

---

# GAN-based Priors for Quantifying Uncertainty

---

**Dhruv V. Patel**

University of Southern California  
Los Angeles, CA, USA  
dhruvvp@usc.edu

**Assad A. Oberai**

University of Southern California  
Los Angeles, CA, USA  
aoberai@usc.edu

## Abstract

Bayesian inference is used extensively to quantify the uncertainty in an inferred field given the measurement of a related field when the two are linked by a mathematical model. Despite its many applications, Bayesian inference faces challenges when inferring fields that have discrete representations of large dimension, and/or have prior distributions that are difficult to characterize mathematically. In this work we demonstrate how the approximate distribution learned by a deep generative adversarial network (GAN) may be used as a prior in a Bayesian update to address both these challenges. We demonstrate the efficacy of this approach on two distinct, and remarkably broad, classes of problems. The first class leads to supervised learning algorithms for image classification with superior out of distribution detection and accuracy, and for image inpainting with built-in variance estimation. The second class leads to unsupervised learning algorithms for image denoising and for solving physics-driven inverse problems.

Bayesian inference provides a principled approach for quantifying uncertainty. As shown in the following section, it treats the inferred vector as a multivariate stochastic vector and leads to an expression for its distribution. This expression can be used to estimate the most likely solution (the maximum a-posteriori estimate, or the MAP), the mean, the variance, or any other population parameter of interest. Thus providing a recipe for thoroughly quantifying the uncertainty in an inference problem. For the image recovery problems considered in this paper, Bayesian inference not only provides the best guess of the true image, but also a means to estimate measures of uncertainty such as the pixel-wise variance or the likelihood of an out-of-distribution input.

The knowledge of uncertainty in a prediction can directly influence the downstream action that depends on the inference. Consider an image recovery problem where two distinct inputs lead to similar recovered images: those of a traffic sign with a high speed limit. However, for the first input the predicted variance is small, while for the second input it is large. Further, the set of likely images in the second set also includes images of a Stop Sign. Then the appropriate action for the two inputs, determined after solving the inference problem and quantifying uncertainty, is very different. For the first input, the appropriate action is one of continued motion, whereas for the second input it is to slow down. Similar examples can be drawn from other areas like medical imaging, high frequency trading and autonomous systems (Gal [2016]).

The knowledge of uncertainty can also be useful in determining the optimal location of a sensor. Consider an image recovery problem, where the goal is to infer the signal, and associated uncertainty, using limited amount of measurement data. In this problem a user can leverage information about the spatial distribution of uncertainty to choose the location with maximum uncertainty as next measurement location. This task falls within the fields of active learning and/or design of experiments (DeGroot et al. [1962]) and is particularly useful in applications

## 1 INTRODUCTION

Quantifying uncertainty in an inference problem amounts to making a prediction and quantifying the confidence in that prediction. In the context of an image recovery problem, this may be understood as follows. A typical computer vision algorithm uses a noisy version of an image and prior knowledge to produce the inferred image which can be interpreted as the “best guess” of the original image. Quantifying uncertainty in this context involves generating an estimate of the level of confidence in the best guess, in addition to the guess itself.

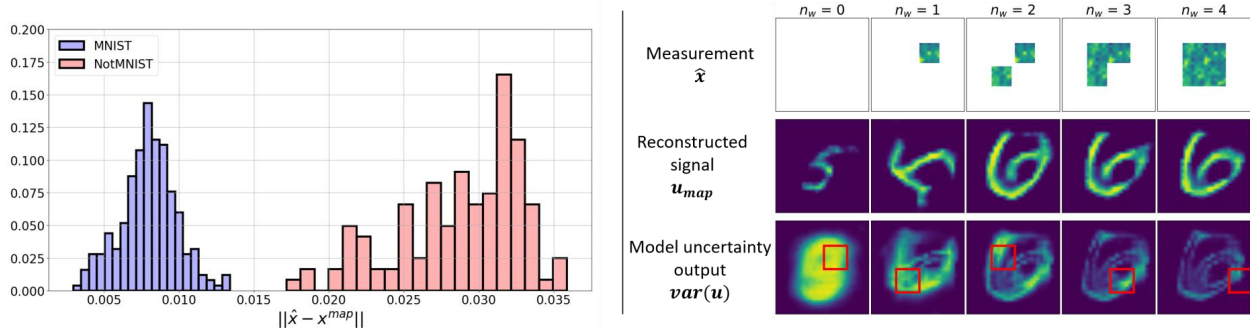


Figure 1: Left panel: Histogram of  $\|\hat{\mathbf{x}} - \mathbf{x}^{\text{map}}\|$ , our measure for out-of-distribution (OOD) data detection, on classification experiments on MNIST. The proposed method is able to successfully distinguish in-distribution (MNIST) and OOD (NotMNIST) test inputs. A large value of this parameter is a warning to the end user to disregard the classification results. Right panel: Estimate of the MAP (2nd row) and pixel-wise variance (3rd row) from the limited view of a noisy image (1st row) using the proposed method for image inpainting with a prior trained on MNIST images. An active learning strategy based on the maximum value of variance is used to determine the location of the subsequent window. An accurate reconstruction of the original image is obtained with just 4 windows.

like satellite imaging, where each measurement requires significant time and/or resources.

In Figure 1, we demonstrate how the proposed GAN-based Bayesian inference algorithm for quantifying uncertainty is useful in the scenarios described above. In the first example it is used to compute a measure that detects, with perfect accuracy, that the input to an image classification algorithm is not from the dataset that was used to train it - the so-called out of distribution (OOD) detection problem. In the second example it computes the pixel-wise variance in an image inpainting task, which is then used to determine the location of the subsequent window to be revealed in an optimal iterative strategy. We return to these applications in greater detail in Section 4.

### 1.1 BAYESIAN INFERENCE

Bayesian inference is a well-established technique for quantifying uncertainties in inference problems (Dashti and Stuart [2016], Kaipio and Somersalo [2006]). It has found applications in diverse fields such as geophysics, climate modeling, chemical kinetics, heat conduction, astrophysics, materials modeling, and the detection and diagnosis of disease. The two critical ingredients of this technique are - an informative prior distribution representing the prior belief about the parameters to be inferred and an efficient method for sampling from the posterior distribution. In this manuscript we describe how deep generative adversarial networks (GANs) can be effectively used in these roles.

We consider the setting where we wish to infer a vector of parameters  $\mathbf{y} \in \mathbb{R}^N$  from the measurement of a re-

lated vector  $\mathbf{x} \in \mathbb{R}^P$ . We allow for two broad classes of problems. In one class (labeled Class 1 in Section 2) the map from  $\mathbf{y}$  to  $\mathbf{x}$  is known through a forward model  $\mathbf{x} = \mathbf{f}(\mathbf{y})$ . These types of problems are often referred to as inverse problems. In the other class (labeled Class 2 in Section 2) this map is not known and must be inferred from prior data. In the discussion that follows, we apply Bayesian inference to Class 1 problems and point out two main challenges. We note that the same challenges apply to problems in Class 2 as well.

A noisy measurement of  $\mathbf{x}$  is denoted by  $\hat{\mathbf{x}} = \mathbf{f}(\mathbf{y}) + \boldsymbol{\eta}$ , where the vector  $\boldsymbol{\eta} \in \mathbb{R}^P$  represents noise. While the forward map  $\mathbf{f}$  is typically well-posed, its inverse is not, and hence to infer  $\mathbf{y}$  from the measurement  $\hat{\mathbf{x}}$  requires techniques that account for this ill-posedness. Classical techniques based on regularization tackle this ill-posedness by using additional information about the sought solution field explicitly or implicitly (Tarantola [2005]). Bayesian inference offers a different approach to this problem by modeling the unknown solution as well as measurements as random variables. This addresses the ill-posedness of the inverse problem, and allows for the characterization of the uncertainty in the inferred solution.

The notion of a prior distribution plays a key role in Bayesian inference. Through multiple observations of the field  $\mathbf{y}$ , denoted by the set  $\mathcal{S} = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(S)}\}$ , we have some prior knowledge of  $\mathbf{y}$  that can be utilized when inferring  $\mathbf{y}$  from  $\hat{\mathbf{x}}$ . This is used to build, or intuit, a prior distribution for  $\mathbf{y}$ , denoted by  $p_Y^{\text{prior}}(\mathbf{y})$ . Some typical examples include Gaussian process prior with specified co-variance kernels, Gaussian Markov random fields, Gaussian priors defined through differential oper-

ators, and hierarchical Gaussian priors. These priors promote smoothness and/or structure in the inferred solution and can be expressed explicitly in an analytical form.

Another key component of Bayesian inference is a distribution that represents the likelihood of  $\mathbf{x}$  given an instance of  $\mathbf{y}$ , denoted by  $p^l(\mathbf{x}|\mathbf{y})$ . This is often determined by the distribution of the error in the model. Given this, the posterior distribution of  $\mathbf{y}$ , determined using Bayes' rule after accounting for the observation  $\hat{\mathbf{x}}$  is given by,

$$p_Y^{\text{post}}(\mathbf{y}|\mathbf{x}) = \frac{1}{\mathbb{Z}} p^l(\mathbf{x}|\mathbf{y}) p_Y^{\text{prior}}(\mathbf{y}) \quad (1)$$

Here,  $\mathbb{Z}$  is the prior-predictive distribution of  $\mathbf{y}$ .

The posterior distribution characterizes the uncertainty in  $\mathbf{y}$ ; however for vectors of large dimension characterizing this distribution explicitly is a challenging task. Consequently the expression above is used to perform tasks that are more manageable. These include determining estimates such as the maximum a-posteriori estimate (MAP), expanding the posterior distribution in terms of other distributions that are simpler to work with (Bui-Thanh et al. [2012]), or using techniques like Markov Chain Monte-Carlo (MCMC) to generate samples that are close to the samples generated by the true posterior distribution (Parno and Marzouk [2018]).

In summary, despite its numerous applications, Bayesian inference faces significant challenges. These include defining a reliable and informative prior distribution for  $\mathbf{x}$  when the set  $\mathcal{S}$  is difficult to characterize analytically, and efficiently sampling from the posterior distribution when the dimension of  $\mathbf{x}$  is large; a typical situation in many practical science and engineering applications.

## 1.2 OUR CONTRIBUTION

The main contribution of this paper are:

1. A novel method for performing Bayesian inference involving complex priors and high dimensional posterior. We utilize the distribution learned by a GAN as a surrogate for the prior distribution and reformulate the inference problem in the low-dimensional latent space of the GAN. Furthermore, we provide a theoretical analysis of the weak convergence of the posterior density learned by the proposed method to the true posterior density.
2. Application of this method to problems where the map from the inferred to the measured vector is known *a-priori*. This leads to novel unsupervised algorithms for physics-based inverse problems and image denoising problems with quantitative measures of uncertainty.
3. Application of this method to problems where the map from the inferred to the measured vector is not

known and is determined from data. This leads to novel algorithms for image classification and image inpainting with quantitative measures of uncertainty.

4. Demonstration of the utility of quantifying uncertainty in the detection of out-of-distribution (OOD) samples and in active learning.

## 1.3 RELATED WORK

The main idea developed in this paper tackles the challenges described above by training a generative adversarial network (GAN) using the sample set  $\mathcal{S}$ , and then using the distribution learned by the GAN as the prior distribution in Bayesian inference. Related work in this area can be organized by considering the two broad classes of problems this idea is applied to.

In one class of problems, which are referred to as inverse problems, the map  $\mathbf{x} = \mathbf{f}(\mathbf{y})$ , that is the map from the inferred field to the measurement is known. The use of sample-based priors for solving inverse problems has a rich history (Calvetti and Somersalo [2005]). As does the idea of reducing the dimension of the parameter space by mapping it to a lower-dimensional space (Marzouk and Najm [2009]). However, the use of learning-based deep generative models like GANs in these tasks is novel. Recently, several authors have considered the use of learning-based methods for solving inverse problems arising in different domains. These include the use of deep convolutional neural networks (CNNs) or recurrent neural networks (RNNs) to solve physics-driven inverse problems (Adler and Öktem [2017], Patel et al. [2019], Pesah et al. [2018]) and use of deep generative models like VAEs and GANs to solve inverse problems arising in computer vision (Kupyn et al. [2018], Zhu et al. [2017], Chang et al.). There is also a growing body of work dedicated to using GANs as a regularizer in solving inverse problems (Lunz et al. [2018] and in compressed sensing (Bora et al. [2018], Kabkab et al. [2018], Shah and Hegde [2018])). However, these approaches differ from ours in that they solve the inverse problem as an optimization problem and do not rely on Bayesian inference; as a result, they add regularization in an ad-hoc manner and do not attempt to quantify the uncertainty in the inferred field. More recently, the approach described in (Adler and Öktem [2018]) utilizes GANs in a Bayesian setting; however the GAN is trained to approximate the posterior distribution and not the prior, as in our case.

In another class of problems the forward map is not known; however in its lieu pair-wise instances of  $\mathbf{y}$  and  $\mathbf{x}$  are available. In this case the algorithms most closely related to our approach are the so-called hybrid methods, where an invertible generator is trained to learn  $p(\mathbf{x})$  and

is linked to another network that is trained to produce  $p(\mathbf{y}|\mathbf{x})$  (Nalisnick et al. [2019], Chen et al. [2019]). This algorithm is then applied to image classification problems, where for a given input image  $\mathbf{x}$ ,  $p(\mathbf{x})$  is used to determine the likelihood of the input and  $p(\mathbf{y}|\mathbf{x})$  is used to infer the probability of the corresponding label. In contrast to this, we train a Wasserstein GAN to learn the joint density  $p(\mathbf{u})$ , where  $\mathbf{u} = [\mathbf{x}, \mathbf{y}]$  and then use this as the prior in a Bayesian update of the posterior density for a given input  $\mathbf{x}$ . The sampling problem for the posterior is reduced to the latent space of the GAN, which is of smaller dimension, and a Markov Chain Monte-Carlo algorithm is trained to generate samples of  $p(\mathbf{u}|\mathbf{x})$  and fully characterize the posterior density.

We note that deep learning based Bayesian networks, where the network weights are stochastic parameters that are determined using Bayesian inference, are another means of quantifying uncertainty (Gal and Ghahramani [2016]), and have recently been applied to semantic image-segmentation and super-resolution (Kendall et al. [2019]).

## 2 PROBLEM FORMULATION

We consider the problem where we wish to infer the vector  $\mathbf{y}$  from the noisy measurement of a related vector  $\mathbf{x}$ . We consider two broad classes of problems of this type. In one class we assume that the forward operator which maps  $\mathbf{y}$  to  $\mathbf{x}$ , that is  $\mathbf{x} = \mathbf{f}(\mathbf{y})$ , is known. While in the other, we assume that  $\mathbf{f}(\mathbf{y})$  is not known and any relation between  $\mathbf{x}$  and  $\mathbf{y}$  must be determined from data.

### 2.1 CLASS 1: THE MAP $\mathbf{f}(\mathbf{y})$ IS KNOWN

These problems are commonly referred to as inverse problems and the map  $\mathbf{x} = \mathbf{f}(\mathbf{y})$  is called the forward map. In this class of problems, this forward map is usually well-defined and is assumed to be known either through physics-based principles or through other modeling paradigms.

A noisy measurement of  $\mathbf{x}$  is denoted by  $\hat{\mathbf{x}} = \mathbf{f}(\mathbf{y}) + \boldsymbol{\eta}$ , where  $\boldsymbol{\eta} \sim p_{\boldsymbol{\eta}}$  is the noise vector. In addition, it is assumed that the sample set  $\mathcal{S} = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(S)}\}$ , which contains multiple realizations of  $\mathbf{y}$  drawn from the distribution  $P_Y$ , is also known. The goal is to use the prior information encoded in  $\mathcal{S}$ , the noisy measurement  $\hat{\mathbf{x}}$ , and the forward map  $\mathbf{f}$  to determine the distribution of the vector  $\mathbf{y}$ .

The prior information for this class of problems is built from the distribution of  $\mathbf{y}$  alone. Thus problems in this class fall into unsupervised learning category, where training only requires instances of one type of data (the

vector to be inferred). The need to have access to pairwise samples of  $\mathbf{x}$  and  $\mathbf{y}$  is circumvented by the knowledge of the forward operator. An example problem in this class is that of image denoising, where  $\mathbf{x}$  represents a noisy image, the operator  $\mathbf{f}$  is the identity, and  $\mathbf{y}$  is the de-noised image. Another example, which is drawn from physics-based inverse problems, is where one wishes to infer the initial temperature field from a measurement of the temperature field at a later time. Here  $\mathbf{x}$  represents the temperature field at a finite time  $T > 0$ ,  $\mathbf{f}$  is the forward in time heat conduction operator, and  $\mathbf{y}$  is the temperature at the initial time. A large number of other physics-based inverse problems can also be cast in this form.

Given  $\hat{\mathbf{x}}$ , using Bayes' rule we may write the posterior distribution of  $\mathbf{y}$  as,

$$\begin{aligned} p_Y^{\text{post}}(\mathbf{y}|\mathbf{x}) &= \frac{1}{\mathbb{Z}} p^{\text{l}}(\mathbf{x}|\mathbf{y}) p_Y^{\text{prior}}(\mathbf{y}) \\ &= \frac{1}{\mathbb{Z}} p_{\boldsymbol{\eta}}(\hat{\mathbf{x}} - \mathbf{f}(\mathbf{y})) p_Y(\mathbf{y}). \end{aligned} \quad (2)$$

where  $\mathbb{Z}$  is the prior-predictive distribution of  $\mathbf{y}$  and ensures that the posterior integrates to one.

In order to efficiently sample the posterior density, we first train a GAN using the set  $\mathcal{S}$  whose elements are sampled from  $P_Y$ . We let  $\mathbf{z} \sim p_Z(\mathbf{z})$  characterize the latent vector space of the GAN, and let  $\mathbf{g}(\mathbf{z})$  and  $d(\mathbf{y})$  denote its generator and discriminator, respectively. In Appendix A, assuming that (a) the stationarity conditions for the adversarial loss function are satisfied and (b) that the set of basis functions obtained by taking the derivative of the discriminator with respect to its weights forms a complete basis in  $L^\infty(\Omega_y)$ , in the limit of infinite capacity, we prove that

$$\mathbb{E}_{\mathbf{y} \sim p_Y} [m(\mathbf{y})] = \mathbb{E}_{\mathbf{z} \sim p_Z} [m(\mathbf{g}(\mathbf{z}))], \quad (3)$$

for sufficiently smooth  $m(\mathbf{y})$ . This equation states that once a GAN is trained using the sample set  $\mathcal{S}$ , it may be used to evaluate any population parameter for  $p_Y$  by sampling from  $p_Z$  and then passing the samples through the generator. Since the dimension of the latent space is much smaller than that of  $\mathbf{y}$ , this represents an efficient means of evaluating population parameters.

In order to turn this into an expression for evaluating a population parameter for the posterior density, we select  $m(\mathbf{y}) = \frac{l(\mathbf{y}) p_{\boldsymbol{\eta}}(\hat{\mathbf{x}} - \mathbf{f}(\mathbf{y}))}{\mathbb{Z}}$ , substitute this expression on both side of (3), and use (2) to arrive at,

$$\mathbb{E}_{\mathbf{y} \sim p_Y^{\text{post}}} [l(\mathbf{y})] = \mathbb{E}_{\mathbf{z} \sim p_Z^{\text{post}}} [l(\mathbf{g}(\mathbf{z}))], \quad (4)$$

where

$$p_Z^{\text{post}}(\mathbf{z}|\mathbf{x}) \equiv \frac{1}{\mathbb{Z}} p_{\boldsymbol{\eta}}(\hat{\mathbf{x}} - \mathbf{f}(\mathbf{g}(\mathbf{z}))) p_Z(\mathbf{z}). \quad (5)$$

The distribution  $p_Z^{\text{post}}$  is the analog of  $p_Y^{\text{post}}$  in the latent vector space. The measurement  $\hat{x}$  updates the prior distribution for  $\mathbf{y}$  to the posterior distribution. Similarly, it updates the prior distribution for  $z$ ,  $p_Z$ , to the posterior distribution,  $p_Z^{\text{post}}$ , defined above.

Equation (4) implies that sampling from the posterior distribution of  $\mathbf{y}$  is equivalent to sampling from the posterior distribution for  $z$  and passing the sample through the generator  $g$ . That is,

$$\mathbf{y} \sim p_Y^{\text{post}}(\mathbf{y}|\hat{x}) \Rightarrow \mathbf{y} = g(z), z \sim p_Z^{\text{post}}(z|\hat{x}). \quad (6)$$

Since the dimension of  $z$  is typically smaller than that of  $\mathbf{y}$ , this represents an efficient approach to sampling from the posterior of  $\mathbf{y}$ .

The left hand side of (4) is an expression for a population parameter of the posterior. The right hand side of this equation describes how this parameter may be evaluated by sampling  $z$  (instead of  $\mathbf{y}$ ) from  $p_Z^{\text{post}}$ . In practise this is accomplished by generating an MCMC approximation,  $p_Z^{\text{mcmc}}(z|\mathbf{x}) \approx p_Z^{\text{post}}(z|\mathbf{x})$  using the definition in (5), and thereafter sampling  $z$  from this distribution. This circumvents the calculation of the prior-predictive distribution of  $\mathbf{y}$  (denoted by  $\mathbb{Z}$ ), which would be necessary when using (5) directly. Then from (4), any desired population parameter for posterior distribution may be approximated as

$$\begin{aligned} \overline{l(\mathbf{y})} &\equiv \mathbb{E}_{\mathbf{y} \sim p_Y^{\text{post}}} [l(\mathbf{y})] \\ &\approx \frac{\sum_{n=1}^{N_{\text{samp}}} l(g(z))}{N_{\text{samp}}}, \quad z \sim p_Z^{\text{mcmc}}(z|\mathbf{x}). \end{aligned} \quad (7)$$

For all the numerical experiments in this paper we have used this approach to evaluate population parameters.

## 2.2 CLASS 2: THE MAP $f(\mathbf{y})$ IS NOT KNOWN

We now consider problems where the relation between  $\mathbf{x}$  and  $\mathbf{y}$  is not known and must be inferred from data. We denote by  $\mathbf{u} = [\mathbf{x}, \mathbf{y}]$  the joint vector and recognize that the measurement has the form  $\hat{x} = \mathbb{1}_x \mathbf{u} + \boldsymbol{\eta}$ , where  $\mathbb{1}_x$  is the indicator function that extracts components of  $\mathbf{x}$  from  $\mathbf{u}$ , and  $\boldsymbol{\eta}$  is the noise vector drawn from the distribution  $p_\eta$ . Further we assume that the sample set  $S = \{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(S)}\}$  contains multiple measurements of  $\mathbf{u}$  drawn from the distribution  $P_U$ . The goal is to use the prior information encoded in  $S$  and the new, noisy measurement  $\hat{x}$  to determine the distribution for the corresponding vector  $\mathbf{y}$ , and perhaps also the de-noised version of  $\mathbf{x}$ .

Since the prior information is built from the joint distribution  $P_U$  this class of problems is one of supervised learning where training requires *pair-wise* instances of  $\mathbf{x}$

and  $\mathbf{y}$ . An example problem in this class is that of image classification, where  $\mathbf{x}$  represents an image and  $\mathbf{y}$  represents the corresponding one-hot encoded label vector. Another example is that of image inpainting, where  $\mathbf{x}$  represents the portion of an image that is revealed and  $\mathbf{y}$  represents the portion that is occluded.

The use of GANs as priors in this class of problems closely parallels the development for problems treated in the previous section. Therefore, rather than repeating the entire development below, we only highlight the important steps and salient differences.

Using Bayes' rule we may write the posterior distribution of  $\mathbf{u} = [\mathbf{x}, \mathbf{y}]$  as,

$$\begin{aligned} p_U^{\text{post}}(\mathbf{u}|\mathbf{x}) &= \frac{1}{\mathbb{Z}} p^l(\mathbf{x}|\mathbf{u}) p_U^{\text{prior}}(\mathbf{u}) \\ &= \frac{1}{\mathbb{Z}} p_\eta(\hat{x} - \mathbb{1}_x(\mathbf{u})) p_U(\mathbf{u}), \end{aligned} \quad (8)$$

where  $\mathbb{Z}$  is the prior-predictive distribution of  $\mathbf{u}$ . In order to efficiently sample the posterior density we train a GAN using the set  $S$  to generate a prior. As before, we let  $z \sim p_Z(z)$  characterize the latent vector space of the GAN, and let  $g(z)$  denote its generator. Under the assumptions of the result derived in Appendix A, we have

$$\mathbb{E}_{\mathbf{u} \sim p_U} [m(\mathbf{u})] = \mathbb{E}_{z \sim p_Z} [m(g(z))], \quad (9)$$

for sufficiently smooth  $m(\mathbf{u})$ . We choose  $m(\mathbf{u}) = \frac{l(\mathbf{u}) p_\eta(\hat{x} - \mathbb{1}_x(\mathbf{u}))}{\mathbb{Z}}$ , substitute it in (9), and make use of (8) to arrive at,

$$\mathbb{E}_{\mathbf{u} \sim p_U^{\text{post}}} [l(\mathbf{u})] = \mathbb{E}_{z \sim p_Z^{\text{post}}} [l(g(z))], \quad (10)$$

where

$$p_Z^{\text{post}}(z|\mathbf{x}) \equiv \frac{1}{\mathbb{Z}} p_\eta(\hat{x} - \mathbb{1}_x(g(z))) p_Z(z). \quad (11)$$

The distribution  $p_Z^{\text{post}}$  is the analog of  $p_U^{\text{post}}$  in the latent vector space. We use (11) to generate an MCMC approximation,  $p_Z^{\text{mcmc}}(z|\mathbf{x}) \approx p_Z^{\text{post}}(z|\mathbf{x})$  of the posterior. Thereafter, from (10) we conclude that any population parameter for the posterior can be approximated as

$$\begin{aligned} \overline{l(\mathbf{u})} &\equiv \mathbb{E}_{\mathbf{u} \sim p_U^{\text{post}}} [l(\mathbf{u})] \\ &\approx \frac{\sum_{n=1}^{N_{\text{samp}}} l(g(z))}{N_{\text{samp}}}, \quad z \sim p_Z^{\text{mcmc}}(z|\mathbf{x}) \end{aligned} \quad (12)$$

We note that this approach allows us to compute population parameters for the entire vector  $\mathbf{u}$ , which includes the vector  $\mathbf{y}$ , which is not observed, as well as the vector  $\mathbf{x}$ , for which a noisy measurement,  $\hat{x}$ , is available. While it is clear that parameters related to  $\mathbf{y}$  are useful, in some

instances it is also useful to estimate population parameters related to  $x$ . A case in point is the image classification problem considered in the following section. In this problem  $\hat{x}$  represents the input image and  $y$  represents its label. Here computing parameters associated with  $y$  provide information about the label for an image. In addition, computing  $x^{\text{map}}$  is useful since large values of the quantity  $\|x^{\text{map}} - \hat{x}\|$ , which measures the distance between the mode of the posterior distribution and the input image, are strongly correlated with input images that lie outside of the range of the prior, thus enabling the detection of out of distribution (OOD) samples.

**Summary** We have described a method for probing the posterior distribution in two broad classes of problems when the prior is defined by a GAN. The steps of our algorithm are: (1) Train a GAN using the sample set  $S$  to learn the prior distribution. (2) Reformulate the posterior distribution in the latent space of the GAN. (3) Run a Markov chain Monte Carlo algorithm to generate samples from this low-dimensional posterior distribution. (4) Use MCMC-generated samples to compute population parameters that quantify the uncertainty in the inference.

In the following section we apply the above algorithm to a broad class of problems where we draw inferences from noisy measurements and quantify uncertainty in these inferences. Wherever possible, we compare our predictions with related methods and/or benchmark solutions and also highlight the role of uncertainty quantification in downstream tasks. In Appendix B, we also derive a computationally efficient approach for estimating the MAP for the posterior density of the latent vector under the assumptions of Gaussian noise and prior.

### 3 EXPERIMENTS

In this section we apply our method to the two broad classes considered earlier. One where the forward map is known and another where it is inferred from data. In each case we apply our method to determine important population parameters that include  $y^{\text{mean}}$ ,  $\text{var}(y)$  and  $\|\hat{x} - x^{\text{map}}\|$ . Thereafter, we use these to answer important questions like: Is the input to the inference problem consistent with the prior data it was trained on? Do we have confidence in the inference? How do we utilize this knowledge in order to design the next measurement.

In all cases we use a Wasserstein GAN-GP (Gulrajani et al. [2017]) to learn the prior density (architecture described in Appendix D). We also ensure that the target images are not chosen from the set used to train the GAN. We sample from the posterior using Hamiltonian Monte Carlo (Brooks et al. [2012]) and implement it using Tensorflow-probability library. We use initial step

Table 1: Comparison of different hybrid models. Arrows indicate which direction is better.

	Configuration / Rejection rule	MNIST		NotMNIST
		Acc $\uparrow$	FPR $\downarrow$	Entropy $\uparrow$
Nalisnick et al. [2019] ( $\lambda = 10.0/D$ )	$\log p(x)$	95.99 %	-	<b>2.300</b>
Chen et al. [2019] ( $\lambda = 1$ )	Coupling	95.42%	-	-
	+ $1 \times 1$ Conv	94.22%	-	-
	Residual	98.69%	-	-
Ours	$\ \hat{x} - x^{\text{map}}\ $	96.81%	0	<b>2.300</b>
	+ $\ \text{var}(y)\ $	<b>99.57%</b>	0.064	<b>2.300</b>

size of 1.0 for HMC and adapt it following (Andrieu and Thoms [2008]) based on the target acceptance probability. We use 64k samples with burn-in period of 0.5. We select these parameters to ensure convergence of chains. Using the HMC sampler we compute the population parameters of interest.

#### 3.1 IMAGE CLASSIFICATION

This problem belongs to the second class, where the forward map is not known and must be learnt from data. The objective of this task is to infer the label  $y$  along with its uncertainty for a given input image  $\hat{x}$ . This predictive uncertainty estimation is crucial in deep learning applications where high-stakes decisions are made based on the output of a model (Kahn et al. [2017]). It has been shown that in real-world scenarios, where a model might encounter inputs that are anomalous to its training data distribution, many models produce overconfident predictions (Lakshminarayanan et al. [2016]) raising serious concerns about AI safety (Amodei et al. [2016]). In this situation, it is desirable that such out-of-distribution (OOD) data points are detected upfront before making any prediction. A useful probabilistic predictive model should therefore flag all OOD data points, maintain high levels of accuracy on in-distribution data points, and provide a measure of confidence in its predictions. In order to achieve this goal, we compute three different quantities:  $\|\hat{x} - x^{\text{map}}\|$  for OOD detection,  $y^{\text{mean}}$  for prediction, and  $\text{var}(y)$  as a measure of confidence in the prediction.

We consider the MNIST database of hand-written digits and use 55k images and the corresponding labels to train a WGAN-GP. Thereafter, we use the MNIST test set to test the performance of our algorithm for in-distribution data, and NotMNIST test set for OOD data. Our approach of learning and inferring the joint distribution is closely related to hybrid models and hence we compare our performance against the most recent hybrid models

in Table 1.

We determine whether a given test image is OOD based on a rejection rule. If this condition is satisfied then following Nalisnick et al. [2019] we set the probability of each label to be equal. We then quantify the performance of the rejection rule by reporting the average entropy of the labels for all test samples from the OOD set and the false positive rate (FPR = # of in-distribution samples rejected as OOD / # of in-distribution samples). Thereafter, for all in-distribution samples that are correctly identified, we report the accuracy of predicting the label, which is determined from  $\mathbf{y}^{\text{mean}}$ . We consider two rejection rules:  $\|\hat{\mathbf{x}} - \mathbf{x}^{\text{map}}\| > c_1$ , and  $\|\hat{\mathbf{x}} - \mathbf{x}^{\text{map}}\| + \|\text{var}(\mathbf{y})\| > c_2$ .

The performance of  $\|\hat{\mathbf{x}} - \mathbf{x}^{\text{map}}\| > c_1$  rejection rule can be discerned in Figure 1, where we observe that it perfectly segregates the in-distribution and OOD samples. This is also apparent in Table 1, where it yields zero FPR and maximum entropy. Its accuracy for the in-distribution samples is also quite high. This accuracy can be further improved by using the combined rejection rule  $\|\hat{\mathbf{x}} - \mathbf{x}^{\text{map}}\| + \|\text{var}(\mathbf{y})\| > c_2$ , since it rejects some incorrectly labeled in-distribution samples with high variance as OOD. However, this comes at the cost of a slightly higher FPR. The usefulness of  $\|\text{var}(\mathbf{y})\|$  as a predictor of accuracy is evident in Figure 2, where we observe that most correctly labeled samples have low variance (avg. value =  $0.0026 \pm 6\text{e-}3$ ) when compared with their incorrectly labeled counterparts (avg. value =  $0.025 \pm 5\text{e-}3$ ).

In Table 1 we compare the performance of our GAN-based approach with two hybrid invertible flow-based models Nalisnick et al. [2019], Chen et al. [2019]. The explicit nature of these models allows the joint density to be decomposed into generative and discriminative components, enabling a way to explore the generative-discriminative trade-off by introducing a weighted likelihood objective with a scaling parameter  $\lambda$ . Values of  $\lambda < 1$  favor discriminative performance, while  $\lambda > 1$  favors generative performance. In this context our approach may be regarded as one where the generative and discriminative components are equally weighted, and is therefore close to the choice  $\lambda = 1$ . Given this, in Table 1, we have compared our approach with hybrid models where  $\lambda \approx 1$ . We note that with the  $\|\hat{\mathbf{x}} - \mathbf{x}^{\text{map}}\| \leq c_1$  rejection rule our model performs competitively with both the scaled (Nalisnick et al. [2019]) and the un-scaled hybrid models (Chen et al. [2019]) for both in-distribution and OOD datasets. With the  $\|\hat{\mathbf{x}} - \mathbf{x}^{\text{map}}\| + \|\text{var}(\mathbf{y})\| > c_2$  rejection rule it outperforms both hybrid models giving maximum accuracy and entropy but with non-zero FPR.

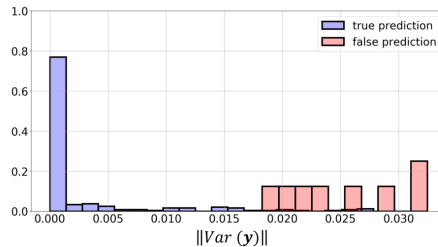


Figure 2: Histogram of  $\|\text{var}(\mathbf{y})\|$  for MNIST dataset.

### 3.2 IMAGE INPAINTING

Image inpainting also belongs to second class of problems, where the forward map is unknown and has to be learnt from data. In this case the quantity to be inferred ( $\mathbf{y}$ ) is the occluded region of an image, and the measurement ( $\hat{\mathbf{x}}$ ) is the noisy version of its visible portion. The goal is to recover the entire image ( $\mathbf{u} = [\mathbf{x}, \mathbf{y}]$ ). While there has been great interest in recent years in developing efficient deep learning-based image inpainting algorithms (Yu et al. [2018]), most of it has focused on deterministic algorithms that lack the ability to quantify uncertainty in a prediction. In contrast, we use the algorithm in Section 2 to perform probabilistic image inpainting.

We consider MNIST dataset and use 55k images to train a WGAN-GP. We generate measurements by selecting an image from the test set, occluding a significant region and then adding Gaussian noise. With this as input we use our algorithm to generate samples from the posterior distribution of the entire image (both occluded and retained regions). From these samples we evaluate the relevant population parameters,  $\mathbf{u}^{\text{map}}$ ,  $\mathbf{u}^{\text{mean}}$ , and  $\text{var}(\mathbf{u})$ . These results are shown in Figure 3 and indicate that the map and mean images are close to the true image, even in the presence of significant occlusion and noise. The image of pixel-wise variance reveals that we are most uncertain along the boundaries of the digits and around the occlusion window.

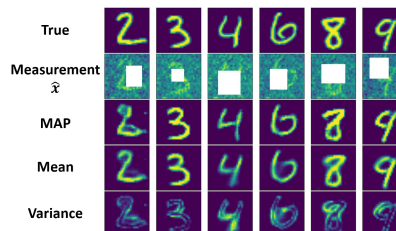


Figure 3: Estimate of the MAP, mean and pixel-wise variance from noisy occluded images using the proposed method. The variance is peaked at the occluded region.

In Figure 1, we demonstrate how uncertainty may be



used in active learning/design of experiment, where the goal is to determine the optimal location for a measurement. We begin with an input where the entire image is occluded and in every subsequent step, we allow for a small  $7 \times 7$  pixel window to be revealed. We select the window with the largest average pixel-wise variance. As the iterations progress, the MAP estimate converges to the true digit, and the variance decreases. In about 4 iterations we arrive at a very good guess for the digit. The performance of this approach is quantified in Figure 4, where we have plotted reconstruction error versus the number windows for this strategy, and a strategy where the subsequent window is selected randomly. The variance-driven strategy consistently performs better. We are not aware of any other methods for computing uncertainty in recovered images that have been applied to drive an active learning task in image inpainting. While methods based on dropout (Kendall et al. [2019]) or variational inference (Kohl et al. [2018]) could be extended to accomplish this, this has not been done thus far.

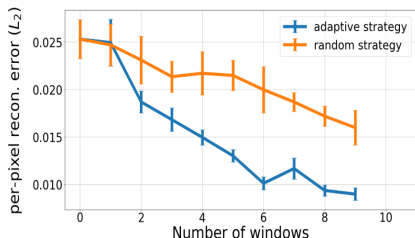


Figure 4: Average reconstruction error (with 95% confidence interval) as a function of number of windows for variance-driven (adaptive) and random sampling strategies.

Results for the variance-based window selection strategy applied to the CelebA dataset are shown in Figure 5. We observe that the algorithm produces realistic images at each iteration; however, the initial variance is large indicating large uncertainty. As more windows are sampled using the active learning strategy, the variance reduces and by the 7th iteration a good approximation of the true image is obtained, even though only a small, noisy portion is revealed. Additional results for this dataset are discussed in Appendix C.2.

### 3.3 A PHYSICS-DRIVEN INFERENCE PROBLEM

We now consider a problem from Class 1, where the forward map is known. In particular, we consider the problem of inverse heat conduction, where the goal is to infer the initial temperature distribution (at  $t = 0$ ) in a domain given a noisy measurement of temperature at later time ( $t = 1$ ) and the thermal conductivity of

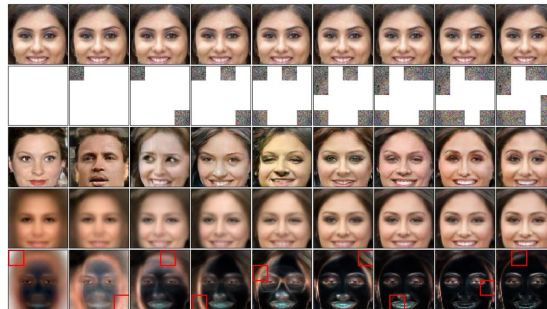


Figure 5: CelebA dataset: Estimate of the MAP, mean and variance from the limited view of a noisy image (2nd row) of a true image (1st row) using the variance-driven adaptive learning strategy.

the material. The forward map is the solution of the time-dependent heat conduction problem with uniform conductivity,  $\kappa = 0.64$ , in a square domain of length  $L = 2\pi$  with Dirichlet boundary conditions. This operator maps the initial temperature field ( $y$ ) to the temperature field at later time ( $x$ ). Its discrete version is obtained by discretizing the time-dependent linear heat conduction equation using the central difference scheme in space and backward Euler scheme in time. Much like a blurring kernel, the forward operator smooths the initial temperature distribution, and the extent of smoothing increases with  $\kappa \times t$ .

We consider a family of initial temperatures where the background is zero, and the temperature on a rectangular sub-domain varies linearly from 2 units on the left edge to 4 units on the right edge. This distribution is parameterized by four parameters,  $\{\xi_i\}_{i=1}^4$ , which are the coordinates of the lower left and upper right corners of the rectangular region. The sample set  $\mathcal{S}$  is created by sampling each parameter from a uniform distribution and is used to train the GAN prior. The posterior distribution is sampled using the HMC sampler.

In the top two rows of Figure 6, we have plotted the true initial condition, the noise-free temperature at  $t = 1$ , and the noisy temperature measurement ( $\sigma_x = 1$ ) used as input in the GAN-based prior approach. The corresponding MAP, mean and pixel-wise variance estimated by the MCMC approximation are shown next. We observe that the MAP is very close to the true initial temperature distribution and the variance is concentrated along the edges of the rectangle where the uncertainty is the largest. In the following columns we have plotted the MAP estimate obtained assuming  $L_2$  and  $H^1$  Gaussian priors, which are often used to solve these types of problems, and are clearly much less accurate.

For this problem the “true” posterior can be reduced to



the 4-dimensional space of parameters, and sampled by generating initial conditions corresponding to the values of these parameters. A simple Monte-Carlo approximation can be performed to compute the mean and the pixel-wise variance for the true posterior (last two columns of Figure 6). By comparing these with the mean and the pixel-wise variance (columns 5 & 6) estimated by the GAN-based prior, we conclude that GAN-based posterior has converged to the true posterior.

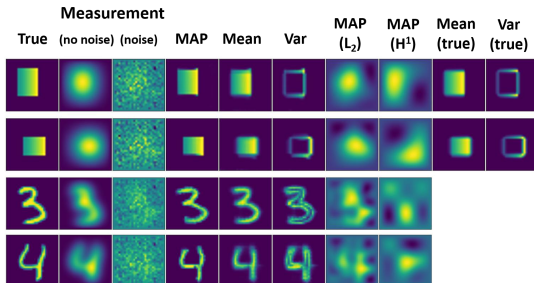


Figure 6: From left to right: (1) true initial temperature, (2) temperature at  $t = 1$ , (3) noisy version temperature used as measurement, (4), (5), (6) MAP, mean and pixel-wise variance estimates using GAN priors, (7) and (8) MAP estimates using  $L_2$  and  $H^1$  Gaussian priors, (9) & (10) true MAP and variance.

In the bottom rows of Figure 6, we plot similar results for initial conditions generated from the MNIST database when the GAN prior was also trained on the MNIST database. The measurement is made at  $t = 0.2$ . Once again we observe that mean and the map estimated by our approach is very close to the true initial condition, while the MAP solutions obtained from the  $L_2$  and  $H^1$  priors are inaccurate. The pixel-wise variance illustrates the uncertainty in determining the boundary of the digits.

### 3.4 IMAGE DENOISING

As another example of a problem where the forward map is known (Class 1 problem) we consider image denoising. Here the forward map is the identity, the measured data is the noisy image and the inferred field is its denoised version. We consider the MNIST dataset and use 55k images to train the GAN. We add Gaussian noise with zero mean and specified variance ( $\sigma_x$ ) to the test images and use these as measurements to recover the distribution of likely images using the MCMC approach. In Figure 7, we have plotted the noisy input image, the MAP estimate, and the pixel-wise mean and variance. For low and medium noise levels ( $\sigma_x = 0.1, 1$ ), we are able to recover the original image with good accuracy, the pixel-wise variance is small overall, and is largest around the boundary of the recovered digit. For the high-

est noise level ( $\sigma_x = 10$ ), however, the image recovered by the MAP is incorrect in 2/3 cases, and would be misleading if viewed by itself. However, when viewed in conjunction with the estimated variance, which is large, it is clear that the confidence in the inference is small and the inferred image ought not be trusted for downstream tasks. The dependence of the average per-pixel variance in the recovered image on the variance of noise in the measured image is shown in Figure 7, and it increases with noise. Additional results for CelebA are discussed in Appendix C.2

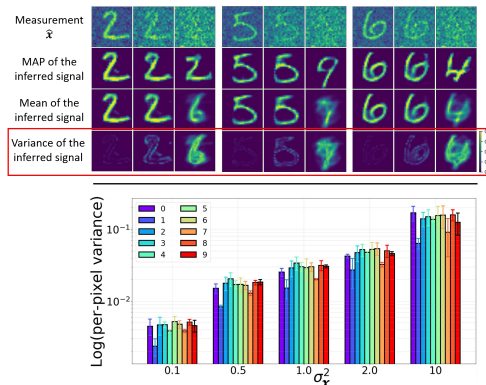


Figure 7: Top panel: Estimate of the MAP, mean and pixel-wise variance from a noisy image using the proposed method. In the first three panels  $\sigma_x = 0.1, 1, \& 10$ , when moving from left to right. Bottom panel: Average variance per pixel ( $\text{var}(\mathbf{y})$ ) in a reconstructed image as a function of variance of noise for 10 MNIST digits (along with 95% confidence interval).

## 4 CONCLUSIONS

The ability to quantify the uncertainty in an inference problem is useful in developing confidence in that inference, identifying measurements that are outliers, and in designing strategies to improve the confidence. In this paper we have described how this may be accomplished when solving a Bayesian inference problem by using GANs as priors. Since GANs can learn complex distributions of a wide variety of fields from their samples, this approach can be applied to a range of problems in computer vision and physics-driven inference. This includes those where the operator that maps the inferred field to measurement is known (so-called inverse problems) and those where this map is not known and must be inferred from data. It derives its efficiency by mapping the posterior distribution to the latent space, whose dimension is often much smaller than that of the inferred field. We have presented applications of this approach to image classification, image inpainting, image denoising and physics-driven inverse problems.

## References

- Jonas Adler and Ozan Öktem. Solving ill-posed inverse problems using iterative deep neural networks. *Inverse Problems*, 33(12):124007, dec 2017. ISSN 0266-5611. doi: 10.1088/1361-6420/aa9581. URL <http://stacks.iop.org/0266-5611/33/i=12/a=124007?key=crossref.65c4fa88a47e07d4789aa10592f2090c>.
- Jonas Adler and Ozan Öktem. Deep bayesian inversion. *arXiv preprint arXiv:1811.05910*, 2018.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete Problems in AI Safety. jun 2016. URL <http://arxiv.org/abs/1606.06565>.
- Christophe Andrieu and Johannes Thoms. A tutorial on adaptive MCMC. *Statistics and Computing*, 18(4): 343–373, dec 2008. ISSN 09603174. doi: 10.1007/s11222-008-9110-y.
- Ashish Bora, Eric Price, and Alexandros G Dimakis. Ambientgan: Generative models from lossy measurements. *ICLR*, 2:5, 2018.
- Steve Brooks, Andrew Gelman, Galin Jones, Xiao-Li Meng, and Radford M Neal. MCMC using Hamiltonian dynamics. Technical report, 2012. URL <https://arxiv.org/pdf/1206.1901.pdf>.
- Tan Bui-Thanh, Carsten Burstedde, Omar Ghattas, James Martin, Georg Stadler, and Lucas C. Wilcox. Extreme-scale UQ for Bayesian inverse problems governed by PDEs. In *2012 International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–11. IEEE, nov 2012. ISBN 978-1-4673-0805-2. doi: 10.1109/SC.2012.56. URL <http://ieeexplore.ieee.org/document/6468442/>.
- Daniela Calvetti and Erkki Somersalo. Priorconditioners for linear systems. *Inverse Problems*, 21(4):1397–1418, aug 2005. ISSN 0266-5611. doi: 10.1088/0266-5611/21/4/014. URL <http://stacks.iop.org/0266-5611/21/i=4/a=014?key=crossref.ab419ffb66111e3db21bf3d9fd3836f7>.
- J H Rick Chang, Chun-Liang Li, Barnab Barnabás, P Oczos, B V K Vijaya Kumar, and Aswin C Sankaranarayanan. One Network to Solve Them All-Solving Linear Inverse Problems using Deep Projection Models. Technical report. URL <https://arxiv.org/pdf/1703.09912.pdf>.
- Ricky T. Q. Chen, Jens Behrmann, David Duvenaud, and Jörn-Henrik Jacobsen. Residual Flows for Invertible Generative Modeling. jun 2019. URL <http://arxiv.org/abs/1906.02735>.
- Masoumeh Dashti and Andrew M Stuart. The bayesian approach to inverse problems. *Handbook of Uncertainty Quantification*, pages 1–118, 2016.
- Morris H DeGroot et al. Uncertainty, information, and sequential experiments. *The Annals of Mathematical Statistics*, 33(2):404–419, 1962.
- Yarin Gal. *Uncertainty in deep learning*. PhD thesis, PhD thesis, University of Cambridge, 2016.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059, 2016.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017.
- Maya Kabkab, Pouya Samangouei, and Rama Chellappa. Task-aware compressed sensing with generative adversarial networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Gregory Kahn, Adam Villaflor, Vitchyr Pong, Pieter Abbeel, and Sergey Levine. Uncertainty-Aware Reinforcement Learning for Collision Avoidance. feb 2017. URL <http://arxiv.org/abs/1702.01182>.
- Jari Kaipio and Erkki Somersalo. *Statistical and computational inverse problems*, volume 160. Springer Science & Business Media, 2006.
- Alex Kendall, Vijay Badrinarayanan, and Roberto Cipolla. Bayesian SegNet: Model Uncertainty in Deep Convolutional Encoder-Decoder Architectures for Scene Understanding. British Machine Vision Association and Society for Pattern Recognition, may 2019. doi: 10.5244/c.31.57.
- Simon A. A. Kohl, Bernardino Romera-Paredes, Clemens Meyer, Jeffrey De Fauw, Joseph R. Ledsam, Klaus H. Maier-Hein, S. M. Ali Eslami, Danilo Jimenez Rezende, and Olaf Ronneberger. A Probabilistic U-Net for Segmentation of Ambiguous Images. jun 2018. URL <http://arxiv.org/abs/1806.05034>.
- Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jii Matas. DeblurGAN: Blind Motion Deblurring Using Conditional Adversarial Networks. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 8183–8192. IEEE Computer Society, dec 2018. ISBN 9781538664209. doi: 10.1109/CVPR.2018.00854. URL <http://arxiv.org/abs/1711.07064>.

- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. *Advances in Neural Information Processing Systems*, 2017-December:6403–6414, dec 2016. URL <http://arxiv.org/abs/1612.01474>.
- Sebastian Lunz, Ozan Öktem, and Carola-Bibiane Schönlieb. Adversarial regularizers in inverse problems. In *Advances in Neural Information Processing Systems*, pages 8507–8516, 2018.
- Youssef M. Marzouk and Habib N. Najm. Dimensionality reduction and polynomial chaos acceleration of Bayesian inference in inverse problems. *Journal of Computational Physics*, 228(6):1862–1902, apr 2009. ISSN 0021-9991. doi: 10.1016/J.JCP.2008.11.024. URL <https://www.sciencedirect.com/science/article/pii/S0021999108006062>.
- Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Hybrid Models with Deep and Invertible Features. *36th International Conference on Machine Learning, ICML 2019*, 2019-June:8295–8304, feb 2019. URL <http://arxiv.org/abs/1902.02767>.
- Matthew D Parno and Youssef M Marzouk. Transport map accelerated markov chain monte carlo. *SIAM/ASA Journal on Uncertainty Quantification*, 6(2):645–682, 2018.
- Dhruv Patel, Raghav Tibrewala, Adriana Vega, Li Dong, Nicholas Hugenberg, and Assad A Oberai. Circumventing the solution of inverse problems in mechanics through deep learning: Application to elasticity imaging. *Computer Methods in Applied Mechanics and Engineering*, 353:448–466, 2019.
- Arthur Pesah, Antoine Wehenkel, and Gilles Louppe. Recurrent machines for likelihood-free inference. nov 2018. URL <http://arxiv.org/abs/1811.12932>.
- Viraj Shah and Chinmay Hegde. Solving linear inverse problems using gan priors: An algorithm with provable guarantees. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4609–4613. IEEE, 2018.
- Albert Tarantola. *Inverse problem theory and methods for model parameter estimation*, volume 89. siam, 2005.
- Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S. Huang. Generative image inpainting with contextual attention. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017.

## A Weak convergence of the prior density

Let the generator of the Wasserstein GAN be given by  $\mathbf{g}(\mathbf{z}; \boldsymbol{\theta})$ , where  $\mathbf{z} \in \mathbb{R}^M$  is the latent vector, and  $\boldsymbol{\theta} \in \mathbb{R}^{N_\theta}$  is the vector of weights. The vector  $\mathbf{z}$  is selected from the distribution  $p_Z(\mathbf{z})$ . Note that  $\mathbf{g} : \mathbb{R}^M \rightarrow \mathbb{R}^N$ .

Let the discriminator of the GAN be given by  $d(\mathbf{y}; \boldsymbol{\phi})$ , where  $\mathbf{y} \in \mathbb{R}^N$ , and  $\boldsymbol{\phi} \in \mathbb{R}^{N_\phi}$  be the vector of weights. Note that  $d : \mathbb{R}^N \rightarrow \mathbb{R}(0, 1)$ .

Assume that the GAN is trained using a set of samples of  $\mathbf{y}$ , drawn from  $p_Y(\mathbf{y})$ .

Then under the following assumptions:

1. The stationarity conditions for the adversarial loss function are satisfied.
2. In the limit of infinite capacity ( $N_\phi \rightarrow \infty$ ), the set of basis functions obtained by taking the derivative of the discriminator with respect to its weights forms a complete basis in  $L^\infty(\Omega_y)$

For a sufficiently smooth  $m(\mathbf{y})$ , we prove that

$$\left| \mathbb{E}_{\mathbf{y} \sim p_Y} [m(\mathbf{y})] \right| = \left| \mathbb{E}_{\mathbf{z} \sim p_Z} [m(\mathbf{g}(\mathbf{z}))] \right|. \quad (13)$$

That is we establish the weak convergence of the density obtained by using a GAN as a prior to the true prior density.

**Proof.** For the loss function, consider

$$\begin{aligned} L(\boldsymbol{\theta}, \boldsymbol{\phi}) &= \mathbb{E}_{\mathbf{y} \sim p_Y} [\rho(1 - d(\mathbf{y}; \boldsymbol{\phi}))] \\ &\quad + \mathbb{E}_{\mathbf{z} \sim p_Z} [\rho(d(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}); \boldsymbol{\phi}))]. \end{aligned} \quad (14)$$

Here  $\rho$  is a monotone real-valued function which defines the GAN family being analyzed. For example, for the Wasserstein GAN,  $\rho(\xi) = \xi$ .

The optimal values of the weights are given by

$$\boldsymbol{\theta}^*, \boldsymbol{\phi}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} (\underset{\boldsymbol{\phi}}{\operatorname{argmin}} (L(\boldsymbol{\theta}, \boldsymbol{\phi}))). \quad (15)$$

The necessary conditions for these optimal values are

$$\frac{\partial L(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*)}{\partial \boldsymbol{\phi}} = \mathbf{0} \quad (16)$$

$$\frac{\partial L(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*)}{\partial \boldsymbol{\theta}} = \mathbf{0}. \quad (17)$$

Using the definition of the loss function (14) in (16), we have

$$\begin{aligned} \mathbb{E}_{\mathbf{y} \sim p_Y} [\rho'(1 - d(\mathbf{y}; \boldsymbol{\phi}^*)) \frac{\partial d}{\partial \boldsymbol{\phi}}(\mathbf{y}; \boldsymbol{\phi}^*)] &= \\ \mathbb{E}_{\mathbf{z} \sim p_Z} [\rho'(d(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}^*); \boldsymbol{\phi}^*)) \frac{\partial d}{\partial \boldsymbol{\phi}}(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}^*); \boldsymbol{\phi}^*)]. \end{aligned} \quad (18)$$

Similarly, using (14) in (17), we have

$$\begin{aligned} \mathbb{E}_{\mathbf{z} \sim p_Z} [\rho'(d(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}^*); \boldsymbol{\phi}^*)) \frac{\partial d}{\partial \boldsymbol{\theta}}(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}^*); \boldsymbol{\phi}^*) \cdot \\ \frac{\partial \mathbf{g}}{\partial \boldsymbol{\theta}}(\mathbf{z}; \boldsymbol{\theta}^*)] = \mathbf{0}. \end{aligned} \quad (19)$$

For the Wasserstein GAN,  $\rho(\xi) = \xi$  and  $\rho'(\xi) = 1$ . As a result (18) reduces to

$$\mathbb{E}_{\mathbf{y} \sim p_Y} [\frac{\partial d}{\partial \boldsymbol{\phi}}(\mathbf{y}; \boldsymbol{\phi}^*)] = \mathbb{E}_{\mathbf{z} \sim p_Z} [\frac{\partial d}{\partial \boldsymbol{\phi}}(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}^*); \boldsymbol{\phi}^*)] \quad (20)$$

and (19) reduces to,

$$\mathbb{E}_{\mathbf{z} \sim p_Z} [\frac{\partial d}{\partial \mathbf{y}}(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}^*); \boldsymbol{\phi}^*) \cdot \frac{\partial \mathbf{g}}{\partial \boldsymbol{\theta}}(\mathbf{z}; \boldsymbol{\theta}^*)] = \mathbf{0}. \quad (21)$$

Let  $w_a(\mathbf{y}) \equiv \frac{\partial d}{\partial \phi_a}(\mathbf{y}; \boldsymbol{\phi}^*)$ , then for  $a = 1, \dots, N_\phi$ . (20) implies

$$\mathbb{E}_{\mathbf{y} \sim p_Y} [w_a(\mathbf{y})] = \mathbb{E}_{\mathbf{z} \sim p_Z} [w_a(\mathbf{g}(\mathbf{z}; \boldsymbol{\theta}^*))]. \quad (22)$$

As  $N_\phi \rightarrow \infty$ , this equation implies that the push forward of the measure in the latent space under the function  $\mathbf{g}(\mathbf{z})$  weakly converges to the measure associated with distribution of  $\mathbf{y}$ . We note with increasing number of weights in the discriminator, the relation above is required to hold for an increasing number of test functions,  $w_a$ . In addition, we have implicitly assumed that the generator is rich enough, that is it has enough weights/layers, such that this relation can actually be satisfied. To make this clear consider the extreme case of a generator with a single weight; in this case there is no way that (22) will be satisfied for a large number  $N_\phi$ . Thus in order for this relation to hold for a large  $N_\phi$ , we must also provide the generator with a large  $N_\theta$ .

Let  $\mathcal{V} \equiv \operatorname{span}\{w_a(\mathbf{y}), a = 1, \dots, N_\phi\}$ . Then from (22) for any  $v \in \mathcal{V}$ , we have

$$\mathbb{E}_{\mathbf{y} \sim p_Y} [v(\mathbf{y})] = \mathbb{E}_{\mathbf{z} \sim p_Z} [v(\mathbf{g}(\mathbf{z}))]. \quad (23)$$

In the equation above, and hereafter, we have suppressed the arguments  $\boldsymbol{\theta}^*$  and  $\boldsymbol{\phi}^*$  for ease of notation, implicitly assuming that the relations hold at the optimal values of weights. Now consider a sufficiently smooth function  $m(\mathbf{y})$  which defines the point estimate we wish to compute, and let  $\bar{m}(\mathbf{y})$  be its  $L^\infty$  projection on to  $\mathcal{V}$ . That is,

$$\bar{m}(\mathbf{y}) = \underset{v \in \mathcal{V}}{\operatorname{argmin}} \|m(\mathbf{y}) - v(\mathbf{y})\|_{L^\infty(\Omega_y)}. \quad (24)$$

and let  $\epsilon = \|m(\mathbf{y}) - \bar{m}(\mathbf{y})\|_{L^\infty(\Omega_y)}$ . Given the assumption that the functions  $w_a$  form a complete basis in  $L^\infty(\Omega_y)$ , we note that as  $N_\phi \rightarrow \infty$ ,  $\epsilon \rightarrow 0$ .

Now consider the difference,

$$\begin{aligned}
& \left| \mathbb{E}_{\mathbf{y} \sim p_Y} [m(\mathbf{y})] - \mathbb{E}_{\mathbf{z} \sim p_Z} [m(\mathbf{g}(\mathbf{z}))] \right| \\
& \leq \left| \mathbb{E}_{\mathbf{y} \sim p_Y} [m(\mathbf{y})] - \mathbb{E}_{\mathbf{y} \sim p_Y} [\bar{m}(\mathbf{y})] \right| \\
& \quad + \left| \mathbb{E}_{\mathbf{z} \sim p_Z} [\bar{m}(\mathbf{g}(\mathbf{z}))] - \mathbb{E}_{\mathbf{z} \sim p_Z} [m(\mathbf{g}(\mathbf{z}))] \right| \\
& \leq \left| \mathbb{E}_{\mathbf{y} \sim p_Y} [m(\mathbf{y}) - \bar{m}(\mathbf{y})] \right| \\
& \quad + \left| \mathbb{E}_{\mathbf{z} \sim p_Z} [m(\mathbf{g}(\mathbf{z})) - \bar{m}(\mathbf{g}(\mathbf{z}))] \right| \\
& \leq \mathbb{E}_{\mathbf{y} \sim p_Y} [\epsilon] + \mathbb{E}_{\mathbf{z} \sim p_Z} [\epsilon] \\
& = 2\epsilon. \tag{25}
\end{aligned}$$

In the equation above, the first inequality is obtained by recognizing that  $\bar{m}(\mathbf{y}) \in \mathcal{V}$  and using (23), the second inequality is a consequence of the triangle inequality and the third is due to the definition of  $\epsilon$ . Now in the limit  $N_\phi \rightarrow \infty$ ,  $\mathcal{V}$  tends to a complete basis, therefore  $\epsilon \rightarrow 0$  and we have the desired result.

## B Expression for the maximum a-posteriori estimate

The techniques described in Section 2.1 focus on sampling from the posterior distribution and computing approximations to population parameters. These techniques can be applied in conjunction with any distribution used to model noise and the latent space vector; that is, any choice of  $p_\eta$  (likelihood) and  $p_Z$  (prior). In this section we consider the special case when Gaussian models are used for noise and the latent vector. In this case, we can derive a simple optimization algorithm to determine the maximum a-posteriori estimate (MAP) for  $p_Z^{\text{post}}(\mathbf{z}|\mathbf{x})$ . This point is denoted by  $\mathbf{z}^{\text{map}}$  in the latent vector space and represents the most likely value of the latent vector in the posterior distribution. It is likely that the operation of the generator on  $\mathbf{z}^{\text{map}}$ , that is  $\mathbf{g}(\mathbf{z}^{\text{map}})$ , will yield a value that is close to  $\mathbf{y}^{\text{map}}$ , and may be considered as a likely solution to the inference problem.

We consider the case when the components of the latent vector are iid with a normal distribution with zero mean and unit variance. This is often the case in many typical applications of GANs. Further, we assume that the components of noise vector are defined by a normal distribution with zero mean and a covariance matrix  $\Sigma$ . Using these assumptions in (5), we have

$$p_Z^{\text{post}}(\mathbf{z}|\mathbf{x}) \propto \exp\left(-\frac{1}{2} \overbrace{(\|\Sigma^{-1/2}(\hat{\mathbf{x}} - \mathbf{f}(\mathbf{g}(\mathbf{z}))\|)^2 + |\mathbf{z}|^2)}^{\equiv r(\mathbf{z})}\right). \tag{26}$$

The MAP estimate for this distribution is obtained by minimizing the negative of the argument of the exponen-

tial. That is

$$\mathbf{z}^{\text{map}} = \underset{\mathbf{z}}{\operatorname{argmin}} r(\mathbf{z}). \tag{27}$$

This minimization problem may be solved using any gradient-based optimization algorithm. The input to this algorithm is the gradient of the functional  $r$  with respect to  $\mathbf{z}$ , which is given by

$$\frac{\partial r}{\partial \mathbf{z}} = \mathbf{H}^T \Sigma^{-1} (\mathbf{f}(\mathbf{g}(\mathbf{z})) - \hat{\mathbf{x}}) + \mathbf{z}, \tag{28}$$

where the matrix  $\mathbf{H}$  is defined as

$$\mathbf{H} \equiv \frac{\partial \mathbf{f}(\mathbf{g}(\mathbf{z}))}{\partial \mathbf{z}} = \frac{\partial \mathbf{f}}{\partial \mathbf{y}} \frac{\partial \mathbf{g}}{\partial \mathbf{z}}. \tag{29}$$

Here  $\frac{\partial \mathbf{f}}{\partial \mathbf{y}}$  is the derivative of the forward map  $\mathbf{f}$  with respect to its input  $\mathbf{x}$ , and  $\frac{\partial \mathbf{g}}{\partial \mathbf{z}}$  is the derivative of the generator output with respect to the latent vector. In evaluating the gradient above we need to evaluate the operation of the matrices  $\frac{\partial \mathbf{f}}{\partial \mathbf{y}}$  and  $\frac{\partial \mathbf{g}}{\partial \mathbf{z}}$  on a vector, and not the matrices themselves. The operation of  $\frac{\partial \mathbf{g}}{\partial \mathbf{z}}$  on a vector can be determined using a back-propagation algorithm within the GAN; while the operation of  $\frac{\partial \mathbf{f}}{\partial \mathbf{y}}$  can be determined by making use of the adjoint of the linearization of the forward operator.

Once  $\mathbf{z}^{\text{map}}$  is determined, one may evaluate  $\mathbf{g}(\mathbf{z}^{\text{map}})$  by using the GAN generator. This represents the value of the field we wish to infer at the most likely value value of latent vector. Note that this is not the same as the MAP estimate of  $p_Y^{\text{post}}(\mathbf{y}|\mathbf{x})$ .

We note that the result derived above applies to Class 1 problems, that is inverse problems. A similar result can also be derived for problems from Class 2, by using (11) as a starting point and repeating the steps outlined above. In this case  $\mathbf{z}^{\text{map}}$  is the minimizer of

$$r \equiv \|\Sigma^{-1/2}(\hat{\mathbf{x}} - \mathbf{R}(\mathbf{g}(\mathbf{z}))\|)^2 + |\mathbf{z}|^2, \tag{30}$$

where  $\mathbf{R}$  is the restriction operator. The gradient for this optimization problem is given by

$$\frac{\partial r}{\partial \mathbf{z}} = \mathbf{H}^T \Sigma^{-1} (\mathbf{R}(\mathbf{g}(\mathbf{z})) - \hat{\mathbf{x}}) + \mathbf{z}, \tag{31}$$

where the matrix  $\mathbf{H}$  is defined as

$$\mathbf{H} = \mathbf{R} \frac{\partial \mathbf{g}}{\partial \mathbf{z}}. \tag{32}$$

## C Additional results

In this section we provide additional results for both MNIST and CelebA dataset for different tasks discussed in the main paper.

## C.1 MNIST

First we provide additional examples in Figure 8 for variance-based adaptive measurement window selection procedure described in Section 3.2.

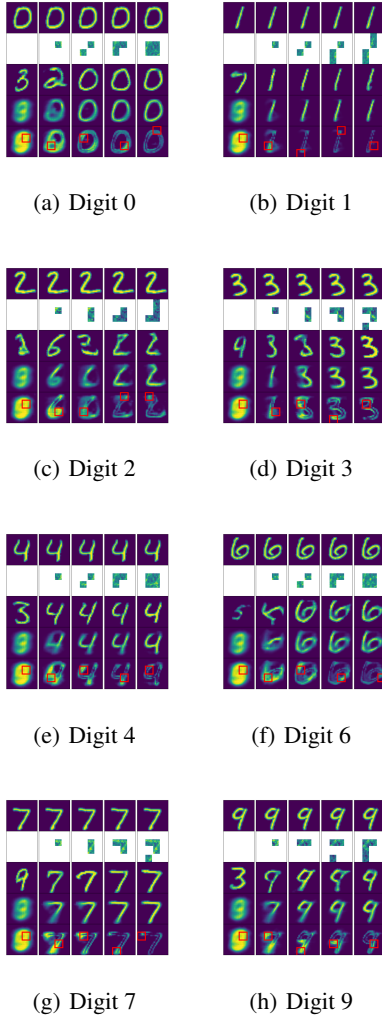


Figure 8: Estimate of the MAP (3rd row), mean (4th row) and variance (5th row) from the limited view of a noisy image (2nd row) using the proposed method. The window to be revealed at a given iteration (shown in red box) is selected using a variance-driven strategy. Top row indicates ground truth. For all digits  $\sigma_y = 1$ .

Figure 9 shows additional results for the inpainting + denoising task, where an MNIST digit is occluded with masks of different sizes at different locations. Note that the variance is high where the occlusion mask is located indicating lower confidence in reconstructed image in that location.

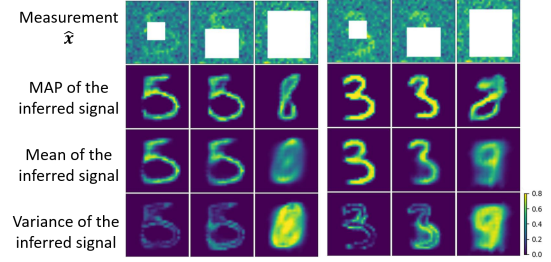


Figure 9: Estimate of the MAP (2nd row), mean (3rd row) and variance (4th row) from a noisy image (1st row) using the proposed method. Note that all variance images are plotted on the same color scale and it highlights increasing level of uncertainty as more and more portion of an image is occluded.

## C.2 CelebA

For the CelebA dataset, we trained WGAN-GP model using more than 200k celebrity facial images and perform inference using remaining test set images. The input images were cropped to a  $64 \times 64$  RGB image and were normalized between  $[-1, 1]$ .

Once the GAN was trained, the HMC algorithm was used for posterior sampling and inference on a complementary set of images (not used for training). In Figure 10 we show some additional results for variance-based adaptive measurement window selection procedure for CelebA dataset.

Next, in figure 11 we show some additional results for image recovery task for CelebA dataset. Once again we note that the MAP estimate and the mean is close to the true image. On the other hand, the closest image from the training set (in an  $L_2$  sense) is not as accurate. This points to the utility of using the GAN as an interpolant in the latent vector space.

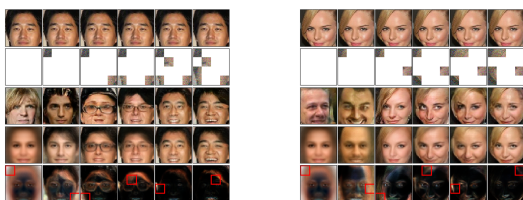
## D Architecture and training details

We use the WGAN-GP model for learning prior density. The tuned value of hyper-parameters is shown in Table 2.

We use the same generator and discriminator architecture for the MNIST and the physics-based inference problem; whereas for the CelebA dataset we use a slightly different architecture to accommodate different input image size. The layout of both these architecture is shown in Figure 12 and 13. Some notes regarding nomenclature used in Figure 12 and 13.

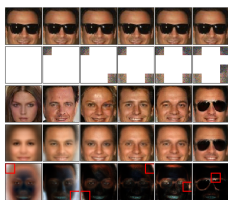
- Conv ( $H \times W \times C \mid s=n$ ) indicates convolutional



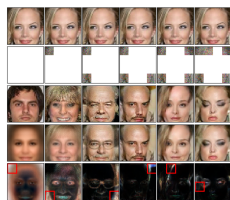


(a)

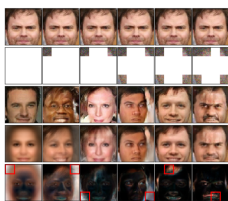
(b)



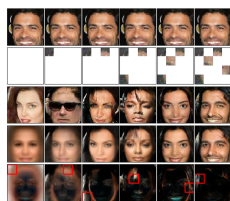
(c)



(d)



(e)

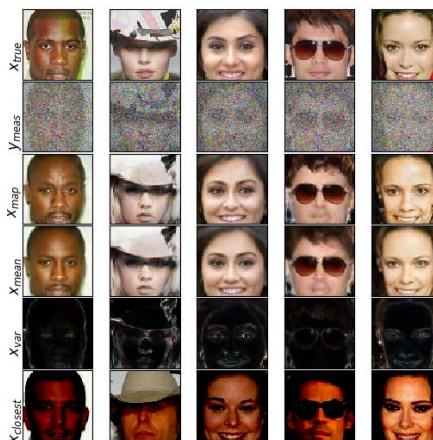


(f)

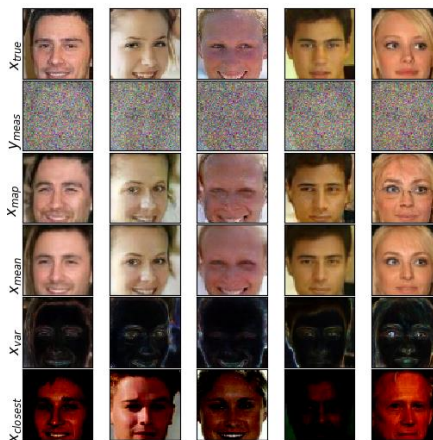
Figure 10: Estimate of the MAP (3rd row), mean (4th row) and variance (5th row) from the limited view of a noisy image (2nd row) using the proposed adaptive method. The window to be revealed at a given iteration (shown in red box) is selected using a variance-driven strategy. Top row indicates ground truth. For all images  $\sigma_y = 1$ .

layer with filter size of  $H \times W$  and number of filters= $C$  with stride= $n$ .

- FC (x,y) indicates fully connected layer with x neurons in input layer and y neurons in output layer.
- BN = Batch norm, LN = Layer norm.
- TrConv = Transposed Convolution.
- LReLU = Leaky ReLU with  $\alpha=0.2$ .



(a)



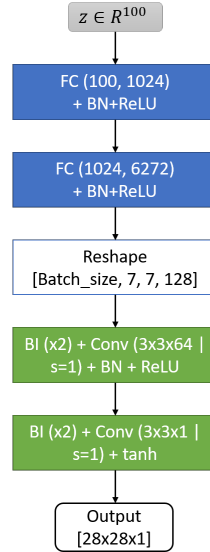
(b)

Figure 11: Estimate of the MAP (3rd row), mean (4th row) and variance (5th row) from a noisy image (2nd row) using the proposed method. Top row shows the ground truth. The last row shows the closest example in training set (by the  $L_2$  measure). For all images  $\sigma_y = 1$ .

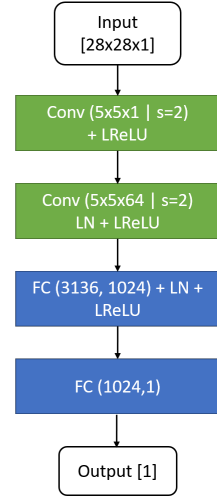
Table 2: Hyper-parameters for WGAN-GP model

Task	CLASS 1 PROBLEMS			CLASS 2 PROBLEMS		
	Image denoising	Physics-based inversion	synthetic	Image classification	Image inpainting and active learning	
Dataset	MNIST	CelebA		MNIST/NotMNIST	MNIST	CelebA
Epochs	1000	500	200	100	1000	500
Learning rate	0.0002	0.0001	0.0002	0.0002	0.0002	0.0001
Batch size	64	64	64	64	64	64
$n_{critic}/n_{gen}$	5	5	1	2	5	5
Momentum params. ( $\beta_1, \beta_2$ )	0.5, 0.999	0.5, 0.999	0.5, 0.999	0.5, 0.999	0.5, 0.999	0.5, 0.999

### Generator

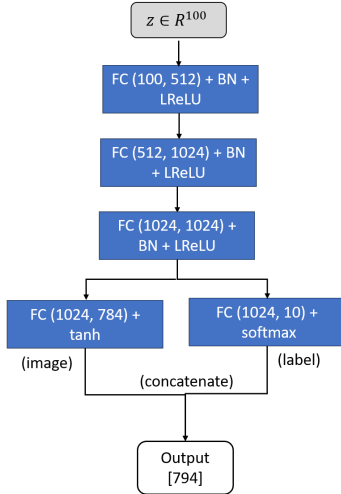


### Discriminator



(a) Architecture for MNIST and synthetic dataset (used in physics-based inference problem)

### Generator



### Discriminator

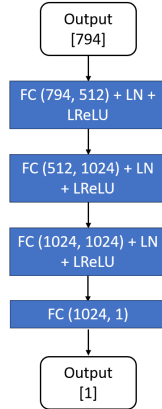
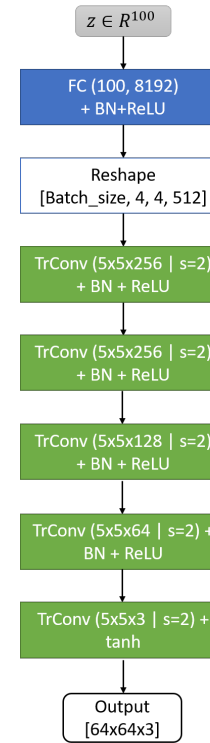
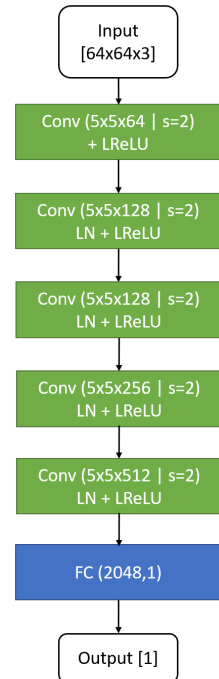


Figure 12: Generator and discriminator architecture used for image classification (hybrid modeling) task for both MNIST and NotMNIST dataset.

### Generator



### Discriminator



(b) Architecture for CelebA dataset

Figure 13: Generator and discriminator architectures for (a) MNIST and synthetic dataset and (b) CelebA dataset used in image denoising, inpainting and physics-driven inversion.