# Pixelated Reconstruction of Foreground Density and Background Surface Brightness in Gravitational Lensing Systems using Recurrent Inference Machines

Alexandre Adam,[1,2,3] Laurence Perreault-Levasseur,[1,2,3,4] Yashar Hezaveh,[1,2,3,4] and Max Welling[5]

[1]*Department of Physics, Université de Montréal, Montréal, Canada*
[2]*Mila - Quebec Artificial Intelligence Institute, Montréal, Canada*
[3]*Ciela - Montreal Institute for Astrophysical Data Analysis and Machine Learning, Montréal, Canada*
[4]*Center for Computational Astrophysics, Flatiron Institute, 162 5th Avenue, 10010, New York, NY, USA*
[5]*Microsoft Research AI4Science*

## ABSTRACT

Modeling strong gravitational lenses in order to quantify the distortions in the images of background sources and to reconstruct the mass density in the foreground lenses has been a difficult computational challenge. As the quality of gravitational lens images increases, the task of fully exploiting the information they contain becomes computationally and algorithmically more difficult. In this work, we use a neural network based on the Recurrent Inference Machine (RIM) to simultaneously reconstruct an undistorted image of the background source and the lens mass density distribution as pixelated maps. The method iteratively reconstructs the model parameters (the image of the source and a pixelated density map) by learning the process of optimizing the likelihood given the data using the physical model (a ray-tracing simulation), regularized by a prior implicitly learned by the neural network through its training data. When compared to more traditional parametric models, the proposed method is significantly more expressive and can reconstruct complex mass distributions, which we demonstrate by using realistic lensing galaxies taken from the IllustrisTNG cosmological hydrodynamic simulation.

## 1. INTRODUCTION

Strong gravitational lensing is a natural phenomenon through which multiple, distorted images of luminous background sources are formed by the gravity of massive foreground objects along the line of sight (e.g., Vieira et al. 2013; Marrone et al. 2018; Rizzo et al. 2020; Sun et al. 2021). These distortions are tracers of the distribution of mass in foreground structures, irrespective of their light emission properties. As such, this phenomenon offers a powerful probe of the distribution of dark matter (e.g., Dalal & Kochanek 2002; Treu & Koopmans 2004; Hezaveh et al. 2016; Gilman et al. 2020, 2021).

Lens modeling is the process through which the parameters describing both the mass distribution in the foreground lens and the undistorted image of the background source are inferred. This has traditionally been done through explicit likelihood-based modeling methods, a time- and resource-consuming procedure. A common practice to model strong lenses is to model the light profile of the background source with a Sérsic (1963)

profile and the density of the foreground lens with a power law function, $\rho \propto r^{-\gamma}$. These simple profiles allow for the exploration of their low-dimensional parameter space with non-linear samplers such as Markov Chain Monte Carlo (MCMC) methods (e.g., Koopmans et al. 2006; Barnabè et al. 2009; Auger et al. 2010) and generally provide a good fit to low-resolution data. However, as high-resolution and high signal-to-noise ratio (SNR) images become available, lensing analysis with simple models requires the introduction of additional parameters representing the true complexity of the mass distribution in lensing galaxies and the complexity of surface brightness in the background sources (e.g., Sluse et al. 2017; Wong et al. 2017; Birrer et al. 2019; Rusu et al. 2020, 2017; Li et al. 2021). This approach becomes intractable as the complexity of the mass distribution and the quality of images increases (e.g., Schmidt et al. 2022). For example, no simple parametric model of the *Hubble Space Telescope* (*HST*) images of the Cosmic Horseshoe (J1148+1930) — initially discovered by Belokurov et al. (2007) — has been able to model the

fine features of the extended arc (e.g., Bellagamba et al. 2016; Cheng et al. 2019; Schuldt et al. 2019).

Free-form methods attempt to relax the assumptions about the smoothness and symmetries of these parametric profiles using more expressive families like regular (or adaptive) grid representations and meshfree representations (Saha & Williams 1997; Abdelsalam et al. 1998a,b; Diego et al. 2005; Birrer et al. 2015; Merten 2016; Galan et al. 2022). These methods strive to model the signal contained in lensed images in a data-agnostic way, in order to place better constraints on the morphology of the source brightness or the projected mass density of the lens. However, most free-form parametrization choices make the inference problem under-constrained, meaning that imposing a prior on the reconstructed parameters becomes essential to penalize unphysical solutions and avoid overfitting the data.

In the context of traditional likelihood-based modeling, there exists a number of commonly used priors for the inference of high dimensional representations of background sources. For example, the strategy to impose a quadratic-log prior for linear inversion of pixelated-source models was developed by Warren & Dye (2003); Suyu et al. (2006). Other methods include iteratively specified priors for shapelets (Birrer et al. 2015; Birrer & Amara 2018; Nightingale et al. 2018) and a sparsity prior for wavelets representations (Galan et al. 2021).

On the other hand, for lens mass reconstruction, the issue of specifying an appropriate prior is still unsolved. This subject has been studied extensively in the context of cluster-scale strong lensing (Bartelmann et al. 1996; Seitz et al. 1998; Abdelsalam et al. 1998a,b; Bradač et al. 2005; Diego et al. 2005; Cacciato et al. 2006; Diego et al. 2007; Liesenborgs et al. 2006, 2007; Jee et al. 2007; Coe et al. 2008; Merten et al. 2009; Deb et al. 2012; Merten 2016; Ghosh et al. 2020; Torres-Ballesteros & Castañeda 2022). Free-form approaches in the context of strong galaxy-galaxy lenses have been comparatively less studied (see however Saha & Williams (1997, 2004); Birrer et al. (2015); Coles et al. (2014); Galan et al. (2022)).

Another major challenge for these models is the issue of optimizing or sampling these high dimensional posteriors. Given the non-linear nature of the model and the existence of multiple local optima, non-linear global optimizers and samplers are needed, which often results in extremely expensive computational procedures. The high computational cost of these methods also limits the extent to which they can be thoroughly tested and validated to identify and characterize potential systematics.

Over the recent years, deep learning methods have proven extremely successful at accurate modeling of strong lensing systems (Hezaveh et al. 2017; Perreault Levasseur et al. 2017; Morningstar et al. 2018; Coogan et al. 2020; Park et al. 2021; Legin et al. 2021, 2022; Wagner-Carena et al. 2021; Schuldt et al. 2022; Wagner-Carena et al. 2022; Karchev et al. 2022; Anau Montel et al. 2022; Mishra-Sharma & Yang 2022; Schuldt et al. 2022). More specifically, Morningstar et al. (2019) demonstrated that recurrent convolutional neural networks can implicitly learn complex prior distributions from their training data to successfully reconstruct pixelated undistorted images of strongly lensed sources, circumventing the need to explicitly specify a prior distribution over those parameters. Motivated by this, we propose a method that extends this framework to solve the full lensing problem and simultaneously reconstruct a pixelated mass map and a pixelated image of the undistorted background source.

The method we propose here is based on the Recurrent Inference Machine (RIM), originally developed by Putzky & Welling (2017). In its original version, this method proposed to solve inverse problems using a Recurrent Neural Network as a metalearner to learn the iterative process of the optimization of a likelihood. RIMs have been trained on a range of linear inverse problems both within and outside of astrophysics (Lønning et al. 2019). In Modi et al. (2021), this method was generalized to non-linear inference problems while using a U-net architecture (Ronneberger et al. 2015) to separate the dynamics of different scales.

In the present paper, we leverage this framework to learn an optimization process over the highly non-convex strong lensing likelihood, and implicitly learn a data-driven prior, which allows for the reconstruction of complex mass distributions representative of realistic galaxies taken from the IllustrisTNG (Nelson et al. 2019) hydrodynamical simulations. We also introduce a fine-tuning procedure, which allows us to directly exploit the prior encoded in the neural network parameters in order to further optimize the posterior down to noise levels. We apply this to the reconstruction of high signal-to-noise galaxy-galaxy lensing systems simulated using IllustrisTNG (Nelson et al. 2019) projected density maps and background galaxy images collected from the COSMOS survey (Koekemoer et al. 2007; Scoville et al. 2007).

The paper is organised as follows. Section 2 details the inference pipeline. In Section 3, we present the data production and preprocessing for the training of the RIM and the generative models used in this paper. In Section 4, we report on the training strategies used. In Section 5, we discuss our results on a test set of gravitational lenses. We conclude in Section 6.

## 2. METHODS

Our goal is to predict pixelated maps of both the undistorted image of the background source and the projected density in the foreground lens from noisy lensed images. Our model consists of a Recurrent Inference Machine that predicts these variables of interest. Training this model requires a large number of training data, which we produce using a Variational Autoencoder (VAE) trained on density maps from the IllustrisTNG simulation and background sources from the Cosmos dataset (Section 4.1).

In this section, we present the structure of the lensing inference problem and provide information about our analysis method. We begin with a general introduction to maximum a posteriori (MAP) inference in Section 2.1. We describe the lensing simulation pipeline in Section 2.2. In Section 2.3, we motivate the use of the Recurrent Inference Machine and describe its computational graph. The architecture of the neural network is described in Section 2.4. Finally, we describe the fine-tuning procedure and the transfer learning technique applied to achieve noise-level reconstructions in Section 2.5.

### 2.1. Maximum a posteriori inference

We consider the task of reconstructing a vector of parameters of interest $\mathbf{x} \in \mathcal{X}$ given a vector of noisy observed data $\mathbf{y} \in \mathcal{Y}$, a known forward (or physical) model $F$, and an additive noise vector $\boldsymbol{\eta}$. In what follows, we assume this vector to be sampled from a Gaussian distribution with known covariance matrix $C$, such that we can write

$$\begin{aligned} \mathbf{y} &= F(\mathbf{x}) + \boldsymbol{\eta}\,; \\ \boldsymbol{\eta} &\sim \mathcal{N}(0, C)\,. \end{aligned} \tag{1}$$

In our case study, $F$ is a many-to-one non-linear mapping between the parameter space $\mathcal{X}$ and the data space $\mathcal{Y}$. Finding physically allowed solutions for this ill-posed inverse problem requires strong priors. The maximum a posteriori (MAP) solution maximizes the product of the likelihood $p(\mathbf{y} \mid \mathbf{x})$ and the prior $p(\mathbf{x})$:

$$\hat{\mathbf{x}}_{\mathrm{MAP}} = \underset{\mathbf{x} \in \mathcal{X}}{\mathrm{argmax}} \; \log p(\mathbf{y} \mid \mathbf{x}) + \log p(\mathbf{x})\,. \tag{2}$$

Assuming a Gaussian noise model for $\boldsymbol{\eta}$, the log-likelihood can be written analytically as

$$\log p(\mathbf{y} \mid \mathbf{x}) \propto -\big(\mathbf{y} - F(\mathbf{x})\big)^T C^{-1}\big(\mathbf{y} - F(\mathbf{x})\big)\,. \tag{3}$$

The prior distribution, however, is problem-dependent and encodes expert knowledge of the model domain. As such, it is typically harder to write explicitly.

### 2.2. The Forward Model

The forward model, $F$, is a simulation pipeline that receives a map of the surface brightness in the background source and a map of the projected density in the foreground lens to produce distorted images of background galaxies. This pipeline uses ray tracing to calculate the deflection angles, $\boldsymbol{\alpha}$, and maps the observed coordinates, $\boldsymbol{\theta}$, into the coordinates of the background plane, $\boldsymbol{\beta}$, through the lens equation

$$\boldsymbol{\beta} = \boldsymbol{\theta} - \boldsymbol{\alpha}(\boldsymbol{\theta})\,. \tag{4}$$

The deflection angles are obtained using the projected surface density field $\kappa$ — also referred to as convergence — through the integral

$$\boldsymbol{\alpha}(\boldsymbol{\theta}) = \frac{1}{\pi} \int_{\mathbb{R}^2} \kappa(\boldsymbol{\theta}') \frac{\boldsymbol{\theta} - \boldsymbol{\theta}'}{\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|^2} d^2\boldsymbol{\theta}'\,. \tag{5}$$

Since we also use a discrete representation for the convergence, we approximate this integral by a discrete global convolution. Taking advantage of the convolution theorem, this operation can be computed in near-linear time using the Fast Fourier Transform algorithm (FFT).

Assuming the observation has $M^2$ pixels, the convolution kernel would have $(2M + 1)^2$ pixels. Both the convergence map and the kernel are zero-padded to a size of $(4M + 1)^2$ pixels in order to approximate a linear convolution and significantly reduce aliasing.

To produce simulated images, the intensity of an image pixel is obtained through bilinear interpolation of the source brightness distribution at the coordinate $\boldsymbol{\beta}$. A blurring operator — convolution by a point spread function (PSF) — is then applied to the lensed image to replicate the response of an optical system. This operator is implemented as a GPU-accelerated matrix operation since the blurring kernels used in this paper have a significant proportion of their energy distribution encircled inside a small pixel radius. Gaussian noise is then applied to the images, as described in more details in section 3.3.

### 2.3. Recurrent Inference Machine

Instead of handcrafting a prior distribution to solve the inverse problem (1), we build an inference pipeline with a data-driven implicit prior encoded in a deep neural network architecture (Bengio 2009). The RIM (Putzky & Welling 2017) is a form of learnt gradient-based inference algorithm, intended to solve inverse problems of the form (1). This framework has mainly been applied in the context of linear under-constrained inverse problems — i.e. where $F(\mathbf{x})$ can be represented as a matrix product $A\mathbf{x}$ — for which the prior on the

parameters $\mathbf{x}$, $p(\mathbf{x})$, is either intractable or hard to compute (Morningstar et al. 2018, 2019; Lønning et al. 2019). The use of the RIM to solve non-linear inverse problems was first investigated in (Modi et al. 2021). In our case, the function representing the physical model $F$ encodes the lens equation (4), which is highly non-linear.

The RIM is made up of a recurrent unit, which, given an observation $\mathbf{y}$, solves (1) for $\mathbf{x}$ through the governing equation

$$\hat{\mathbf{x}}^{(t+1)} = \hat{\mathbf{x}}^{(t)} + g_\varphi\big(\hat{\mathbf{x}}^{(t)}, \mathbf{y}, \boldsymbol{\nabla}_{\hat{\mathbf{x}}^{(t)}} \log p(\mathbf{y} \mid \hat{\mathbf{x}}^{(t)})\big), \quad (6)$$

where $\hat{\mathbf{x}}^{(t)}$ is the estimate of the parameters of interest at time $t$ of the recursion (here, the pixel values of the image of the undistorted background source and of the density field $\kappa$) and $g_\varphi$ is a neural network. In the text, we will often use the shorthand notation $\boldsymbol{\nabla}_{\mathbf{y}|\mathbf{x}}$ to refer to $\boldsymbol{\nabla}_{\mathbf{x}} \log p(\mathbf{y} \mid \mathbf{x})$, the gradient of the likelihood evaluated at $\mathbf{x}$. By minimizing a weighted mean squared loss backpropagated through time,

$$\mathcal{L}_\varphi(\mathbf{x}, \mathbf{y}) = \frac{1}{T} \sum_{t=1}^{T} \sum_{i=1}^{M} \mathbf{w}_i (\hat{\mathbf{x}}_i^{(t)} - \mathbf{x}_i)^2, \quad (7)$$

where $T$ is the total number of time steps in the recursion, the index $i$ labels the pixels of the reconstructions, $\mathbf{w}_i$ is the per-pixel weight, and $M$ is the total number of pixels in the reconstructions, the neural network $g_\varphi$ learns to optimize the parameters $\mathbf{x}$ given a likelihood function. The converged parameters of the neural network given the training set $\mathcal{D}$, $\varphi_\mathcal{D}^\star$, are those that minimize the cost — or empirical risk — which is defined as the expectation of the loss over $\mathcal{D}$

$$\varphi_\mathcal{D}^\star = \underset{\varphi}{\mathrm{argmin}} \; \mathbb{E}_{(\mathbf{x},\mathbf{y})\sim\mathcal{D}}\big[\mathcal{L}_\varphi(\mathbf{x}, \mathbf{y})\big]. \quad (8)$$

Unlike previous works (Andrychowicz et al. 2016; Putzky & Welling 2017; Morningstar et al. 2018, 2019; Lønning et al. 2019), the data vector $\mathbf{y}$ containing the observations is fed to the neural network in order to learn a better initialization of the parameters, $\mathbf{x}^{(0)} = g_\varphi(0, \mathbf{y}, 0)$, in addition to their optimization process. Empirically, we found that this significantly improves the performance of the model for our problem and avoids situations where the model would get stuck in local minima at test time due to poor initialization.

We follow previous works in setting a uniform weight over the time steps ($\mathbf{w}^{(t)} = \frac{\mathbf{w}}{T}$). The choice of the pixel weights $\mathbf{w}_i$ is informed by our empirical observations when training the network. Details are reported in appendix C.

In Figure 1, we show the rolled computational graph of the RIM. During training of the neural network $g_\varphi$,
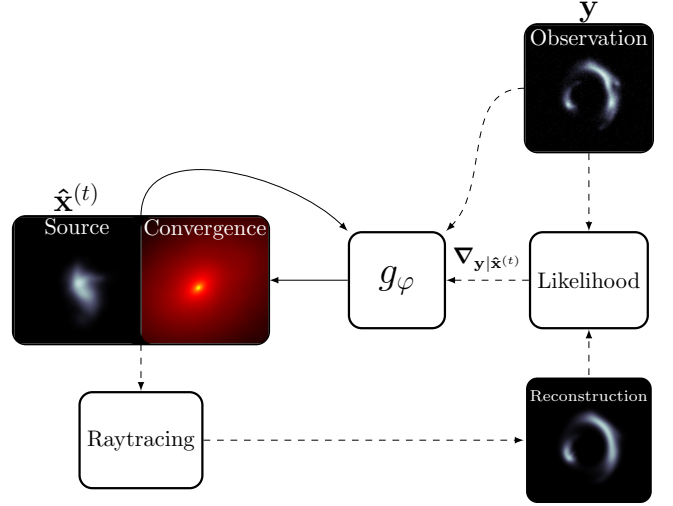


**Figure 1.** Rolled computational graph of the RIM. Dashed arrows represent operations not recorded for BPTT.

operations along the solid arrows are being recorded for backpropagation through time. The recording is stopped along the dashed arrow since these operations are part of the forward modelling process and contain no trainable parameters.

The gradient of the likelihood is computed using automatic differentiation. Following (Modi et al. 2021), we preprocess the gradients using the Adam algorithm (Kingma & Ba 2014). For clarity, we only illustrate this step in Figure 2.

### 2.4. The Neural Network

The neural network architecture is illustrated in Figure 2, which shows a single time step of the unrolled computation graph of the RIM. We use a U-net (Ronneberger et al. 2015) architecture with Gated Recurrent Units (GRU: Cho et al. 2014) placed in each skip connection.

Each GRU cell has its own memory tensor that is updated through time at each iteration of equation 6. The shape of a memory tensor is set to match the feature tensor fed into it from the parent layer in the network graph. Instead of learning a compressed representation like in the hourglass architecture (or autoencoder), the U-net architecture naturally separates the spatial frequency components of the signal into its levels (neural network layers at the same height in Figure 2). The first level generally encodes high frequency features while the lower (or deeper) levels encode low frequency features (due to downsampling of the feature maps). Adding an independent memory unit at each level preserves this property.
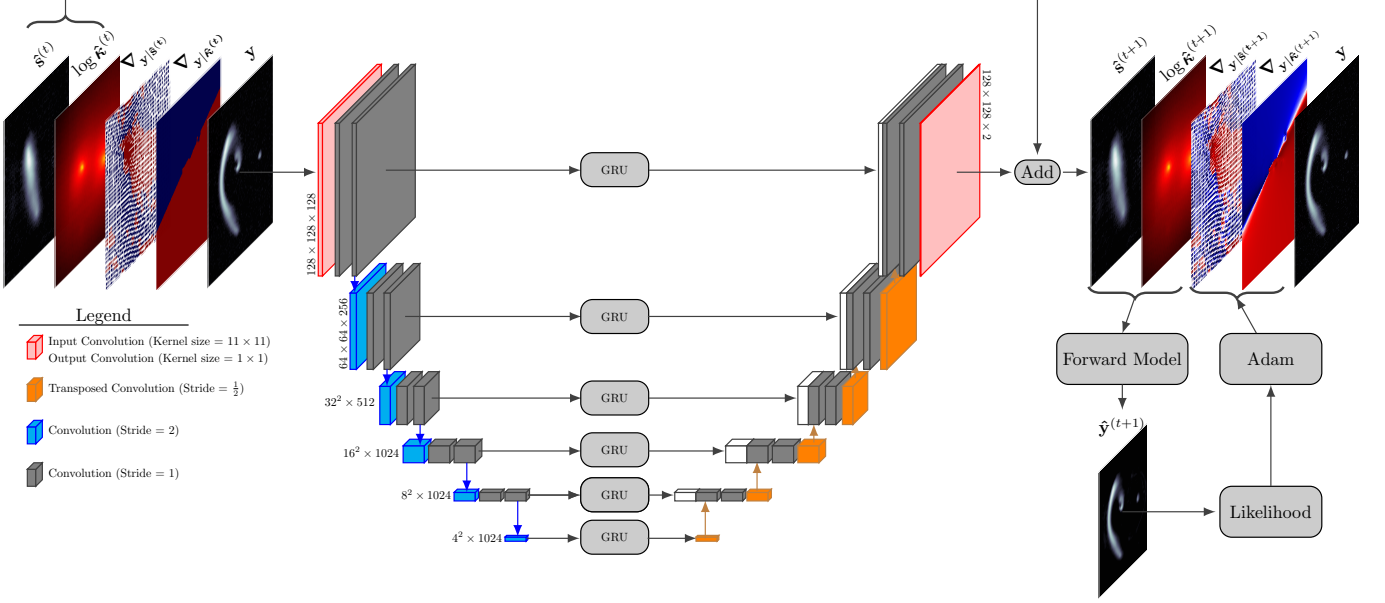
**Figure 2.** A single time step of the unrolled computation graph of the RIM. GRU units are placed in the skip connections to guide the reconstruction of the source and convergence. A schematic of the steps to compute the likelihood gradients is shown in the bottom right of the figure, including the Adam processing step of the likelihood gradient. See appendix C for more details.

Convolutional layers with a stride of 2 are used for downsampling and stride of $\frac{1}{2}$ for upsampling of the feature maps (identified in blue and orange respectively in figure 2). Half-stride convolutions are implemented in practice with the transposed convolution layers from `Tensorflow` (Abadi et al. 2015). Most layers use a kernel size of $3 \times 3$, except the first and last layer. The first layer has larger receptive field ($11 \times 11$) in order to capture more details in the input tensor. The last layer has kernels of size $1 \times 1$. A tanh activation function is used for each convolutional layer, including strided convolutions, except for the output layer. The U-net outputs an image tensor with two channels, one dedicated for the update of the source and the other for the update of the convergence (see figure 2).

### 2.5. Fine-Tuning

#### 2.5.1. Objective function

Once trained, the RIM produces a baseline (point) estimate of the parameters $\mathbf{x}$ given a noisy observation $\mathbf{y}$, a PSF and a noise covariance matrix. We now concern ourselves with a strategy to improve this estimate. This is important for observations with high SNR, for which the estimate must be very accurate to model all the fine features present in the arcs. The metric for the goodness of fit is the reduced chi squared $\chi_\nu^2 = \frac{\chi^2}{\nu}$, where $\nu$ is the total number of degrees of freedom which here corresponds to the total number of pixels in $\mathbf{y}$. Generally, our goal will be to reach $\chi_\nu^2 = 1$, or equivalently $|\chi^2 - \nu| = 0$, which suggests that the RIM's estimate has modeled all

the signal to be recovered from the observations. We note that this metric overestimates the number of degrees of freedom, which cannot be computed reliably for our particular model choice (pixelated source and convergence maps). In practice, the $\chi_\nu^2$ might be biased low.

To improve this figue of merit, we can optimize the log-likelihood directly w.r.t. the network weights given an appropriate prior on those weights (to avoid forgetting the implicit priors that have been learned during training, see section 2.5.2). The new objective function is given by

$$\hat{\varphi}_{\mathrm{MAP}} = \underset{\varphi}{\operatorname{argmax}} \; \frac{1}{T} \sum_{t=1}^{T} \log p(\mathbf{y} \mid \hat{\mathbf{x}}^{(t)}) + \log p(\varphi), \quad (9)$$

where $\varphi$ are the network weights, $\log p(\mathbf{y} \mid \hat{\mathbf{x}}^{(t)})$ is the log-likelihood, and $\log p(\varphi)$ is the log prior over the network weights. Unlike the loss in equation (7), this objective function makes no use of labels ($\mathbf{x}$). This allows us to use equation (9) at test time in order to fine-tune the RIM's weights to a specific test example.

#### 2.5.2. Transfer Learning

We now address the issue of transferring knowledge from the training task defined by the loss function in equation (8), to a test task specific to an observation, as defined by the loss given in equation (9). The reader might refer to reviews on transfer learning (Pan & Yang 2010; Zhuang et al. 2019) for a broad overview of the
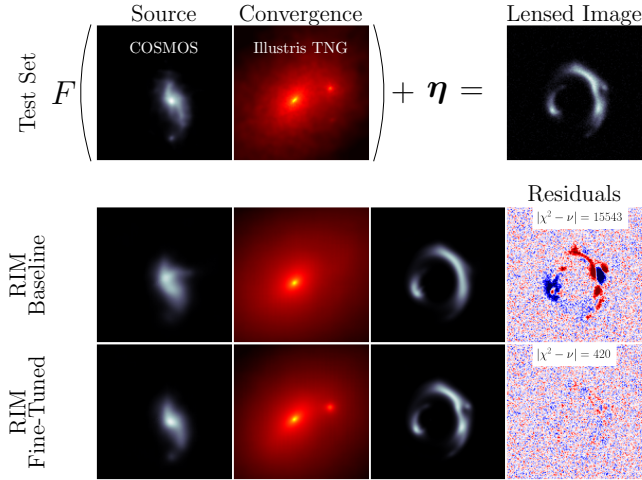
**Figure 3.** Example of a simulated lensed image in the test set that exhibits a large deflection in its eastern arc which indicates the presence of a massive object — in this case a dark matter subhalo. The fine-tuning procedure is able to recover this subhalo because of its strong signal in the lensed image and reduces the residuals to noise level.

field. The strategy we outline falls within the category of inductive transfer learning.

Optimizing the log-likelihood alone without a prior term over the weights (i.e. just the first term from the r.h.s. in (9)) by initializing the weights at $\varphi_{\mathcal{D}}^{\star}$ is not strong enough to preserve the knowledge learned from the training task. This has long been observed in the literature and was coined as the catastrophic interference phenomenon in connectionist networks (McCloskey & Cohen 1989; Ratcliff 1990). In summary, a sequential learning problem exhibits catastrophic forgetting of old knowledge when confronted with new examples (possibly from a different distribution or process), in a manner

1. proportional to the amount of learning;

2. strongly dependent on the disruption of the parameters involved in representing the old knowledge.

While introducing an early stopping condition could potentially alleviate the former issue, the latter could still remain a problem.

We therefore follow the work of Kirkpatrick et al. (2016) to define a prior distribution over $\varphi$ that address this issue

$$\log p(\varphi) \propto -\frac{\lambda}{2} \sum_{j} \mathrm{diag}(\mathcal{I}(\varphi_{\mathcal{D}}^{\star}))_j (\varphi_j - [\varphi_{\mathcal{D}}^{\star}]_j)^2 \, . \quad (10)$$

where $\mathrm{diag}(\mathcal{I}(\varphi_{\mathcal{D}}^{\star}))$ is the diagonal of the Fisher information matrix encoding the amount of information that
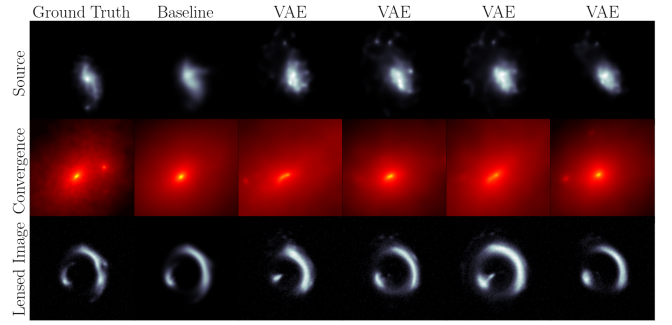


**Figure 4.** Examples similar to the test task, also shown in Figure 3. The first column shows the ground truth used to simulate the lensed image. The second column shows the baseline prediction that is then encoded in the latent space of the VAE in order to sample the next 4 columns.

some set of gravitational lensing systems from the training set, and similar to the observed test task, carries about the baseline RIM weights $\varphi_{\mathcal{D}}^{\star}$ — the parameters that minimize the empirical risk (equation 8). We can also understand this prior using the Cramér-Rao lower bound (Rao 1945; Cramér 1946). The prior can thus be framed as a multivariate Gaussian distribution characterised by a diagonal covariance matrix with $\mathrm{diag}(\mathcal{I})$ as its inverse and by $\varphi_{\mathcal{D}}^{\star}$ as its first moment. Within this view, the Lagrange multiplier is tuning our estimated uncertainty about the neural network weights for the particular task at hand. We have included a derivation of this term in the appendix A.

Examples are drawn from the set of training examples similar to the test task by sampling the latent space of two variational autoencoders (VAE) that model a distribution over the background sources and the convergence maps respectively (as described in Section 3.1 and 3.2) near the baseline prediction of the RIM. In practice, we choose an isotropic Gaussian distribution centered around $\hat{\mathbf{z}}^{(T)}$ — the latent code of the baseline prediction — as a sampling distribution. While we leave the possibility of improving this choice to future work, it is sufficient for our goals. Figure 4 illustrates examples of what is meant here by *similar*.

## 3. DATA

### 3.1. *COSMOS*

The maps of surface brightness of background sources are taken from the *Hubble Space Telescope* (*HST*) Advanced Camera for Surveys Wide Field Channel COSMOS field (Koekemoer et al. 2007; Scoville et al. 2007), a $1.64 \deg^2$ contiguous survey acquired in the F814W filter. A dataset of magnitude limited (F814W $< 23.5$) deblended galaxy postage stamps (Leauthaud et al. 2007) was compiled as part of the GREAT3 challenge (Mandelbaum et al. 2014). The data is publicly available
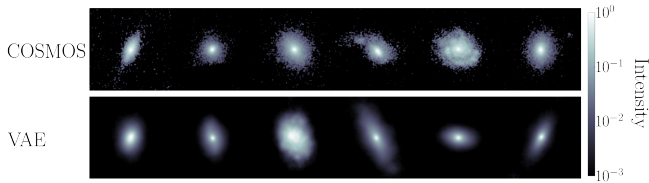
**Figure 5.** Examples of COSMOS galaxy images (top row) and VAE generated samples (bottom row) used as labels in $\mathcal{D}$.



**Figure 6.** Examples of smoothed Illustris TNG100 convergence map (top row) and VAE generated samples (bottom row) used as labels in $\mathcal{D}$.

(Mandelbaum et al. 2012), and the preprocessing is done through the open-source software `GALSIM` (Rowe et al. 2015).

We apply the `marginal` selection criteria (see the `COSMOSCatalog` class) and impose a flux per image greater than 50 photons cm$^{-2}$ s$^{-1}$. This final set has a total of 13 321 individual images. Each image is saved as a postage stamp of $158^2$ pixels. We then subtract the background from each image, apply a random shift, rotate them by an angle multiple of 90°, crop them down to $128^2$ pixels, and finally normalize them to pixel intensities in the range $[0, 1]$. We then train an autoencoder to denoise the galaxy images (Vincent et al. 2008, 2010). More specifically, we use the informational bottleneck principle (Tishby et al. 2000) to learn a lossy lower-dimensional representation of the data. For a generic CNN autoencoder, this amounts to learning a low-pass frequency filter on the COSMOS dataset. Indeed, CNNs are known to exhibit a spectral bias in their learning phase (Rahaman et al. 2018), which we exploit to our advantage in order to filter pixel noise from the galaxy surface brightness. Furthermore, using an expressive CNN autoencoder produces much fewer artifacts than a naive implementation of such a low-pass filter — e.g. by masking Fourier modes.

We split the galaxies into a training set (90%) and a test set (10%). The augmented training set ($\sim$ 50 000 images) is used to train a VAE, as described in Section 4.1, and produce simulated observations to train the RIM.

## 3.2. *IllustrisTNG*

### 3.2.1. *Smooth Particle Lensing*

To compute convergence maps from an N-body simulation, we use Kernel Density Estimation to produce smooth densities on a regular grid from discrete simulation particles. This reduces the particle noise affecting all important lensing quantities. At the same time, the choice of the kernel size is important to preserve substructures in the lens that we might potentially be interested in. Following Aubert et al. (2007); Rau et al. (2013), we use Gaussian smoothing with an adaptive kernel size determined by the distance of the 64$^{\text{th}}$ near-
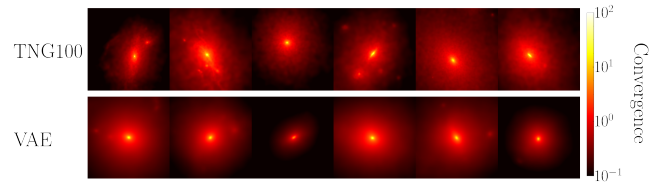
est neighbours, $D_{64,i}$, of a given particle of mass $m_i$ and projected position $\mathbf{x}_i$

$$\kappa(\mathbf{x}) = \frac{1}{\Sigma_{\text{crit}}} \sum_{i=1}^{N_{\text{part}}} \frac{m_i}{2\pi \hat{\ell}_i^2} \exp\left(-\frac{1}{2} \frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{\hat{\ell}_i^2}\right)$$
$$\hat{\ell}_i = \sqrt{\frac{103}{1024}} D_{64,i}. \tag{11}$$

The nearest neighbours are found by fitting a k-d tree — implemented in `scikit-learn` (Pedregosa et al. 2011) — to the $N_{\text{part}}$ particles in a cylinder centered on the centre of mass of the halo of interest. The critical surface density is defined as

$$\Sigma_{\text{crit}} = \frac{c^2 D_s}{4\pi G D_{\ell s} D_\ell} \tag{12}$$

where $D_\ell$, $D_s$ and $D_{\ell s}$ are angular diameter distances to the lens, source and between the lens and the source respectively, $G$ is the gravitational constant, and $c$ the speed of light.

### 3.2.2. *Preprocessing*

The projected surface density maps (convergence) of lensing galaxies are made using the redshift $z = 0$ snapshot of the IllustrisTNG-100 simulation (Nelson et al. 2019) in order to produce physically realistic realizations of density maps containing dark and baryonic matter. We select 1604 halos with the criteria that they have a total dark matter mass of at least $9 \times 10^{11} M_\odot$. We then collect all dark matter, gas, stars and black holes particles from the data in the vicinity of the halo to create a smooth projected surface density map as prescribed in section 3.2.1.

We adopt the $\Lambda$CDM cosmology from Planck Collaboration (2020) with $h = 0.68$ to compute angular diameter distances. We also fix the source redshift to $z_s = 1.5$ and the deflector redshift to $z_\ell = 0.5$. We note that changing the redshifts or the cosmology only amount in a rescaling of the $\kappa$ map by a global scalar. Thus, this choice does not change the generality of our method. The smoothed density maps from equation (11) are rendered into a regular grid of $188^2$ pixels with a comoving

field of view of 105 kpc/$h$. To avoid edge effects in the pixelated maps, we include particles outside of the field of view in the sum of equation (11).

Before applying augmentations or considering different projections, our dataset of halos is split into a training set (90%) and a test set (10%), in order to make sure that the test set consists only of convergence maps unseen by the RIM during training. We take 3 different projections ($xy$, $xz$ and $yz$) of each 3D particle distribution, which amounts to a dataset with a total of 4 812 individual convergence maps. Random rotations by an angle multiple of 90° and random shifts to the pixel coordinates are applied to each image. The $\kappa$ maps are then rescaled by a random factor to change their estimated Einstein radius to the range [0.5, 2.5] arcseconds. The Einstein radius is defined as

$$\theta_E = \sqrt{\frac{4GM(\theta_E)}{c^2}\frac{D_{\ell s}}{D_\ell D_s}} \qquad (13)$$

where $M(\theta_E)$ is the mass enclosed inside the Einstein radius. In practice, we estimate this quantity by summing over the mass of pixels with a value greater than the critical density ($\kappa > 1$). For data augmentation purposes, this procedure gives a good enough estimate of the lensed image separation resulting from a given $\kappa$ map. We test multiple scaling factors for each $\kappa$ map, then uniformly sample between those that produce an estimated Einstein radius within the desired range. This step is used to remove any bias in the Einstein radius that might come from the mass function of the simulation.

The final maps are cropped down to $128^2$ pixels. Placed at a redshift $z_\ell = 0.5$, a $\kappa$ map will thus span an angular field of view of 7.69″ with a resolution similar to *HST*. With these augmentation procedures, we create a total of 50 000 maps from the training split to train a VAE, as described in Section 4.1, and produce simulated observations to train the RIM.

### 3.3. *Simulated Observations*

With a given source map and convergence map, we apply the ray tracing simulation described in section 2.2 to produce a lensed image.

For each lensed image, we create a Gaussian PSF with a full width at half maximum (FWHM) randomly generated from a truncated normal distribution. The support of the distribution is truncated below by the angular size of a single pixel and above by the angular size of 4 pixels. White noise with a standard deviation randomly generated from a truncated normal distribution is then added to the convolved lensed image to simulate noisy observations. These noise realizations result in SNRs between

10 and 1000. For simplicity, we define SNR $= \frac{1}{\sigma}$. This definition is equivalent to the peak signal-to-noise ratio.

To ensure that the images are representative of strongly lensed sources, we require a minimum flux magnification of 3. We also make sure that most pixel coordinates in the image plane are mapped inside the source coordinate system through the lens equation (4).

**Table 1.** Physical model parameters.

| Parameter | Distribution/Value |
|---|---|
| Lens redshift $z_\ell$ | 0.5 |
| Source redshift $z_s$ | 1.5 |
| Field of view (″) | 7.69 |
| Source field of view (″) | 3 |
| PSF FWHM (″) | $\mathcal{TN}(0.06, 0.3; 0.08, 0.05)$ [a] |
| Noise amplitude $\sigma$ | $\mathcal{TN}(0.001, 0.1; 0.01, 0.03)$ |

[a] We defined the parameters of the truncated normal in the order $\mathcal{TN}(a, b; \mu, \sigma)$, where $[a, b]$ defines the support of the distribution.

In total, 400 000 training observations are simulated from random pairs of COSMOS sources and IllustrisTNG convergence maps in order to train the RIM. We create an additional 200 000 observations from pairs of COSMOS sources and pixelated SIE convergence maps. The parameters for these $\kappa$ maps are listed in table 2.

**Table 2.** SIE parameters.

| Parameter | Distribution |
|---|---|
| Radial shift (″) | $\mathcal{U}(0, 0.1)$ |
| Azimutal shift | $\mathcal{U}(0, 2\pi)$ |
| Orientation | $\mathcal{U}(0, \pi)$ |
| $\theta_E$ (″) | $\mathcal{U}(0.5, 2.5)$ |
| Ellipticity | $\mathcal{U}(0, 0.6)$ |

We generate 1 600 000 simulated observations from the VAE background sources and convergence maps as part of the training set. We apply some validation checks to each example in order to avoid configurations like a single image of the background source or an Einstein ring cropped by the field of view.

## 4. TRAINING

### 4.1. *VAE*

Here, we describe the training of two VAEs that are used to produce density maps and images of unlensed background galaxies to train and test our inference model. For an introduction to VAEs we refer the reader to Kingma & Welling (2019).

As mentioned in Kingma & Welling (2019), direct optimisation of the ELBO loss can prove difficult because the reconstruction term $\log p_\theta(\mathbf{x} \mid \mathbf{z})$ is relatively weak compared to the Kullback Leibler (KL) divergence term. To alleviate this issue, we follow the work of Bowman et al. (2015) and Kaae Sønderby et al. (2016) in setting a warm-up schedule for the KL term, starting from $\beta = 0.1$ up to $\beta_{\max}$.

Usually, $\beta_{\max} = 1$ is considered optimal since it matches the original ELBO objective derived by Kingma & Welling (2013). However, we are more interested in the sharpness of our samples and accurate inference around small regions of the latent space for fine-tuning. Thus, setting $\beta_{\max} < 1$ allows us to increase the size of the information bottleneck (i.e. latent space) of the VAE and improve the reconstruction cost of the model. This is a variant of the $\beta$-VAE (Higgins et al. 2017), where $\beta > 1$ was found to improve disentangling of the latent space (Burgess et al. 2018).

The value for $\beta_{\max}$ and the steepness of the schedule are grid searched alongside the architecture for the VAE. These values are found in practice by manually looking at the quality of generated samples for different VAE hyperparameters. A similar method is explored and formalized in the InfoVAE framework (Zhao et al. 2017).

A notable element of the VAE architecture is the use of a fully connected layer to reshape the features of the convolutional layer into the chosen latent space dimension. Following the work of Lanusse et al. (2021), we introduce an $\ell_2$ penalty between the input and output of the bottleneck dense layers to encourage an identity mapping. This regularisation term is slowly removed during training.

### 4.2. RIM

The architecture of the neural network is grid searched on a smaller dataset ($\lesssim 10\,000$ examples) in order to quickly identify a small set of valid hyperparameters. The best hyperparameters are then identified using a two-stage training process. In the first stage, we train 24 different architectures from this small hyperparameter set for approximately 4 days (wall time using a single Nvidia A100 GPU). Different architectures have a training time much longer than others, which is factored in the architecture selection process. For example, adding more time steps ($T$) to the recurrent relation (6) yields better generalisation on the test set, but this comes at great costs to training time until convergence. Thus, $T < 10$ is preferred.

Following this first stage, 4 architectures are deemed efficient enough to be trained for an additional 6 days.

We only report the results for the best architectures out of these 4.

Each reconstruction is performed by fine-tuning the baseline model on a test task composed of an observation, a PSF, and a noise covariance. In practice, fine-tuning predictions on the test set of 3 000 examples can be accomplished in parallel so as to be done in at most a few days by spreading the computation on $\sim 10$ Nvidia A100 GPUs. Each reconstruction uses at most 2000 steps, correspondling to approximately 20 minutes (wall-time) per reconstruction. Early stopping is applied when the $\chi^2$ reaches noise level. The hyperparameters for this procedure are reported in Table 3.

**Table 3.** Hyperparameters for fine-tuning the RIM.

| Parameter | Value |
|---|---|
| Optimizer | RMSProp |
| Learning rate | $10^{-6}$ |
| Maximum number of steps | 2 000 |
| $\lambda$ | $2 \times 10^5$ |
| $\ell_2$ | 0 |
| Number of samples from VAE | 200 |
| Latent space distribution | $\mathcal{N}(\mathbf{z}^{(T)}, \sigma = 0.3)$ [a] |

[a]  $\mathbf{z}^{(T)}$ is the latent code of the RIM baseline source or convergence.

## 5. RESULTS

In this section, we present the performance of our model on the held out test set. A sample of 3000 reconstruction problems is generated from the held-out *HST* and IllustrisTNG data with noise levels and PSFs similar to the training set.

### 5.1. Goodness of Fit

Figure 7 shows a sample of reconstructions for high SNR data with a wide range of lensing configurations from the test set. We select examples representative of all levels of reconstruction performance (covering the entire range of goodness of fit) for data with complex structures in their convergence map to showcase the expressivity of the approach. We also show a randomly selected sample from the test set in Figure 13.

Figure 8 shows a comparison between the goodness of fit of the baseline model and the fine-tuned prediction.

**Table 4.** $\log_{10}$-normal moments of the loss on the test set

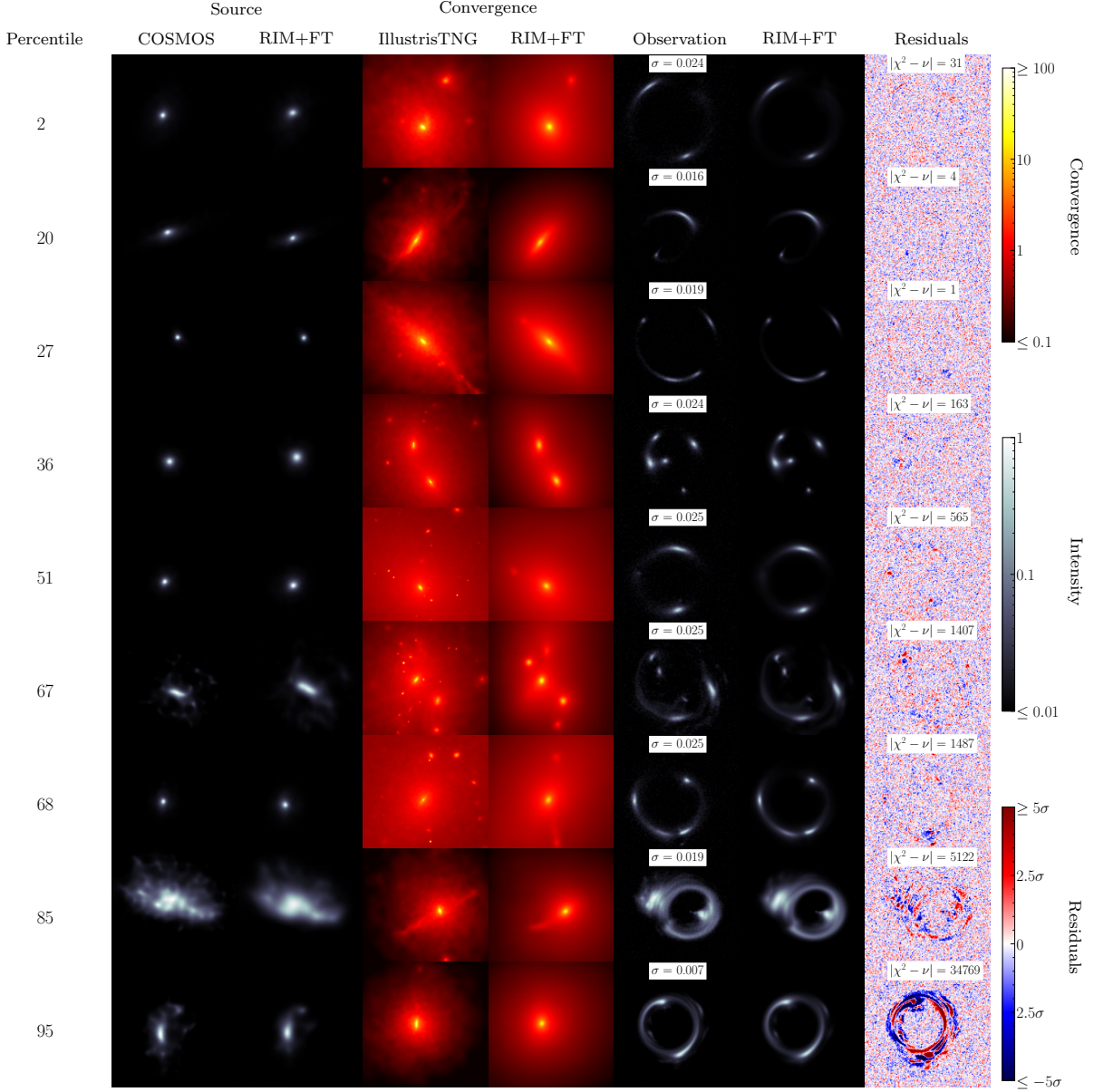| Model | $\mu(\log \mathcal{L}_\varphi)$ | $\sigma(\log \mathcal{L}_\varphi)$ |
|---|---|---|
| Baseline ($\varphi_\mathcal{D}^\star$) | -1.96 | 0.36 |
| Fine-tuned ($\hat{\varphi}_{\mathrm{MAP}}$) | -2.02 | 0.37 |

**Figure 7.** Sample of the fine-tuned RIM reconstructions on a test set of 3000 examples. Examples are ordered from the best $\chi^2$ (top) to the worst (bottom). The percentile rank of each example is in the leftmost column. The last example shown has SNR above the threshold defined in Figure 9.

Since we empirically observe that the distribution of the loss on the test set (and the training set) follows a log-normal distribution, we find that it is more informative to look at the log-loss distribution to extract information about the fine-tuning procedure. The left panel of Figure 8 shows the distribution of the log-loss difference between the fine-tuned prediction and the baseline model. This distribution shows that the fine-tuning procedure loss is constrained within $\sim 1$ order of magnitude of the original loss with a probability $> 99.73\,\%$. We find that the log-loss difference has a scatter of $\sigma = 0.28$, which is smaller than the scatter of the baseline log-loss over

the entire test set $\sigma(\log \mathcal{L}_{\varphi_{\mathcal{D}}^\star}) = 0.36$ reported in Table 4. We note that the loss is not optimized during fine-tuning, still we notice that the fine-tuning procedure does not significantly deteriorate or improve the loss of the baseline prediction on average. We report the first 2 moments of the loss log-normal distribution for the baseline and the fine-tuned reconstructions in Table 4 in order to explicitly compare them. As can be seen in this table, there is no significant difference between the two distributions. This statement can be proven for the measured mean values — $\mu(\log \mathcal{L}_{\hat{\varphi}_{\mathrm{MAP}}}) = \mu(\log \mathcal{L}_{\varphi_{\mathcal{D}}^\star})$ — using the two-sided normal p-value test (Casella &

**Figure 8.** Distribution of the goodness of fit for the baseline and fine-tuned network (right panel), as well as log-loss difference between the two network for a given example in the test set (left panel).
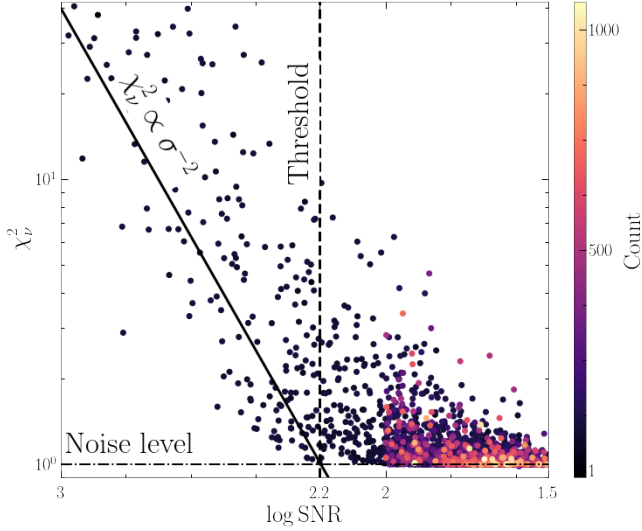


**Figure 9.** Goodness of fit as a function of SNR shows a threshold behavior where our method reaches its limit.

Berger 2001), which we find satisfy the null hypothesis with $p = 0.87$ ($Z = -0.16$). All those observations support our claim that EWC regularisation preserves the prior learned during pretraining, or at least that it preserves the surrogate measures of the prior we reported.

The right panel of Figure 8 shows the distribution of $\chi^2$ for the test set before and after the fine-tuning procedure and the theoretical $\chi^2$ distribution corresponding to $\nu = 128^2$ degrees of freedom. We observe that the fine-tuning procedure significantly improves our $\chi^2$, bringing their distribution closer to that of the expected $\chi^2$ distribution (black curve). However, the improved distribution is still far from the theoretical expectation,

implying that there are statistically significant residuals in a subset of the reconstructions.

In figure 9, we explore how the goodness of fit of the fine-tuned RIM changes as a function of SNR over the examples in the test set. Two behaviors can be identified. For SNR below a certain threshold, the goodness of fit of the fine-tuned model is essentially flat, with a certain scatter, around the noise level. This scatter increases as a function of SNR, which reflects the fact that above a certain SNR threshold (vertical dashed line in Figure 9), our reconstructions are dominated by systematics in the inference algorithm. For SNR above the threshold, the goodness of fit follows the trend $\chi^2 \propto \sigma^{-2}$ (the solid line in Figure 9), which means the reconstructions have stopped improving on par with the SNR.

This behavior is exhibited in a few examples of reconstructions taken from the test set in Figure 10, where we order reconstructions with increasing SNR from top to bottom and plot the surface brightness and foreground densities in log scale. As can been seen, the amplitude of the residual increases significantly as we increase the SNR. Above the SNR threshold ($\sim 220$), the reconstructions are dominated by systematics.

### 5.2. *Quality of the Reconstructions*

In addition to a visual inspection of the reconstructed sources and convergences, we compute the coherence spectrum to quantitatively assess the quality of the reconstructions

$$\gamma(k) = \frac{P_{12}(k)}{\sqrt{P_{11}(k)P_{22}(k)}} \, . \tag{14}$$

Here, $P_{12}(k)$ is the cross power spectrum between a reconstructed and a true image at the wavenumber $k$.
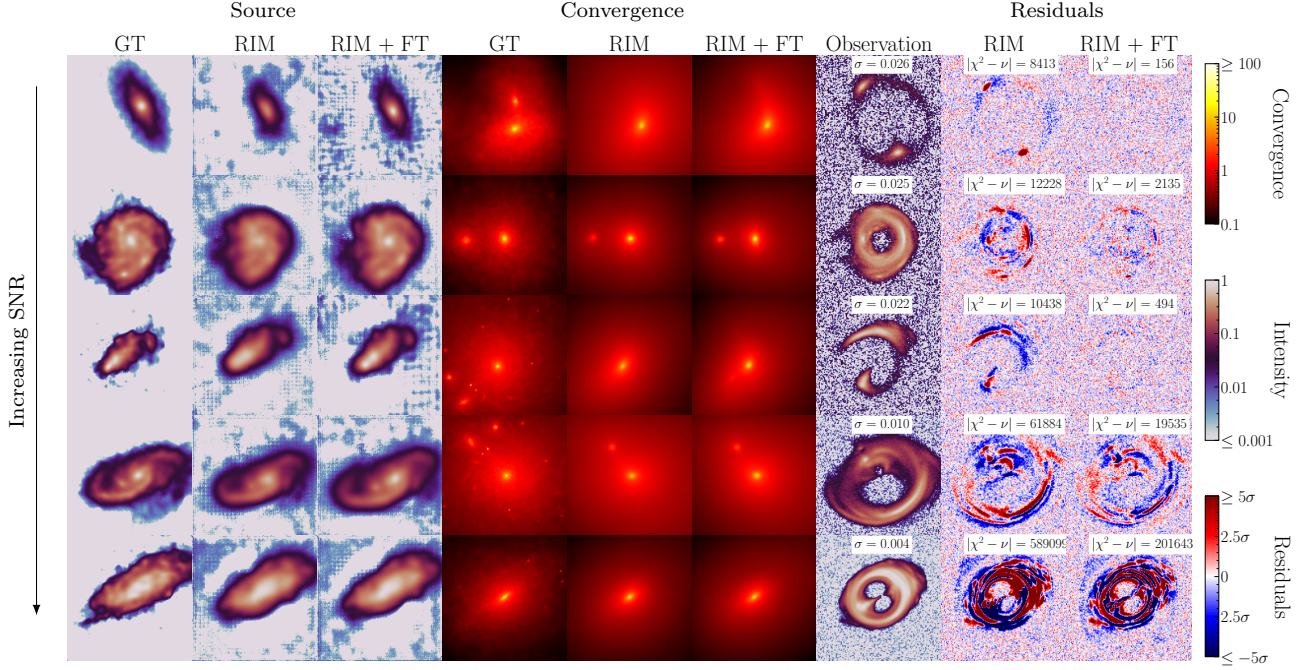
**Figure 10.** Comparison between baseline (RIM) and fine-tuned (RIM+FT) reconstructions for gravitational lensing systems from the test set (GT). From top to bottom, we increase SNR.
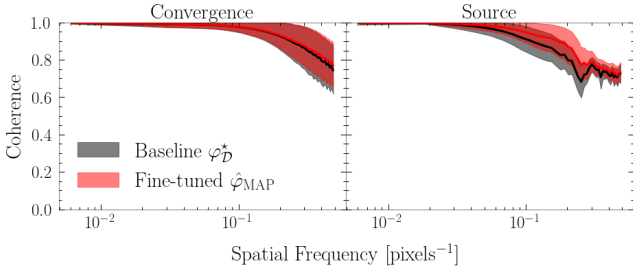


**Figure 11.** Statistics of the coherence spectrum on the test set. The solid line is the average coherence. The transparent region is the 68% confidence interval. The fine-tuning procedure yields a noticeable improvement on the coherence of the source at all frequencies.

Figure 11 shows the coherence of the source and convergence maps for the entire test set of 3000 examples, with the mean value and the 68% inclusion interval of $\gamma(k)$ reported in the solid line and shaded area respectively. The fine-tuning procedure, shown in red, is able to significantly improve the coherence of the baseline background source, shown in black, at all scales. The coherence spectrum of the convergence sees a slight improvement due to the fine-tuning procedure. Still, we note that many examples in the dataset exhibit significant improvement, which we illustrate in Figure 3.

## 6. CONCLUSION

The results obtained here demonstrate the effectiveness of machine learning methods, specifically a recurrent inference machine, for inferring pixelated maps of the distribution of mass in lensing galaxies and the distribution of surface brightness in the background galaxies. Since this is a heavily under-constrained problem, stringent priors are needed to avoid overfitting the data, a task that has traditionally been difficult to accomplish with traditional statistical models (e.g., Saha & Williams 1997). The model proposed here can implicitly learn these priors from a set of training data.

The fine-tuning step that we propose in this work is a general procedure (i.e. not specific to our model or problem), which enables us to exploit a diagonal second-order Laplace approximation of the implicit prior learned by a baseline estimator during pre-training. We use fine-tuning in order to significantly improve this baseline estimator (i.e., a better MAP estimate), by using the likelihood of the data and the EWC prior. In the context of our work, we find that fine-tuning has a limiting — or threshold — behavior, which we speculate is due to the limited expressivity of the neural network and its inductive biases learned during pre-training.

The flexible and expressive form of the reconstructions shown in this work means that, in principle, any lensing system (e.g., a single simple galaxy or a group of complex galaxies) could be analyzed by this model, without

any need for pre-determining the model parameterization. This is of high value given the diversity of observed lensing systems, and their relevance for constraining astrophysical and cosmological parameters.

Perhaps the most important limitation of the method is the fact that, in its current form, the model only provides point estimates of the parameters of interest. Quantifying the posteriors of such high-dimensional data will require an efficient and accurate generative process (e.g., see Adam et al. 2022), which we plan to explore and develop in future works.

## SOFTWARE AND DATA

The source code, as well as the various scripts and parameters used to produce the model and results is available as open-source software under the package Censai[1]. The model parameters, as well as convergence maps used to train these models and the test set examples and reconstructions results are also available as open-source datasets hosted by Zenodo[2]. This research made use of Tensorflow (Abadi et al. 2015), Tensorflow-Probability (Dillon et al. 2017), Numpy (Harris et al. 2020), Scipy (Virtanen et al. 2020), Matplotlib (Hunter 2007), Scikit-image (Van der Walt et al. 2014), IPython (Pérez & Granger 2007), Pandas (Wes McKinney 2010; pandas development team 2020), Scikit-learn (Pedregosa et al. 2011), Astropy (Astropy Collaboration et al. 2013, 2018) and GalSim (Rowe et al. 2015).

## ACKNOWLEDGEMENTS

---

[1] https://github.com/AlexandreAdam/Censai
[2] https://doi.org/10.5281/zenodo.6555463

REFERENCES

Abadi, M., Agarwal, A., Barham, P., et al. 2015, TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. https://www.tensorflow.org/

Abdelsalam, H. M., Saha, P., & Williams, L. L. R. 1998a, AJ, 116, 1541, doi: 10.1086/300546

—. 1998b, MNRAS, 294, 734, doi: 10.1046/j.1365-8711.1998.01356.x

Adam, A., Coogan, A., Malkin, N., et al. 2022, arXiv e-prints, arXiv:2211.03812. https://arxiv.org/abs/2211.03812

Anau Montel, N., Coogan, A., Correa, C., Karchev, K., & Weniger, C. 2022, arXiv e-prints, arXiv:2205.09126. https://arxiv.org/abs/2205.09126

Andrychowicz, M., Denil, M., Gomez, S., et al. 2016, arXiv e-prints, arXiv:1606.04474. https://arxiv.org/abs/1606.04474

Astropy Collaboration, Robitaille, T. P., Tollerud, E. J., et al. 2013, A&A, 558, A33, doi: 10.1051/0004-6361/201322068

Astropy Collaboration, Price-Whelan, A. M., Sipőcz, B. M., et al. 2018, AJ, 156, 123, doi: 10.3847/1538-3881/aabc4f

Aubert, D., Amara, A., & Benton Metcalf, R. 2007, Monthly Notices of the Royal Astronomical Society, 376, 113, doi: 10.1111/j.1365-2966.2006.11296.x

Auger, M. W., Treu, T., Bolton, A. S., et al. 2010, ApJ, 724, 511, doi: 10.1088/0004-637X/724/1/511

Barnabè, M., Czoske, O., Koopmans, L. V. E., et al. 2009, MNRAS, 399, 21, doi: 10.1111/j.1365-2966.2009.14941.x

Bartelmann, M., Narayan, R., Seitz, S., & Schneider, P. 1996, ApJL, 464, L115, doi: 10.1086/310114

Bellagamba, F., Tessore, N., & Metcalf, R. B. 2016, Monthly Notices of the Royal Astronomical Society, 464, 4823, doi: 10.1093/mnras/stw2726

Belokurov, V., Evans, N. W., Moiseev, A., et al. 2007, ApJL, 671, L9, doi: 10.1086/524948

Bengio, Y. 2009, Found. Trends Mach. Learn., 2, 1–127, doi: 10.1561/2200000006

Birrer, S., & Amara, A. 2018, Physics of the Dark Universe, 22, 189, doi: 10.1016/j.dark.2018.11.002

Birrer, S., Amara, A., & Refregier, A. 2015, ApJ, 813, 102, doi: 10.1088/0004-637X/813/2/102

Birrer, S., Treu, T., Rusu, C. E., et al. 2019, MNRAS, 484, 4726, doi: 10.1093/mnras/stz200

Bowman, S. R., Vilnis, L., Vinyals, O., et al. 2015, arXiv e-prints, arXiv:1511.06349. https://arxiv.org/abs/1511.06349

Bradač, M., Schneider, P., Lombardi, M., & Erben, T. 2005, A&A, 437, 39, doi: 10.1051/0004-6361:20042233

Burgess, C. P., Higgins, I., Pal, A., et al. 2018, arXiv e-prints, arXiv:1804.03599. https://arxiv.org/abs/1804.03599

Cacciato, M., Bartelmann, M., Meneghetti, M., & Moscardini, L. 2006, A&A, 458, 349, doi: 10.1051/0004-6361:20054582

Casella, G., & Berger, R. 2001, Statistical Inference (Duxbury Resource Center)

Cheng, J., Wiesner, M. P., Peng, E.-H., et al. 2019, ApJ, 872, 185, doi: 10.3847/1538-4357/ab0029

Cho, K., van Merrienboer, B., Gulcehre, C., et al. 2014, arXiv e-prints, arXiv:1406.1078. https://arxiv.org/abs/1406.1078

Coe, D., Fuselier, E., Benítez, N., et al. 2008, ApJ, 681, 814, doi: 10.1086/588250

Coles, J. P., Read, J. I., & Saha, P. 2014, MNRAS, 445, 2181, doi: 10.1093/mnras/stu1781

Coogan, A., Karchev, K., & Weniger, C. 2020, arXiv e-prints, arXiv:2010.07032. https://arxiv.org/abs/2010.07032

Cramér, H. 1946, Mathematical methods of statistics, Vol. 9 (Princeton University Press, Princeton, NJ)

Dalal, N., & Kochanek, C. S. 2002, ApJ, 572, 25, doi: 10.1086/340303

Deb, S., Morandi, A., Pedersen, K., et al. 2012, arXiv e-prints, arXiv:1201.3636. https://arxiv.org/abs/1201.3636

Diego, J. M., Protopapas, P., Sandvik, H. B., & Tegmark, M. 2005, MNRAS, 360, 477, doi: 10.1111/j.1365-2966.2005.09021.x

Diego, J. M., Tegmark, M., Protopapas, P., & Sandvik, H. B. 2007, MNRAS, 375, 958, doi: 10.1111/j.1365-2966.2007.11380.x

Dillon, J. V., Langmore, I., Tran, D., et al. 2017, arXiv e-prints, arXiv:1711.10604. https://arxiv.org/abs/1711.10604

Galan, A., Peel, A., Joseph, R., Courbin, F., & Starck, J. L. 2021, A&A, 647, A176, doi: 10.1051/0004-6361/202039363

Galan, A., Vernardos, G., Peel, A., Courbin, F., & Starck, J. L. 2022, A&A, 668, A155, doi: 10.1051/0004-6361/202244464

Ghosh, A., Williams, L. L. R., & Liesenborgs, J. 2020, MNRAS, 494, 3998, doi: 10.1093/mnras/staa962

Gilman, D., Birrer, S., Nierenberg, A., et al. 2020, MNRAS, 491, 6077, doi: 10.1093/mnras/stz3480

Gilman, D., Bovy, J., Treu, T., et al. 2021, MNRAS, 507, 2432, doi: 10.1093/mnras/stab2335

Harris, C. R., Millman, K. J., van der Walt, S. J., et al. 2020, Nature, 585, 357, doi: 10.1038/s41586-020-2649-2

Hezaveh, Y. D., Perreault Levasseur, L., & Marshall, P. J. 2017, Nature, 548, 555, doi: 10.1038/nature23463

Hezaveh, Y. D., Dalal, N., Marrone, D. P., et al. 2016, ApJ, 823, 37, doi: 10.3847/0004-637X/823/1/37

Higgins, I., Matthey, L., Pal, A., et al. 2017, in ICLR

Hunter, J. D. 2007, Computing in Science & Engineering, 9, 90, doi: 10.1109/MCSE.2007.55

Jee, M. J., Ford, H. C., Illingworth, G. D., et al. 2007, ApJ, 661, 728, doi: 10.1086/517498

Kaae Sønderby, C., Raiko, T., Maaløe, L., Kaae Sønderby, S., & Winther, O. 2016, arXiv e-prints, arXiv:1602.02282. https://arxiv.org/abs/1602.02282

Karchev, K., Coogan, A., & Weniger, C. 2022, MNRAS, 512, 661, doi: 10.1093/mnras/stac311

Kingma, D. P., & Ba, J. 2014, arXiv e-prints, arXiv:1412.6980. https://arxiv.org/abs/1412.6980

Kingma, D. P., & Welling, M. 2013, arXiv e-prints, arXiv:1312.6114. https://arxiv.org/abs/1312.6114

—. 2019, arXiv e-prints, arXiv:1906.02691. https://arxiv.org/abs/1906.02691

Kirkpatrick, J., Pascanu, R., Rabinowitz, N., et al. 2016, arXiv e-prints, arXiv:1612.00796. https://arxiv.org/abs/1612.00796

Koekemoer, A. M., Aussel, H., Calzetti, D., et al. 2007, The Astrophysical Journal Supplement Series, 172, 196, doi: 10.1086/520086

Koopmans, L. V. E., Treu, T., Bolton, A. S., Burles, S., & Moustakas, L. A. 2006, ApJ, 649, 599, doi: 10.1086/505696

Lanusse, F., Mandelbaum, R., Ravanbakhsh, S., et al. 2021, MNRAS, 504, 5543, doi: 10.1093/mnras/stab1214

Leauthaud, A., Massey, R., Kneib, J.-P., et al. 2007, The Astrophysical Journal Supplement Series, 172, 219, doi: 10.1086/516598

Legin, R., Hezaveh, Y., Perreault Levasseur, L., & Wandelt, B. 2021, arXiv e-prints, arXiv:2112.05278. https://arxiv.org/abs/2112.05278

Legin, R., Stone, C., Hezaveh, Y., & Perreault-Levasseur, L. 2022, arXiv e-prints, arXiv:2207.04123. https://arxiv.org/abs/2207.04123

Li, N., Becker, C., & Dye, S. 2021, MNRAS, 504, 2224, doi: 10.1093/mnras/stab984

Liesenborgs, J., De Rijcke, S., & Dejonghe, H. 2006, MNRAS, 367, 1209, doi: 10.1111/j.1365-2966.2006.10040.x

Liesenborgs, J., de Rijcke, S., Dejonghe, H., & Bekaert, P. 2007, MNRAS, 380, 1729, doi: 10.1111/j.1365-2966.2007.12236.x

Lønning, K., Putzky, P., Sonke, J. J., et al. 2019, Medical Image Analysis, 53, 64, doi: 10.1016/j.media.2019.01.005

Mandelbaum, R., Lackner, C., Leauthaud, A., & Rowe, B. 2012, Zenodo. https://zenodo.org/record/3242143

Mandelbaum, R., Rowe, B., Bosch, J., et al. 2014, The Astrophysical Journal Supplement Series, 212, 5, doi: 10.1088/0067-0049/212/1/5

Marrone, D. P., Spilker, J. S., Hayward, C. C., et al. 2018, Nature, 553, 51, doi: 10.1038/nature24629

McCloskey, M., & Cohen, N. J. 1989in (Academic Press), 109–165, doi: 10.1016/S0079-7421(08)60536-8

Merten, J. 2016, MNRAS, 461, 2328, doi: 10.1093/mnras/stw1413

Merten, J., Cacciato, M., Meneghetti, M., Mignone, C., & Bartelmann, M. 2009, A&A, 500, 681, doi: 10.1051/0004-6361/200810372

Mishra-Sharma, S., & Yang, G. 2022, Strong Lensing Source Reconstruction Using Continuous Neural Fields, arXiv, doi: 10.48550/ARXIV.2206.14820

Modi, C., Lanusse, F., Seljak, U., Spergel, D. N., & Perreault-Levasseur, L. 2021, arXiv e-prints, arXiv:2104.12864. https://arxiv.org/abs/2104.12864

Morningstar, W. R., Hezaveh, Y. D., Levasseur, L. P., et al. 2018, arXiv e-prints. https://arxiv.org/abs/1808.00011v1

Morningstar, W. R., Levasseur, L. P., Hezaveh, Y. D., et al. 2019, The Astrophysical Journal, 883, 14, doi: 10.3847/1538-4357/ab35d7

Nelson, D., Springel, V., Pillepich, A., et al. 2019, MNRAS, 6, doi: 10.1186/s40668-019-0028-x

Nightingale, J. W., Dye, S., & Massey, R. J. 2018, MNRAS, 478, 4738, doi: 10.1093/mnras/sty1264

Pan, S. J., & Yang, Q. 2010, IEEE Transactions on Knowledge and Data Engineering, 22, 1345

pandas development team, T. 2020, pandas-dev/pandas: Pandas, latest, Zenodo, doi: 10.5281/zenodo.3509134

Park, J. W., Wagner-Carena, S., Birrer, S., et al. 2021, ApJ, 910, 39, doi: 10.3847/1538-4357/abdfc4

Pedregosa, F., Varoquaux, G., Gramfort, A., et al. 2011, Journal of Machine Learning Research, 12, 2825

Pérez, F., & Granger, B. E. 2007, Computing in Science and Engineering, 9, 21, doi: 10.1109/MCSE.2007.53

Perreault Levasseur, L., Hezaveh, Y. D., & Wechsler, R. H. 2017, ApJL, 850, L7, doi: 10.3847/2041-8213/aa9704

Planck Collaboration. 2020, A&A, 641, A6, doi: 10.1051/0004-6361/201833910

Putzky, P., & Welling, M. 2017, arXiv e-prints. https://arxiv.org/abs/1706.04008

Rahaman, N., Baratin, A., Arpit, D., et al. 2018, arXiv e-prints, arXiv:1806.08734. https://arxiv.org/abs/1806.08734

16

Rao, C. 1945, Bulletin fo the Calcutta Mathematical Society

Ratcliff, R. 1990, Psychological Review, 97, 285, doi: 10.1037/0033-295X.97.2.285

Rau, S., Vegetti, S., & White, S. D. 2013, Monthly Notices of the Royal Astronomical Society, 430, 2232, doi: 10.1093/mnras/stt043

Rizzo, F., Vegetti, S., Powell, D., et al. 2020, Nature, 584, 201, doi: 10.1038/s41586-020-2572-6

Ronneberger, O., Fischer, P., & Brox, T. 2015, arXiv e-prints, arXiv:1505.04597. https://arxiv.org/abs/1505.04597

Rowe, B. T., Jarvis, M., Mandelbaum, R., et al. 2015, Astronomy and Computing, 10, 121, doi: 10.1016/j.ascom.2015.02.002

Rowe, B. T. P., Jarvis, M., Mandelbaum, R., et al. 2015, Astronomy and Computing, 10, 121, doi: 10.1016/j.ascom.2015.02.002

Rusu, C. E., Fassnacht, C. D., Sluse, D., et al. 2017, MNRAS, 467, 4220, doi: 10.1093/mnras/stx285

Rusu, C. E., Wong, K. C., Bonvin, V., et al. 2020, MNRAS, 498, 1440, doi: 10.1093/mnras/stz3451

Saha, P., & Williams, L. L. R. 1997, MNRAS, 292, 148, doi: 10.1093/mnras/292.1.148

—. 2004, AJ, 127, 2604, doi: 10.1086/383544

Schmidt, T., Treu, T., Birrer, S., et al. 2022, arXiv e-prints, arXiv:2206.04696. https://arxiv.org/abs/2206.04696

Schuldt, S., Chirivì, G., Suyu, S. H., et al. 2019, A&A, 631, A40, doi: 10.1051/0004-6361/201935042

Schuldt, S., Suyu, S. H., Canameras, R., et al. 2022, arXiv e-prints, arXiv:2207.10124. https://arxiv.org/abs/2207.10124

Scoville, N., Aussel, H., Brusa, M., et al. 2007, The Astrophysical Journal Supplement Series, 172, 1, doi: 10.1086/516585

Seitz, S., Schneider, P., & Bartelmann, M. 1998, A&A, 337, 325. https://arxiv.org/abs/astro-ph/9803038

Sérsic, J. L. 1963, Boletin de la Asociacion Argentina de Astronomia La Plata Argentina, 6, 41

Sluse, D., Sonnenfeld, A., Rumbaugh, N., et al. 2017, MNRAS, 470, 4838, doi: 10.1093/mnras/stx1484

Sun, F., Egami, E., Pérez-González, P. G., et al. 2021, ApJ, 922, 114, doi: 10.3847/1538-4357/ac2578

Suyu, S. H., Marshall, P. J., Hobson, M. P., & Blandford, R. D. 2006, MNRAS, 371, 983, doi: 10.1111/j.1365-2966.2006.10733.x

Tishby, N., Pereira, F. C., & Bialek, W. 2000, arXiv e-prints, physics/0004057. https://arxiv.org/abs/physics/0004057

Torres-Ballesteros, D. A., & Castañeda, L. 2022, arXiv e-prints, arXiv:2201.10076. https://arxiv.org/abs/2201.10076

Treu, T., & Koopmans, L. V. E. 2004, ApJ, 611, 739, doi: 10.1086/422245

Van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., et al. 2014, PeerJ, 2, e453

Vieira, J. D., Marrone, D. P., Chapman, S. C., et al. 2013, Nature, 495, 344, doi: 10.1038/nature12001

Vincent, P., Larochelle, H., Bengio, Y., & Manzagol, P.-A. 2008, in Proceedings of the 25th International Conference on Machine Learning, ICML '08 (New York, NY, USA: Association for Computing Machinery), 1096–1103, doi: 10.1145/1390156.1390294

Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., & Manzagol, P.-A. 2010, J. Mach. Learn. Res., 11, 3371–3408

Virtanen, P., Gommers, R., Oliphant, T. E., et al. 2020, Nature Methods, 17, 261, doi: 10.1038/s41592-019-0686-2

Wagner-Carena, S., Aalbers, J., Birrer, S., et al. 2022, arXiv e-prints, arXiv:2203.00690. https://arxiv.org/abs/2203.00690

Wagner-Carena, S., Park, J. W., Birrer, S., et al. 2021, ApJ, 909, 187, doi: 10.3847/1538-4357/abdf59

Warren, S. J., & Dye, S. 2003, ApJ, 590, 673, doi: 10.1086/375132

Wes McKinney. 2010, in Proceedings of the 9th Python in Science Conference, ed. Stéfan van der Walt & Jarrod Millman, 56 – 61, doi: 10.25080/Majora-92bf1922-00a

Wong, K. C., Suyu, S. H., Auger, M. W., et al. 2017, MNRAS, 465, 4895, doi: 10.1093/mnras/stw3077

Zhao, S., Song, J., & Ermon, S. 2017, arXiv e-prints, arXiv:1706.02262. https://arxiv.org/abs/1706.02262

Zhuang, F., Qi, Z., Duan, K., et al. 2019, arXiv e-prints, arXiv:1911.02685. https://arxiv.org/abs/1911.02685

APPENDIX

## A. ELASTIC WEIGHT CONSOLIDATION

Suppose we are given a training set $\mathcal{D}$ and a test task $\mathcal{T}$. The posterior of the RIM parameters $\varphi$ can be rewritten using the Bayes rule as

$$p(\varphi \mid \mathcal{D}, \mathcal{T}) = \frac{p(\mathcal{T} \mid \mathcal{D}, \varphi)p(\varphi \mid \mathcal{D})}{p(\mathcal{T} \mid \mathcal{D})}. \tag{A1}$$

We suppose that $\varphi$ encode information about $\mathcal{D}$, while $\mathcal{T}$ is unseen by $\varphi$. It follows that $\mathcal{T}$ and $\mathcal{D}$ are conditionally independent when given $\varphi$. We do not make the stronger assumption that $\mathcal{D}$ and $\mathcal{T}$ are completely independent. In fact, such an assumption would contradict the premise of our work that building a dataset $\mathcal{D}$ can inform a machine (RIM) about task $\mathcal{T}$ — or that, more broadly, $\mathcal{D}$ contains information about $\mathcal{T}$.

We rewrite the marginal $p(\mathcal{T} \mid \mathcal{D})$ using the Bayes rule in order to extract $p(\mathcal{D} \mid \mathcal{T})$, the sampling distribution used to compute the Fisher diagonal elements

$$p(\varphi \mid \mathcal{D}, \mathcal{T}) = \frac{p(\mathcal{T} \mid \varphi)p(\varphi \mid \mathcal{D})}{p(\mathcal{D} \mid \mathcal{T})} \frac{p(\mathcal{D})}{p(\mathcal{T})}. \tag{A2}$$

The log-likelihood $\log p(\mathcal{T} \mid \varphi)$ is equivalent to the negative of the loss function for the particular task at hand. In this work, we assign a uniform probability density to $p(\mathcal{T})$ and $p(\mathcal{D})$ in order to ignore them.

We now turn to the prior $p(\varphi \mid \mathcal{D})$, which appears as a conditional relative to the training dataset. We use the Laplace approximation around the maximum $\varphi_{\mathcal{D}}^\star$ to evaluate the prior, where $\varphi_{\mathcal{D}}^\star$ are the trained parameters of the RIM that minimize the empirical risk (equation (8)). The Taylor expansion of the prior around this maximum yields

$$\log p(\varphi \mid \mathcal{D}) \approx \log p(\varphi_{\mathcal{D}}^\star \mid \mathcal{D}) + \frac{1}{2}(\varphi - \varphi_{\mathcal{D}}^\star)^T \underbrace{\left( \left. \frac{\partial^2 \log p(\varphi \mid \mathcal{D})}{\partial^2 \varphi} \right|_{\varphi_{\mathcal{D}}^\star} \right)}_{\mathbf{H}(\varphi_{\mathcal{D}}^\star)} (\varphi - \varphi_{\mathcal{D}}^\star). \tag{A3}$$

Since $\varphi_{\mathcal{D}}^\star$ is an extremum of the prior, the linear term vanishes. The empirical estimate of the negative hessian matrix is the observed Fisher information matrix which can be written as

$$\mathcal{I}(\varphi_{\mathcal{D}}^\star) = -\mathbb{E}_{\mathcal{D}\mid\mathcal{T}}[\mathbf{H}(\varphi_{\mathcal{D}}^\star)] = \mathbb{E}_{\mathcal{D}\mid\mathcal{T}}\left[ \left. \left( \left( \frac{\partial \log p(\varphi \mid \mathcal{D})}{\partial \varphi} \right) \left( \frac{\partial \log p(\varphi \mid \mathcal{D})}{\partial \varphi} \right)^T \right) \right|_{\varphi_{\mathcal{D}}^\star} \right]. \tag{A4}$$

The expectation is taken over the sample space $p(\mathcal{D} \mid \mathcal{T})$ since the network parameters are held fixed during sampling. In order to compute the Fisher score, we apply the Bayes rule to the prior to extract a loss function, which we take to be proportional to the training loss (equation (7)) and the $\chi^2$:

$$\log p(\varphi \mid (\mathbf{x}, \mathbf{y}) = \mathcal{D}) \propto -\mathcal{L}_\varphi(\mathbf{x}, \mathbf{y}) + \frac{1}{T}\sum_{t=1}^{T} \log p(\mathbf{y} \mid \hat{\mathbf{x}}^{(t)}) - \frac{\ell_2}{2}\|\varphi\|_2^2 \tag{A5}$$

We find in practice the $\ell_2$ term has little effect on the Fisher diagonal and our results. Thus, we set $\ell_2 = 0$.

Since the full Fisher matrix is intractable for a neural network, we approximate the quadratic term of the prior with the diagonal of the Fisher matrix following Kirkpatrick et al. (2016). For an optimisation problem, the first term of (A3) is constant. Thus, the posterior becomes proportional to

$$\log p(\varphi \mid \mathcal{D}, \mathcal{T}) \propto \log p(\mathcal{T} \mid \varphi) - \frac{\lambda}{2}\sum_j \mathrm{diag}(\mathcal{I}(\varphi_{\mathcal{D}}^\star))_j (\varphi_j - [\varphi_{\mathcal{D}}^\star]_j)^2. \tag{A6}$$

The Lagrange multiplier $\lambda$ is introduced to tune our uncertainty about the network parameters during fine-tuning.

## B. VAE ARCHITECTURE AND OPTIMISATION

For the following architectures, we employ the notion of *level* to mean layers in the encoder and the decoder with the same resolution. In each level, we place a block of convolutional layers before downsampling (encoder) or after upsampling (decoder). These operations are done with strided convolutions like in the U-net architecture of the RIM.

**Table 5.** Hyperparameters for the background source VAE.

| Parameter | Value |
|---|---|
| Input preprocessing | $\mathbb{1}$ |
| *Architecture* | |
| Levels (encoder and decoder) | 3 |
| Convolutional layer per level | 2 |
| Latent space dimension | 32 |
| Hidden Activations | Leaky ReLU |
| Output Activation | Sigmoid |
| Filters (first level) | 16 |
| Filters scaling factor (per level) | 2 |
| Number of parameters | 3 567 361 |
| *Optimization* | |
| Optimizer | Adam |
| Initial learning rate | $10^{-4}$ |
| Learning rate schedule | Exponential Decay |
| Decay rate | 0.5 |
| Decay steps | 30 000 |
| Number of steps | 500 000 |
| $\beta_{\mathrm{max}}$ | 0.1 |
| Batch size | 20 |

**Table 6.** Hyperparameters for the convergence VAE.

| Parameter | Value |
|---|---|
| Input preprocessing | $\log_{10}$ |
| *Architecture* | |
| Levels (encoder and decoder) | 4 |
| Convolutional layer per level | 1 |
| Latent space dimension | 16 |
| Hidden Activations | Leaky ReLU |
| Output Activation | $\mathbb{1}$ |
| Filters (first level) | 16 |
| Filters scaling factor (per level) | 2 |
| Number of parameters | 1 980 033 |
| *Optimization* | |
| Optimizer | Adam |
| Initial learning rate | $10^{-4}$ |
| Learning rate schedule | Exponential Decay |
| Decay rate | 0.7 |
| Decay steps | 20 000 |
| Number of steps | 155 000 |
| $\beta_{\mathrm{max}}$ | 0.2 |
| Batch size | 32 |

## C. RIM ARCHITECTURE AND OPTIMISATION

The notion of a link function $\Psi : \Xi \to \mathcal{X}$, introduced by Putzky & Welling (2017), is an invertible transformation between the network prediction space $\boldsymbol{\xi} \in \Xi$ and the forward modelling space $\mathbf{x} \in \mathcal{X}$. This is a different notion from preprocessing, discussed in section 3, because this transformation is applied inside the recurrent relation 6 as opposed to before training. In the case where the forward model has some restricted support or it is found that some transformation helps the training, then the link function chosen must be implemented as part of the network architecture as shown in the unrolled computational graph in Figure 12. Also, the loss $\mathcal{L}_{\varphi}$ must be computed in the $\Xi$ space in order to avoid gradient vanishing problems when $\Psi$ is a non-linear mapping, which happens if the non-linear link function is applied in an operation recorded for backpropagation through time (BPTT). For the convergence, we use an exponential link function with base 10: $\hat{\boldsymbol{\kappa}} = \Psi(\boldsymbol{\xi}) = 10^{\boldsymbol{\xi}}$. This $\Psi$ encodes the non-negativity of the convergence. Furthermore, it is a power transformation that leaves the linked pixel values $\boldsymbol{\xi}_i$ normally distributed, thus improving the learning through the non-linearities in the neural network.

The pixel weights $\mathbf{w}_i$ in the loss function (7) are chosen to encode the fact that the pixels with critical mass density ($\boldsymbol{\kappa}_i > 1$) have a stronger effect on the lensing configuration than other pixels. We find in practice that the weights

$$\mathbf{w}_i = \frac{\sqrt{\boldsymbol{\kappa}_i}}{\sum_i \boldsymbol{\kappa}_i}, \tag{C7}$$

encode this knowledge in the loss function and improve both the empirical risk and the goodness of fit of the baseline model on early test runs.

For the source, we found that we do not need a link function — the identity is generally better compared to other link functions we tried like sigmoid and power transforms — and we found that the pixel weights can be taken to be uniform, i.e. $\mathbf{w}_i = \frac{1}{M}$.
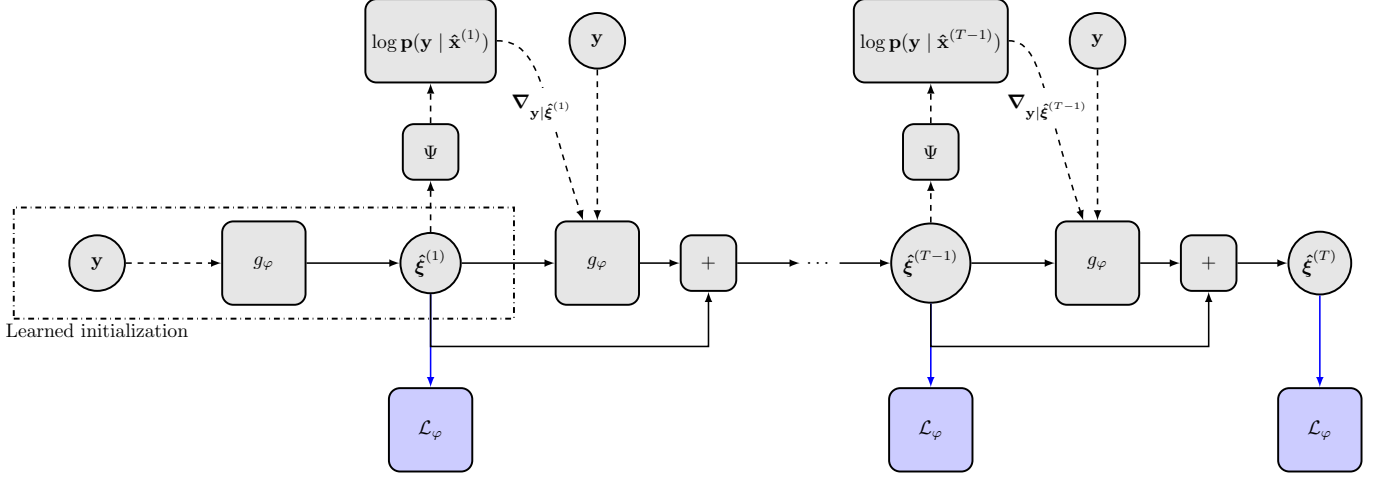


**Figure 12.** Unrolled computational graph of the RIM. Operations along solid arrows are being recorded for BPTT, while operations along dashed arrows are not. The blue arrows are only used for optimisation during training. During fine-tuning or testing, the loss is computed only as an oracle metric to validate that our methods can recover the ground truth.

**Table 7.** Hyperparameters for the RIM.

| Parameter | Value |
|---|---|
| Source link function | $\mathbb{1}$ |
| $\kappa$ link function | $10^{\boldsymbol{\xi}}$ |
| | |
| *Architecture* | Figure 2 |
| Recurrent steps $(T)$ | 8 |
| Number of parameters | 348 546 818 |
| | |
| *First Stage Optimisation* | |
| Optimizer | Adamax |
| Initial learning rate | $10^{-4}$ |
| Learning rate schedule | Exponential Decay |
| Decay rate | 0.95 |
| Decay steps | 100 000 |
| Number of steps | 610 000 |
| Batch size | 1 |
| | |
| *Second Stage Optimisation* | |
| Optimizer | Adamax |
| Initial learning rate | $6 \times 10^{-5}$ |
| Learning rate schedule | Exponential Decay |
| Decay rate | 0.9 |
| Decay steps | 100 000 |
| Number of steps | 870 000 |
| Batch size | 1 |

COSMOS   RIM+FT   IllustrisTNG   RIM+FT   Observation   RIM+FT   COSMOS   RIM+FT   IllustrisTNG   RIM+FT   Observation   RIM+FT
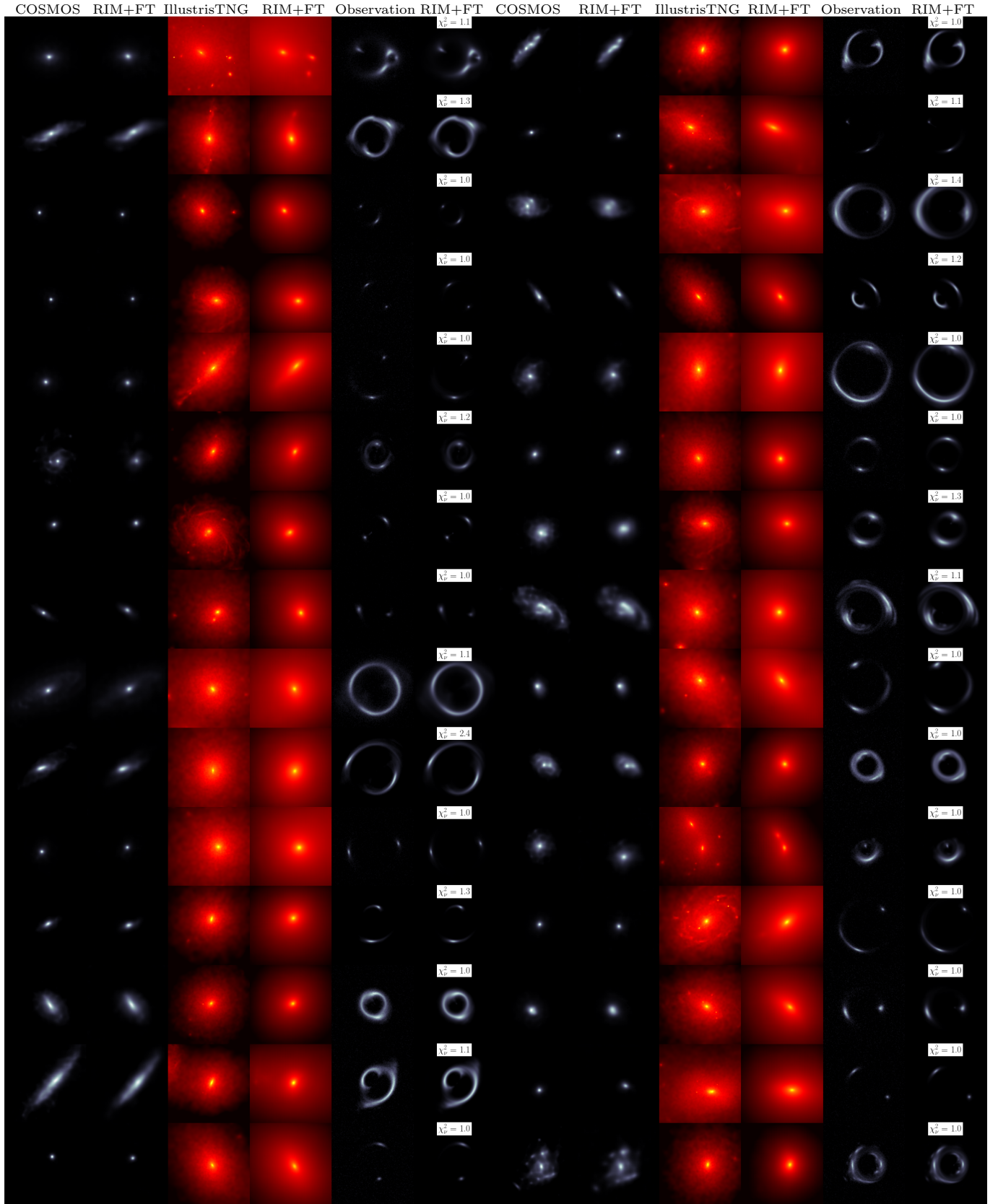
**Figure 13.** 30 reconstructions taken at random from the test set of 3000 examples simulated from COSMOS and IllustrisTNG data at high SNR. The colorscale are the same as in Figure 7.