

Differentiable Strong Lensing: Uniting Gravity and Neural Nets through Differentiable Probabilistic Programming

Marco Chianese,^{1*} Adam Coogan,^{1†} Paul Hofma^{1‡} Sydney Otten^{1,2§}
and Christoph Weniger^{1¶}

¹*Gravitation Astroparticle Physics Amsterdam (GRAPPA), Institute for Theoretical Physics Amsterdam and Delta Institute for Theoretical Physics, University of Amsterdam, Science Park 904, 1098 XH Amsterdam, The Netherlands*

²*Institute for Mathematics, Astrophysics and Particle Physics (IMAPP), Radboud University, Heyendaalseweg 135, 6525 AJ Nijmegen, The Netherlands*

Accepted XXX. Received YYY; in original form ZZZ

ABSTRACT

Since upcoming telescopes will observe thousands of strong lensing systems, creating fully-automated analysis pipelines for these images becomes increasingly important. In this work, we make a step towards that direction by developing the first end-to-end differentiable strong lensing pipeline. Our approach leverages and combines three important computer science developments: (a) convolutional neural networks, (b) efficient gradient-based sampling techniques, and (c) deep probabilistic programming languages. The latter automatize parameter inference and enable the combination of generative deep neural networks and physics components in a single model. In the current work, we demonstrate that it is possible to combine a convolutional neural network trained on galaxy images as a source model with a fully-differentiable and exact implementation of gravitational lensing physics in a single probabilistic model. This does away with hyperparameter tuning for the source model, enables the simultaneous optimization of nearly one hundred source and lens parameters with gradient-based methods, and allows the use of efficient gradient-based posterior sampling techniques. These features make this automated inference pipeline potentially suitable for processing a large amount of data. By analyzing mock lensing systems with different signal-to-noise ratios, we show that lensing parameters are reconstructed with percent-level accuracy. More generally, we consider this work as one of the first steps in establishing *differentiable probabilistic programming* techniques in the particle astrophysics community, which have the potential to significantly accelerate and improve many complex data analysis tasks.

Key words: gravitational lensing: strong – galaxies: structure – dark matter.

1 INTRODUCTION

Strong lensing is a gravitational effect through which an astrophysical light source is observed in distorted, multiple images in the sky due to the deflection of its light by matter distributed along the line of sight (Treu 2010). It has become one of the main ways to probe the small-scale structure of dark matter halos, since subhalos with mass below $\sim 10^8 M_\odot$ do not host stars and are thus invisible (Fitts et al. 2017; Read et al. 2017). The detection (or non-detection) of

these subhalos is a critical tool for discriminating among different paradigms of dark matter (DM) (Bertone et al. 2005; Bertone & Tait 2018). In the standard Λ CDM cosmological model, the large-scale structures of the Universe form through the collapse of primordial density fluctuations. The matter content of the Universe is dominated by non-relativistic and almost collisionless substance dubbed cold dark matter (CDM). In this scenario, an abundance of small DM substructures is formed, as confirmed by ab initio N -body cosmological simulations (Kuhlen et al. 2012). On the other hand, alternative well-motivated particle DM scenarios such as warm dark matter (WDM) (Bode et al. 2001; Lovell et al. 2014) and self-interacting DM (Vogelberger et al. 2018; Kahlhoefer et al. 2019) predict a lower abundance of low-mass DM substructures. Recent analyses

* m.chianese@uva.nl

† a.m.coogan@uva.nl

‡ paul.hofma@student.uva.nl

§ s.m.m.otten@uva.nl

¶ c.weniger@uva.nl

of strong gravitational lensing of extended sources (Vegetti et al. 2010; Vegetti et al. 2012; Hezaveh et al. 2016; Vegetti et al. 2018; Ritondale et al. 2019) and quasars (Fadely & Keeton 2012; Gilman et al. 2019b,a) have already demonstrated sensitivity to DM haloes with masses larger than $10^8 M_{\odot}$. In the near future, new observatories like DES (Abbott et al. 2016), LSST (LSST Science Collaboration et al. 2009; Drlica-Wagner et al. 2019; Verma et al. 2019), Euclid (Refregier et al. 2010), and next-generation observatories like ELT (Simon et al. 2019) will observe thousands of strong lensing systems with very high precision, pushing the sensitivity of lensing probes of DM substructures to even lower masses. Moreover, these observations will be dominated by lenses at high redshift, increasing the likelihood of detecting small DM haloes along the line-of-sight (Despali et al. 2018).

In analyzing a galaxy-galaxy strong lens, the observed lensed image is reconstructed by simultaneously modeling the surface brightness of the source and the matter distribution of the lens galaxy. This requires parametrizing both of these components. While N -body simulations show that the density profiles of galactic DM halos are well-described by various analytic profiles (Navarro et al. 1996), the distribution of the source’s light is more complicated. The Sérsic brightness profile (Sérsic 1963) is a common choice (see e.g. Brewer et al. (2011)), but is inadequate for high-resolution observations and modeling high-redshift source galaxies, which generally have more complex morphologies than low-redshift ones. Another class of methods computes the source brightness profile on a grid by linearly inverting the observed lensed image given a fixed lens model (Warren & Dye 2003; Suyu et al. 2006). This requires using specific prior with a particular form to regularize the source, depending on two-point quantities calculated between pairs of pixels as well as hyperparameters. These methods can be cast in a fully Bayesian framework and performed on an adaptive grid in the source plane (Koopmans 2005; Vegetti & Koopmans 2009). The public code PyAutoLens implements this analysis strategy (Nightingale et al. 2018, 2019). Extensions of these methods include grid-free approaches using radial basis functions centered on image pixels ray-traced back to the source plane (Merten 2016) and methods decomposing the source as a sum of shapelets (Birrer et al. 2015). These methods are available in the SaWLens2 (Merten 2017) and Lenstronomy (Birrer & Amara 2018; Birrer 2019) software packages. Depending on the pipeline’s fitting scheme, the choice of priors for the lensing system’s parameters may be restricted, making it challenging to perform sensitivity analysis. These analysis pipelines also generally require dedicated, time-consuming hyperparameter optimization efforts for fitting each strong lensing system. On the other hand, fully-automated lensing pipelines which do not require these interventions will become increasingly useful for analyzing large strong lensing datasets, and for subsequently characterizing the underlying dark matter subhalo mass function. To this aim, it is of paramount importance to develop automated lens modeling techniques. In this paper, we demonstrate that this is possible through the use of new computational tools from the field of deep learning such as *automatic differentiation* (AD). This enables the use of powerful gradient descent based methods which are crucial to speed up and automatize the fitting procedure of high-dimensional pa-

rameter space. In particular, this paper focuses on how deep generative models can be combined with known physics using AD to create a new pipeline for analyzing images of galaxy-galaxy strong lensing systems.

Deep learning has advanced dramatically over the past decade, with accurate image classifiers (Krizhevsky et al. 2012) and a host of generative methods such as variational autoencoders (Kingma & Welling 2013; Jimenez Rezende et al. 2014), generative adversarial networks (Goodfellow et al. 2014; Brock et al. 2018; Karras et al. 2018) and flow-based models (Kingma & Dhariwal 2018) capable of producing novel, realistic images counting among its successes. These methods have also been applied to the physical sciences (Carleo et al. 2019), with topics including deblending galaxy images (Reiman & Göhre 2019), generating weak gravitational lens convergence maps (Mustafa et al. 2019), and classifying LHC jet events (Larkoski et al. 2017). However, there has been considerably less scientific investigation into leveraging automatic differentiation (Baydin et al. 2015), the core technology enabling training of neural networks with millions of parameters using gradient descent. AD libraries (Paszke et al. 2017; Revels et al. 2016; Innes et al. 2019) make it possible to take exact derivatives of arbitrary computable functions by using the chain rule to compose the gradients of individual operations.

The approach of creating automatically-differentiable mechanistic models that can be combined with deep neural networks is known as *differentiable programming* (Innes et al. 2019). Differentiable programming has recently been applied to engineering problems, with demonstrated benefits for challenging optimization problems in various domains. Constructing a differentiable ray-tracer (Li et al. 2018b) and image-processing algorithms (Li et al. 2018a; V. Sitzmann 2018) simplifies the inverse problem of fitting parameters describing the lighting, materials and objects in a scene from a photograph, as well as the forward problem of optimizing image-processing pipelines. Differentiable rigid-body physics engines make it possible to train deep learning controllers for robots in accurate environments (Degraeve et al. 2016). These simulators can also be combined with neural networks to model physical systems based on video and predict their future behavior (de Avila Belbute-Peres et al. 2018).

Recently, deep learning methods have been brought to bear on strong lensing analyses. Hezaveh et al. (2017) and Perreault Levasseur et al. (2017) applied a convolutional neural network (CNN) (LeCun et al. 1989) to infer the parameters of singular isothermal ellipsoid lenses (Keeton & Kochanek 1998; Kormann et al. 1994) with unprecedented speed. This CNN was later coupled to an optimizer controlled by another CNN that performs a linear inversion to recover the source’s pixelated brightness profile without requiring regularization hyperparameters (Morningstar et al. 2019). Both CNNs were trained using supervised learning on mock lensing datasets. Very recently, supervised CNNs have been trained using toy models of lensing systems containing substructure to differentiate between different DM models (Alexander et al. 2019) and infer parameters in the subhalo mass function (Brehmer et al. 2019). Moreover, in Diaz Rivero & Dvorkin (2019) a supervised CNN trained to classify whether or not a lensing system observation contained detectable substructure, thus identifying images worthy of follow-up observations and analyses.

In this work we construct the first differentiable programming-based strong lensing analysis pipeline, consisting of a deep generative model for the source galaxy brightness profile and a physics model for the lens. This approach removes the need for manual modeling or hyperparameter optimization and speeds up inference using gradient-based techniques. In particular, we use a variational autoencoder (VAE) (Kingma & Welling 2013; Jimenez Rezende et al. 2014) to learn a parametric description of the source galaxy’s light. In contrast with the aforementioned deep learning strategies for lensing, the VAE is trained in an unsupervised manner on unlensed galaxies; a similar VAE was previously constructed in Ravanbakhsh et al. (2016). The lensing physics is implemented by solving the Poisson equation for analytical models of the lens and external shear. Since our analysis pipeline is modularized, it is possible to change parameters’ priors or include additional lensing effects such as multiple sources and lenses, and line-of-sight halos and subhalos without having to retrain the source VAE. The source model can also be improved independently from the rest of the code.

Our pipeline is implemented in the PyTorch machine learning framework (Paszke et al. 2017), which contains an automatic differentiation engine and graphical processing unit (GPU)-accelerated functions for array computations. As with the differentiable programming examples mentioned above, pervasive AD makes it straightforward to fit the high-dimensional parameter vector of our lensing model using gradient-based methods. We also exploit the Pyro (Bingham et al. 2018) probabilistic programming language to characterize the uncertainties of our parameter fits. Pyro is capable of sampling model parameters, computing likelihoods, and automatically performing inference with tools such as Markov Chain Monte Carlo (MCMC) and variational inference (Blei et al. 2016; Hoffman et al. 2012) on arbitrary probabilistic models, even those with stochastic control flow. While probabilistic programming has existed for a long time (Lunn et al. 2000, 2009), it has only recently been integrated with automatic differentiation, making parameter inference possible even for models including neural networks (Tran et al. 2016; Carpenter et al. 2017; Burrone et al. 2018; Bingham et al. 2018; Ge et al. 2018; Charnock et al. 2019). In the context of high energy physics, recent work has reframed the Sherpa event generator using probabilistic programming, enabling more efficient simulation of rare events (Baydin et al. 2018; Güneş Baydin et al. 2019). In our present work, we again leverage the differentiability of our pipeline by employing Hamiltonian Monte Carlo (HMC) (Neal 2012; Betancourt 2017), a gradient-based MCMC method, to efficiently sample from the high-dimension posterior for the lens and source parameters. The unique meshing of an exact physics model with a deep generative source model in an AD-compatible, probabilistic programming framework makes the current paper one of the first differentiable probabilistic programs to our knowledge.

We begin this paper by discussing the variational autoencoder model for source galaxies in Section 2, and describing the training procedure and validation tests. In Section 3 we review the physics of strong gravitational lensing, and introduce all the ingredients required to generate lensed images. In particular, we define the lens and the source models considered in the present paper. Section 4 delineates our

lensing inference pipeline. In Section 5 we test the pipeline on two mock galaxy-galaxy lensing systems, and discuss parameter fitting and posterior analysis. We draw our conclusions in Section 6.

2 DEEP GENERATIVE MODELS FOR SOURCE GALAXIES

This section describes how we use a variational autoencoder (VAE) to construct a parametric model for source galaxy light. After reviewing VAEs and explaining the architecture we selected, we describe how ours is trained and present tests validating the performance of the parametric model we use it to construct.

2.1 Variational Autoencoders

Natural images (such as those of galaxies) lie on or near a low-dimensional submanifold in the space of all possible images (Belkin 2003). This motivates the concept of *probabilistic latent variable models* (Bishop 1998), where a datum \mathbf{x} (such as a galaxy image) is related to a latent variable \mathbf{z} through a conditional probability density $p(\mathbf{x}|\mathbf{z})$, and the latent variables are assumed to have a prior density $p(\mathbf{z})$. We will use this to construct a parametric model for galaxy images, where \mathbf{z} maps onto the mean of the decoding distribution and has prior $p(\mathbf{z})$.

The variational autoencoder (VAE) was constructed to efficiently approximate models of this form (Kingma & Welling 2013; Jimenez Rezende et al. 2014). It approximates the conditional density $p(\mathbf{x}|\mathbf{z})$ with a decoder $d_\theta(\mathbf{x}|\mathbf{z})$ whose parameters are functions represented by a neural network; θ represents the network’s parameters. The VAE also includes an encoder, $e_\phi(\mathbf{z}|\mathbf{x})$, that similarly approximates the conditional distribution $p(\mathbf{z}|\mathbf{x})$ using a neural network with parameters ϕ . Figure 1 illustrates this structure. The decoder and encoder are typically both taken to be diagonal Gaussians:

$$d_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mu_d(\mathbf{z}), \sigma_d) \quad (1)$$

$$e_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mu_e(\mathbf{x}), \sigma_e(\mathbf{x})). \quad (2)$$

In our notation $\mathcal{N}(x|\mu, \sigma)$ denotes that x follows a normal distribution with mean μ and standard deviation σ . The functions μ_d , μ_e and σ_e are implemented using neural networks, and the decoder’s standard deviation σ_d is a constant hyperparameter. We identify σ_d with the approximate standard deviation of the Gaussian noise in our training dataset.

Training VAE requires a dataset $\{\mathbf{x}\}$ as well as selecting the latent space’s dimensionality and prior $p(\mathbf{z})$, which is typically taken to be $\mathcal{N}(0, I)$.¹ Ideally, the VAE would be trained by maximizing the marginal likelihood of the training dataset, which requires computing the integral

$$p_\theta(\mathbf{x}) = \int d\mathbf{z} d_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \quad (3)$$

for each training point. Note that the marginal likelihood

¹ Several works have studied more complex prior distributions, including learnable ones (Chen et al. 2016; Dilokthanakul et al. 2016; Tomczak & Welling 2017; Alemi et al. 2017; Goyal et al. 2017).

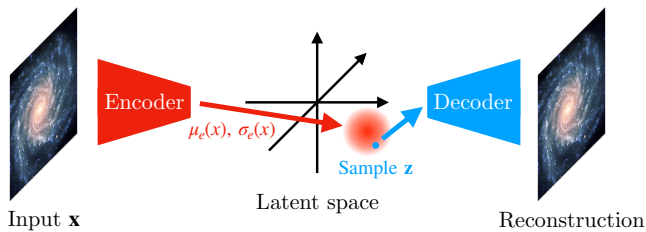


Figure 1. Diagram of a variational autoencoder. The diagram shows how an input image \mathbf{x} is passed through the encoder (red trapezoid) to yield an encoding distribution (red blob) with mean and standard deviation $\mu_e(\mathbf{x})$ and $\sigma_e(\mathbf{x})$. A point is then sampled from this distribution (blue dot) and passed through the decoder (blue trapezoid) to yield a reconstructed image.

only depends on the decoder. However, this integral is generally intractable. The difficulty is circumvented by defining an alternative objective function called the evidence lower bound (ELBO), whose derivation is reviewed in Appendix B. The reason the encoder network was introduced is because it is required to compute the ELBO, whose value for each training point is given by²

$$\begin{aligned} \log p_\theta(\mathbf{x}) &\geq \sum_i \text{ELBO}(\mathbf{x}; \theta, \phi) \\ &\equiv \mathbb{E}_{e_\phi(\mathbf{z}|\mathbf{x})} [\log d_\theta(\mathbf{x}|\mathbf{z})] - D_{KL} [e_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]. \end{aligned} \quad (4)$$

The first term is related to the quality of the reconstruction obtained by passing an image through the encoder followed by the decoder: maximizing this term improves reconstruction quality. The function $D_{KL}[\cdot||\cdot]$ in the second term is the Kullback-Leibler divergence, which measures the difference between its two argument probability distributions. This term’s maximization drives the averaged encoding distribution to look more like the prior $p(\mathbf{z})$. The VAE is trained to maximize the ELBO with stochastic gradient descent by taking random minibatches of training images. The second term can often be computed analytically, which makes computing its gradient straightforward. A Monte Carlo estimate of the first term can be computed. Appendix C explains how the derivative of this estimate can be computed.

Our VAE architecture choice is influenced by three widespread trends in deep convolutional neural network design (Radford et al. 2015): replacing pooling functions with strided convolutions, relying on convolutional layers rather than fully-connected ones when possible, and interleaving batch normalization layers (Ioffe & Szegedy 2015). In more detail, the encoder uses five blocks made up of strided convolutions, batch normalization layers and LeakyReLU (leaky rectified linear unit) activation functions (Maas et al. 2013). The output from these blocks is processed by two separate dense layers, which give the $\mu_e(\mathbf{x})$ and $\sigma_e(\mathbf{x})$, the parameters of the encoding distribution. The decoder is similarly built from five blocks consisting of a transposed convolution, batch normalization layers and ReLU (rectified linear unit) (Hahnloser et al. 2000) activation function. The output of last block is passed through a tanh activation function to

produce $\mu_d(\mathbf{x})$, the pixel values of which are thus restricted to lie within $[-1, 1]$. The standard deviation of the decoding distribution was set to $\sigma_d = 1/50$, which is the approximate standard deviation of the noise in our training dataset. We also studied the impact of making σ_d a trainable parameter (as recommended in Dai & Wipf (2019)), in which case we find it converges to approximately this value. We used 64 latent-space dimensions. Complete details of the architecture and weight initializations can be found in Appendix D, where we also describe other architectures with which we experimented.

After training the VAE as described in the next subsection, we obtain a parametric model for galaxy images where \mathbf{z} corresponds to the image $\mu_d(\mathbf{z})$, and has prior $p(\mathbf{z})$. However, it is a well-known and difficult-to-solve problem that the assumed prior $p(\mathbf{z})$ does not actually match the “aggregate prior” obtained by encoding the full training dataset, (Jimenez Rezende et al. 2014; Alemi et al. 2017)

$$q_\phi(\mathbf{z}) \equiv \frac{1}{N} \sum_{i=1}^N e_\phi(\mathbf{z}|\mathbf{x}^{(i)}). \quad (5)$$

This mismatch can cause problems when performing maximum a posteriori (MAP) estimation of \mathbf{z} for a given galaxy image. The assumed prior can drag \mathbf{z} into unrealistic regions of parameter space far from the location preferred by the likelihood, where the corresponding decoded image $\mu_d(\mathbf{z})$ does not look like a galaxy.

A variety of methods have been proposed to address this problem, including constructing simple priors using $q_\phi(\mathbf{z})$ (Tomczak & Welling 2017; Otten et al. 2019), fitting $q_\phi(\mathbf{z})$ with a second VAE after training the primary one (Dai & Wipf 2019), or using normalizing flows (Jimenez Rezende & Mohamed 2015; Kingma et al. 2016; Papamakarios et al. 2017; van den Berg et al. 2018). Here we follow the simpler approach of creating a weakly-informative prior for \mathbf{z} by fitting a multivariate normal distribution to the set of encoded means of the training data $\{\mu_e(\mathbf{x}^{(i)})\}_{i=1}^N$ and rescaling its covariance matrix by a factor of 9. This prior roughly confines the latent variable to realistic regions of the latent space while remaining diffuse enough that the likelihood drives MAP estimates of \mathbf{z} from observations.

2.2 Training

We construct a dataset for training the VAE starting from 56,062 images of galaxies in the COSMOS field taken by the Hubble Space Telescope (Scoville et al. 2007a,b; Koekemoer et al. 2007). These were used for the GREAT3 weak gravitational lensing challenge (Mandelbaum et al. 2014; Team 2019). An estimate of the pixel noise is included for each image, which is assumed to be Gaussian and uncorrelated. Parameters for Sersic profiles fit to each of the galaxies are also provided. Details about the image processing and parameter fits can be found in Appendix E of Mandelbaum et al. (2014).

We discard the small number of galaxies for which the Sersic fits failed or gave unphysical parameters. The dimensions of the images in the dataset vary, so we remove any that are smaller than 64×64 pixels. Images with unequal width and height are cropped into squares, and all are then down-

² The objective eq. (4) is equivalent to the β -VAE (Higgins et al. 2017) objective obtained by taking $\beta = \sigma^2$ and setting the decoder’s standard deviation to 1.

scaled to 64×64 pixels. Finally, many of the galaxies are extremely bright and small, and some are nearly indistinguishable from the pixel noise. In our experiments both of these degraded the quality of the VAE’s reconstructions, in the former case by biasing it towards producing only compact, bright galaxies and in the later case by degrading the fidelity of the reconstructions. We find a useful, very heuristic way of removing these is to make a cut on the image “signal-to-noise” quantity $\max(\mathcal{I}_{1/4})/\sigma_{\text{noise}}$, where σ_{noise} is the standard deviation of the pixel noise and $\mathcal{I}_{1/4}$ is the image downsampled by a factor of 4. Restricting this quantity to lie between 15 and 50 reduces the dataset to 17,543 images. We then split these into a training, test and validation datasets consisting of 15,500, 500 and 1,543 images respectively. This is a fairly small training dataset by industrial machine learning standards (Deng et al. 2009), but could be greatly augmented by future astronomical observations.

The training and validation sets are preprocessed by dividing each image by its maximum pixel value. The VAE’s parameters are optimized to minimize the ELBO of the training using the Adam optimizer (Kingma & Ba 2015), with a learning rate of 10^{-6} , minibatch size of 32 and momentum parameters $(\beta_1, \beta_2) = (0.5, 0.999)$. The mean ELBO for the images in the validation set is monitored during training. While this decreases at first, it inevitably increases as the VAE starts overfitting the training data. We terminate training at this point (~ 450 epochs in practice). The test-set images are not seen by the VAE during training.

2.3 Validation tests

To assess the quality of our parametric model $\mu_d(\mathbf{z})$ for galaxy images, we present the reconstructions of galaxies from the test set that the model has not seen during training in Fig. 2. We can observe that the reconstructions of our VAE model are denoised versions of the original images that reliably contain galaxy substructure. This indicates the latent space learned by the VAE contains a point corresponding to each of these galaxies, even though they have complex and varied morphologies.

It is also important that the VAE’s reconstructions are equivariant under transformations such as rotations, since the parameter inferences by the analysis pipeline should be as well. In Fig. 3, a galaxy image from the test set is rotated by various angles. Each of the rotated images is then derotated for comparison. The reconstructed images are once again denoised versions of the input images, and it can be seen by inspection that the derotated reconstructions are nearly identical. From this we conclude that our generative model has (approximately) learned rotational equivariance, even though the training dataset was not augmented with rotated images to teach this explicitly.

3 STRONG LENSING

Here we review the physics of strong gravitational lensing, describe how we model the main lens and external shear, and explain how we construct the source model using the variational autoencoder from the previous section.

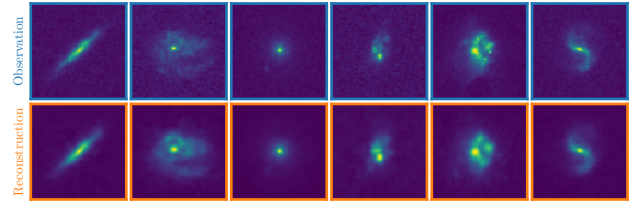


Figure 2. Reconstructions of galaxies from the test set. The input images are shown in the top row and the reconstructions in the bottom row. The reconstructions were obtained by passing the observations \mathbf{x} through the encoder, sampling $\mathbf{z} \sim \mathcal{N}(\mu_e(\mathbf{x}), \sigma_e(\mathbf{x}))$ from the encoded distribution, and taking the mean of the decoded distribution $\hat{\mathbf{x}} = \mu_d(\mathbf{z})$.

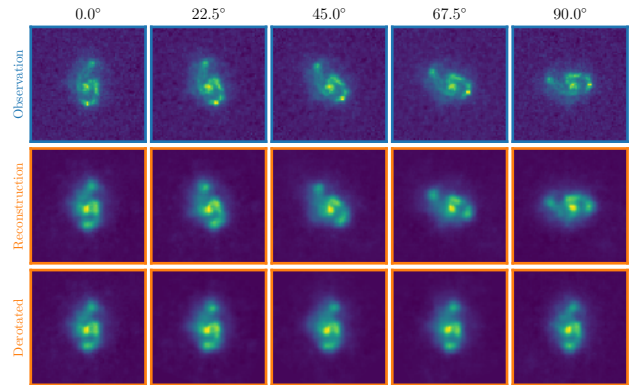


Figure 3. Approximate rotational invariance of reconstructions. The first row shows a galaxy image from the test set and versions rotated counterclockwise by the amount indicated above each column, with new pixels filled using the known noise distribution. The second row shows the corresponding VAE reconstructions, with derotated versions shown in the third row to make comparison simpler. The same color scale is used in each subplot.

3.1 Strong lensing physics

A gravitational galaxy-galaxy lensing system mainly consists of a background galaxy, playing the role of the source, and a foreground galaxy, acting as the main lens that deflects the light through its gravitational potential. In the thin lens approximation, the relation between the two-dimensional angular coordinates of the lens plane θ and the ones in the source plane β is encoded by the lens equation (Kormann et al. 1994; Treu 2010)

$$\beta = \theta - \alpha(\theta), \quad (6)$$

where α is the displacement field, which defines the deflection experienced by the light ray. The geometry of the system is displayed in Fig. 4. The displacement field depends on the Newtonian gravitational potential related to the mass distribution of the foreground galaxy. By means of the Poisson equation, it can be expressed as

$$\alpha = \frac{4G_N}{c^2} \frac{D_{\text{OL}}D_{\text{LS}}}{D_{\text{OS}}} \int \Sigma(\theta') \frac{\theta - \theta'}{|\theta - \theta'|^2} d^2\theta', \quad (7)$$

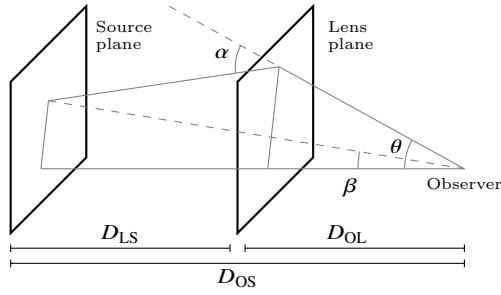


Figure 4. Diagram of a gravitational lensing system. The simplest galaxy-galaxy lensing system is represented by a background and a foreground galaxy, which define the source (S) and the lens plane (L), respectively. The quantities D are the angular diameter distances between the different planes and the observer (O). The two-dimensional angular coordinates, β and θ , are related through the displacement field α by the lens equation (6).

where Σ is the surface mass density of the lens, and the quantities D_{OL} , D_{OS} , D_{LS} are the angular diameter distances between the observer and the lens, the observer and the source, the lens and the source, respectively. Moreover, G_N is the gravitational constant while c is the speed of light.

Since the lens equation (eq. (6)) preserves the surface brightness (photon flux density per unit angular area), the image of the system in the lens plane, denoted as $\mathcal{I}_{\text{lens}}$, is simply obtained by evaluating the source light distribution \mathcal{I}_{src} on the lens plane. Hence, we have³

$$\mathcal{I}_{\text{lens}}(\theta) = \mathcal{I}_{\text{src}}(\theta - \alpha(\theta)). \quad (8)$$

This equation is solved on a square pixel grid defined in the lens plane according to the observed image. In particular, we consider a grid of 256×256 pixels with angular size $10 \text{ arcsec} \times 10 \text{ arcsec}$, corresponding to a pixel size of 0.04 arcsec .⁴ Then, the predicted lensed image $\mathcal{I}_{\text{pred}}$ is obtained by taking into account the Point Spread Function (PSF), which defines how a point-like source (a pixel) is spread due to atmospheric distortions and defects in the optics, and the noise from instrumental and astrophysical backgrounds. In the present paper, we consider a symmetric two-dimensional Gaussian PSF with a standard deviation of 0.05 arcsec . Moreover, to each pixel we add uncorrelated Gaussian noise with $\sigma_{\text{noise}} = 0.333$, $\sigma_{\text{noise}} = 0.1$ and $\sigma_{\text{noise}} = 0.0333$ for mock images with low, medium and high signal-to-noise (S/N) ratio, respectively. Hence, we have

$$\mathcal{I}_{\text{pred}} = \mathcal{N}(\text{PSF} * \mathcal{I}_{\text{lens}}, \sigma_{\text{noise}}), \quad (9)$$

where the symbol $*$ stands for the mathematical convolution between the Point Spread Function and the image in the lens plane.

The two main ingredients describing a gravitational lensing system are, therefore, the total mass distribution of the lens Σ (the lens model) and the surface brightness profile of the background source \mathcal{I}_{src} (the source model).

³ Note that we do not consider the light distribution of the foreground galaxy since it is in general subtracted in real data analyses.

⁴ For reference, this is the pixel size of the optical/UV CCDs of WFC3 Hubble (Dressel 2012).

3.2 Lens model

In a typical lensing system, the dominant contributions to the total displacement field α come from the smooth main halo (mh) of the foreground galaxy and the external shear (ext), namely

$$\alpha = \alpha_{\text{mh}} + \alpha_{\text{ext}}. \quad (10)$$

For the first component, we consider the so-called Singular Power-Law Ellipsoid (SPLE) profile to model the surface mass distribution of the main halo. Such profiles can fit gravitational potentials of lenses in images of galaxy-scale strong lensing systems at the percent level (see e.g. Suyu et al. (2009)). In this case, we have (Kassiola & Kovner 1993; Barkana 1998)

$$\Sigma_{\text{mh}} = \Sigma_{\text{cr}} \frac{3 - \gamma}{2} \left[\frac{\rho(\theta', q)}{r_{\text{Ein}}} \right]^{1-\gamma}, \quad (11)$$

where $\Sigma_{\text{cr}} = c^2 D_{OS} / (4\pi G D_{OL} D_{LS})$ is the critical surface mass density, γ is the slope, r_{Ein} denotes the Einstein radius, and ρ encodes the dependence on the position. In the special case $\gamma = 2$ this distribution reduces to the one for a singular isothermal ellipsoid (Kormann et al. 1994). In the coordinates system $\theta'(\theta - \theta_{\text{lens}}, \phi)$ with origin in the centroid of the foreground galaxy (denoted as θ_{lens}) and axes aligned to the minor and major axes of the elliptical galaxy (after a rotation of an angle ϕ), we have

$$\rho(\theta', q) = \theta_1'^2 + \theta_2'^2 / q^2, \quad (12)$$

with q being the minor to major axis ratio of the elliptical contours of equal surface mass density. Hence, the displacement field α_{mh} induced by the main halo mass distribution is simply given by plugging eq. (11) into eq. (7).

Since numerical integration is difficult to implement in an automatically-differentiable manner, we instead compute α_{mh} by interpolating over a precomputed grid, as described in detail in Appendix A. This enables computation of gradients of α_{mh} since the interpolation function is itself automatically differentiable.

The external shear contribution represents the additional angular structure provided by additional matter distribution in the cluster where the galaxy is located. The corresponding displacement field is parametrized as

$$\alpha_{\text{ext}} = \begin{pmatrix} \gamma_1 & \gamma_2 \\ \gamma_2 & -\gamma_1 \end{pmatrix} \theta. \quad (13)$$

Our lens model is thus completely defined by a set of 8 parameters: $\Theta_{\text{lens}} \equiv \{\gamma_1, \gamma_2, \phi, q, \gamma, r_{\text{Ein}}, \theta_{\text{lens},1}, \theta_{\text{lens},2}\}$.

3.3 Source model

The VAE's decoder $d_{\theta}(\mathbf{z})$ provides a parametrized model for 64×64 -pixel galaxy images, which is the basis for our source model. One of its parameters is the 64-dimensional vector \mathbf{z} specifying a point in the VAE's latent space. As described earlier, the prior for \mathbf{z} is defined by fitting a multivariate normal distribution to the means of the encoded distributions of the training points and increasing its covariance by a factor of 9. We introduce four other parameters to complete the model: $\theta_{\text{src},1}$, $\theta_{\text{src},2}$, s and t . The first two are the position of the center of the decoded image $d(\mathbf{z})$. The second specifies the spatial scale of the image and the third is the

normalization of the pixel intensities. In this work we adopt the priors

$$\theta_{\text{src},1}, \theta_{\text{src},2} \sim \mathcal{N}(0, 0.1) \quad (14)$$

$$s \sim \mathcal{N}(5, 1) \quad (15)$$

$$\iota \sim \mathcal{N}(1, 0.5) \quad (16)$$

for our mock data generation and analysis. Our overall model for the surface brightness of the source galaxy at a position $\theta = (\theta_1, \theta_2)$ is thus

$$\mathcal{I}_{\text{src}}(\theta|\Theta_{\text{src}}) = \iota d(\mathbf{z}) \left(\frac{\theta_1 - \theta_{\text{src},1}}{s}, \frac{\theta_2 - \theta_{\text{src},2}}{s} \right), \quad (17)$$

where $\Theta_{\text{src}} \equiv \{\mathbf{z}, \theta_{\text{src},1}, \theta_{\text{src},2}, s, \iota\}$. We use bilinear interpolation to allow the right-hand side of this expression to be evaluated between adjacent pixels in the 64×64 output image from the decoder.

4 THE LENSING PIPELINE

Our lensing pipeline is unique since it combines the VAE-based source and physical lens models detailed in the previous section in a fully-differentiable manner. As sketched in Fig. 5, the pipeline consists of a *forward* flow (black solid lines) and a *backward* one (red dashed lines).

In the forward flow, the predicted lensed image $\mathcal{I}_{\text{pred}}(\Theta_{\text{lens}}, \Theta_{\text{src}})$ is obtained once the displacement field $\alpha(\theta|\Theta_{\text{lens}})$ and the source surface brightness $\mathcal{I}_{\text{src}}(\theta|\Theta_{\text{src}})$ have been computed in the Lens Model (see Section 3.2) and the Source Model (see Section 3.3), respectively. This image is then compared with the observed one to estimate the likelihood function given the parameters of the lens and the source models. Thanks to the differentiable programming framework, it is then possible to compute the derivatives of the likelihood function with respect to all the parameters, Θ_{lens} and Θ_{src} . This step represents the backward flow of the whole pipeline.

The pipeline is implemented in a differentiable probabilistic programming framework comprised of the PyTorch (Paszke et al. 2017) machine learning library and Pyro (Bingham et al. 2018) probabilistic programming library. PyTorch provides an automatic differentiation engine, enabling the backward flow of the pipeline. The likelihood calculations in the forward flow are made straightforward by Pyro. The VAE Source Model is constructed from neural network layers contained in PyTorch and trained using the variational inference module in Pyro. We employ Pyro’s variational inference and Hamiltonian Monte Carlo modules for parameter fitting and posterior sampling our mock data analysis in the following section.

We stress that automatic differentiability is automatically guaranteed if all the pipeline is written in the differentiable programming language. Moreover, we note that the lensing pipeline is implemented so that the lens and the source models are independent building blocks. In the present paper, the former is fully based on physical models while the latter is provided by the VAE’s decoder. However, thanks to the modularity of the pipeline, both can be easily modified or substituted, as can any of the priors on the lensing parameters. This fundamental feature allows one to generalize the present lensing pipeline to analyze more realistic systems and to include the gravitational effect of dark

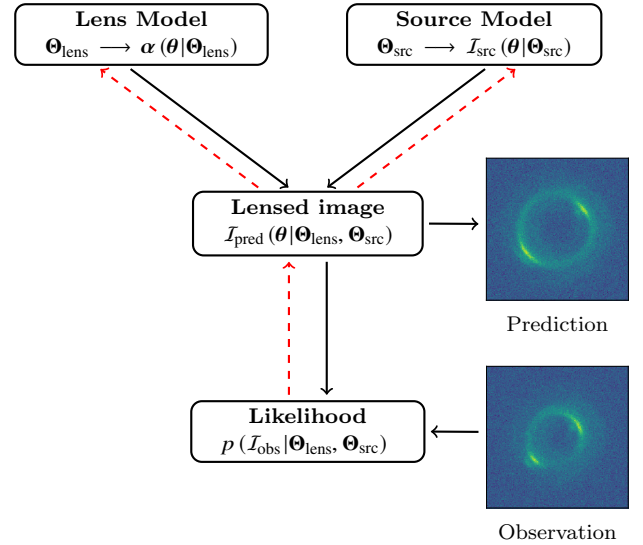


Figure 5. Lensing Pipeline. The forward flow (black solid lines) estimates the likelihood by comparing the observation with the predicted lensed image. This is obtained through the Lens Model and the Source Model. The backward flow (red dashed lines) computes the derivatives of the likelihood with respect to all the parameters of the models, Θ_{lens} and Θ_{src} . Each box represents an independent module.

matter substructures in the lensing physics. This is left for future investigation.

5 RESULTS

In this section, we test the lensing pipeline on mock lensing system observations. We describe the generation of mock data, and discuss the results obtained by the parameter fitting and the posterior analysis for two different mock lensing systems.

5.1 Mock data

To test our pipeline, we generate mock lensing system observations. We create mock sources by first denoising images from the test dataset using the non-local means algorithm (Buades et al. 2005, 2011). These images are then rescaled to fill roughly the central third of the source image plane, which ensures the lensed image pixel intensities drop to zero along the image boundaries. The mock lensing parameters are determined randomly by drawing from the following distributions:

$$\gamma_1, \gamma_2 \sim \mathcal{N}(0, 1) \quad (18)$$

$$\phi \sim \mathcal{N}(0, 1 \text{ rad}) \quad (19)$$

$$q \sim \mathcal{N}(0.5, 0.5) \quad (20)$$

$$r_{\text{Ein}} \sim \mathcal{N}(1.5 \text{ arcsec}, 0.5 \text{ arcsec}) \quad (21)$$

$$\gamma \sim \mathcal{N}(2, 0.5) \quad (22)$$

$$\theta_{\text{lens},1}, \theta_{\text{lens},2} \sim \mathcal{N}(0, 0.5). \quad (23)$$

The mock lensed images are then produced using the formalism from the previous section on a 256×256 pixel grid

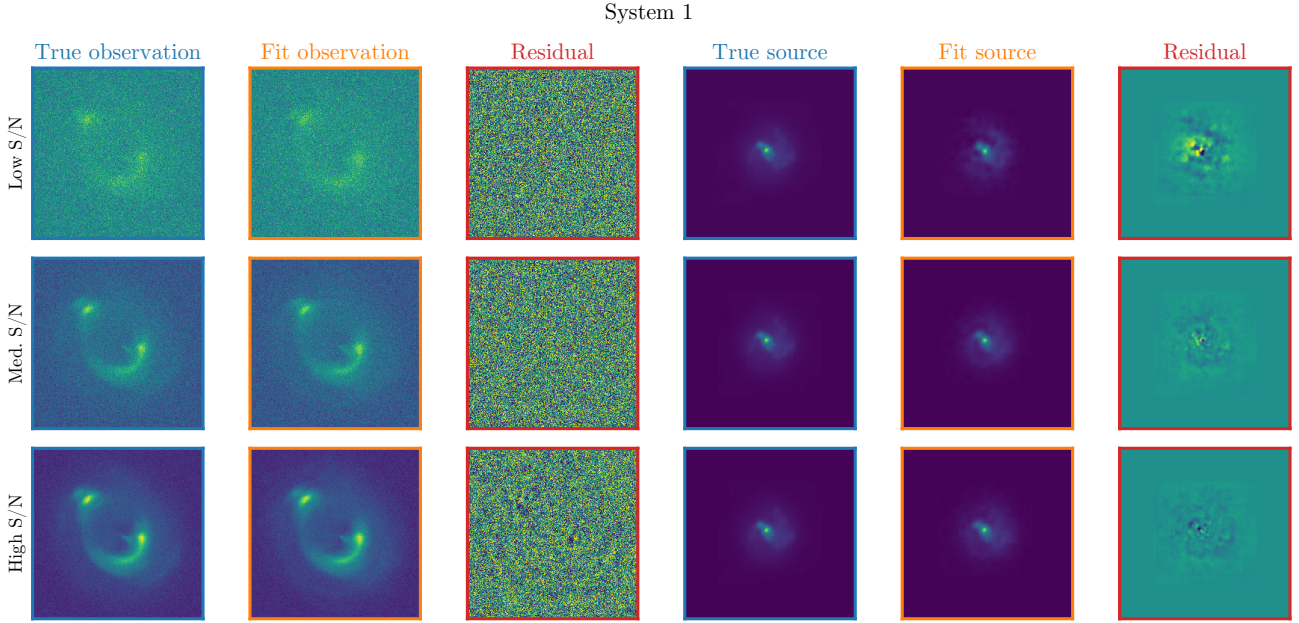


Figure 6. Results of MAP fit of mock lensing system 1. The different columns show the true and fit observations, the residual between these, the true and fit sources and the residuals between those. The rows correspond to different observed signal-to-noise values. Within each row, the color scales for the true and fit observations are the same, as are the scales for the true and fit sources. The observation residuals are normalized by dividing by the standard deviation of the observation noise. The source residuals are normalized by dividing by 10% of the maximum value of the source. The color scale ranges from dark blue to bright yellow.

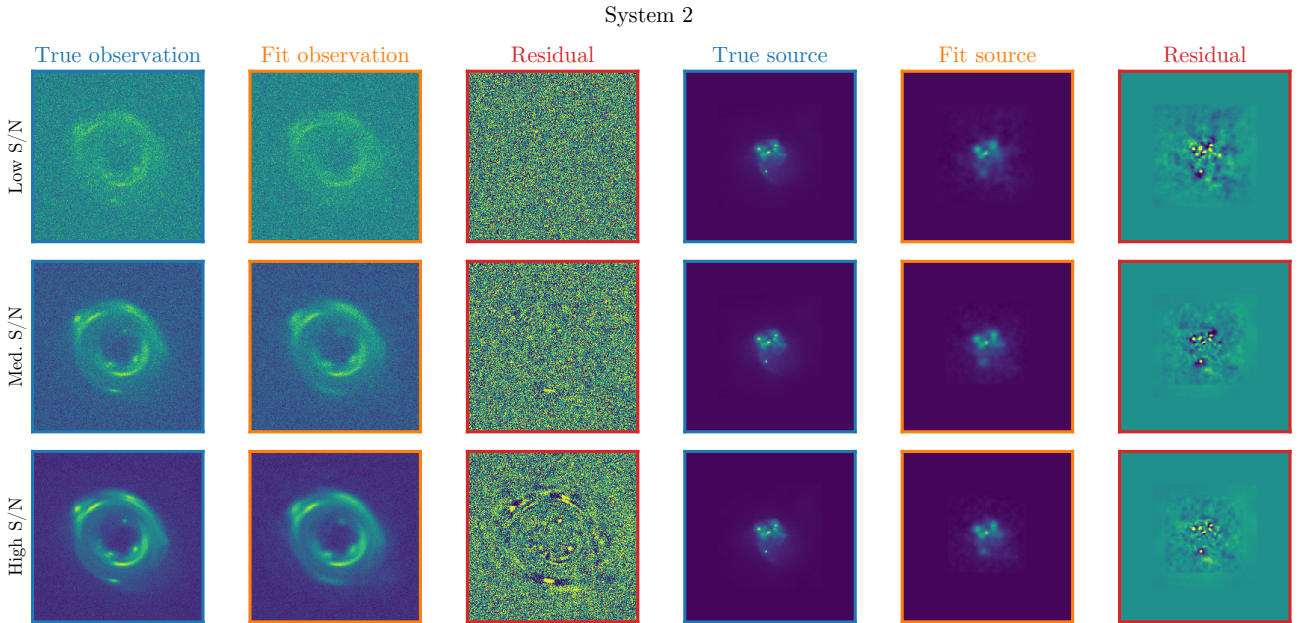


Figure 7. Results of MAP fit of mock lensing system 2. The subplots and scales are explained in the caption of Fig. 6.

	Low S/N	Med. S/N	High S/N
σ_{noise}	0.333	0.1	0.0333
System 1	66	272	757
System 2	77	313	849

Table 1. Signal-to-noise ratios for the systems considered in this work. The second row shows the pixel noise level while the third and fourth show the S/N values for the two systems.

with angular size $10 \text{ arcsec} \times 10 \text{ arcsec}$. Finally, we convolve with the PSF and add Gaussian pixel noise.

To test how data quality impacts our pipeline’s performance, we fix the pixel noise level σ_{noise} to three different values to obtain observations with low, medium and high signal-to-noise ratios. This ratio is defined by creating a mask m to select only the pixels belonging to the galaxy in the observed image \mathcal{I}_{obs} and computing the following quantity (see e.g. O’Riordan et al. (2019)):

$$S/N(\sigma_{\text{noise}}) = \frac{\sum_{i,j} m_{ij} \mathcal{I}_{\text{obs}ij}}{\sigma_{\text{noise}} \sqrt{\sum_{i,j} m_{ij}}}. \quad (24)$$

The mask is constructed by first convolving \mathcal{I}_{obs} with a Gaussian G with standard deviation 0.16 arcsec (four pixels) and then thresholding the blurred image using the noise level:

$$m_{ij} = \begin{cases} 1 & (G * \mathcal{I}_{\text{obs}})_{ij} \geq \sigma_{\text{noise}} \\ 0 & (G * \mathcal{I}_{\text{obs}})_{ij} < \sigma_{\text{noise}} \end{cases}. \quad (25)$$

Changing the threshold and/or width of the Gaussian G does not substantially change the value of S/N . The pixel noise levels and corresponding S/N values for each system are shown in Tab. 1 for completeness, though their specific values are not important for our study.

In the rest of this paper, we focus on two particular mock systems, hereafter referred to as system 1 and system 2. Source 1 has a fairly simple spiral morphology. Source 2 has a significant amount of substructure, making it representative of complex, high-redshift source galaxies.

5.2 Parameter fitting

We first test our pipeline by finding the best-fit source and lens parameter values. For a given image \mathcal{I}_{obs} we compute the maximum a posteriori parameter (MAP) estimates:

$$\hat{\Theta}_{\text{lens}}, \hat{\Theta}_{\text{src}} = \max_{\Theta_{\text{lens}}, \Theta_{\text{src}}} p(\Theta_{\text{lens}}, \Theta_{\text{src}} | \mathcal{I}_{\text{obs}}), \quad (26)$$

$$p(\Theta_{\text{lens}}, \Theta_{\text{src}} | \mathcal{I}_{\text{obs}}) \propto p(\mathcal{I}_{\text{obs}} | \Theta_{\text{lens}}, \Theta_{\text{src}}) p(\Theta_{\text{lens}}) p(\Theta_{\text{src}}).$$

The first term on the right-hand side of the posterior is the likelihood of the observed image for fixed source and lens parameters, which is Gaussian due to our noise assumption:

$$p(\mathcal{I}_{\text{obs}} | \Theta_{\text{lens}}, \Theta_{\text{src}}) = \mathcal{N}(\mathcal{I}_{\text{obs}} | \mathcal{I}_{\text{pred}}(\Theta_{\text{lens}}, \Theta_{\text{src}}), \sigma_{\text{obs}}).$$

The second term is the prior on the lens parameters, for which we adopt eqs. (18)-(23). The source priors are specified in eqs. (14)-(16), and the prior on \mathbf{z} is described in Section 2.

Since it is possible to differentiate through the full lensing pipeline, we obtain the best-fit parameters using the gradient-based Adam optimizer (Kingma & Ba 2015). Adam has been empirically shown to outperform other gradient-based methods in nonconvex optimization problems with

large numbers of parameters. These fits converge after $\sim 10^4$ iterations, which takes approximately 20 (7.5) minutes on CPU (GPU).

The MAP reconstructions of the mock lensing systems 1 and 2 are shown in Figs. 6 and 7. We report the true and fit images, and the residuals between these, in the lens and in the source planes, for three values of signal-to-noise ratio. In case of system 1, the residuals in the lens plane are at the noise level even for the smallest pixel noise level (highest S/N) considered. This is related to the fact that the reconstruction of the source improves as the signal-to-noise ratio increases, as can clearly be seen in the last column in Fig. 6. This does not occur for the complex source galaxy of system 2. Indeed, as shown in Fig. 7, even for the highest S/N value, the VAE’s decoder is not able to reproduce all the substructures exhibited in the source surface brightness. This directly affects the image reconstruction in the lens plane, and the corresponding residuals (third column in Fig. 7) increase and show more structure as the signal-to-noise ratio increases.

5.3 Posterior analysis

To study our pipeline’s parameter inference capabilities, we ran Hamiltonian Monte Carlo (Duane et al. 1987; Neal 2012; Betancourt 2017) to sample from the parameters’ posteriors. Hamiltonian Monte Carlo (HMC) is a Markov-chain Monte Carlo (MCMC) procedure that uses Hamiltonian dynamics based on the gradient of the posterior to efficiently traverse parameter space. HMC can take larger steps than other MCMC procedures such as Metropolis-Hastings while keeping the acceptance probabilities high, leading to more efficient exploration of parameter space. After 50 steps during which the internal HMC parameters are calibrated, we use it to sample 500 times from the posterior starting from the MAP parameter estimates. This takes about 10 (6.5) hours on a CPU (GPU).

The HMC results for the two systems are reported in Figs. 8 and 9, which show the marginalized posteriors of the lens parameters $(r_{\text{Ein}}, \gamma, \phi)$. As expected, for both systems the posteriors shrink as the S/N ratio increases. In particular, the statistical error on the parameter estimates moves from $\sim 3\%$ to $\lesssim 1\%$ for the lowest and highest S/N ratio, respectively. However, especially for the system 1, we note that our lensing pipeline gives biased estimates typically at the level of 1%. We hypothesize this is because the VAE generally produces slightly blurred and hazy reconstructions (most clearly visible in the last row of residual plots in Fig. 7). Improving the fidelity of reconstructions and samples from VAEs is an active area of research in machine learning (Zhao et al. 2017; Rezende & Viola 2018; Dai & Wipf 2019). In addition to improving the VAE architecture and training procedure, our source model would benefit from a larger, higher-resolution training dataset, as will be made available by future astronomical surveys. Surprisingly, the level of bias is smaller for the system 2, even though the source-plane residuals are larger due to the presence of fine substructure.

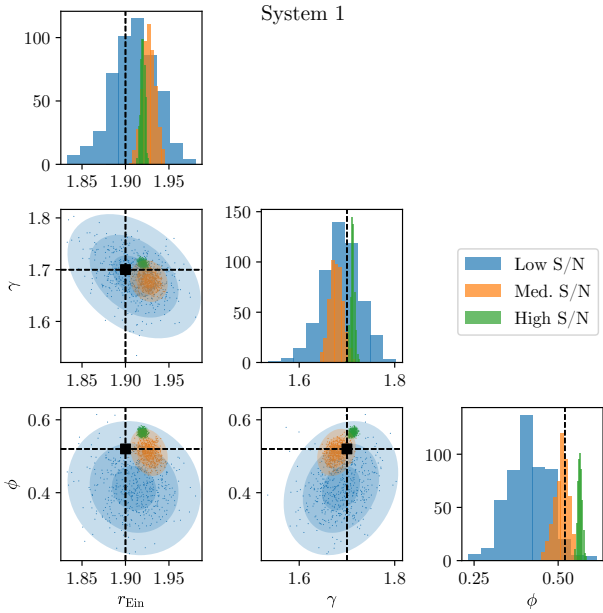


Figure 8. Results of HMC parameter posteriors of mock lensing system 1. The different panels show the marginalized one-dimensional and two-dimensional posteriors for a subset of lens parameter, ($r_{\text{Ein}}, \gamma, \phi$). The different colors (blue, orange, green) refer to the different signal-to-noise ratios while the shading in the ellipses corresponds to 68%, 95% and 99% confidence levels. The sampled points are also plotted. The black squares and dashed lines represent the true values of parameters.

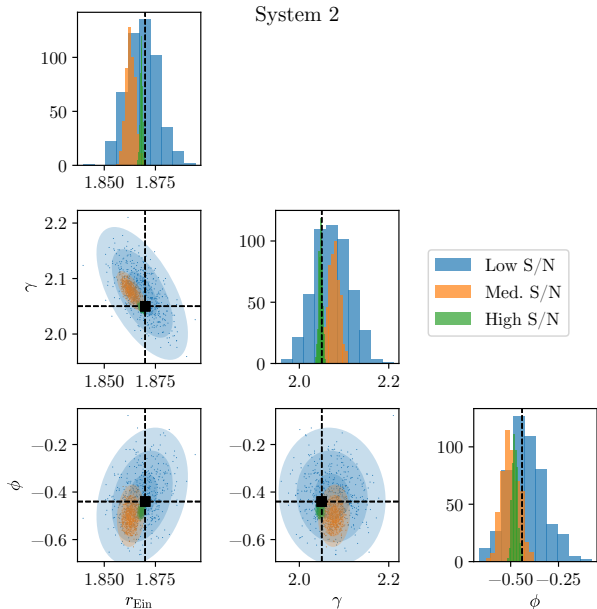


Figure 9. Results of HMC parameter posteriors of mock lensing system 2. The panels are explained in the caption of Fig. 8. Note the substantial overlap between the green ellipses and black squares.

6 CONCLUSIONS

We have presented the first step towards a new inference pipeline to analyze present and future strong gravitational lensing systems. The main novelty of our approach is the use of a differentiable probabilistic programming framework in which all operations are automatically differentiable with respect to the input parameters. This powerful approach makes Bayesian inference feasible for complex models with hundreds (or thousands) of free parameters thanks to efficient sampling techniques utilizing gradient descent.

Our lensing pipeline, shown in Fig. 5, consists of two independent blocks describing the surface brightness of the source galaxy (Source Model) and the mass distribution of the lens galaxy (Lens Model). They are combined to generate the lensed image, which is used to estimate the likelihood and perform Bayesian inference. The advantages of this strategy are:

- Exact gradients of the pipeline’s output can be computed with respect to its inputs using automatic differentiation. This makes it possible to use efficient gradient-based fitting and posterior sampling procedures.
- Using a differentiable probabilistic programming framework allows us to integrate the variational autoencoder source model learned from unlensed galaxy images directly with a physical lensing model. Learning the source model directly from data removes the need for tuning hyperparameters to regularize the source model.
- We fully automatize the inference step using probabilistic programming. The lens and source parameters are fit and sampled simultaneously and can have arbitrary priors.
- By implementing our pipeline in the PyTorch framework, we automatically gain the ability to perform computations on graphical processing units.

From a quantitative perspective, the best-fit lens parameter values we obtained in our mock data tests were within $\sim 1\%$ of the true values, albeit with some bias at very high signal-to-noise ratios.

Since our fitting and posterior sampling is gradient-based, the computational cost of adding parameters to increase the realism of our model is low. For example, a prerequisite for applying our pipeline to real data is to model the light from the lens galaxy. We anticipate this could be done by adding an additional unlensed light source again modeled by our variational autoencoder. In future work, we will also explore more realistic lens models with additional components such as dark matter subhalos, a baryonic disk, line-of-sight halos, or a more complicated main lens model. For these very high-dimensional models, automated gradient-based inference techniques with favorable scaling behavior such as variational inference (Hoffman et al. 2012; Blei et al. 2016) could enable analysis of posterior distributions for hundreds or even thousands of parameters. Lastly, we expect that our pipeline’s efficiency in analyzing large datasets could be improved by making technical changes so it can operate on batches of images in parallel, as is done when training convolutional neural networks.

While our pipeline is an early example of differentiable probabilistic programming, we anticipate this approach will enable other challenging and exciting data analyses in the

future by leveraging the advantages of deep learning and physics modeling.

ACKNOWLEDGEMENTS

We would like to thank Simona Vegetti and Rajat Thomas for helpful conversations. This work was carried out on the Dutch national e-infrastructure with the support of SURF Cooperative. We acknowledge funding from the Netherlands Organization for Scientific Research (NWO) through the VIDI research program “Probing the Genesis of Dark Matter” (680-47-532).

REFERENCES

- Abbott T., et al., 2016, *Mon. Not. Roy. Astron. Soc.*, 460, 1270
- Alemi A. A., Poole B., Fischer I., Dillon J. V., Saurous R. A., Murphy K., 2017, CoRR, abs/1711.00464
- Alexander S., Gleyzer S., McDonough E., Toomey M. W., Usai E., 2019
- Barkana R., 1998, *Astrophys. J.*, 502, 531
- Baydin A. G., Pearlmutter B. A., Radul A. A., 2015, CoRR, abs/1502.05767
- Baydin A. G., et al., 2018, CoRR, abs/1807.07706
- Belkin M., 2003, PhD thesis, The University Of Chicago
- Bertone G., Tait Tim M. P., 2018, *Nature*, 562, 51
- Bertone G., Hooper D., Silk J., 2005, *Phys. Rept.*, 405, 279
- Betancourt M., 2017, arXiv e-prints,
- Bingham E., et al., 2018, Journal of Machine Learning Research
- Birrer S., 2019, sibirrer/lenstronomy: multi-purpose lens modeling software package, <https://github.com/sibirrer/lenstronomy>
- Birrer S., Amara A., 2018,] 10.1016/j.dark.2018.11.002
- Birrer S., Amara A., Refregier A., 2015, *Astrophys. J.*, 813, 102
- Bishop C. M., 1998, Latent Variable Models. Springer Netherlands, Dordrecht, pp 371–403, doi:10.1007/978-94-011-5014-9_13, https://doi.org/10.1007/978-94-011-5014-9_13
- Blei D. M., Kucukelbir A., McAuliffe J. D., 2016, arXiv e-prints, p. arXiv:1601.00670
- Bode P., Ostriker J. P., Turok N., 2001, *Astrophys. J.*, 556, 93
- Brehmer J., Mishra-Sharma S., Hermans J., Louppe G., Cranmer K., 2019
- Brewer B. J., Lewis G. F., Belokurov V., Irwin M. J., Bridges T. J., Evans N. W., 2011, *Monthly Notices of the Royal Astronomical Society*, 412, 2521
- Brock A., Donahue J., Simonyan K., 2018, CoRR, abs/1809.11096
- Buades A., Coll B., Morel J.-M., 2005, in Proceedings of the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05) - Volume 2 - Volume 02. CVPR ’05. IEEE Computer Society, Washington, DC, USA, pp 60–65, doi:10.1109/CVPR.2005.38, <http://dx.doi.org/10.1109/CVPR.2005.38>
- Buades A., Coll B., Morel J.-M., 2011, *Image Processing On Line*, 1, 208
- Burroni J., Baudart G., Mandel L., Hirzel M., Shinnar A., 2018, CoRR, abs/1810.00873
- Carleo G., Cirac I., Cranmer K., Daudet L., Schuld M., Tishby N., Vogt-Maranto L., Zdeborova L., 2019
- Carpenter B., et al., 2017, *Journal of Statistical Software, Articles*, 76, 1
- Charnock T., Lavaux G., Wandelt B. D., Boruah S. S., Jasche J., Hudson M. J., 2019
- Chen X., Kingma D. P., Salimans T., Duan Y., Dhariwal P., Schulman J., Sutskever I., Abbeel P., 2016, CoRR, abs/1611.02731
- Dai B., Wipf D., 2019, arXiv e-prints, p. arXiv:1903.05789
- Degrave J., Hermans M., Dambre J., Wyffels F., 2016, CoRR, abs/1611.01652
- Deng J., Dong W., Socher R., Li L., Kai Li Li Fei-Fei 2009, in 2009 IEEE Conference on Computer Vision and Pattern Recognition. pp 248–255, doi:10.1109/CVPR.2009.5206848
- Despali G., Vegetti S., White S. D. M., Giocoli C., van den Bosch F. C., 2018, *Mon. Not. Roy. Astron. Soc.*, 475, 5424
- Diaz Rivero A., Dvorkin C., 2019
- Dilokthanakul N., Mediano P. A. M., Garnelo M., Lee M. C. H., Salimbeni H., Arulkumaran K., Shanahan M., 2016, CoRR, abs/1611.02648
- Dressel L., 2012, Wide Field Camera 3 Instrument Handbook for Cycle 21 v. 5.0
- Drlica-Wagner A., et al., 2019
- Duane S., Kennedy A., Pendleton B. J., Roweth D., 1987, *Physics Letters B*, 195, 216
- Fadely R., Keeton C. R., 2012, *MNRAS*, 419, 936
- Fitts A., et al., 2017, *Mon. Not. Roy. Astron. Soc.*, 471, 3547
- Ge H., Xu K., Ghahramani Z., 2018, in International Conference on Artificial Intelligence and Statistics, AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain. pp 1682–1690, <http://proceedings.mlr.press/v84/ge18b.html>
- Gilman D., Birrer S., Nierenberg A., Treu T., Du X., Benson A., 2019a
- Gilman D., Birrer S., Treu T., Nierenberg A., Benson A., 2019b, *Mon. Not. Roy. Astron. Soc.*, 487, 5721
- Goodfellow I. J., Pouget-Abadie J., Mirza M., Xu B., Warde-Farley D., Ozair S., Courville A., Bengio Y., 2014, arXiv e-prints, p. arXiv:1406.2661
- Goyal P., Hu Z., Liang X., Wang C., Xing E. P., 2017, CoRR, abs/1703.07027
- Güneş Baydin A., et al., 2019, arXiv e-prints, p. arXiv:1907.03382
- Hahnloser R. H. R., Sarpeshkar R., Mahowald M. A., Douglas R. J., Seung H. S., 2000, *Nature*, 405, 947
- Hezaveh Y. D., et al., 2016, *Astrophys. J.*, 823, 37
- Hezaveh Y. D., Perreault Levasseur L., Marshall P. J., 2017, *Nature*, 548, 555
- Higgins I., Matthey L., Pal A., Burgess C., Glorot X., Botvinick M., Mohamed S., Lerchner A., 2017, in ICLR.
- Hoffman M., Blei D. M., Wang C., Paisley J., 2012, arXiv e-prints, p. arXiv:1206.7051
- Innes M., Edelman A., Fischer K., Rackauckas C., Saba E., Shah V. B., Tebbutt W., 2019, CoRR, abs/1907.07587
- Ioffe S., Szegedy C., 2015, arXiv e-prints, p. arXiv:1502.03167
- Jimenez Rezende D., Mohamed S., 2015, arXiv e-prints, p. arXiv:1505.05770
- Jimenez Rezende D., Mohamed S., Wierstra D., 2014, arXiv e-prints, p. arXiv:1401.4082
- Kahlhoefer F., Kaplinghat M., Slatyer T. R., Wu C.-L., 2019
- Karras T., Laine S., Aila T., 2018, CoRR, abs/1812.04948
- Kassiola A., Kovner I., 1993, *ApJ*, 417, 450
- Keeton C. R., Kochanek C. S., 1998, *ApJ*, 495, 157
- Kingma D. P., Ba J., 2015, in 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings. <http://arxiv.org/abs/1412.6980>
- Kingma D. P., Dhariwal P., 2018, in Bengio S., Wallach H., Larochelle H., Grauman K., Cesa-Bianchi N., Garnett R., eds, , Advances in Neural Information Processing Systems 31. Curran Associates, Inc., pp 10215–10224, <http://papers.nips.cc/paper/8224-glow-generative-flow-with-invertible-1x1-convolutions.pdf>
- Kingma D. P., Welling M., 2013, arXiv e-prints, p. arXiv:1312.6114
- Kingma D. P., Salimans T., Jozefowicz R., Chen X., Sutskever I.,

- Welling M., 2016, arXiv e-prints, p. [arXiv:1606.04934](https://arxiv.org/abs/1606.04934)
- Koekemoer A. M., et al., 2007, *Astrophys. J. Suppl.*, 172, 196
- Koopmans L. V. E., 2005, *Mon. Not. Roy. Astron. Soc.*, 363, 1136
- Kormann R., Schneider P., Bartelmann M., 1994, *Astronomy and Astrophysics*, 284, 285
- Krizhevsky A., Sutskever I., Hinton G. E., 2012, in Pereira F., Burges C. J. C., Bottou L., Weinberger K. Q., eds., *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc., pp 1097–1105, <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks>
- Kuhlen M., Vogelsberger M., Angulo R., 2012, *Phys. Dark Univ.*, 1, 50
- LSST Science Collaboration Abell P. A., et al., 2009, arXiv e-prints, p. [arXiv:0912.0201](https://arxiv.org/abs/0912.0201)
- Larkoski A. J., Moulton I., Nachman B., 2017
- LeCun Y., Boser B., Denker J. S., Henderson D., Howard R. E., Hubbard W., Jackel L. D., 1989, *Neural Computation*, 1, 541
- Li T.-M., Gharbi M., Adams A., Durand F., Ragan-Kelley J., 2018a, *ACM Trans. Graph. (Proc. SIGGRAPH)*, 37, 139:1
- Li T.-M., Aittala M., Durand F., Lehtinen J., 2018b, *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 37, 222:1
- Lovell M. R., Frenk C. S., Eke V. R., Jenkins A., Gao L., Theuns T., 2014, *Mon. Not. Roy. Astron. Soc.*, 439, 300
- Lunn D. J., Thomas A., Best N., Spiegelhalter D., 2000, *Statistics and Computing*, 10, 325
- Lunn D., Spiegelhalter D., Thomas A., Best N., 2009, *Statistics in Medicine*, 28, 3049
- Maas A. L., Hannun A. Y., Ng A. Y., 2013, in *ICML Workshop on Deep Learning for Audio, Speech and Language Processing*.
- Mandelbaum R., et al., 2014, *Astrophys. J. Suppl.*, 212, 5
- Merten J., 2016, *Mon. Not. Roy. Astron. Soc.*, 461, 2328
- Merten J., 2017, *SaWLens2 Wiki*, <https://bitbucket.org/jmerten82/libsa2/wiki/sawlens2>
- Morningstar W. R., et al., 2019
- Mustafa M., Bard D., Bhimi W., Lukić Z., Al-Rfou R., Kratochvil J. M., 2019, *Computational Astrophysics and Cosmology*, 6, 1
- Navarro J. F., Frenk C. S., White S. D. M., 1996, *Astrophys. J.*, 462, 563
- Neal R. M., 2012, arXiv e-prints, p. [arXiv:1206.1901](https://arxiv.org/abs/1206.1901)
- Nightingale J., Dye S., Massey R., 2018, *Mon. Not. Roy. Astron. Soc.*, 478, 4738
- Nightingale J., Hayes R., Kelly A., Etherington A., Can X., He Q., Li N., 2019, *Jammy2211/PyAutoLens: PyAutoLens: Automated Strong Gravitational Lens Modeling*, <https://github.com/Jammy2211/PyAutoLens>
- O’Riordan C. M., Warren S. J., Mortlock D. J., 2019, *MNRAS*, 487, 5143
- Otten S., et al., 2019
- Papamakarios G., Pavlakou T., Murray I., 2017, arXiv e-prints, p. [arXiv:1705.07057](https://arxiv.org/abs/1705.07057)
- Paszke A., et al., 2017.
- Perreault Levasseur L., Hezaveh Y. D., Wechsler R. H., 2017, *Astrophys. J.*, 850, L7
- Radford A., Metz L., Chintala S., 2015, arXiv e-prints, p. [arXiv:1511.06434](https://arxiv.org/abs/1511.06434)
- Ravanbakhsh S., Lanusse F., Mandelbaum R., Schneider J., Poczoz B., 2016
- Read J. I., Iorio G., Agertz O., Fraternali F., 2017, *Mon. Not. Roy. Astron. Soc.*, 467, 2019
- Refregier A., Amara A., Kitching T. D., Rassat A., Scaramella R., Weller J., 2010, arXiv e-prints, p. [arXiv:1001.0061](https://arxiv.org/abs/1001.0061)
- Reiman D. M., Göhre B. E., 2019, *Monthly Notices of the Royal Astronomical Society*, 485, 2617–2627
- Revels J., Lubin M., Papamarkou T., 2016, arXiv:1607.07892 [cs.MS]
- Rezende D. J., Viola F., 2018, *Taming VAEs* ([arXiv:1810.00597](https://arxiv.org/abs/1810.00597))
- Ritondale E., Vegetti S., Despali G., Auger M. W., Koopmans L. V. E., McKean J. P., 2019, *Mon. Not. Roy. Astron. Soc.*, 485, 2179
- Schramm T., 1990, *Astronomy and Astrophysics*, 231, 19
- Scoville N., et al., 2007a, *Astrophys. J. Suppl.*, 172, 1
- Scoville N., et al., 2007b, *Astrophys. J. Suppl.*, 172, 38
- Sérsic J. L., 1963, *Boletín de la Asociación Argentina de Astronomía La Plata Argentina*, 6, 41
- Simonovic B., et al., 2019
- Suyu S. H., Marshall P. J., Hobson M. P., Blandford R. D., 2006, *Mon. Not. Roy. Astron. Soc.*, 371, 983
- Suyu S. H., Marshall P. J., Blandford R. D., Fassnacht C. D., Koopmans L. V. E., McKean J. P., Treu T., 2009, *Astrophys. J.*, 691, 277
- Team G. C., 2019, *GREAT3 | The third GRavitational lEnsing Accuracy Testing challenge*, <http://great3challenge.info/>
- Tomczak J. M., Welling M., 2017, *CoRR*, abs/1705.07120
- Tran D., Kucukelbir A., Dieng A. B., Rudolph M., Liang D., Blei D. M., 2016, arXiv preprint [arXiv:1610.09787](https://arxiv.org/abs/1610.09787)
- Treu T., 2010, *Ann. Rev. Astron. Astrophys.*, 48, 87
- V. Sitzmann S. Diamond Y. P. X. D. S. B. W. H. F. H. G. W., 2018, *ACM Trans. Graph. (SIGGRAPH)*
- Vegetti S., Koopmans L. V. E., 2009, *Mon. Not. Roy. Astron. Soc.*, 392, 945
- Vegetti S., Koopmans L. V. E., Bolton A., Treu T., Gavazzi R., 2010, *MNRAS*, 408, 1969
- Vegetti S., Lagattuta D. J., McKean J. P., Auger M. W., Fassnacht C. D., Koopmans L. V. E., 2012, *Nature*, 481, 341
- Vegetti S., Despali G., Lovell M. R., Enzi W., 2018, *Mon. Not. Roy. Astron. Soc.*, 481, 3661
- Verma A., Collett T., Smith G. P., *Strong Lensing Science Collaboration the DESC Strong Lensing Science Working Group* 2019, arXiv e-prints, p. [arXiv:1902.05141](https://arxiv.org/abs/1902.05141)
- Vogelsberger M., Zavala J., Schutz K., Slatyer T. R., 2018,] 10.1093/mnras/stz340
- Warren S., Dye S., 2003, *Astrophys. J.*, 590, 673
- Zhao S., Song J., Ermon S., 2017, *CoRR*, abs/1702.08658
- de Avila Belbute-Peres F., Smith K., Allen K., Tenenbaum J., Kolter J. Z., 2018, in Bengio S., Wallach H., Larochelle H., Grauman K., Cesa-Bianchi N., Garnett R., eds., *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., pp 7178–7189, <http://papers.nips.cc/paper/7948-end-to-end-differentiable-physics-for-learning-and-control.pdf>
- van den Berg R., Hasenclever L., Tomczak J. M., Welling M., 2018, arXiv e-prints, p. [arXiv:1803.05649](https://arxiv.org/abs/1803.05649)

APPENDIX A: THE MAIN HALO MODEL

In this appendix, we describe how the contribution to the displacement field due to the main halo, the Singular Power-Law Ellipsoid (SPLE) profile, is implemented in the Lens Model block of the lensing pipeline shown in Fig. 5. It is worth observing that, in the absence of an analytical expression for α_{mh} , one has to numerically compute the integral in eq. (7). However, this is not feasible in our framework because the numerical integration is not coded up in an autodifferentiable way. For this reason, the displacement field is instead determined by means of an interpolation of a precomputed numerical table of the corresponding integral, eq. (7).

In case of a surface mass density profile with elliptical contours (like for example the SPLE profile), the two-dimensional integral of eq. (7) can be reduce to a simpler

one-dimensional integral (Schramm 1990; Barkana 1998)

$$\alpha_1(\theta_1, \theta_2) = 2\theta_1 q \int_0^{\rho(\theta_1, \theta_2)} \frac{\rho' \kappa(\rho') \omega}{\theta_1^2 + \omega^4 \theta_2^2} d\rho', \quad (\text{A1})$$

$$\alpha_2(\theta_1, \theta_2) = 2\theta_2 q \int_0^{\rho(\theta_1, \theta_2)} \frac{\rho' \kappa(\rho') \omega^3}{\theta_1^2 + \omega^4 \theta_2^2} d\rho', \quad (\text{A2})$$

with

$$\omega^2 = \frac{\Delta + r^2 + \rho'^2(1 - q^2)}{\Delta + r^2 - \rho'^2(1 - q^2)}, \quad (\text{A3})$$

$$\Delta^2 = \left[\rho'^2(1 - q^2) + \theta_2^2 - \theta_1^2 \right]^2 + 4\theta_1^2 \theta_2^2. \quad (\text{A4})$$

where $r^2 = \theta_1^2 + \theta_2^2$, $\rho(\theta_1, \theta_2)$ is defined in eq. (13), and the quantity κ is the surface mass density Σ in units of the critical density Σ_{cr} . For a given profile, the integrals in eqs. (A1) and (A2) can be tabulated for different values of a subset of the lens parameters Θ_{lens} . In particular, in case of the SPLE lens, the displacement field has been evaluated on a unit circle ($r = 1$), for which the coordinates are $\theta_1 = \cos \eta$ and $\theta_2 = \sin \eta$ with the angle η being in the first quadrant ($0 \leq \eta \leq \pi/2$), for different values of the axis ratio q in the interval $[0, 1]$ and the slope γ . The Einstein radius is instead fixed to be $r_{\text{Ein}} = 1$. Such a procedure provided a three-dimensional table in the variables $\{\eta, q, \gamma\}$. This numerical table is then interpolated to compute the SPLE displacement field for any values of q and γ , and at any position (θ_1, θ_2) . Each component of the displacement field is indeed given by

$$\alpha_i(\theta_1, \theta_2) = \alpha_i(|\theta_1|, |\theta_2|) \text{sign}(\theta_i), \quad (\text{A5})$$

$$\alpha_i(\theta_1, \theta_2)|_r = r^{2-\gamma} \alpha_i(\theta_1, \theta_2)|_{r=1}. \quad (\text{A6})$$

Moreover, in case of the SPLE profile, these two components show a simple scaling relation as a function of the Einstein radius. We have

$$\alpha_i(\theta_1, \theta_2)|_{r_{\text{Ein}}} = \left(\frac{1}{r_{\text{Ein}}} \right)^{1-\gamma} \alpha_i(\theta_1, \theta_2)|_{r_{\text{Ein}}=1} \quad (\text{A7})$$

It is worth noticing that the accuracy in computing the displacement field by means of this procedure depends on the size of the interpolation table, which can be defined without any constraint.

APPENDIX B: ELBO DERIVATION

This appendix demonstrates one possible derivation of eq. (4). Consider a latent variable model defined by the joint probability distribution $p(x, z)$, where x and z are the observed and latent variables, respectively. We start by rewriting an expression for the log of the evidence:

$$\log p(x) = \log \int_z p(x, z) dz \quad (\text{B1})$$

$$= \log \int_z p(z|x) \frac{p(x, z)}{p(z|x)} dz, \quad (\text{B2})$$

This can be recognized as an expectation value over $p(z|x)$:

$$\log p(x) = \log \mathbb{E}_{p(z|x)} \left[\frac{p(x, z)}{p(z|x)} \right]. \quad (\text{B3})$$

By Jensen's inequality, which relates the expectation value of a convex function to that function applied to an expectation

value, we have

$$\log \mathbb{E}_{p(z|x)} \left[\frac{p(x, z)}{p(z|x)} \right] \geq \mathbb{E}_{p(z|x)} \left[\log \left(\frac{p(x, z)}{p(z|x)} \right) \right], \quad (\text{B4})$$

and thus:

$$\log p(x) \geq \mathbb{E}_{p(z|x)} \left[\log \left(\frac{p(x, z)}{p(z|x)} \right) \right]. \quad (\text{B5})$$

Using the definition of the Kullback-Leibler divergence for continuous random variables

$$D_{KL}[p||q] = -\mathbb{E}_{p(x)} \left[\log \frac{q(x)}{p(x)} \right], \quad (\text{B6})$$

this can be manipulated to yield

$$\log p(x) \geq \mathbb{E}_{p(z|x)} \left[\log \frac{p(x|z)p(z)}{p(z|x)} \right] \quad (\text{B7})$$

$$= \mathbb{E}_{p(z|x)} \left[\log p(x|z) + \log \frac{p(z)}{p(z|x)} \right] \quad (\text{B8})$$

$$= \mathbb{E}_{p(z|x)} [\log p(x|z)] - D_{KL}[p(z|x)||p(z)] \quad (\text{B9})$$

$$\equiv \text{ELBO}(\theta, \phi; x). \quad (\text{B10})$$

The VAE training objective is obtained by substituting the approximate distributions $p(x|z) \rightarrow d_\theta(x|z)$ and $p(z|x) \rightarrow e_\phi(z|x)$ for the true ones.

APPENDIX C: OPTIMIZING THE ELBO

Training the variational autoencoder requires taking the gradient of the ELBO for (batches of) training images $\{x^{(i)}\}_{i=1}^N$ with respect to the encoder and decoder's parameters θ and ϕ :

$$\begin{aligned} \nabla_{\theta, \phi} \text{ELBO}(\theta, \phi; x^{(i)}) \\ = \nabla_{\theta, \phi} \mathbb{E}_{e_\phi(z|x^{(i)})} \left[\log d_\theta(x^{(i)}|z) \right] \\ - \nabla_{\theta, \phi} D_{KL} \left[e_\phi(z|x^{(i)}) || p(z) \right]. \end{aligned} \quad (\text{C1})$$

For the normal encoding distribution and latent space prior adopted in this work, the KL divergence term can be integrated analytically, which makes it simple to compute the second term on above. The first term is more challenging: while a Monte Carlo estimate of the derivative with respect to θ can be performed by sampling $\{z^{(j)} \sim e_\phi(z|x^{(i)})\}_{j=1}^M$, it is not obvious how to compute the derivatives of the sampled latent variable values with respect to ϕ .

The solution is the reparameterization trick introduced in the two original papers on variational autoencoders (Kingma & Welling 2013; Jimenez Rezende et al. 2014). The insight is that (assuming a normal encoding distribution) the randomness and ϕ -dependent parts of the sampling process can be factored, allowing the sampled values to be written as

$$z^{(j)} = \mu_e(x^{(i)}; \phi) + \epsilon^{(j)} \sigma_e(x^{(i)}; \phi), \quad (\text{C2})$$

with $\epsilon^{(j)} \sim \mathcal{N}(0, I)$. These can be used to construct the following Monte Carlo gradient estimator by treating the $\epsilon^{(j)}$

Conv2d(1, 64, 4, 2, 1)	
LeakyReLU(0.2)	
Conv2d(64, 128, 4, 2, 1)	
BatchNorm2d(128)	
LeakyReLU(0.2)	
Conv2d(128, 256, 4, 2, 1)	
BatchNorm2d(256)	
LeakyReLU(0.2)	
Conv2d(256, 512, 4, 2, 1)	
BatchNorm2d(512)	
LeakyReLU(0.2)	
Conv2d(512, 4096, 4, 1, 0)	
LeakyReLU(0.2)	
Linear(4096, 64)	Linear(4096, 64)
	Exp
$\mu_e(x)$	$\sigma_e(x)$

Table D1. Encoder neural network architecture. The notation uses the same conventions as pytorch. The arguments of Conv2d indicate the number of input channels, number of output channels, kernel size, stride and zero padding; all convolutions are unbiased. The LeakyReLU argument is slope for inputs less than 0. The BatchNorm2d argument is the number of input channels. The output of the last convolutional block is flattened before being passed to the two separate Linear layers to produce the mean and standard deviation of the encoding distribution. Linear layers' arguments show the number of input and output channels.

values as constants for the optimization epoch:

$$\begin{aligned} \nabla_{\theta, \phi} \mathbb{E}_{e_\phi(z|x^{(i)})} \left[\log d_\theta \left(x^{(i)} | z \right) \right] \\ \approx \frac{1}{M} \sum_{j=1}^M \log d_\theta \left(x^{(i)} | z^{(j)} \right). \end{aligned} \quad (\text{C3})$$

This estimator is stable; we set $M = 10$ in our work by using batches of 32 training images.

APPENDIX D: VARIATIONAL AUTOENCODER ARCHITECTURE

The architectural details of the encoder and decoder networks of our variational autoencoder are presented in Tables D1 and D2, respectively. All weights were initialized to 0.02. The biases in the final linear layers of the encoder were initialized to 0.

We tried several experiments to see whether we could improve upon this VAE design. For example, since the decoder can have trouble saturating the final tanh activation, we tried exchanging this for a LeakyReLU, as well as removing it altogether. We also tested 32 and 128 latent-space dimensions. The former lead to blurry reconstructions while the later yielded little improvement relative to 64 latent-space dimension. Dai & Wipf (2019) analytically demonstrated that making the hyperparameter σ_d a trainable parameter should lead to sharper reconstructions. We did not find this to be the case, and instead found that σ_d converged to roughly the value we selected by hand based on the signal-to-noise ratio of the training data. Finally, we experimented with a residual network-based architecture, as was used in Dai & Wipf (2019), which did not yield any improvement relative to the architecture we eventually selected.

ConvTranspose2d(64, 512, 4, 1, 0)	
BatchNorm2d(512)	
ReLU	
ConvTranspose2d(512, 256, 4, 2, 1)	
BatchNorm2d(256)	
ReLU	
ConvTranspose2d(256, 128, 4, 2, 1)	
BatchNorm2d(128)	
ReLU	
ConvTranspose2d(128, 64, 4, 2, 1)	
BatchNorm2d(64)	
ReLU	
ConvTranspose2d(64, 1, 4, 2, 1)	
Tanh	
$\mu_d(z)$	

Table D2. Decoder neural network architecture. The notation is described in the caption of Table D1, and is the same for Conv2d and ConvTranspose2d. The input vector \mathbf{z} is reshaped to have 64 channels and spatial dimensions equal to 1 along both axes.

This paper has been typeset from a $\text{\TeX}/\text{\LaTeX}$ file prepared by the author.