
DeepGEM: Generalized Expectation-Maximization for Inverse Problems with Model Mismatch

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of jointly solving an inverse problem coupled with correcting
2 for model mismatch. Typically, inversion algorithms assume that a forward model,
3 which relates a source to its resulting measurements, is known and fixed. Using
4 collected indirect measurements and the forward model, the goal becomes to
5 recover the source. When the forward model is unknown, or imperfect, artifacts
6 due to model mismatch occur in the recovery of the source. We aim to recover
7 the desired source with imperfect knowledge of the forward model. In this paper
8 we propose DeepGEM, a variational Expectation-Maximization (EM) framework
9 that can be used to solve for the unknown parameters of the forward model in
10 an unsupervised manner. DeepGEM makes use of a normalizing flow generative
11 network to efficiently capture complex posterior distributions, which leads to more
12 accurate evaluation of the source’s posterior distribution used in EM. We showcase
13 the effectiveness of our DeepGEM approach by achieving strong performance
14 on the challenging problem of blind seismic tomography, where we significantly
15 outperform the standard method used in seismology. We also demonstrate the
16 generality of DeepGEM by applying it to blind deconvolution.

17 1 Introduction

18 Physics-based inversion methods typically recover an unknown source from indirect measurements
19 by assuming that the source and measurements are related via a known forward model [14, 7, 29]. For
20 example, non-blind deconvolution algorithms often assume that a measured blurry image is related to
21 its true sharp image via a known spatially-invariant blur kernel [9]; and traditional seismic inversion
22 methods assume that the spatially-varying velocity of the Earth’s interior is known *a priori* when
23 solving for an earthquake’s hypocenter [29]. However, these “known” forward models are generally
24 idealized and ignore intricacies of the systems that are either hard to capture or simply unknown.
25 Inversion algorithms with forward *model mismatch* result in biased reconstructions. For instance, bias
26 is regularly seen in non-blind deconvolution results, where reconstruction artifacts are often present
27 due to the use of an incorrect blur kernel [9].

28 Reducing model mismatch is key to reducing inversion bias and eliminating artifacts in recovered
29 sources. When considering how learning can help, a natural first idea might be to learn a direct map
30 from measurements to the desired source via supervised learning. However, such a “model-free”
31 approach is generally not practical, due to the lack of available ground truth training data. For
32 example, the blur kernel caused by handheld camera shake cannot be reproduced to get a training set
33 of sharp-blurry pairs to train a deconvolution approach. In seismic tomography, synthetic earthquakes
34 cannot be placed densely throughout the interior of the Earth to measure the ground response as a
35 function of hypocenter location. An alternative approach, which we adopt, is to develop unsupervised
36 methods that treat the true forward model as something that is unobserved and must be inferred.
37 An additional consideration is that we must solve for the true forward model from a single dataset,

38 without knowledge of the true source that produces the measurements – this problem setting is
 39 commonly referred to as *blind inversion*.

40 In this paper, we propose Deep Generalized Expectation-Maximization (DeepGEM) for solving
 41 inverse problems that are plagued by model mismatch. Using the indirect measurements as input,
 42 DeepGEM jointly estimates the source and forward model that together produce the observed
 43 measurements. DeepGEM is a variational inference based framework that makes use of deep learning
 44 machinery to easily capture and optimize complex probabilistic distributions that cannot be easily
 45 integrated in analytic Expectation-Maximization (EM) solutions. Our proposed framework is generic
 46 and can be applied to blind inversion problems described by differentiable forward models. In
 47 Section 4 we showcase the effectiveness of our DeepGEM approach on the challenging problem of
 48 blind seismic tomography, where we significantly outperform methods used in seismology. We also
 49 demonstrate the generality of DeepGEM by applying it to blind deconvolution in Section 5.

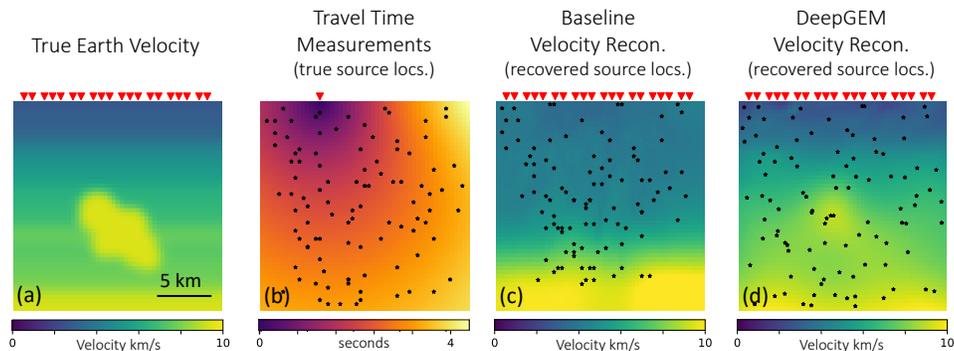


Figure 1: **DeepGEM applied to blind seismic tomography.** (a) A simulated cross section of the Earth’s interior (velocity structure), along with the locations of receivers on the surface (red triangles) that collect measurements. (b) The time it takes for a wave traveling from a source below the surface to reach the specified receiver is visualized for each location in the region of interest. The overlaid dots represent the true locations of simulated earthquakes and indicate the measured travel times that constrain optimization. (c) The subsurface velocity reconstruction obtained using a baseline approach optimized with the help of a seismologist. Note that the bright anomaly is missing from this reconstruction. Overlaid dots represent the inferred earthquake locations. (d) DeepGEM reconstructed subsurface velocity and inferred earthquake locations. Note that DeepGEM is able to accurately recover the gradient of the velocity field as well as partially recover the central anomaly.

50 2 Background and related work

51 The joint optimization of a forward model with source recovery is a very challenging ill-posed
 52 problem, leading to many possible solutions that are hard to disambiguate. For example, in blind
 53 deconvolution the blurry image observed can be equivalently explained by a sharp source image
 54 convolved with an extended kernel or a blurry source image convolved with a impulse kernel; to
 55 prefer one solution over another, additional information, such as image priors, must be considered.
 56 Previous work on inversion in poorly characterized systems (e.g., model mismatch) focused on
 57 limited contexts, such as spatially-invariant blind deconvolution [20, 21, 13, 10] and CT with a simple
 58 rotational error [33]. These methods tend to be highly specialized to each application domain, and
 59 cannot be easily generalized. In contrast, we propose a flexible model-based Bayesian framework
 60 that can be applied across multiple differentiable inversion problems subject to model mismatch.

61 2.1 Model-based Bayesian inversion

62 In model-based inversion, unobserved sources x and observed measurements y are related through a
 63 forward model: $y = f(x)$. When model mismatch is suspected, one can parameterize the forward
 64 model as $f_\theta(x)$ and then solve for the true model parameters θ^* . A common approach is to solve a
 65 *maximum a posteriori* (MAP) objective: either $\text{MAP}_{\theta,x}$ or MAP_θ .

66 $\text{MAP}_{\theta,x}$ solves for the optimal point estimate of the pair $\{\hat{\theta}, \hat{x}\}$ that maximizes a joint objective:

$$\{\hat{\theta}, \hat{x}\} = \arg \max_{\theta,x} [\log p(\theta, x|y)] = \arg \max_{\theta,x} [\log p(y|\theta, x) + \log p(\theta) + \log p(x)]. \quad (1)$$

67 In practice, formal probabilistic definitions of $p(\theta)$ and $p(x)$ are often unknown and replaced with
 68 regularization terms (e.g., total variation, sparsity) [21, 20]. Although $\text{MAP}_{\theta,x}$ provides a straightfor-
 69 ward approach to solve for θ , the energy landscape (1) is typically poorly behaved due to the ill-posed
 70 nature of the problem [21]. As a result, optimization is likely to get stuck in a (bad) local minimum.

71 **MAP $_{\theta}$** attempts to smooth the energy landscape by reducing the number of parameters that must be
 72 optimized. This is done by solving for parameters of the forward model, θ , that perform best under
 73 the full volume of possible x interpretations:

$$\hat{\theta} = \arg \max_{\theta} [\log p(\theta|y)] = \arg \max_{\theta} \left[\log \left(\int_x p(y|\theta, x)p(x)dx \right) + \log p(\theta) \right]. \quad (2)$$

74 Since this marginalization integral is often intractable, Expectation-Maximization (EM) algorithms
 75 have long been used for solving MAP_{θ} efficiently [11]. EM is an iterative algorithm that alternates
 76 between: “E”-Step) calculating the posterior of x conditioned on the current estimated forward model
 77 parameters $\theta^{(t-1)}$; and M-Step) updating θ to maximize the expected value of the log likelihood:

$$\theta^{(t)} = \arg \max_{\theta} \left[\mathbb{E}_{x \sim p(x|y, \theta^{(t-1)})} [\log p(y|\theta, x)] + \log p(\theta) \right]. \quad (3)$$

78 The advantage of MAP_{θ} over $\text{MAP}_{\theta,x}$ optimization for the blind deconvolution problem was de-
 79 scribed in [21] and its success demonstrated via EM optimization in [20]. However, it is important
 80 to note that evaluating the expectation in (3) over complex distributions is often intractable. For
 81 instance, the authors of [20] were forced to restrict the posterior distribution to a Gaussian distri-
 82 bution. Alternatively, stochastic EM methods [5, 8] bypasses the need to evaluate the expectation
 83 directly, approximating it by sampling the distribution. In this paper, we solve MAP_{θ} using complex
 84 distributions parameterized by deep neural networks.

85 **Forward model parameterization** Model-based inversion requires that the parametric form of
 86 $f_{\theta}(x)$ is well matched with the true forward model, which is not always known. Alternatively, neural
 87 networks can be used approximate the forward model, where θ represents the network weights. This
 88 setup is flexible, in that it can be used to approximate arbitrarily complex forward models, with
 89 the downside that it is often not interpretable when the parameters have no physical meaning. In
 90 this work, we develop and make use of interpretable, physically-motivated deep neural networks to
 91 parameterize $f_{\theta}(\cdot)$ for the problems of blind seismic tomography and blind deconvolution.

92 2.2 Deep variational distributions

93 We are interested in solving blind inverse problems using the EM algorithm to optimize the MAP_{θ}
 94 objective. Optimizing in this fashion requires the use of the posterior distribution $p(x|y, \theta^{(t-1)})$ in
 95 evaluating Eq. 3. As inverse problems are often ill-posed, we expect that the posterior distribution of
 96 source x conditioned on the forward model parameters θ is likely to be complex and perhaps even
 97 multi-modal. Therefore, it is important to be able to parameterize a flexible family of distributions to
 98 best estimate this conditional posterior.

99 Deep Probabilistic Imaging (DPI) [31] uses a normalizing flow-based generative model, $G_{\phi}(\cdot)$, to
 100 solve for the uncertainty of an unknown source x given a fixed forward model $f(x)$ and measurements
 101 y . DPI solves the variational objective:

$$\hat{\phi} = \arg \min_{\phi} \text{KL}(q_{\phi}(x)||p(x|y)) = \arg \min_{\phi} E_{x \sim q_{\phi}(x)} [-\log p(y|x) - \log p(x) + \log q_{\phi}(x)], \quad (4)$$

102 where $q_{\phi}(x)$ is the implicit distribution defined by $G_{\phi}(z)$ for $z \in \mathcal{R}^{|x|} \sim \mathcal{N}(0, 1)$ and $\log p(y|x) \propto$
 103 $\|y - f(x)\|_2 + c$ when there exists *i.i.d.* Gaussian noise on the measurements, y . After inference,
 104 $q_{\phi}(x)$ can be efficiently sampled by evaluating $G_{\phi}(z)$ for a Gaussian sample z .

105 The DPI variational objective is equivalent to the Variational Autoencoder [16, 27] objective, except
 106 with a fixed decoder, $f(x)$. In practice, vanilla VAEs constrain the posterior to be a Gaussian
 107 distribution, relying on the reparameterization trick for tractable optimization. Alternatively, DPI
 108 uses flow-based networks to efficiently sample and directly evaluate $q_{\phi}(x)$ [17, 12]. DPI’s use of a
 109 flow-based network allows for complicated and multi-modal posterior distributions constrained only
 110 by the space of possible invertible network architectures. Our proposed DeepGEM approach utilizes
 111 similar tools to model flexible distributions over x , while simultaneously learning the forward model
 112 parameters θ .

113 **3 Methods**

114 We propose a deep variational EM approach (DeepGEM) that optimizes the MAP_θ objective in Eq. 2
 115 to recover the parameters of a forward model $f_\theta(x)$ using only measurements y . Once learned, the
 116 updated forward model can then be used to estimate the posterior distribution of the unknown source,
 117 x . DeepGEM iterates between two stages that are inspired by the standard EM algorithm: (1) an
 118 ‘‘E’’-step that learns a variational distribution, $q_{\phi^{(t)}}(x)$, to approximate the posterior distribution of
 119 x given the current forward model parameters $\theta^{(t-1)}$, and (2) an M-step (refer to Eq. 3) that solves
 120 for $\theta^{(t)}$ that maximizes the expected value of the log likelihood function of θ , with respect to the
 121 posterior distribution $q_{\phi^{(t)}}(x)$ estimated in the prior ‘‘E’’-step. Each step is alternated and solved to
 122 convergence.

123 **3.1 ‘‘Expectation’’ step (‘‘E’’-step)**

124 The goal of DeepGEM’s ‘‘E’’-step is to solve for the posterior distribution $p(x|y, \theta^{(t-1)})$ that facilitates
 125 optimizing Eq. 3. Because this posterior distribution can be very complex, and even multi-modal, we
 126 propose to use a flexible variational approach to learn the parameters ϕ of an approximate posterior
 127 distribution $q_\phi(x)$. The variational distribution $q_\phi(x)$ can then be used to evaluate Eq. 3 via efficient
 128 sampling.

129 Using DPI, we solve for a flexible variational distribution, $q_\phi(x)$ that well approximates the posterior
 130 distribution $p(x|y, \theta^{(t-1)})$. DPI uses a normalizing flow network, $G_\phi(z)$, with input $z \in \mathbb{R}^{|x|}$, where
 131 $x = G_\phi(z) \sim q_\phi(x)$ when $z \sim \mathcal{N}(0, \mathbb{1})$. Normalizing flow networks allow for exact computation of
 132 the log-likelihood $\log q_\phi(x)$ needed to solve

$$\begin{aligned} \phi^{(t+1)} &= \arg \min_{\phi} \text{KL}(q_\phi(x) || p(x|y, \theta^{(t-1)})) \\ &\approx \arg \min_{\phi} \frac{1}{N} \sum_{n=1}^N [-\log p(y|\theta^{(t-1)}, x_n) - \log p(x_n) + \log q_\phi(x_n)] \\ &\quad \text{for } x_n = G_\phi(z_n), \quad z_n \sim \mathcal{N}(0, \mathbf{1}), \end{aligned} \quad (5)$$

133 (as derived from Eq. 4) where $\log p(x)$ is a prior on the source and $\log p(y|x, \theta^{(t)})$ is the data
 134 likelihood. When assuming the measurements y experience *i.i.d* additive Gaussian noise with
 135 standard deviation σ_y , $\log p(y|\theta^{(t)}, x_n) = \frac{1}{2\sigma_y^2} \|y - f_{\theta^{(t)}}(x_n)\|^2 + c$.

136 **3.2 Maximization step (M-step)**

137 The goal of DeepGEM’s M-step is to use the parameterized approximate posterior distribution,
 138 $q_{\phi^{(t)}}(x)$, from the prior ‘‘E’’-step to update θ , the parameters of the unknown forward model $f_\theta(\cdot)$.
 139 This is achieved by sampling from the learned normalizing flow network, $G_{\phi^{(t)}}(\cdot)$, to stochastically
 140 solve Eq. 3 :

$$\theta^{(t)} \approx \arg \max_{\theta} \left[\frac{1}{N} \sum_{n=1}^N [\log p(y|\theta, x_n)] + \log p(\theta) \right] \quad \text{for } x_n = G_{\phi^{(t)}}(z_n), \quad z_n \sim \mathcal{N}(0, \mathbb{1}), \quad (6)$$

141 where $p(\theta)$ is a prior on the forward model. This prior can be used to encourage the forward model
 142 parameters to remain close to an initial model $\tilde{\theta}$ by defining $\log p(\theta) \propto \|\theta - \tilde{\theta}\|_2 + c$.

143 **4 Case study: blind seismic tomography**

144 Two fundamental seismic inverse problems are spatially localizing an earthquake’s hypocenter (also
 145 referred to as source localization) and tomographic reconstruction of the Earth’s subsurface [29].
 146 These problems are interconnected: source localization relies on knowing how fast waves propagate
 147 through different regions of the Earth’s interior, referred to as the Earth’s *velocity* structure. However,
 148 in standard seismological practice, due to difficulties in solving these problems jointly, they are
 149 generally treated separately: source localization is performed initially using oversimplified velocity
 150 models [29], and then the tomography problem is performed by taking those best-fitting hypocenters

151 as ground truth [25]. This approach typically results in the need to over-regularize the inverse
 152 problem by smoothing out high-frequency information [2], and can only be improved by carefully
 153 incorporating other forms of information such as waveform-derived quantities [15, 22, 4, 23] or
 154 by performing costly experiments such as controlled explosions. In contrast, we demonstrate our
 155 DeepGEM approach on blind seismic tomography, solving for the subsurface velocity (parameterized
 156 by θ) when the source hypocenters, x , are unknown. Measurements, y , used to constrain the inverse
 157 problem are the time it takes for the first wave to propagate from its source to a receiver on the Earth’s
 158 surface, referred to as a *travel time* measurement (refer to Fig. 1).

159 4.1 Seismic tomography background

160 **Physics of earthquake source localization:** The earthquake source location, also called the hypocen-
 161 ter, is the location where the earthquake nucleates [29]. The source location can be triangulated using
 162 travel times from multiple receivers. However, when there are very few receivers (< 3 in 3D, < 2 in
 163 2D), the source localization is ill-posed and there exists multiple equally optimal solutions.

164 **Physics of travel time tomography:** Travel time tomography is a method for reconstructing the
 165 velocity structure using the absolute arrival times of earthquake waves from the earthquake to the
 166 receiver [29]. Given perfect knowledge of the earthquake locations x and receivers r at every
 167 position in the ground, the exact velocity can be computed by solving the Eikonal equation $V(r) =$
 168 $\|\nabla_r T(x, r)\|_2^{-1}$, where $T(x, r)$ is the travel time from an earthquake to a receiver. The solution to
 169 the Eikonal equation is physically sound; however, seismologists often use simplifications, such as
 170 straight ray tomography, for efficiency [2].

171 **Deep Learning for travel time tomography:** EikoNet [30] implicitly solves the Eikonal equation
 172 [32] by learning a mapping from a source-receiver pair (x, r) to its associated travel time, learning θ
 173 such that $f_\theta(x, r) \approx T(x, r)$ where $T(x, r)$ is the true travel time. The velocity structure can then be
 174 extracted from the learned EikoNet by solving the Eikonal equation through automatic differentiation.
 175 In particular, the computed velocity $V(s)$ at location s is $\|\nabla_s f_\theta(x, s)\|_2^{-1}$. Note that this computed
 176 velocity depends on the source location x . Ideally, the velocity should be invariant to the source
 177 location; however, in practice, this is only true when EikoNet is trained with densely-sampled (x, r)
 178 pairs. Additionally, the travel time from a source to a receiver, $T(x, r)$, should be identical to the
 179 travel time from a receiver to a source, $T(r, x)$, but remains unconstrained by EikoNet.

180 4.2 DeepGEM setup for blind seismic tomography

181 For blind seismic tomography, we parameterize the forward model using a modification of EikoNet,
 182 $f_\theta(x, r)$, with unknown source location x and known receiver location r as inputs and $y \approx T(x, r)$,
 183 the absolute travel time between x and r , as output. In order to solve Eqs. 5 and 6 for an updated
 184 forward model, we must define priors $p(x)$ and $p(\theta)$. The prior over source locations, $p(x)$, is often
 185 well defined, typically a Gaussian distribution $\mathcal{N}(\bar{x}, \sigma_x)$ with a standard deviation of $\sigma_x \sim 2$ km.

186 We construct a prior over the forward model, $p(\theta)$, that encourages EikoNet to learn a velocity close
 187 to $\tilde{V}(s)$. Additionally, as discussed above, there are constraints specific to seismic tomography that
 188 are not explicitly enforced through EikoNet’s architecture: (1) velocity reconstruction invariance
 189 with respect to the source location, and (2) travel time symmetry between sources and receivers. We

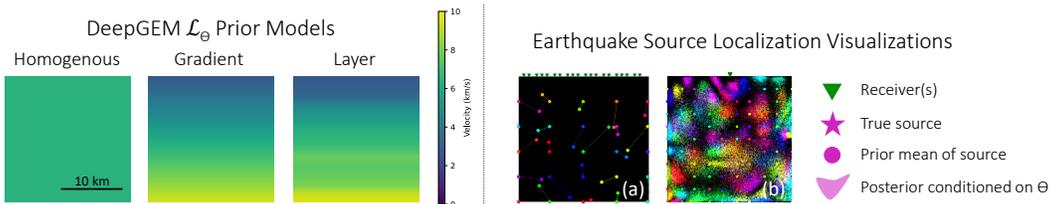


Figure 2: (Left) A visualization of the homogeneous, gradient, and layer models used in $p(\theta)$ ’s \mathcal{L}_θ term. (Right) Visualizations used to describe source configurations and earthquake source posteriors. (a) Visualizations of the true and initialized source locations are plotted as stars (true) connected to circles, which indicate the expected source locations according to $p(x)$. Note that the expected source locations deviate significantly from the true locations. (b) Visualizations of the learned posterior distribution, $q_\phi(x|y, \theta)$, for each source are plotted as colored histograms and overlaid with stars (of the same color) indicating the true source location.

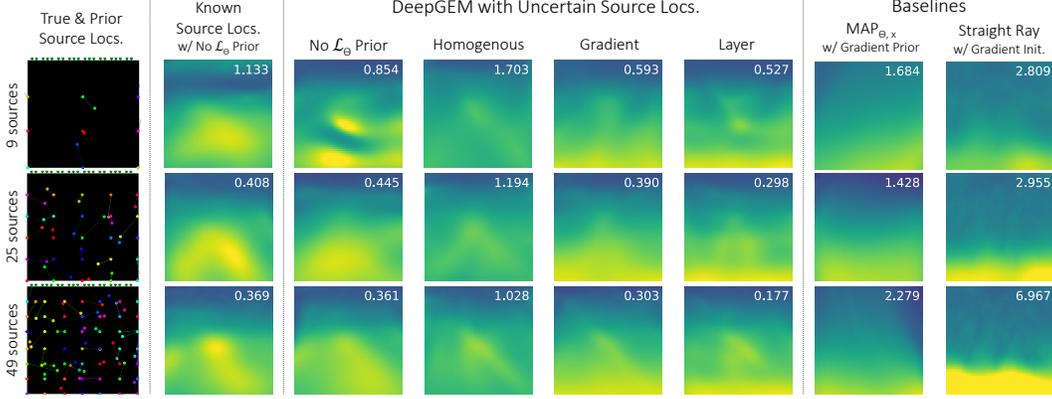


Figure 3: **DeepGEM reconstructions significantly outperform baselines, and improve with more sources and better \mathcal{L}_θ prior models.** Each row corresponds to the reconstructions obtained from a single noisy observation of travel times, where the true Earth velocity is shown in Fig. 1(a). Results shown are simulated using 20 surface receivers and a varying number of sources (9, 25, and 49) in a uniform grid. Columns 3-6 show DeepGEM results obtained using different \mathcal{L}_θ priors. Note that results improve as the \mathcal{L}_θ prior becomes closer to the true velocity structure, and as the number of sources increases. As a reference, column 2 shows the velocity reconstruction obtained using DeepGEM under fixed, perfectly known source locations. Columns 7-8 show results obtained by the baseline approaches. The velocity reconstruction MSE is included in the top right of each reconstruction. DeepGEM substantially outperforms both straight ray and $\text{MAP}_{\theta,x}$ baselines.

190 augment the prior $p(\theta)$ to include these constraints, implemented as soft constraints:

$$\log p(\theta) = \lambda_\theta \underbrace{\sum_{\substack{r \in \mathcal{R}, \\ s \in \mathcal{S}}} \|\tilde{V}(s) - V_r(s)\|_2}_{\mathcal{L}_\theta} + \lambda_V \underbrace{\sum_{\substack{r_i, r_j \in \mathcal{R}, \\ s \in \mathcal{S}}} \|V_{r_i}(s) - V_{r_j}(s)\|_2}_{\mathcal{L}_V} + \lambda_T \underbrace{\sum_{\substack{r \in \mathcal{R}, \\ s \in \mathcal{S}}} \|T(s, r) - T(r, s)\|_2}_{\mathcal{L}_T}. \quad (7)$$

191 The velocity constraint is represented through \mathcal{L}_V and travel time constraint through \mathcal{L}_T , with
 192 corresponding hyperparameters λ_V and λ_T . \mathcal{S} is a set of points sampled uniformly from the region
 193 of interest, \mathcal{R} is the set of all receiver locations, and $V_r(s) = \|\nabla_s f_\theta(r, s)\|_2^{-1}$.

194 **Implementation details:** In our experiments we define a realistic prior over unknown source
 195 locations as $p(x) = \mathcal{N}(\bar{x}, \sigma_x)$, where $\bar{x} \sim \mathcal{N}(x, \sigma_x)$ and $\sigma_x = 2$ km. We assume the measurements
 196 y are sampled from a Gaussian distribution with mean \tilde{y} – true travel times computed using the
 197 package `eikonal_fm` [28, 32] – and a realistic standard deviation of $\sigma_y = 0.2$ seconds. To simulate
 198 real world passive tomography, we assume receivers are located only along the surface (the top edge
 199 of the 2D image) of the region of interest, which is 20 km \times 20 km, and sources can be anywhere
 200 within the region of interest.

201 The posterior distribution of the source locations, x , is estimated using a Real-NVP network $G_\phi(\cdot)$
 202 with 16 affine coupling layers. An updated EikoNet (described in the supplemental material) has been
 203 modified to parameterize $f_\theta(x)$. This EikoNet is pretrained with samples from the prior $p(x)$ as input
 204 and the simulated travel time measurements as output. We use Adam as the optimizer [34] with a
 205 batch size of 32 and an E-step learning rate of 1e-3 and M-step learning rate of 5e-5. Hyperparameters
 206 $(\lambda_T, \lambda_V, \lambda_\theta)$ were empirically chosen by inspecting the loss on a grid search over hyperparameters.
 207 Results presented have been run with 10 EM iterations, each with 800 “E”-step epochs and 2000
 208 M-step epochs. Each DeepGEM model takes ~ 6 hours on a NVIDIA Quatro RTX 5000, for a total
 209 of ~ 500 hours of development time.

210 4.3 Results

211 **Comparison to Baseline Approaches:** We compare results from DeepGEM to results obtained
 212 using a baseline run by a seismologist. This iterative baseline alternates between source inversion
 213 and straight ray tomography, and is the standard approach used for blind tomography. Further detail
 214 on this baseline is provided in the supplemental material. The gradient model shown in Fig. 2 is used
 215 to perform the initial source inversion for this baseline; nonetheless, we find that the solution quickly
 216 diverges. Therefore, we rely on the expertise of the domain expert to decide when to terminate the

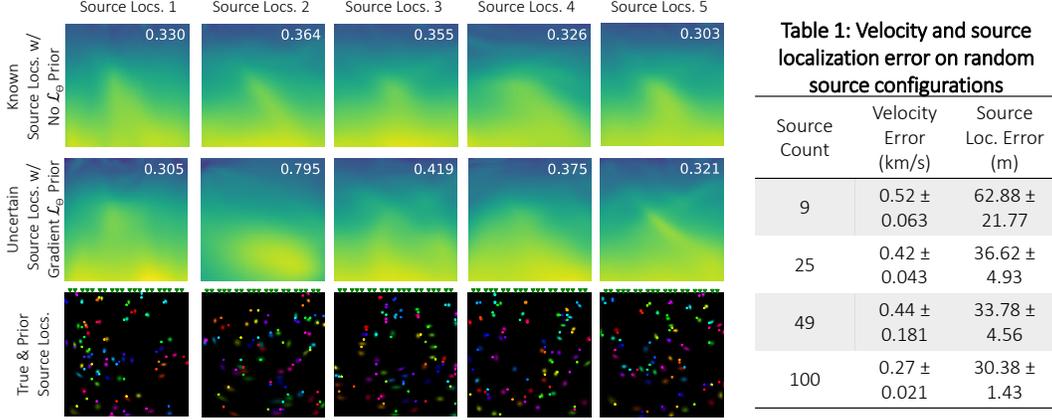


Figure 4: **DeepGEM consistently recovers prominent features across various source configurations.** Row 2 shows DeepGEM results obtained using different random source configurations, where the true Earth velocity is shown in Fig. 1(a). As a reference, row 1 shows the velocity reconstruction obtained using DeepGEM under fixed, perfectly known source locations. As can be seen by the table, both velocity and localization error decrease with an increasing number of sources in the region of interest. Each mean and standard deviation is computed using 5 random realizations of true source configurations.

217 optimization. As seen in Fig. 3, DeepGEM consistently outperforms this human-in-the-loop baseline
 218 across all source configurations. In Fig. 3 we also compare with a $\text{MAP}_{\theta,x}$ solution. $\text{MAP}_{\theta,x}$ is
 219 consistently outperformed by the DeepGEM MAP_{θ} approach across all source configurations.

220 **Sensitivity to \mathcal{L}_{θ} Prior Choice:** We evaluate DeepGEM’s recovery of the true velocity structure
 221 shown in Fig. 1 using one of three different \mathcal{L}_{θ} priors shown in Fig. 2: homogeneous, gradient, and
 222 layer, as well as $\mathcal{L}_{\theta} = 0$. The homogeneous model takes on value of 6.419 km/s, the average velocity
 223 value of the true velocity structure. The gradient captures the increasing velocity as depth increases,
 224 and the layer model represents the true model without the added anomaly. As shown in Fig. 3, the
 225 mean squared error tends to decrease with the gradient and layer \mathcal{L}_{θ} prior models, which are closer to
 226 the true velocity structure.

227 **Sensitivity to Source Configuration:** We evaluate results with sources that are both uniformly
 228 and randomly spaced throughout the region of interest. Improved performance is expected when
 229 the number of sources is increased and/or when sources are well distributed. To better understand
 230 DeepGEM’s performance we introduce an ablation, where the velocity structure is learned by training
 231 $f_{\theta}(\cdot)$ with access to perfect source locations x (i.e., $p(x)$ is a delta function). As is shown in Fig 3,
 232 even when training with true source locations, the anomaly is not well resolved until ~ 49 sources.

233 Fig. 3 shows one realization of DeepGEM results obtained from different counts of uniformly spaced
 234 sources. As expected, the MSE between the reconstructed and true velocity structure tends to
 235 decrease as the number of sources increases. Fig. 4 shows results obtained using five randomly
 236 generated source configurations, with 49 sources each. These results demonstrate that, although
 237 the reconstructed velocity structure is somewhat sensitive to the underlying source configuration,
 238 the primary features of the true velocity can still be recovered in all cases. Velocity and source
 239 localization error obtained for different random configurations are shown in the accompanying table.

240 **Sensitivity to Number of Receivers:** In the case of a single receiver, there exists an entire ring of
 241 source locations that result in the same travel time measurement. Fig. 5 contains results from this
 242 challenging one-receiver setting using DeepGEM with/without measurement noise and with/without
 243 known source locations. Since there is only one receiver, the velocity model is able to easily
 244 overfit. However, perhaps surprisingly, artifacts are more severe when source locations are perfectly
 245 constrained. These artifacts are caused by the velocity model overfitting to noise in the travel time
 246 measurements, and are substantially reduced when noise-free travel time measurements are used.
 247 Note that the recovered source location posterior $q_{\phi}(x)$ obtained by the “E”-step is non-Gaussian.

248 **Sensitivity to Velocity Structure:** In Fig. 6, DeepGEM is tested on randomly generated velocity
 249 fields, each generated from a Gaussian random field (GRF) described in the supplemental material.
 250 These results show that DeepGEM works well at recovering the primary features across a variety

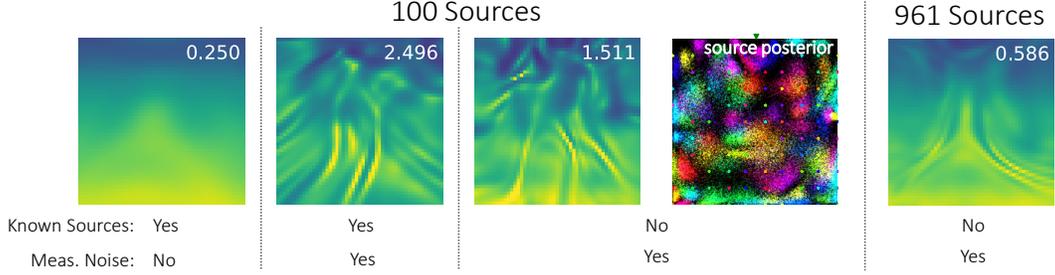


Figure 5: **DeepGEM recovery with a single receiver.** Velocity reconstructions shown in columns 3 and 5 demonstrate that DeepGEM is able to learn some of the true velocity features (see Fig. 1(a)), even when limited to measurements from a single receiver. However, these reconstructions show clear signs of overfitting to the measurement data. This is demonstrated by observing that the DeepGEM reconstruction with perfectly known source locations is significantly better with a single receiver when no measurement error is included on the measurements (comparing columns 1 and 2).

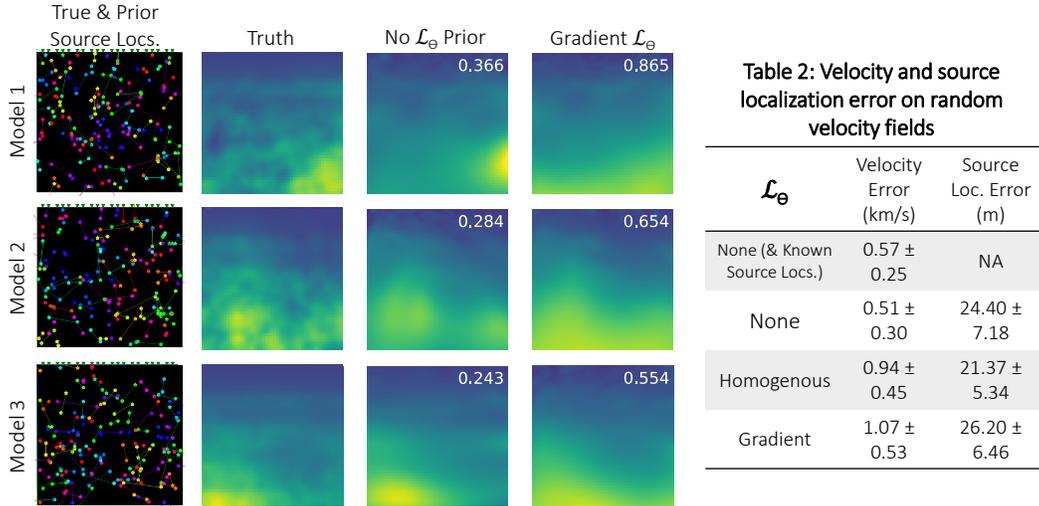


Figure 6: **Performance of DeepGEM recovery on random velocity fields.** Ten random velocity fields were drawn from a GRF-based distribution and used to simulate travel time measurements with 20 receivers and 100 randomly placed sources. (left) Reconstructions obtained for 3 of these configurations are shown. (right) A table lists the mean and standard deviation of velocity and source error obtained across the ten models, each recovered using different \mathcal{L}_θ priors.

251 of different velocity structures. As compared to when $\mathcal{L}_\theta = 0$, the \mathcal{L}_θ gradient prior biases the
 252 reconstruction towards smoother velocity structures. The accompanying table contains error statistics
 253 for ten different randomly generated velocity fields.

254 5 Case study: blind deconvolution

255 We apply DeepGEM to the problem of blind deconvolution to further demonstrate the generality of
 256 our approach. Blind deconvolution is a classic ill-posed imaging problem that aims to reconstruct
 257 a sharp image from a blurry image with an unknown PSF [18, 13, 24, 20, 21, 19, 6, 1, 26]. Blurry
 258 images, caused by handheld camera shake, can be modeled using a single spatially-invariant blur
 259 kernel:

$$y = x * k_\theta + \varepsilon \text{ for } \varepsilon \sim \mathcal{N}(0, \sigma), \quad (8)$$

260 where y is the blurry image, $*$ represents a 2D convolution, x is the true sharp image, k_θ is the
 261 spatially invariant blur kernel, and ε is additive Gaussian noise.

262 5.1 DeepGEM setup for blind deconvolution

263 For blind deconvolution, we parameterize the forward model $f_\theta(\cdot)$ using a deep network consisting of
 264 multiple convolution layers without non-linear activation, as proposed in [3]. Multiple convolutional
 265 layers without activation simply overparameterizes a linear blur kernel, which has been empirically
 266 shown to produce multiple good minima that are easier to converge to. To ensure the blur kernel is

267 non-negative and volume preserving, we use a Softmax layer and normalize the kernel. For an n
 268 layer network with weights θ_i for $i = 1, \dots, n$, the resulting parameterized forward model is:

$$f_{\theta}(x) = x * \hat{k}_{\theta} = x * \left[\frac{\text{Softmax}(\theta_1 * \theta_2 * \dots * \theta_n)}{\|\text{Softmax}(\theta_1 * \theta_2 * \dots * \theta_n)\|_1} \right] \quad (9)$$

269 **Implementation details:** We demonstrate DeepGEM using simple Total Variation (TV) regulariza-
 270 tion in place of $\log p(x)$. We assume a Gaussian prior on the noise on the blurry measurements where
 271 $\varepsilon \sim \mathcal{N}(0, 0.01)$ as well as a sparsity prior on the reconstructed kernel through an $\ell_{0.8}$ soft constraint.
 272 The posterior distribution of the sharp image, x , is estimated using a Real-NVP network $G_{\phi}(\cdot)$ with
 273 16 affine coupling layers. We use 5 convolution layers to parameterize k_{θ} . We use Adam as the
 274 optimizer [34] with a batch size of 64 and an E-step learning rate of $5e-4$ and M-step learning rate of
 275 $1e-4$. Hyperparameters, weights used for sparsity and TV priors, were empirically chosen by a grid
 276 search over hyperparameters. Results presented have been run with 10 EM iterations, each with 400
 277 “E”-step epochs and 400 M-step epochs, which takes ~ 1 hour on a NVIDIA Tesla V100.

278 5.2 Results

279 In Fig. 7 we show results from DeepGEM on two different blurry images [21]. The blurry image
 280 in Fig. 7(a) exhibits artifacts from the blur kernel in the form of repeated features. This is due to
 281 the blur kernel being similar to the sum of two delta point spread functions. The recovered source
 282 image is much sharper, containing clear text with minimal ringing artifacts. The reconstructed kernel
 283 roughly matches the true kernel’s shape with two lobes along the same diagonal. In Fig. 7(b), ringing
 284 in the blurry image is also reduced. Although the reconstructed kernel is not located spatially at the
 285 same location as the true kernel, this does not significantly harm reconstruction; since the kernels
 286 are shift-invariant, the reconstructed image and learned kernel can both be shifted such that they
 287 reconstruct the same blurry image. Please refer to the supplemental material for more results.

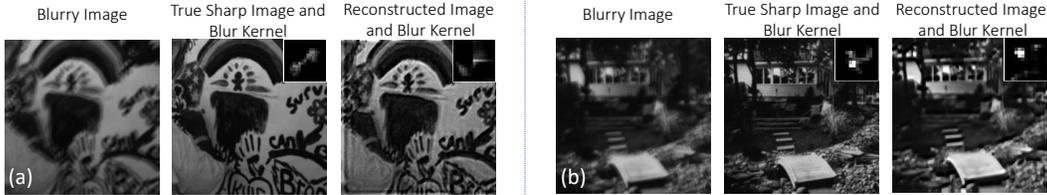


Figure 7: Blurry measured images (columns 1 and 4) generated using the true sharp image and blur kernel (shown in columns 2 and 5). Reconstructed sharp image and the corresponding inferred blur kernel from DeepGEM (columns 3 and 6).

288 6 Conclusions

289 In this paper we present DeepGEM, a deep probabilistic framework for tackling blind inverse
 290 problems through estimation of the forward model. DeepGEM achieves strong performance in the
 291 task of joint seismic tomography and earthquake source localization, substantially outperforming
 292 standard approaches currently being used in seismology on synthetic data. The proposed framework
 293 is flexible and can be applied to different applications that require estimation or fine tuning of forward
 294 model parameters. We demonstrate this flexibility by also applying the approach to the simple, but
 295 challenging, blind deconvolution problem. Future work includes applying this method to real seismic
 296 data, extending to other applications, and incorporating data-driven priors. Our results highlight the
 297 benefit of blending physically sound model-based techniques with learning machinery for inversion
 298 with model mismatch.

299 **Broader Impacts.** DeepGEM can be used to solve for a system’s model mismatch, which can then
 300 help improve our understanding of complex physical systems. However, this tool is not trustworthy
 301 enough for safety-critical systems. Nonetheless, this approach can benefit society through a better
 302 understanding of fundamental science and advanced earthquake prediction models via seismic
 303 imaging.

304 **References**

- 305 [1] Raied Aljadaan, Dipan K Pal, and Marios Savvides. Douglas-rachford networks: Learning both the image
 306 prior and data fidelity terms for blind image deconvolution. In *Proceedings of the IEEE/CVF Conference*
 307 *on Computer Vision and Pattern Recognition*, pages 10235–10244, 2019.
- 308 [2] M. P. Barmin, M. H. Ritzwoller, and A. L. Levshin. *A Fast and Reliable Method for Surface Wave*
 309 *Tomography*, pages 1351–1375. Birkhäuser Basel, Basel, 2001.
- 310 [3] Sefi Bell-Kligler, Assaf Shocher, and Michal Irani. Blind super-resolution kernel estimation using an
 311 internal-gan. *CoRR*, abs/1909.06581, 2019.
- 312 [4] Ebru Bozdağ, Jeannot Trampert, and Jeroen Tromp. Misfit functions for full waveform inversion based on
 313 instantaneous phase and envelope measurements. *Geophysical Journal International*, 185(2):845–870,
 314 2011.
- 315 [5] Olivier Cappé and Eric Moulines. On-line expectation-maximization algorithm for latent data models.
 316 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):593–613, Jun 2009.
- 317 [6] Ayan Chakrabarti. A neural approach to blind motion deblurring. In *European conference on computer*
 318 *vision*, pages 221–235. Springer, 2016.
- 319 [7] Chengqian Che, Fujun Luan, Shuang Zhao, Kavita Bala, and Ioannis Gkioulekas. Inverse transport
 320 networks. *arXiv preprint arXiv:1809.10820*, 2018.
- 321 [8] Jianfei Chen, Jun Zhu, Yee Whye Teh, and Tong Zhang. Stochastic expectation maximization with variance
 322 reduction. In *NeurIPS*, pages 7978–7988, 2018.
- 323 [9] S. Cho, Jue Wang, and S. Lee. Handling outliers in non-blind image deconvolution. In *2011 International*
 324 *Conference on Computer Vision*, pages 495–502, 2011.
- 325 [10] Sunghyun Cho and Seungyong Lee. Fast motion deblurring. In *ACM SIGGRAPH Asia 2009 papers*, pages
 326 1–8. 2009.
- 327 [11] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the
 328 em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- 329 [12] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp, 2017.
- 330 [13] Rob Fergus, Barun Singh, Aaron Hertzmann, Sam T Roweis, and William T Freeman. Removing camera
 331 shake from a single photograph. In *ACM SIGGRAPH 2006 Papers*, pages 787–794. 2006.
- 332 [14] Jeffrey A Fessler. Model-based image reconstruction for mri. *IEEE signal processing magazine*, 27(4):81–
 333 89, 2010.
- 334 [15] Andreas Fichtner, Brian LN Kennett, Heiner Igel, and Hans-Peter Bunge. Theoretical background for
 335 continental-and global-scale full-waveform inversion in the time–frequency domain. *Geophysical Journal*
 336 *International*, 175(2):665–685, 2008.
- 337 [16] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*,
 338 2013.
- 339 [17] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In
 340 *Advances in neural information processing systems*, pages 10215–10224, 2018.
- 341 [18] D. Kundur and D. Hatzinakos. Blind image deconvolution. *IEEE Signal Processing Magazine*, 13(3):43–64,
 342 1996.
- 343 [19] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiří Matas. Deblurgan:
 344 Blind motion deblurring using conditional adversarial networks. In *Proceedings of the IEEE conference on*
 345 *computer vision and pattern recognition*, pages 8183–8192, 2018.
- 346 [20] A. Levin, Y. Weiss, F. Durand, and W. T. Freeman. Efficient marginal likelihood optimization in blind
 347 deconvolution. In *CVPR 2011*, pages 2657–2664, 2011.
- 348 [21] Anat Levin, Yair Weiss, Fredo Durand, and William T Freeman. Understanding and evaluating blind
 349 deconvolution algorithms. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages
 350 1964–1971. IEEE, 2009.
- 351 [22] Jingrui Luo and Ru-Shan Wu. Seismic envelope inversion: reduction of local minima and noise resistance.
 352 *Geophysical Prospecting*, 63(3):597–614, 2015.
- 353 [23] Yi Luo and Gerard T Schuster. Wave-equation travelttime inversion. *Geophysics*, 56(5):645–653, 1991.
- 354 [24] D. Perrone and P. Favaro. Total variation blind deconvolution: The devil is in the details. In *2014 IEEE*
 355 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2909–2916, Los Alamitos, CA,
 356 USA, jun 2014. IEEE Computer Society.
- 357 [25] Nicholas Rawlinson, S Pozgay, and S Fishwick. Seismic tomography: a window into deep earth. *Physics*
 358 *of the Earth and Planetary Interiors*, 178(3-4):101–135, 2010.

- 359 [26] Dongwei Ren, Kai Zhang, Qilong Wang, Qinghua Hu, and Wangmeng Zuo. Neural blind deconvolution using deep priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3341–3350, 2020.
360
361
- 362 [27] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.
363
364
- 365 [28] J. A. Sethian. Fast marching methods. *SIAM Review*, 41(2):199–235, 1999.
- 366 [29] Peter M Shearer. *Introduction to seismology*. Cambridge university press, 2019.
- 367 [30] Jonathan D. Smith, Kamyar Azizzadenesheli, and Zachary E. Ross. Eikonet: Solving the eikonal equation with deep neural networks, 2020.
368
- 369 [31] He Sun and Katherine L. Bouman. Deep probabilistic imaging: Uncertainty quantification and multi-modal solution characterization for computational imaging, 2020.
370
- 371 [32] Eran Treister and Eldad Haber. A fast marching algorithm for the factored eikonal equation. *Journal of Computational Physics*, 324:210–225, 11 2016.
372
- 373 [33] Mingyang Xie, Yu Sun, Jiaming Liu, Brendt Wohlberg, and Ulugbek S. Kamilov. Joint reconstruction and calibration using regularization by denoising, 2020.
374
- 375 [34] Aston Zhang, Zachary C. Lipton, Mu Li, and Alexander J. Smola. *Dive into Deep Learning*. 2020. <https://d21.ai>.
376

377 Checklist

- 378 1. For all authors...
- 379 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
380 contributions and scope? [Yes]
- 381 (b) Did you describe the limitations of your work? [Yes]
- 382 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 383 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
384 them? [Yes]
- 385 2. If you are including theoretical results...
- 386 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 387 (b) Did you include complete proofs of all theoretical results? [N/A]
- 388 3. If you ran experiments...
- 389 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
390 mental results (either in the supplemental material or as a URL)? [Yes]
- 391 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
392 were chosen)? [Yes]
- 393 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
394 ments multiple times)? [Yes]
- 395 (d) Did you include the total amount of compute and the type of resources used (e.g., type
396 of GPUs, internal cluster, or cloud provider)? [Yes]
- 397 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 398 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 399 (b) Did you mention the license of the assets? [N/A]
- 400 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 401 (d) Did you discuss whether and how consent was obtained from people whose data you’re
402 using/curating? [N/A]
- 403 (e) Did you discuss whether the data you are using/curating contains personally identifiable
404 information or offensive content? [N/A]
- 405 5. If you used crowdsourcing or conducted research with human subjects...
- 406 (a) Did you include the full text of instructions given to participants and screenshots, if
407 applicable? [N/A]

408
409
410
411

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]