presented in Coogan et al. (2022, Figure 1). Our posterior ‘distribution’ is a bit wider than the inferred and analytical posterior distribution presented there. However, as our method is not tailored to the inference of a single subhalo and we are not using truncation, we still consider our inference results as high-quality.
- With our method of pixelation, the precision of the method can only be as good as the resolution of the observation. However, the width of the high-probability regions around the true subhalo positions spans multiple pixels. This means that the resolution of the observation, and therefore the resolution of the prediction, does not a bottleneck, as the posterior spans multiple pixels.
- When inferring the position and mass of multiple subhalos, the network is able to localize regions around high-mass subhalos, as presented in Figure 5.3. Furthermore, the network is able to exclude high-mass regions that do not contain a subhalo by giving these pixels a very low prob-

ability. The network returns values around the prior for low-mass channels. Although a direct comparison with [Ostdiek et al. \(2022b\)](#) is not possible because of different modeling and detection threshold, we conclude that our results are in line with what they found.

- The method is able to predict stable results, already without calibration, and the predictions can correctly be interpreted as probabilities. A calibration by fitting the empirical pixel posterior probability can only improve the inference quality.
- The combination of the pixel posterior probability can correctly be interpreted as a subhalo density estimation. We do not force this, but through our definition of the pixel posterior probability, we have shown that this works.

This work uses MNRE, similar to [Anau Montel et al. \(2022\)](#); [Coogan et al. \(2022\)](#), and fills within the larger pipeline to measure dark matter subhalos. This work can be applied as an initial exploration of the subhalo parameter space, before tailored techniques as [Coogan et al. \(2022\)](#) can be applied to measure single subhalos more efficiently and precisely. Furthermore, our probabilistic predictions of the subhalo parameter space can be used to determine the posterior of dark matter properties.

Finally, we are optimistic how MNRE, with the improvements discussed in the next section, will be able to measure dark matter subhalos with strong gravitational lensed images. The combination of cutting-edge machine learning techniques and the vast amount of future observations makes it an exciting time to shed light upon the identity of dark matter, making our understanding of the Universe undark.

6.2 Discussion and future work

Pixel posterior probability estimation The presented normalization and calibration processes have shown to produce reliable and consistent results. However, the procedure could be built in a statistically more robust manner. We will discuss the most important possible improvements to the procedure below.

It is actually quite remarkable that the normalization step gives us values that can be interpreted as probabilities. During training, the network is never informed that it should estimate ratios that can be translated to probabilities. During the development of the method, we have performed tests where this normalization is built within the network, such that the network ‘knows’ that it should estimate probabilities. This procedure is discussed in Appendix [A.2](#), but we have not been able to produce stable results.

The posterior is calculated by multiplying the output of U-Net with the pixel prior, through Equation [\(4.10\)](#). There is, in contrast with other likelihood-free methods ([Anau Montel et al., 2022](#); [Coogan et al., 2022](#)), no sampling used to produce the posterior. We take directly the output of a single prediction and multiply the likelihood ratio with the prior. By doing this, we do not take the variability of the free parameters into account. Sampling the posterior from multiple predictions of an observation could lead to more robust results.

The empirical pixel posterior probability is calculated through the division of two histograms. Using histograms is always delicate because of the binning. The chosen bins where the histogram is computed have a large effect on the actual outcome of the empirical pixel posterior probability. If wider bins would have been chosen, we could have decided that calibration is not needed, because the empirical pixel posterior probability lies already well enough on the diagonal. The opposite is true as well: smaller bins could result in a wrong calibration. This could be the case already now. Let us discuss, for example, the overconfidence for very low and very high predictions that we see in the top middle panel of Figure [5.2](#). It could be very much the case that overconfident predictions are mainly there because we are running

out of statistics since the number of pixels that contain a subhalo is not high enough. Furthermore, a proper definition of error bars for the empirical pixel posterior probability should be developed. There is currently no formal threshold if and when the calibration is working correctly, when the empirical pixel posterior probability is ‘diagonal enough’.

The validation and calibration steps are now performed with aggregated statistics. All the pixels, no matter their location in space or mass, are used to calculate the empirical pixel posterior probability. However, not all pixels with the same predicted value are the same. Pixels in the lower mass channels have a value around the prior because the network does not learn to pick up any signal from these low-mass subhalos, while a pixel in the higher mass channels could have the same value because it is on the border between a region where there is a subhalo and there is no subhalo. These two pixels could need a different calibration. This has been investigated with Figure 5.5, but this could be done in more detail. One could, for example, do the same procedure, but in the spatial direction instead of the mass direction.

Simulating strong lensing This work has mainly been a proof-of-concept of the presented subhalo density estimation applied to gravitational lensing. We have not discussed or investigated the lensing model in great detail. However, to apply the pipeline to existing data, the simulator has to be improved by including more physical known processes. Further improvements to the simulator are a never-ending task, one can always add more complexity from the latest astrophysical theories. However, there are some improvements that could directly be built within the simulator. We will discuss the most important improvements below.

The most immediate point of interest would be the subhalo modeling. Currently, we are modeling the mass and number of subhalos uniformly. A SHMF, such as Equation (2.16), is needed to realistically model multiple subhalos with varying mass. The current choice of uniform mass sampling is mainly based upon the fact that understanding the relation between training data and the results is easier. If one would use a proper SHMF, there would be many more low-mass subhalos compared to high-mass subhalos. However, when one would infer the position of a heavy subhalo, the trained network could return low-confident results. This would not be because high-mass subhalos are hard to measure, but because there have not been many high-mass subhalo training simulations presented to the network. Increasing the size of the training data would not necessarily improve the results, because the training data would mainly contain low-mass subhalos. Furthermore, [Ostdiek et al. \(2022b\)](#) highlighted that this could result in biases of the SHMF. We leave a proper implementation of mass and number sampling of subhalos for future work.

Further improvements on the modeling of the lensing system are including LOS halos and varying the lens- and source parameters, similar to [Anau Montel et al. \(2022\)](#); [Coogan et al. \(2022\)](#). Currently, we are only marginalizing over the noise, so calling our method ‘marginal’ NRE is already debatable. However, it is straightforward to vary the lens and source parameters within the current pipeline. Truncation, which will be discussed in the next subsection, has been shown to efficiently target the lens- and source parameters while inferring the subhalo parameters. Therefore, we expect that, besides wider posteriors, not very different results. [Ostdiek et al. \(2022b\)](#) has also shown that U-Nets are able to handle varying lens and source parameters. However, including LOS halos takes more care. [Coogan et al. \(2022\)](#) has shown that including those can significantly bias the position inference of a single subhalo. Varying sub- and LOS halo parameters, as [Wagner-Carena et al. \(2022\)](#) does, would make the model more flexible to apply on real data.

Besides the lensing system, there are also some simplifications in the modeling of the observation itself. Before the pipeline could be applied to real data, one should include correct modeling of the PSF. Modeling the PSF is dependent on the telescope itself. As a final suggestion to model observations more

accurately, we suggest modeling the lens galaxy light. This light is currently completely neglected, as also is done in similar works. However, telescopes observe the gravitationally lensed light in multiple frequency bands. Including these frequency bands in the simulator would help the network to differentiate between lens and source light, as the lens and source have a different redshifts, and therefore frequencies.

Machine learning In this work we have used MNRE, where we marginalize over the free parameters and apply NRE to determine the posteriors of the parameters of interest. However, if there are many free parameters and/or the prior ranges are wide, this procedure can become inefficient. The loss function can get a lot of examples presented that are by far not what the observation looks like. By truncating, the training data is regenerated in the regions of the parameters space that are the most relevant when analyzing a particular (mock) observation. These truncation steps are done in multiple rounds until there is no improvement after a truncation. See [Coogan et al. \(2022, Figure 1\)](#) for a visualization of the truncation of the subhalo position, and [Anau Montel et al. \(2022, Figure 6\)](#) for the lens and source parameters. In practice, truncation happens by selecting a uniform region around an intermediate posterior.²² The width of the uniform region, the truncated prior, is determined by a user-defined threshold of the intermediate prior. See [Cole et al. \(2021, Figure 2\)](#) for a general visualization of this, and [Anau Montel et al. \(2022, Figure 8\)](#) how this works in practice for the lens and source parameters.

Applying truncation on the subhalo sampling when calculating the pixel posterior probability would take more caution, as the practical implementation of truncation does not generalize to this procedure. This is because our method does produce a pixel posterior probability of a subhalo being in that pixel, instead of a continuous probability distribution of a subhalo parameter. Therefore, one can not define a truncated prior from the intermediate posterior distribution. Instead, one would isolate multiple small regions of high-valued predictions, which would typically be in the higher mass channels. Additionally, one would also keep a large region of pixels that got a prediction around the value of the prior, spanning most pixels in the lower mass channels. In the high-probability regions, a heavy subhalo would be individually modeled. In the lower-probability region, we would simulate multiple subhalos, similar to how we sampled subhalos in the initial round.

Besides truncation, there are a few other improvements on the machine learning that could directly be explored. The used U-Net has just some minor improvements on the original architecture of [Ronneberger et al. \(2015\)](#). Significant improvements on the U-Net are for example nested ([Zhou et al., 2018](#)) and probabilistic ([Kohl et al., 2019](#)) U-Nets. The original architecture of the U-Net was designed for ‘local segmentation’, segmentation where the imprint of a segmented feature is mainly present in its own pixel and its direct neighbors. This is not the case for gravitational lensing. The effect of subhalos reaches well outside of the pixel where it is localized, and generally has an effect on all the light on the Einstein ring. Further research on other designs of networks is needed to decide whether this would improve the inference results.

Other applications The presented technique, the estimation of an ‘object’ density with MNRE is not limited to subhalos from strongly lensed images. Nothing is stopping us to apply this to other problems in astroparticle physics and other disciplines, where we have an *unknown number of interchangeable objects* to estimate. One of the examples where our technique would be useful, is point source detection in Fermi-LAT data to characterize the Galactic Center gamma-ray excess. A subset of point sources could be dark matter halos. Similar work as been done by [Mishra-Sharma & Cranmer \(2022\)](#), where they used neural SBI with normalizing flows to characterize summary statistics of the Galactic Center.

²²‘Intermediate’ in the sense that it is not the final posterior, but the posterior found after a non-final round of truncation.

With our method, we would infer the longitudinal position, the lateral position, and the J -factor of dark matter halos on the sky.²³ With this we could discriminate dark matter halos from other point sources. Our method could work better in this context than gravitational lensing because the imprint of the halos is very localized, and could therefore be treated as point sources.

²³The J -factor describes the dark matter halo distribution. It would have the same role as the subhalo mass in the lensing context.

Appendices

A.1 Additional subhalo density predictions

Below are additional subhalo density estimations shown, similar as Figure 5.3 but for other subhalo configurations. The pixel posterior probability (ppp) is inferred from systems with multiple subhalos varying their position and mass. See Section 5.3 for an in-depth discussion.

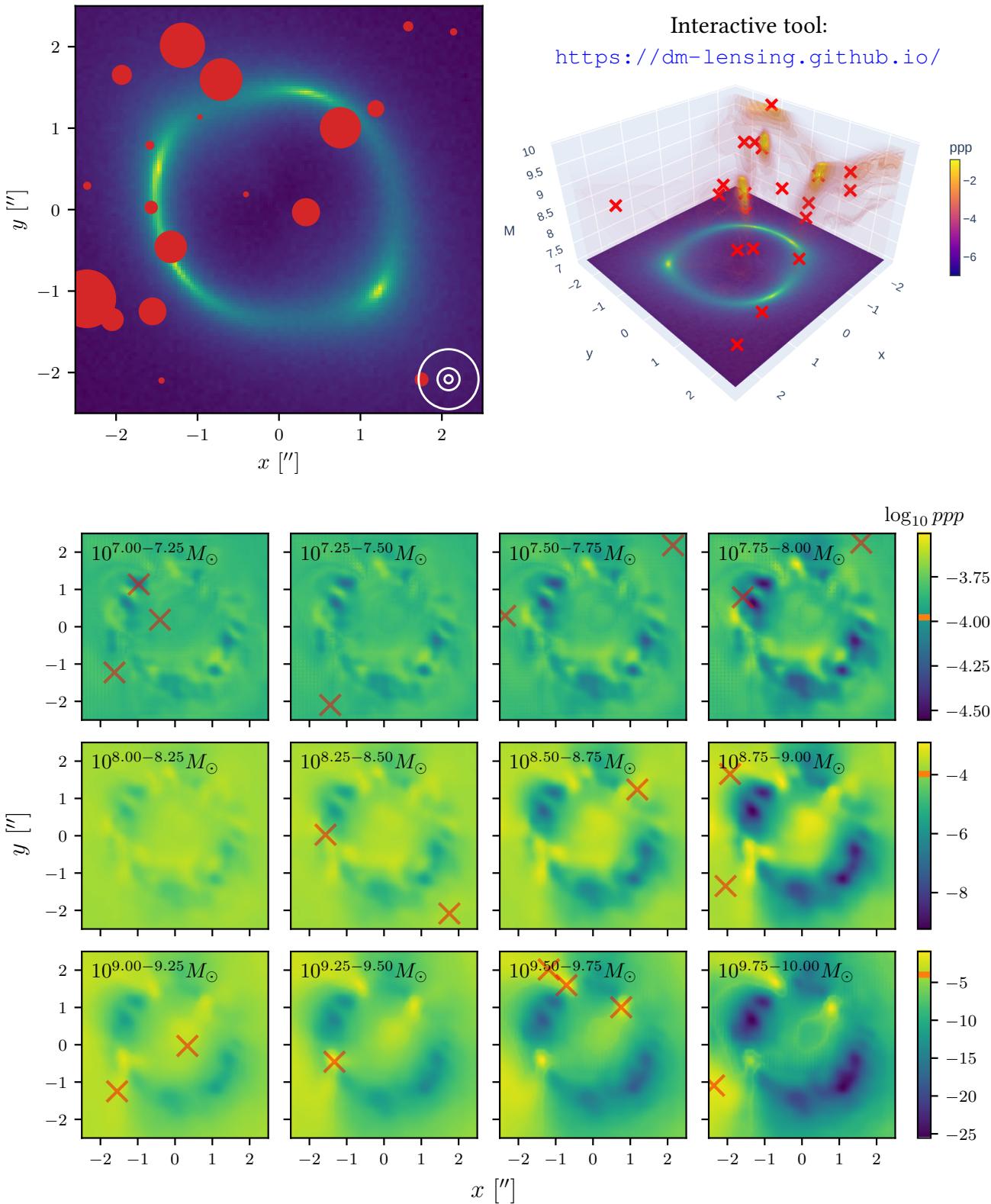


Figure A.1.1: Inference of a mock observation. See main text and Section 5.3 for a discussion. **(Top left)** Mock observation with multiple subhalos modeled. The size of the red dots indicates the mass of the subhalo, with the white circles as a reference for $m_{\text{sub}} = \{10^8, 10^9, 10^{10}\} M_{\odot}$. **(Top right)** The ‘lavalamp’ figure shows the pixel posterior probabilities in three dimensions. The true subhalo positions and masses are indicated by the red crosses. **(Bottom)** The pixel posterior probabilities split out per mass-channel. Note the different colorbar ranges per row. The orange bar indicates the value of the prior.

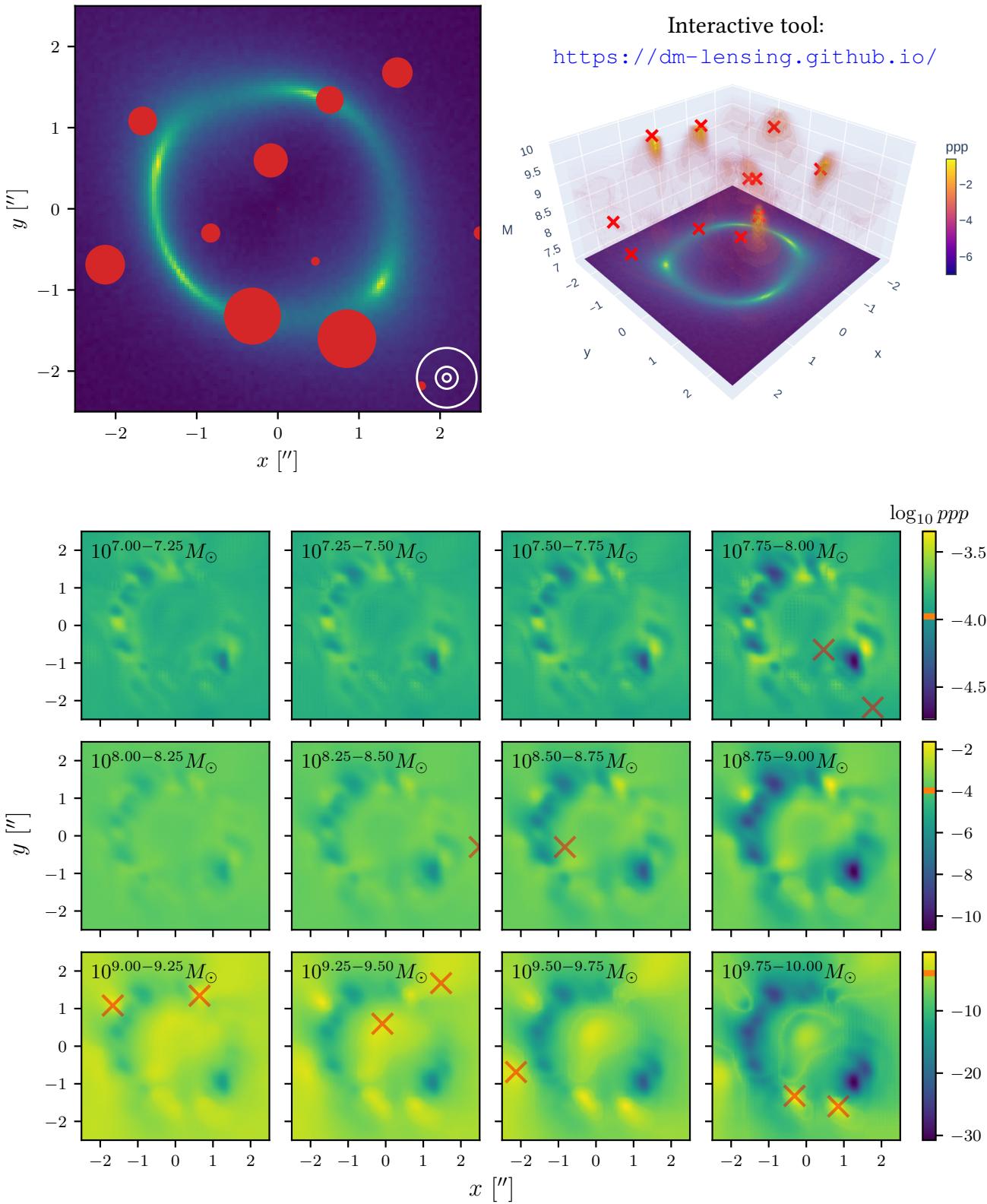


Figure A.1.2: Inference of a mock observation. See main text and Section 5.3 for a discussion. **(Top left)** Mock observation with multiple subhalos modeled. The size of the red dots indicates the mass of the subhalo, with the white circles as a reference for $m_{\text{sub}} = \{10^8, 10^9, 10^{10}\} M_\odot$. **(Top right)** The ‘lavalamp’ figure shows the pixel posterior probabilities in three dimensions. The true subhalo positions and masses are indicated by the red crosses. **(Bottom)** The pixel posterior probabilities split out per mass-channel. Note the different colorbar ranges per row. The orange bar indicates the value of the prior.

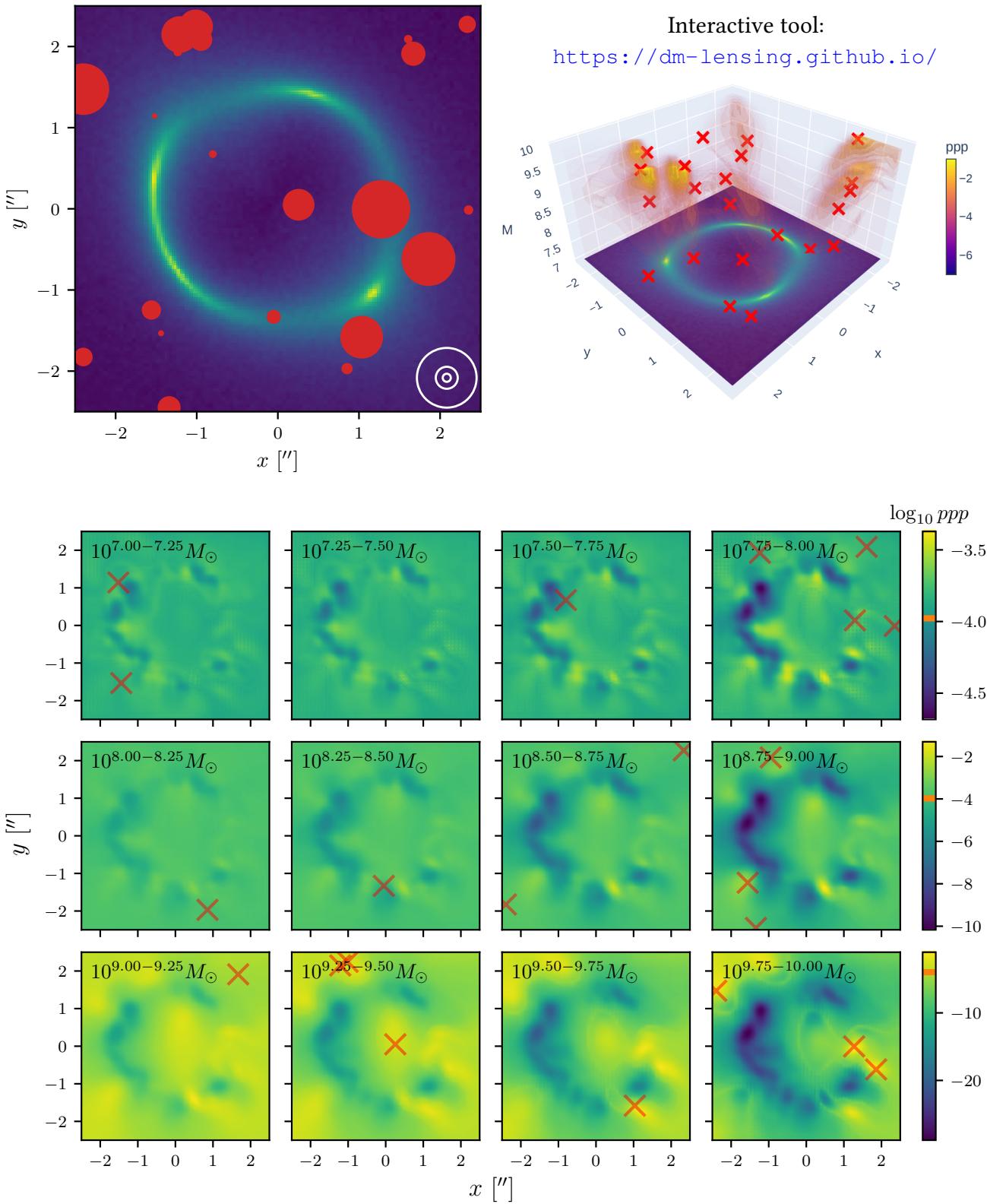


Figure A.1.3: Inference of a mock observation. See main text and Section 5.3 for a discussion. **(Top left)** Mock observation with multiple subhalos modeled. The size of the red dots indicates the mass of the subhalo, with the white circles as a reference for $m_{\text{sub}} = \{10^8, 10^9, 10^{10}\} M_{\odot}$. **(Top right)** The ‘lavalamp’ figure shows the pixel posterior probabilities in three dimensions. The true subhalo positions and masses are indicated by the red crosses. **(Bottom)** The pixel posterior probabilities split out per mass-channel. Note the different colorbar ranges per row. The orange bar indicates the value of the prior.

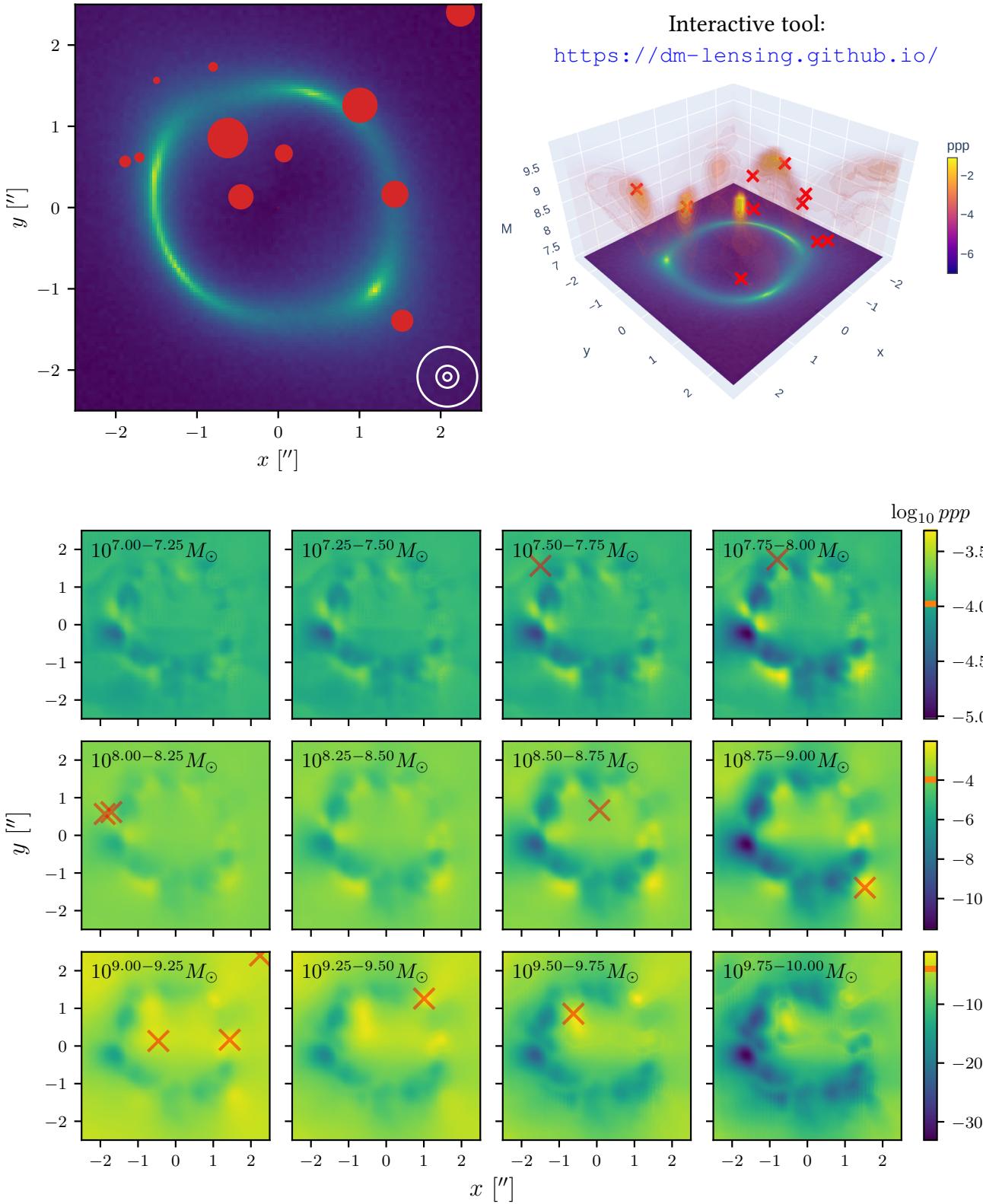


Figure A.1.4: Inference of a mock observation. See main text and Section 5.3 for a discussion. **(Top left)** Mock observation with multiple subhalos modeled. The size of the red dots indicates the mass of the subhalo, with the white circles as a reference for $m_{\text{sub}} = \{10^8, 10^9, 10^{10}\} M_{\odot}$. **(Top right)** The ‘lavalamp’ figure shows the pixel posterior probabilities in three dimensions. The true subhalo positions and masses are indicated by the red crosses. **(Bottom)** The pixel posterior probabilities split out per mass-channel. Note the different colorbar ranges per row. The orange bar indicates the value of the prior.

A.2 Normalizing during training

We have explored the feasibility of normalization during training without success. This normalization step would be performed directly with the output of the U-Net, before the mapping. The pixelating (Section 4.2.2) would be a bit different. Instead of predicting two ratios (there is a subhalo and there is not a subhalo) for each pixel, we would make the network only output the ratio that there is a subhalo. Instead of $2 \times N_{\text{msc}}$ output layers, we would have N_{msc} layers and Equation 4.6 would be replaced by

$$\left\{ f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}^1) \right\} = \text{Network}(\mathbf{x}). \quad (\text{A.2.1})$$

This would be the input for the normalization algorithm, as shown in Algorithm 1. The normalization as described in Section 4.3.1, would effectively be replaced by lines 2 and 3. After line 4, the mapping procedure (Section 4.2.3) would be continued.

Algorithm 1 Constructing the ratio of a pixel containing no subhalo from the ratio of a pixel containing a subhalo. This would be done between the UNet and mapping procedure.

- 1: $p(\tilde{\boldsymbol{\vartheta}}^1 | \mathbf{x}_{\text{obs}}) = \sigma(f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}^1))$ ▷ Set between [0,1] so we can interpret them as probabilities.
 - 2: $p(\tilde{\boldsymbol{\vartheta}}^0 | \mathbf{x}_{\text{obs}}) = 1 - p(\tilde{\boldsymbol{\vartheta}}^1 | \mathbf{x}_{\text{obs}})$
 - 3: $f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}^c) = p(\tilde{\boldsymbol{\vartheta}}^c | \mathbf{x}_{\text{obs}}) / p(\tilde{\boldsymbol{\vartheta}}^0)$ ▷ With $c = 0, 1$ and $p(\tilde{\boldsymbol{\vartheta}}^c)$ the prior.
 - 4: $f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}^c) = \log f_\phi(\mathbf{x}, \tilde{\boldsymbol{\vartheta}}^c)$ ▷ To pass logarithmic values to the loss.
-

Algorithm 1 has not produced stable results. Therefore, we do not include it in our method and leave the implementation of it for future work.

Bibliography

- Alexander S., Gleyzer S., McDonough E., Toomey M. W., Usai E., 2020, *The Astrophysical Journal*, 893, 15
- Alexander S., Gleyzer S., Parul H., Reddy P., Toomey M. W., Usai E., Von Klar R., 2021, Decoding Dark Matter Substructure without Supervision ([arXiv:2008.12731](https://arxiv.org/abs/2008.12731)), doi:10.48550/arXiv.2008.12731
- Alsing J., Wandelt B., Feeney S., 2018, *Monthly Notices of the Royal Astronomical Society*, 477, 2874
- Alsing J., Charnock T., Feeney S., Wandelt B., 2019, *Monthly Notices of the Royal Astronomical Society*, 488, 4440
- Amorisco N. C., et al., 2022, *Monthly Notices of the Royal Astronomical Society*, 510, 2464
- Anau Montel N., Coogan A., Correa C., Karchev K., Weniger C., 2022, Estimating the Warm Dark Matter Mass from Strong Lensing Images with Truncated Marginal Neural Ratio Estimation ([arXiv:2205.09126](https://arxiv.org/abs/2205.09126)), doi:10.48550/arXiv.2205.09126
- Anderson L., et al., 2014, *Monthly Notices of the Royal Astronomical Society*, 441, 24
- Baltz E. A., Marshall P., Oguri M., 2009, *Journal of Cosmology and Astroparticle Physics*, 2009, 015
- Banik N., Bovy J., Bertone G., Erkal D., de Boer T. J. L., 2021, *Journal of Cosmology and Astroparticle Physics*, 2021, 043
- Baur J., Palanque-Delabrouille N., Yèche C., Magneville C., Viel M., 2016, *Journal of Cosmology and Astroparticle Physics*, 2016, 012
- Bayer D., Chatterjee S., Koopmans L. V. E., Vegetti S., McKean J. P., Treu T., Fassnacht C. D., 2018, Observational Constraints on the Sub-Galactic Matter-Power Spectrum from Galaxy-Galaxy Strong Gravitational Lensing ([arXiv:1803.05952](https://arxiv.org/abs/1803.05952)), doi:10.48550/arXiv.1803.05952
- Bertone G., Hooper D., 2018, *Reviews of Modern Physics*, 90, 045002
- Birrer S., Amara A., Refregier A., 2017, *Journal of Cosmology and Astroparticle Physics*, 2017, 037
- Bode P., Ostriker J. P., Turok N., 2001, *The Astrophysical Journal*, 556, 93
- Bolton A. S., Burles S., Koopmans L. V. E., Treu T., Gavazzi R., Moustakas L. A., Wayth R., Schlegel D. J., 2008, *The Astrophysical Journal*, 682, 964
- Bond J. R., Szalay A. S., Turner M. S., 1982, *Physical Review Letters*, 48, 1636
- Boyarsky A., Drewes M., Lasserre T., Mertens S., Ruchayskiy O., 2019, *Progress in Particle and Nuclear Physics*, 104, 1
- Brehmer J., Mishra-Sharma S., Hermans J., Louppe G., Cranmer K., 2019, *The Astrophysical Journal*, 886, 49
- Brewer B. J., Huijser D., Lewis G. F., 2015, Trans-Dimensional Bayesian Inference for Gravitational Lens Substructures ([arXiv:1508.00662](https://arxiv.org/abs/1508.00662)), doi:10.48550/arXiv.1508.00662
- Bullock J. S., 2010, arXiv:1009.4505 [astro-ph]
- Buscicchio R., Roeber E., Goldstein J. M., Moore C. J., 2019, *Physical Review D*, 100, 084041
- Celeux G., 1998, in Payne R., Green P., eds, COMPSTAT. Physica-Verlag HD, Heidelberg, pp 227–232, doi:10.1007/978-3-662-01131-7_26
- Chatterjee S., Koopmans L. V. E., 2018, *Monthly Notices of the Royal Astronomical Society*, 474, 1762
- Chianese M., Coogan A., Hofma P., Otten S., Weniger C., 2020, *Monthly Notices of the Royal Astronomical Society*, 496, 381
- Ciotti L., Bertin G., 1999, *Astronomy and Astrophysics*, 352, 447
- Cole A., Miller B. K., Witte S. J., Cai M. X., Grootes M. W., Nattino F., Weniger C., 2021, Fast and Credible Likelihood-Free Cosmology with Truncated Marginal Neural Ratio Estimation ([arXiv:2111.08030](https://arxiv.org/abs/2111.08030)), doi:10.48550/arXiv.2111.08030
- Collett T. E., 2015, *The Astrophysical Journal*, 811, 20
- Coogan A., Correa C., Karchev K., Anau Montel N., Weniger C., 2022, Pulled from overleaf
- Cordts M., et al., 2016, The Cityscapes Dataset for Semantic Urban Scene Un-

- derstanding ([arXiv:1604.01685](https://arxiv.org/abs/1604.01685)), doi:[10.48550/arXiv.1604.01685](https://doi.org/10.48550/arXiv.1604.01685)
- Cornachione M. A., et al., 2018, *The Astrophysical Journal*, 853, 148
- Cranmer K., Pavez J., Louppe G., 2016, Approximating Likelihood Ratios with Calibrated Discriminative Classifiers ([arXiv:1506.02169](https://arxiv.org/abs/1506.02169)), doi:[10.48550/arXiv.1506.02169](https://doi.org/10.48550/arXiv.1506.02169)
- Cranmer K., Brehmer J., Louppe G., 2020, *Proc. Nat. Acad. Sci.*, 117, 30055
- Cyburt R. H., 2004, *Physical Review D*, 70, 023505
- Cyr-Racine F.-Y., Keeton C. R., Moustakas L. A., 2019, *Physical Review D*, 100, 023013
- Dalal N., Kochanek C. S., 2002, *The Astrophysical Journal*, 572, 25
- Dax M., Green S. R., Gair J., Macke J. H., Buonanno A., Schölkopf B., 2021, *Physical Review Letters*, 127, 241103
- Daylan T., Cyr-Racine F.-Y., Rivero A. D., Dvorkin C., Finkbeiner D. P., 2018, *The Astrophysical Journal*, 854, 141
- Deane R. P., Rawlings S., Garrett M. A., Heywood I., Jarvis M. J., Klöckner H.-R., Marshall P. J., McKean J. P., 2013, *Monthly Notices of the Royal Astronomical Society*, 434, 3322
- Delaunoy A., Wehenkel A., Hinderer T., Nisanke S., Weniger C., Williamson A. R., Louppe G., 2020, Lightning-Fast Gravitational Wave Parameter Inference through Neural Amortization ([arXiv:2010.12931](https://arxiv.org/abs/2010.12931)), doi:[10.48550/arXiv.2010.12931](https://doi.org/10.48550/arXiv.2010.12931)
- Despali G., Vegetti S., 2017, *Monthly Notices of the Royal Astronomical Society*, 469, 1997
- Despali G., Giocoli C., Angulo R. E., Tormen G., Sheth R. K., Baso G., Moscardini L., 2016, *Monthly Notices of the Royal Astronomical Society*, 456, 2486
- Diebolt J., Robert C. P., 1994, *Journal of the Royal Statistical Society. Series B (Methodological)*, 56, 363
- Dodelson S., Schmidt F., 2021, Modern Cosmology, second edition edn. Academic Press, an imprint of Elsevier, London, United Kingdom ; San Diego, CA
- Duffy L. D., van Bibber K., 2009, *New Journal of Physics*, 11, 105008
- Einstein A., 1916, *Annalen der Physik*, 354, 769
- Falcon W., Cho K., 2020, A Framework For Contrastive Self-Supervised Learning And Designing A New Approach ([arXiv:2009.00104](https://arxiv.org/abs/2009.00104)), doi:[10.48550/arXiv.2009.00104](https://doi.org/10.48550/arXiv.2009.00104)
- Fassnacht C. D., et al., 1999, *The Astronomical Journal*, 117, 658
- Fermi LAT Collaboration et al., 2017, *The Astrophysical Journal*, 840, 43
- Fong M., Bowyer R., Whitehead A., Lee B., King L., Applegate D., McCarthy I., 2018, *Monthly Notices of the Royal Astronomical Society*, 478, 5366
- Gardner J. P., et al., 2006, *Space Science Reviews*, 123, 485
- Gilman D., Birrer S., Treu T., Keeton C. R., Nierenberg A., 2018, *Monthly Notices of the Royal Astronomical Society*, 481, 819
- Gilman D., Birrer S., Nierenberg A., Treu T., Du X., Benson A., 2020, *Monthly Notices of the Royal Astronomical Society*, 491, 6077
- Giocoli C., Tormen G., Sheth R. K., van den Bosch F. C., 2010, *Monthly Notices of the Royal Astronomical Society*
- Harris C. R., et al., 2020, *Nature*, 585, 357
- He Q., et al., 2022, *Monthly Notices of the Royal Astronomical Society*, 511, 3046
- Hermans J., Begy V., Louppe G., 2020, Likelihood-Free MCMC with Amortized Approximate Ratio Estimators ([arXiv:1903.04057](https://arxiv.org/abs/1903.04057)), doi:[10.48550/arXiv.1903.04057](https://doi.org/10.48550/arXiv.1903.04057)
- Hermans J., Delaunoy A., Rozet F., Wehenkel A., Louppe G., 2021, Averting A Crisis In Simulation-Based Inference ([arXiv:2110.06581](https://arxiv.org/abs/2110.06581)), doi:[10.48550/arXiv.2110.06581](https://doi.org/10.48550/arXiv.2110.06581)
- Heymans C., et al., 2012, *Monthly Notices of the Royal Astronomical Society*, 427, 146
- Hezaveh Y. D., et al., 2016, *The Astrophysical Journal*, 823, 37
- Hezaveh Y. D., Levasseur L. P., Marshall P. J., 2017, *Nature*, 548, 555
- Hu W., Barkana R., Gruzinov A., 2000, *Physical Review Letters*, 85, 1158
- Hui L., Ostriker J. P., Tremaine S., Witten E., 2017, *Physical Review D*, 95, 043541
- Hunter J. D., 2007, *Computing in Science & Engineering*, 9, 90
- Ioffe S., Szegedy C., 2015, Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift ([arXiv:1502.03167](https://arxiv.org/abs/1502.03167)), doi:[10.48550/arXiv.1502.03167](https://doi.org/10.48550/arXiv.1502.03167)
- Jasra A., Holmes C. C., Stephens D. A., 2005, *Statistical Science*, 20, 50
- Karchev K., Coogan A., Weniger C., 2022, *Monthly Notices of the Royal Astronomical Society*, 512, 661
- Kavanagh B. J., Nichols D. A., Bertone G., Gaggero D., 2020, *Physical Review D*, 102, 083006
- Keskar N. S., Mudigere D., Nocedal J., Smelyanskiy M., Tang P. T. P., 2017, On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima ([arXiv:1609.04836](https://arxiv.org/abs/1609.04836)), doi:[10.48550/arXiv.1609.04836](https://doi.org/10.48550/arXiv.1609.04836)

- Kingma D. P., Ba J., 2014, Adam: A Method for Stochastic Optimization
- Kluyver T., et al., 2016, in Loizides F., Schmidt B., eds, 20th International Conference on Electronic Publishing (01/01/16). IOS Press, pp 87–90, doi:10.3233/978-1-61499-649-1-87
- Klypin A., Kravtsov A. V., Valenzuela O., Prada F., 1999, *The Astrophysical Journal*, 522, 82
- Kohl S. A. A., et al., 2019, arXiv:1806.05034 [cs, stat]
- Koopmans L. V. E., 2006, *EAS Publications Series*, 20, 161
- Kravtsov A., 2009, *Advances in Astronomy*, 2010, e281913
- LSST Science Collaboration et al., 2009, LSST Science Book, Version 2.0 (arXiv:0912.0201), doi:10.48550/arXiv.0912.0201
- Legin R., Hezaveh Y., Levasseur L. P., Wandelt B., 2021, arXiv:2112.05278 [astro-ph]
- Levasseur L. P., Hezaveh Y. D., Wechsler R. H., 2017, *The Astrophysical Journal*, 850, L7
- Liebes S., 1964, *Physical Review*, 133, B835
- Lin H., et al., 2009, *The Astrophysical Journal*, 699, 1242
- Lin J. Y.-Y., Yu H., Morningstar W., Peng J., Holder G., 2020, arXiv:2010.12960 [astro-ph, physics:physics]
- Lovell M. R., Frenk C. S., Eke V. R., Jenkins A., Gao L., Theuns T., 2014, *Monthly Notices of the Royal Astronomical Society*, 439, 300
- Lueckmann J.-M., Goncalves P. J., Bassetto G., Öcal K., Nonnenmacher M., Macke J. H., 2017, Flexible Statistical Inference for Mechanistic Models of Neural Dynamics (arXiv:1711.01861), doi:10.48550/arXiv.1711.01861
- Mao S., Schneider P., 1998, *Monthly Notices of the Royal Astronomical Society*, 295, 587
- Markevitch M., Gonzalez A. H., Clowe D., Vikhlinin A., Forman W., Jones C., Murray S., Tucker W., 2004, *The Astrophysical Journal*, 606, 819
- Meneghetti M., 2021, Introduction to Gravitational Lensing: With Python Examples. No. volume 956 in Lecture Notes in Physics, Springer, Cham, Switzerland
- Miller B. K., Cole A., Louppe G., Weniger C., 2020, Simulation-Efficient Marginal Posterior Estimation with Swyft: Stop Wasting Your Precious Time (arXiv:2011.13951), doi:10.48550/arXiv.2011.13951
- Miller B. K., Cole A., Forré P., Louppe G., Weniger C., 2021, arXiv:2107.01214 [astro-ph, physics:hep-ph, stat 10.5281/zenodo.5043706
- Minaee S., Boykov Y., Porikli F., Plaza A., Kehtarnavaz N., Terzopoulos D., 2020, Image Segmentation Using Deep Learning: A Survey (arXiv:2001.05566), doi:10.48550/arXiv.2001.05566
- Mishra-Sharma S., Cranmer K., 2022, *Physical Review D*, 105, 063017
- Mondino C., Taki A.-M., Van Tilburg K., Weiner N., 2020, *Physical Review Letters*, 125, 111101
- Moore B., Ghigna S., Governato F., Lake G., Quinn T., Stadel J., Tozzi P., 1999, *The Astrophysical Journal*, 524, L19
- Morningstar W. R., Hezaveh Y. D., Levasseur L. P., Blandford R. D., Marshall P. J., Putzky P., Wechsler R. H., 2018, Analyzing Interferometric Observations of Strong Gravitational Lenses with Recurrent and Convolutional Neural Networks (arXiv:1808.00011), doi:10.48550/arXiv.1808.00011
- Morningstar W. R., et al., 2019, *The Astrophysical Journal*, 883, 14
- Nair V., Hinton G. E., 2010, Omnipress, pp 807–814
- Navarro J. F., Frenk C. S., White S. D. M., 1996, *The Astrophysical Journal*, 462, 563
- Navarro J. F., Frenk C. S., White S. D. M., 1997, *The Astrophysical Journal*, 490, 493
- Niculescu-Mizil A., Caruana R., 2005, in Proceedings of the 22nd International Conference on Machine Learning. ICML '05. Association for Computing Machinery, New York, NY, USA, pp 625–632, doi:10.1145/1102351.1102430
- Nuesch P. E., 1991, *Journal of Applied Econometrics*, 6, 105
- O'Riordan C. M., Warren S. J., Mortlock D. J., 2020, *Monthly Notices of the Royal Astronomical Society*, 496, 3424
- Ostdiek B., Rivero A. D., Dvorkin C., 2022a, *Astronomy & Astrophysics*, 657, L14
- Ostdiek B., Rivero A. D., Dvorkin C., 2022b, *The Astrophysical Journal*, 927, 83
- Papamakarios G., Sterratt D., Murray I., 2019, in Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics. PMLR, pp 837–848
- Paszke A., et al., 2019, PyTorch: An Imperative Style, High-Performance Deep Learning Library (arXiv:1912.01703), doi:10.48550/arXiv.1912.01703
- Pearson J., Maresca J., Li N., Dye S., 2021, *Monthly Notices of the Royal Astronomical Society*, 505, 4362
- Peccei R. D., Quinn H. R., 1977a, *Physical Review D*, 16, 1791
- Peccei R. D., Quinn H. R., 1977b, *Physical Review Letters*, 38, 1440

- Pedregosa F., et al., 2018, Scikit-Learn: Machine Learning in Python ([arXiv:1201.0490](https://arxiv.org/abs/1201.0490)), doi:[10.48550/arXiv.1201.0490](https://doi.org/10.48550/arXiv.1201.0490)
- Planck Collaboration et al., 2016, *Astronomy and Astrophysics*, 594, A13
- Plotly~Technologies~Inc. 2015, Collaborative Data Science, Plotly Technologies Inc.
- Refregier A., Amara A., Kitching T. D., Rassat A., Scaramella R., Weller J., 2010, Euclid Imaging Consortium Science Book ([arXiv:1001.0061](https://arxiv.org/abs/1001.0061)), doi:[10.48550/arXiv.1001.0061](https://doi.org/10.48550/arXiv.1001.0061)
- Refsdal S., 1964, *Monthly Notices of the Royal Astronomical Society*, 128, 307
- Refsdal S., Bondi H., 1964, *Monthly Notices of the Royal Astronomical Society*, 128, 295
- Richings J., Frenk C., Jenkins A., Robertson A., Schaller M., 2021, *Monthly Notices of the Royal Astronomical Society*, 501, 4657
- Rivero A. D., Dvorkin C., 2020, *Physical Review D*, 101, 023515
- Roberts M. S., 1966, *The Astrophysical Journal*, 144, 639
- Robertson T., Wright F. T., Dykstra R., 1988, Order Restricted Statistical Inference. Wiley Series in Probability and Mathematical Statistics, Wiley, Chichester ; New York
- Ronneberger O., Fischer P., Brox T., 2015, U-Net: Convolutional Networks for Biomedical Image Segmentation ([arXiv:1505.04597](https://arxiv.org/abs/1505.04597)), doi:[10.48550/arXiv.1505.04597](https://doi.org/10.48550/arXiv.1505.04597)
- Roszkowski L., Sessolo E. M., Trojanowski S., 2018, *Reports on Progress in Physics*, 81, 066201
- Rubin V. C., Ford Jr. W. K., 1970, *The Astrophysical Journal*, 159, 379
- Saharia C., et al., 2022, Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding ([arXiv:2205.11487](https://arxiv.org/abs/2205.11487)), doi:[10.48550/arXiv.2205.11487](https://doi.org/10.48550/arXiv.2205.11487)
- Schneider P., 2006, in Schneider P., Kochanek C. S., Wambsganss J., eds, Saas-Fee Advanced Courses, Gravitational Lensing: Strong, Weak and Micro. Springer, Berlin, Heidelberg, pp 269–451, doi:[10.1007/978-3-540-30310-7_3](https://doi.org/10.1007/978-3-540-30310-7_3)
- Schneider P., Ehlers J., Falco E. E., 1992, Gravitational Lenses, first edn. Astronomy and Astrophysics Library, Springer Berlin, Heidelberg, doi:[10.1007/978-3-662-03758-4](https://doi.org/10.1007/978-3-662-03758-4)
- Schneider A., Smith R. E., Macciò A. V., Moore B., 2012, *Monthly Notices of the Royal Astronomical Society*, 424, 684
- Sérsic J. L., 1963, Boletín de la Asociación Argentina de Astronomía La Plata Argentina, 6, 41
- Shu Y., et al., 2017, *The Astrophysical Journal*, 851, 48
- Simon J. D., et al., 2019, Testing the Nature of Dark Matter with Extremely Large Telescopes ([arXiv:1903.04742](https://arxiv.org/abs/1903.04742)), doi:[10.48550/arXiv.1903.04742](https://doi.org/10.48550/arXiv.1903.04742)
- Sisson S. A., Fan Y., Beaumont M. A., 2018, Overview of Approximate Bayesian Computation ([arXiv:1802.09720](https://arxiv.org/abs/1802.09720)), doi:[10.48550/arXiv.1802.09720](https://doi.org/10.48550/arXiv.1802.09720)
- Skilling J., 2004, *AIP Conference Proceedings*, 735, 395
- Skilling J., 2006, *Bayesian Analysis*, 1, 833
- Soldner von J. G., 1801, Berliner Astronomisches Jahrbruch fuer Jahr 1804, 29, 161
- Springel V., et al., 2008, *Monthly Notices of the Royal Astronomical Society*, 391, 1685
- Stephens M., 1997, PhD thesis, Department of Statistics, University of Oxford.
- Stephens M., 2000, *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 62, 795
- Suyu S. H., Marshall P. J., Blandford R. D., Fassnacht C. D., Koopmans L. V. E., McKean J. P., Treu T., 2009, *The Astrophysical Journal*, 691, 277
- Tessore N., Metcalf R. B., 2015, *Astronomy & Astrophysics*, 580, A79
- The Astropy Collaboration et al., 2013, *Astronomy & Astrophysics*, 558, A33
- The Astropy Collaboration et al., 2018, *The Astronomical Journal*, 156, 123
- Treu T., 2010, *Annual Review of Astronomy and Astrophysics*, 48, 87
- Tsapras Y., 2018, *Geosciences*, 8, 365
- Tulin S., Yu H.-B., 2018, *Physics Reports*, 730, 1
- Varma S., Fairbairn M., Figueroa J., 2020, arXiv:2005.05353 [astro-ph]
- Vegetti S., Koopmans L. V. E., 2009a, *Monthly Notices of the Royal Astronomical Society*, 392, 945
- Vegetti S., Koopmans L. V. E., 2009b, *Monthly Notices of the Royal Astronomical Society*, 400, 1583
- Vegetti S., Czoske O., Koopmans L. V. E., 2010a, *Monthly Notices of the Royal Astronomical Society*, 407, 225
- Vegetti S., Koopmans L. V. E., Bolton A., Treu T., Gavazzi R., 2010b, *Monthly Notices of the Royal Astronomical Society*, 408, 1969
- Vegetti S., Lagattuta D. J., McKean J. P., Auger M. W., Fassnacht C. D., Koopmans L. V. E., 2012, *Nature*, 481, 341
- Wagner-Carena S., Park J. W., Birrer S., Marshall P. J., Roodman A., Wechsler R. H., 2021, *The Astrophysical Journal*, 909, 187
- Wagner-Carena S., Aalbers J., Birrer S., Nadler E. O., Darragh-Ford E., Marshall P. J., Wechsler R. H.,

- 2022, From Images to Dark Matter: End-To-End Inference of Substructure From Hundreds of Strong Gravitational Lenses ([arXiv:2203.00690](https://arxiv.org/abs/2203.00690)),
doi:10.48550/arXiv.2203.00690
- Walsh D., Carswell† R. F., Weymann‡ R. J., 1979, *Nature*, 279, 381
- Wechsler R. H., Tinker J. L., 2018, *Annual Review of Astronomy and Astrophysics*, 56, 435
- White S. D. M., Frenk C. S., Davis M., 1983, *The Astrophysical Journal*, 274, L1
- Will C. M., 1998, *American Journal of Physics*, 56, 413
- Winn J. N., Rusin D., Kochanek C. S., 2004, *Nature*, 427, 613
- Wong K. C., et al., 2020, *Monthly Notices of the Royal Astronomical Society*, 498, 1420
- XENON Collaboration et al., 2019, *Physical Review Letters*, 123, 251801
- Yuan X., Shi J., Gu L., 2021, *Expert Systems with Applications*, 169, 114417
- Zadrozny B., Elkan C., 2001, in Proceedings of the Eighteenth International Conference on Machine Learning. ICML '01. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp 609–616
- Zadrozny B., Elkan C., 2002, in Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '02. Association for Computing Machinery, New York, NY, USA, pp 694–699, [doi:10.1145/775047.775151](https://doi.org/10.1145/775047.775151)
- Zelko I. A., Treu T., Abazajian K. N., Gilman D., Benson A. J., Birrer S., Nierenberg A. M., Kusenko A., 2022, Constraints on Sterile Neutrino Models from Strong Gravitational Lensing, Milky Way Satellites, and Lyman-\$\alpha\$ Forest ([arXiv:2205.09777](https://arxiv.org/abs/2205.09777))
- Zhou Z., Siddiquee M. M. R., Tajbakhsh N., Liang J., 2018, UNet++: A Nested U-Net Architecture for Medical Image Segmentation ([arXiv:1807.10165](https://arxiv.org/abs/1807.10165)), [doi:10.48550/arXiv.1807.10165](https://doi.org/10.48550/arXiv.1807.10165)
- Zwicky F., 1933, *Helvetica Physica Acta*, 6, 110
- da Costa-Luis C., et al., 2022, Tqdm: A Fast, Extensible Progress Bar for Python and CLI, Zenodo, [doi:10.5281/zenodo.6412640](https://doi.org/10.5281/zenodo.6412640)