

---

# DEQGAN: Learning the Loss Function for PINNs with Generative Adversarial Networks

---

Blake Bullwinkel<sup>\* 1</sup> Dylan Randle<sup>\* 1 2</sup> Pavlos Protopapas<sup>1</sup> David Sondak<sup>1 3</sup>

## Abstract

Solutions to differential equations are of significant scientific and engineering relevance. Physics-Informed Neural Networks (PINNs) have emerged as a promising method for solving differential equations, but they lack a theoretical justification for the use of any particular loss function. This work presents Differential Equation GAN (DEQGAN), a novel method for solving differential equations using generative adversarial networks to “learn the loss function” for optimizing the neural network. Presenting results on a suite of twelve ordinary and partial differential equations, including the nonlinear Burgers’, Allen-Cahn, Hamilton, and modified Einstein’s gravity equations, we show that DEQGAN<sup>1</sup> can obtain multiple orders of magnitude lower mean squared errors than PINNs that use  $L_2$ ,  $L_1$ , and Huber loss functions. We also show that DEQGAN achieves solution accuracies that are competitive with popular numerical methods. Finally, we present two methods to improve the robustness of DEQGAN to different hyperparameter settings.

## 1. Introduction

In fields such as physics, chemistry, biology, engineering, and economics, differential equations are used to model important and complex phenomena. While numerical methods for solving differential equations perform well and the theory for their stability and convergence is well established, the recent success of deep learning (Krizhevsky et al., 2012; Sutskever et al., 2014; Bahdanau et al., 2015; Vaswani et al., 2017; Mnih et al., 2013; Dabney et al., 2018; Gu et al.,

2017; Silver et al., 2018) has inspired researchers to apply neural networks to solving differential equations, which has given rise to the growing field of Physics-Informed Neural Networks (PINNs) (Raissi et al., 2019; Hagge et al., 2017; Piscopo et al., 2019; Mattheakis et al., 2019; Stevens & Colonius, 2020; Mattheakis et al., 2020; Han et al., 2018; Raissi, 2018; Sirignano & Spiliopoulos, 2018).

In contrast to traditional numerical methods, PINNs: provide solutions that are closed-form (Lagaris et al., 1998), suffer less from the “curse of dimensionality” (Han et al., 2018; Raissi, 2018; Sirignano & Spiliopoulos, 2018; Grohs et al., 2018), provide a more accurate interpolation scheme (Lagaris et al., 1998), and can leverage transfer learning for fast discovery of new solutions (Flamant et al., 2020; Desai et al., 2021). Further, PINNs do not require an underlying grid and offer a meshless approach to solving differential equations. This makes it possible to use trained neural networks, which typically have small memory footprints, to generate solutions over arbitrary grids in a single forward pass.

PINNs have been successfully applied to a wide range of differential equations, but provide no theoretical justification for the use of a particular loss function. In domains outside of differential equations, data following a known noise model (e.g. Gaussian) have clear justification for fitting models with specific loss functions (e.g.  $L_2$ ). In the case of deterministic differential equations, however, there is no noise model and we lack an equivalent justification.

To address this gap in the theory, we propose generative adversarial networks (GANs) (Goodfellow et al., 2014) for solving differential equations in a fully unsupervised manner. Recently, Zeng et al. (2022) showed that adaptively modifying the loss function throughout training can lead to improved solution accuracies. The discriminator network of our GAN-based method, however, can be thought of as “learning the loss function” for optimizing the generator, thereby eliminating the need for a pre-specified loss function and providing even greater flexibility than an adaptive loss. Beyond the context of differential equations, it has also been shown that where classical loss functions struggle to capture complex spatio-temporal dependencies, GANs may be an effective alternative (Larsen et al., 2015; Ledig

<sup>\*</sup>Equal contribution <sup>1</sup>IACS, Harvard University, Cambridge, Massachusetts, USA <sup>2</sup>Amazon Robotics, North Reading, Massachusetts, USA <sup>3</sup>Dassault Systèmes Simulia Inc., Waltham, Massachusetts, USA. Correspondence to: Blake Bullwinkel <jbullwinkel@fas.harvard.edu>.

2nd AI4Science Workshop at the 39<sup>th</sup> International Conference on Machine Learning (ICML), 2022. Copyright 2022 by the author(s).

<sup>1</sup>We provide our PyTorch code at <https://github.com/dylanrandle/denn>

et al., 2016; Karras et al., 2018).

Our contributions in this work are summarized as follows:

- We present Differential Equation GAN (DEQGAN), a novel method for solving differential equations in a fully unsupervised manner using generative adversarial networks.
- We highlight the advantage of “learning the loss function” with a GAN rather than using a pre-specified loss function by showing that PINNs trained using  $L_2$ ,  $L_1$ , and Huber losses have variable performance and fail to solve the modified Einstein’s gravity equations (Chantada et al., 2022).
- We present results on a suite of twelve ordinary differential equations (ODEs) and partial differential equations (PDEs), including highly nonlinear problems, showing that our method produces solutions with multiple orders of magnitude lower mean squared errors than PINNs that use  $L_2$ ,  $L_1$ , and Huber loss functions.
- We show that DEQGAN achieves solution accuracies that are competitive with popular numerical methods, including the fourth-order Runge-Kutta and second-order finite difference methods.
- We present two techniques to improve the training stability of DEQGAN that are applicable to other GAN-based methods and PINN approaches to solving differential equations.

## 2. Related Work

A variety of neural network methods have been developed for solving differential equations. Some of these are supervised and learn the dynamics of real-world systems from data (Raissi et al., 2019; Choudhary et al., 2020; Greidanus et al., 2019; Bertalan et al., 2019). Others are semi-supervised, learning general solutions to a differential equation and extracting a best fit solution based on observational data (Paticchio et al., 2020). Our work falls under the category of *unsupervised* neural network methods, which are trained in a data-free manner that depends solely on the equation residuals. Unsupervised neural networks have been applied to a wide range of ODEs (Lagaris et al., 1998; Flamant et al., 2020; Mattheakis et al., 2020; 2021) and PDEs (Han et al., 2018; Sirignano & Spiliopoulos, 2018; Raissi, 2018; Stevens & Colonius, 2020), primarily use feed-forward architectures, and require the specification of a particular loss function computed over the equation residuals.

Goodfellow et al. (2014) introduced the idea of learning generative models with neural networks and an adversarial training algorithm, called generative adversarial networks

(GANs). To solve issues of GAN training instability, Arjovsky et al. (2017) introduced a formulation of GANs based on the Wasserstein distance, and Gulrajani et al. (2017) added a gradient penalty to approximately enforce a Lipschitz constraint on the discriminator. Miyato et al. (2018) introduced an alternative method for enforcing the Lipschitz constraint with a spectral normalization technique that outperforms the former method on some problems.

Further work has applied GANs to differential equations with solution data used for supervision. Yang et al. (2018) apply GANs to stochastic differential equations by using “snapshots” of ground-truth data for semi-supervised training. A project by students at Stanford (Subramanian et al., 2018) employed GANs to perform “turbulence enrichment” of solution data in a manner akin to that of super-resolution for images proposed by Ledig et al. (2016). Our work distinguishes itself from other GAN-based approaches for solving differential equations by being *fully unsupervised*, and removing the dependence on using supervised training data (i.e. solutions of the equation).

## 3. Background

### 3.1. Unsupervised Neural Networks for Differential Equations

Early work by Dissanayake & Phan-Thien (1994) proposed solving initial value problems in an unsupervised manner with neural networks. In this work, we extend their approach to handle spatial domains and multidimensional problems. In particular, we consider general differential equations of the form

$$F\left(t, \mathbf{x}, \Psi(t, \mathbf{x}), \frac{d\Psi}{dt}, \frac{d^2\Psi}{dt^2}, \dots, \Delta\Psi, \Delta^2\Psi, \dots\right) = 0 \quad (1)$$

where  $\Psi(t, \mathbf{x})$  is the desired solution,  $d\Psi/dt$  and  $d^2\Psi/dt^2$  represent the first and second time derivatives,  $\Delta\Psi$  and  $\Delta^2\Psi$  are the first and second spatial derivatives, and the system is subject to certain initial and boundary conditions. The learning problem can then be formulated as minimizing the sum of squared residuals (i.e., the squared  $L_2$  loss) of the above equation

$$\min_{\theta} \sum_{(t, \mathbf{x}) \in \mathcal{D}} F\left(t, \mathbf{x}, \Psi_{\theta}(t, \mathbf{x}), \frac{d\Psi_{\theta}}{dt}, \frac{d^2\Psi_{\theta}}{dt^2}, \dots, \Delta\Psi_{\theta}, \Delta^2\Psi_{\theta}, \dots\right) \quad (2)$$

where  $\Psi_{\theta}$  is a neural network parameterized by  $\theta$ ,  $\mathcal{D}$  is the domain of the problem, and derivatives are computed with automatic differentiation. This allows backpropagation

(Hecht-Nielsen, 1992) to be used to train the neural network to satisfy the differential equation. We apply this formalism to both initial and boundary value problems, including multidimensional problems, as detailed in Appendix A.2.

### 3.2. Generative Adversarial Networks

Generative adversarial networks (GANs) (Goodfellow et al., 2014) are generative models that use two neural networks to induce a generative distribution  $p(x)$  of the data by formulating the inference problem as a two-player, zero-sum game.

The generative model first samples a latent random variable  $z \sim \mathcal{N}(0, 1)$ , which is used as input into the generator  $G$  (e.g., a neural network). A discriminator  $D$  is trained to classify whether its input was sampled from the generator (i.e., “fake”) or from a reference data set (i.e., “real”).

Informally, the process of training GANs proceeds by optimizing a minimax objective over the generator and discriminator such that the generator attempts to trick the discriminator to classify “fake” samples as “real”. Formally, one optimizes

$$\min_G \max_D V(D, G) = \min_G \max_D \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [1 - \log D(G(z))] \quad (3)$$

where  $x \sim p_{\text{data}}(x)$  denotes samples from the empirical data distribution, and  $p_z \sim \mathcal{N}(0, 1)$  samples in latent space (Goodfellow et al., 2014). In practice, the optimization alternates between gradient ascent and descent steps for  $D$  and  $G$  respectively.

#### 3.2.1. TWO TIME-SCALE UPDATE RULE

Heusel et al. (2017) proposed the two time-scale update rule (TTUR) for training GANs, a method in which the discriminator and generator are trained with separate learning rates. They showed that their method led to improved performance and proved that, in some cases, TTUR ensures convergence to a stable local Nash equilibrium. One intuition for TTUR comes from the potentially different loss surfaces of the discriminator and generator. Allowing learning rates to be tuned to a particular loss surface can enable more efficient gradient-based optimization. We make use of TTUR throughout this paper as an instrumental lever when tuning GANs to reach desired performance.

#### 3.2.2. SPECTRAL NORMALIZATION

Proposed by Miyato et al. (2018), Spectrally Normalized GAN (SN-GAN) is a method for controlling exploding discriminator gradients when optimizing Equation 3 that leverages a novel weight normalization technique. The key idea is to control the Lipschitz constant of the discriminator by

constraining the spectral norm of each layer in the discriminator. Specifically, the authors propose dividing the weight matrices  $W_i$  of each layer  $i$  by their spectral norm  $\sigma(W_i)$

$$W_{SN,i} = \frac{W_i}{\sigma(W_i)}, \quad (4)$$

where

$$\sigma(W_i) = \max_{\|h_i\|_2 \leq 1} \|W_i h_i\|_2 \quad (5)$$

and  $h_i$  denotes the input to layer  $i$ . The authors prove that this normalization technique bounds the Lipschitz constant of the discriminator above by 1, thus strictly enforcing the 1-Lipschitz constraint on the discriminator. In our experiments, adopting the SN-GAN formulation led to even better performance than WGAN-GP (Arjovsky et al., 2017; Gulrajani et al., 2017).

### 3.3. Guaranteeing Initial & Boundary Conditions

Lagaris et al. (1998) showed that it is possible to exactly satisfy initial and boundary conditions by adjusting the output of the neural network. For example, consider adjusting the neural network output  $\Psi_\theta(t, \mathbf{x})$  to satisfy the initial condition  $\Psi_\theta(t, \mathbf{x})|_{t=t_0} = x_0$ . We can apply the re-parameterization

$$\tilde{\Psi}_\theta(t, \mathbf{x}) = x_0 + t\Psi_\theta(t, \mathbf{x}) \quad (6)$$

which exactly satisfies the initial condition. Mattheakis et al. (2020) proposed an augmented re-parameterization

$$\begin{aligned} \tilde{\Psi}_\theta(t, \mathbf{x}) &= \Phi(\Psi_\theta(t, \mathbf{x})) \\ &= x_0 + (1 - e^{-(t-t_0)}) \Psi_\theta(t, \mathbf{x}) \end{aligned} \quad (7)$$

that further improved training convergence. Intuitively, Equation 7 adjusts the output of the neural network  $\Psi_\theta(t, \mathbf{x})$  to be exactly  $x_0$  when  $t = t_0$ , and decays this constraint exponentially in  $t$ . Chen et al. (2020) provide re-parameterizations to satisfy a range of other conditions, including Dirichlet and Neumann boundary conditions, which we employ in our experiments and detail in Appendix A.2.

### 3.4. Residual Connections

He et al. (2015) showed that the addition of residual connections improves deep neural network training. We employ residual connections in our networks, as they allow gradients to flow more easily through the models and thereby reduce numerical instability. Residual connections augment a typical activation with the identity operation.

$$y = \mathcal{F}(x, W_i) + x \quad (8)$$

where  $\mathcal{F}$  is the activation function,  $x$  is the input to the unit,  $W_i$  are the weights and  $y$  is the output of the unit. This acts as a “skip connection”, allowing inputs and gradients to forego the nonlinear component.

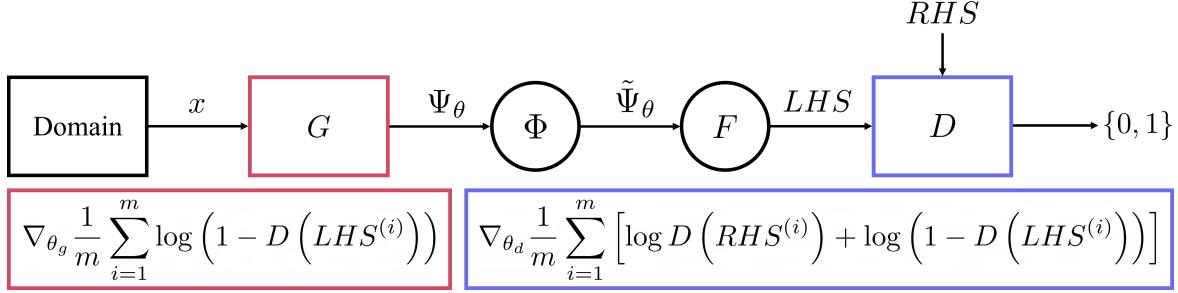


Figure 1. Schematic representation of DEQGAN. We pass input points  $x$  to a generator  $G$ , which produces candidate solutions  $\Psi_\theta$ . Then we analytically adjust these solutions according to  $\Phi$  and apply automatic differentiation to construct  $LHS$  from the differential equation  $F$ .  $RHS$  and  $LHS$  are passed to a discriminator  $D$ , which is trained to classify them as “real” and “fake” respectively.

## 4. Differential Equation GAN

In this section, we present our method, Differential Equation GAN (DEQGAN), which trains a GAN to solve differential equations in a *fully unsupervised* manner. To do this, we rearrange the differential equation so that the left-hand side ( $LHS$ ) contains all the terms which depend on the generator (e.g.  $\Psi$ ,  $d\Psi/dt$ ,  $\Delta\Psi$ , etc.) and the right-hand side ( $RHS$ ) contains only constants (e.g. zero).

During training, we sample points from the domain  $(t, \mathbf{x}) \sim \mathcal{D}$  and use them as input to a generator  $G(x)$ , which produces candidate solutions  $\Psi_\theta$ . We sample points from a noisy grid that spans  $\mathcal{D}$ , which we found reduced interpolation error in comparison to sampling points from a fixed grid. We then adjust  $\Psi_\theta$  for initial or boundary conditions to obtain the re-parameterized output  $\tilde{\Psi}_\theta$ , construct the  $LHS$  from the differential equation  $F$  using automatic differentiation

$$LHS = F \left( t, \mathbf{x}, \tilde{\Psi}_\theta(t, \mathbf{x}), \frac{d\tilde{\Psi}_\theta}{dt}, \frac{d^2\tilde{\Psi}_\theta}{dt^2}, \dots, \Delta\tilde{\Psi}_\theta, \Delta^2\tilde{\Psi}_\theta, \dots \right) \quad (9)$$

and set  $RHS$  to zero. Note that in our procedure, we add Gaussian noise to  $RHS$ , which we found improves training of the discriminator, as described in Section 4.1. Training proceeds in a manner similar to that of traditional GANs. We update the weights of the generator  $G$  and the discriminator  $D$  according to the gradients

$$g_G = \nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left( 1 - D \left( LHS^{(i)} \right) \right), \quad (10)$$

$$g_D = \nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D \left( RHS^{(i)} \right) + \log \left( 1 - D \left( LHS^{(i)} \right) \right) \right] \quad (11)$$

where  $LHS^{(i)}$  is the output of  $G(x^{(i)})$  after adjusting for initial or boundary conditions and constructing the  $LHS$  from  $F$ . Note that we perform stochastic gradient *descent* for  $G$  (gradient steps  $\propto -g_G$ ), and stochastic gradient *ascent* for  $D$  (gradient steps  $\propto g_D$ ). We provide a schematic representation of DEQGAN in Figure 1 and detail the training steps in Algorithm 1.

### Algorithm 1 DEQGAN

**Input:** Differential equation  $F$ , generator  $G(\cdot; \theta_g)$ , discriminator  $D(\cdot; \theta_d)$ , grid  $x$  of  $m$  points with spacing  $\Delta x$ , perturbation precision  $\tau$ , re-parameterization function  $\Phi$ , total steps  $N$ , learning rates  $\eta_G, \eta_D$ , Adam optimizer parameters  $\beta_{G1}, \beta_{G2}, \beta_{D1}, \beta_{D2}$

**for**  $i = 1$  **to**  $N$  **do**

- for**  $j = 1$  **to**  $m$  **do**
- Perturb  $j$ -th point in mesh  $x_s^{(j)} = x^{(j)} + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \frac{\Delta x}{\tau})$
- Forward pass  $\Psi_\theta = G(x_s^{(j)})$
- Analytic re-parameterization  $\tilde{\Psi}_\theta = \Phi(\Psi_\theta)$
- Compute  $LHS^{(j)}$  (Equation 9)
- Set  $RHS^{(j)} = 0$

**end for**

Compute gradients  $g_G, g_D$  (Equation 10 & 11)

Update generator

$$\theta_g \leftarrow \text{Adam}(\theta_g, -g_G, \eta_G, \beta_{G1}, \beta_{G2})$$

Update discriminator

$$\theta_d \leftarrow \text{Adam}(\theta_d, g_D, \eta_D, \beta_{D1}, \beta_{D2})$$

**end for**

**Output:**  $G$

Informally, our algorithm trains a GAN by setting the “fake” component to be the  $LHS$  (in our formulation, the residuals of the equation) and the “real” component to be the  $RHS$  of the equation. This results in a GAN that learns to produce solutions that make  $LHS$  indistinguishable from  $RHS$ , thereby approximately solving the differential equation.

Table 1. Summary of Experiments

Key	Equation	Class	Order	Linear
EXP	$\dot{x}(t) + x(t) = 0$	ODE	1 <sup>st</sup>	Yes
SHO	$\ddot{x}(t) + x(t) = 0$	ODE	2 <sup>nd</sup>	Yes
NLO	$\ddot{x}(t) + 2\beta\dot{x}(t) + \omega^2x(t) + \phi x(t)^2 + \epsilon x(t)^3 = 0$	ODE	2 <sup>nd</sup>	No
COO	$\begin{cases} \dot{x}(t) = -ty \\ \dot{y}(t) = tx \end{cases}$	ODE	1 <sup>st</sup>	Yes
SIR	$\begin{cases} \dot{S}(t) = -\beta I(t)S(t)/N \\ \dot{I}(t) = \beta I(t)S(t)/N - \gamma I(t) \\ \dot{R}(t) = \gamma I(t) \end{cases}$	ODE	1 <sup>st</sup>	No
HAM	$\begin{cases} \dot{x}(t) = p_x \\ \dot{y}(t) = p_y \\ \dot{p}_x(t) = -V_x \\ \dot{p}_y(t) = -V_y \end{cases}$	ODE	1 <sup>st</sup>	No
EIN	$\begin{cases} \dot{x}(z) = \frac{1}{z+1}(-\Omega - 2v + x + 4y + xv + x^2) \\ \dot{y}(z) = \frac{-1}{z+1}(vx\Gamma(r) - xy + 4y - 2yv) \\ \dot{v}(z) = \frac{-v}{z+1}(x\Gamma(r) + 4 - 2v) \\ \dot{\Omega}(z) = \frac{\Omega}{z+1}(-1 + 2v + x) \\ \dot{r}(z) = \frac{-r\Gamma(r)x}{z+1} \end{cases}$	ODE	1 <sup>st</sup>	No
POS	$u_{xx} + u_{yy} = 2x(y-1)(y-2x+xy+2)e^{x-y}$	PDE	2 <sup>nd</sup>	Yes
HEA	$u_t = \kappa u_{xx}$	PDE	2 <sup>nd</sup>	Yes
WAV	$u_{tt} = c^2 u_{xx}$	PDE	2 <sup>nd</sup>	Yes
BUR	$u_t + uu_x - \nu u_{xx} = 0$	PDE	2 <sup>nd</sup>	No
ACA	$u_t - \epsilon u_{xx} - u + u^3 = 0$	PDE	2 <sup>nd</sup>	No

#### 4.1. Instance Noise

While GANs have achieved state of the art results on a wide range of generative modeling tasks, they are often difficult to train. As a result, much recent work on GANs has been dedicated to improving their sensitivity to hyperparameters and training stability (Salimans et al., 2016; Gulrajani et al., 2017; Sønderby et al., 2016; Arjovsky & Bottou, 2017; Karnewar et al., 2019; Kodali et al., 2017; Arjovsky et al., 2017; Berthelot et al., 2017; Mirza & Osindero, 2014; Miyato et al., 2018). In our experiments, we found that DEQGAN could also be sensitive to hyperparameters, such as the Adam optimizer parameters shown in Algorithm 1.

Sønderby et al. (2016) note that the convergence of GANs relies on the existence of a unique optimal discriminator that separates the distribution of “fake” samples  $p_{\text{fake}}$  produced by the generator, and the distribution of the “real” data  $p_{\text{data}}$ . In practice, however, there may be many near-optimal discriminators that pass very different gradients to the generator, depending on their initialization. Arjovsky & Bottou (2017) proved that this problem will arise when there is insufficient overlap between the supports of  $p_{\text{fake}}$  and  $p_{\text{data}}$ . In the DEQGAN training algorithm, setting  $RHS = 0$  constrains  $p_{\text{data}}$  to the Dirac delta function  $\delta(0)$ , and therefore

the distribution of “real” data to a zero-dimensional manifold. This makes it unlikely that  $p_{\text{fake}}$  and  $p_{\text{data}}$  will share support in a high-dimensional space.

The solution proposed by (Sønderby et al., 2016; Arjovsky & Bottou, 2017) is to add “instance noise” to  $p_{\text{fake}}$  and  $p_{\text{data}}$  to encourage their overlap. This amounts to adding noise to the *LHS* and the *RHS*, respectively, at each iteration of Algorithm 1. Because this makes the discriminator’s job more difficult, we add Gaussian noise with standard deviation equal to the difference between the generator and discriminator losses,  $L_g$  and  $L_d$ , i.e.

$$\varepsilon = \mathcal{N}(0, \sigma^2), \quad \sigma = \text{ReLU}(L_g - L_d) \quad (12)$$

As the generator and discriminator reach equilibrium, Equation 12 will naturally converge to zero. We use the ReLU function because when  $L_d > L_g$ , the generator is already able to fool the discriminator, suggesting that additional noise should not be used. In Section 5.2, we conduct an ablation study and find that this improves the ability of DEQGAN to produce accurate solutions across a range of hyperparameter settings.

Table 2. Experimental Results

Key	Mean Squared Error				
	$L_1$	$L_2$	Huber	DEQGAN	Numerical
EXP	$3 \cdot 10^{-3}$	$2 \cdot 10^{-5}$	$1 \cdot 10^{-5}$	$3 \cdot 10^{-16}$	$2 \cdot 10^{-14}$ (RK4)
SHO	$9 \cdot 10^{-6}$	$1 \cdot 10^{-10}$	$6 \cdot 10^{-11}$	$4 \cdot 10^{-13}$	$1 \cdot 10^{-11}$ (RK4)
NLO	$6 \cdot 10^{-2}$	$1 \cdot 10^{-9}$	$9 \cdot 10^{-10}$	$1 \cdot 10^{-12}$	$4 \cdot 10^{-11}$ (RK4)
COO	$5 \cdot 10^{-1}$	$1 \cdot 10^{-7}$	$1 \cdot 10^{-7}$	$1 \cdot 10^{-8}$	$2 \cdot 10^{-9}$ (RK4)
SIR	$7 \cdot 10^{-5}$	$3 \cdot 10^{-9}$	$1 \cdot 10^{-9}$	$1 \cdot 10^{-10}$	$5 \cdot 10^{-13}$ (RK4)
HAM	$1 \cdot 10^{-1}$	$2 \cdot 10^{-7}$	$9 \cdot 10^{-8}$	$1 \cdot 10^{-10}$	$7 \cdot 10^{-14}$ (RK4)
EIN	$6 \cdot 10^{-2}$	$2 \cdot 10^{-2}$	$1 \cdot 10^{-2}$	$3 \cdot 10^{-4}$	$4 \cdot 10^{-7}$ (RK4)
POS	$4 \cdot 10^{-6}$	$1 \cdot 10^{-10}$	$6 \cdot 10^{-11}$	$4 \cdot 10^{-13}$	$3 \cdot 10^{-10}$ (FD)
HEA	$6 \cdot 10^{-3}$	$3 \cdot 10^{-5}$	$1 \cdot 10^{-5}$	$6 \cdot 10^{-10}$	$4 \cdot 10^{-7}$ (FD)
WAV	$6 \cdot 10^{-2}$	$4 \cdot 10^{-5}$	$6 \cdot 10^{-4}$	$1 \cdot 10^{-8}$	$7 \cdot 10^{-5}$ (FD)
BUR	$4 \cdot 10^{-3}$	$2 \cdot 10^{-4}$	$1 \cdot 10^{-4}$	$4 \cdot 10^{-6}$	$1 \cdot 10^{-3}$ (FD)
ACA	$6 \cdot 10^{-2}$	$9 \cdot 10^{-3}$	$4 \cdot 10^{-3}$	$3 \cdot 10^{-3}$	$2 \cdot 10^{-4}$ (FD)

## 4.2. Residual Monitoring

One of the attractive properties of Algorithm 1 is that the “fake”  $LHS$  vector of equation residuals gives a direct measure of solution quality at each training iteration. We observe that when DEQGAN training becomes unstable, the  $LHS$  tends to oscillate wildly, while it decreases steadily throughout training for successful runs. By monitoring the  $L_1$  norm of the  $LHS$  in the first 25% of training iterations, we are able to easily detect and terminate poor-performing runs if the variance of these values exceeds some threshold. We provide further details on this method in Appendix A.6 and experimentally demonstrate that it is able to distinguish between DEQGAN runs that end in high and low mean squared errors in Section 5.2.

Table 3. Ablation Study Results

% Runs with High MSE ( $\geq 10^{-5}$ )		
	Original	Residual Monitoring
Original	12.4	0.4
Instance Noise	8.0	0.0

## 5. Experiments

We conducted experiments on a suite of twelve differential equations (Table 1), including highly nonlinear PDEs and systems of ODEs, comparing DEQGAN to classical unsupervised PINNs that use (squared)  $L_2$ ,  $L_1$ , and Huber (1964) loss functions. We also report results obtained by the fourth-order Runge-Kutta (RK4) and second-order finite difference (FD) numerical methods for initial and boundary value problems, respectively. The numerical solutions were computed over meshes containing the same number of points that were used to train the neural network methods.

Details for each experiment, including exact problem specifications and hyperparameters, are provided in Appendix A.2 and A.4.

## 5.1. DEQGAN vs. Classical PINNs

We report the mean squared error of the solution obtained by each method, computed against known solutions obtained either analytically or with high-quality numerical solvers (Virtanen et al., 2020; Brunton & Kutz, 2019). We added residual connections between neighboring layers of all models, applied spectral normalization to the discriminator, added instance noise to the  $p_{\text{fake}}$  and  $p_{\text{real}}$ , and used residual monitoring to terminate poor-performing runs in the first 25% of training iterations. Results were obtained with hyperparameters tuned for DEQGAN. In Appendix A.5, we tuned each classical PINN method for comparison, but did not observe a significant difference.

Table 2 reports the lowest mean squared error obtained by each method across ten different model weight initializations. We see that DEQGAN obtains lower mean squared errors than classical PINNs that use  $L_2$ ,  $L_1$ , and Huber loss functions for all twelve problems, often by several orders of magnitude. DEQGAN also achieves solution accuracies that are competitive with the RK4 and FD numerical methods.

Figure 2 plots the mean squared error vs. training iteration for six challenging equations and highlights multiple advantages of using DEQGAN over a pre-specified loss function (equivalent plots for the other six problems are provided in Appendix A.3). In particular, there is considerable variation in the quality of the solutions obtained by the classical PINNs. For example, while Huber performs better than  $L_2$  on the Allen-Cahn PDE, it is outperformed by  $L_2$  on the wave equation. Furthermore, Figure 2(f) shows that the  $L_2$ ,  $L_1$  and Huber losses all fail to converge to an accurate

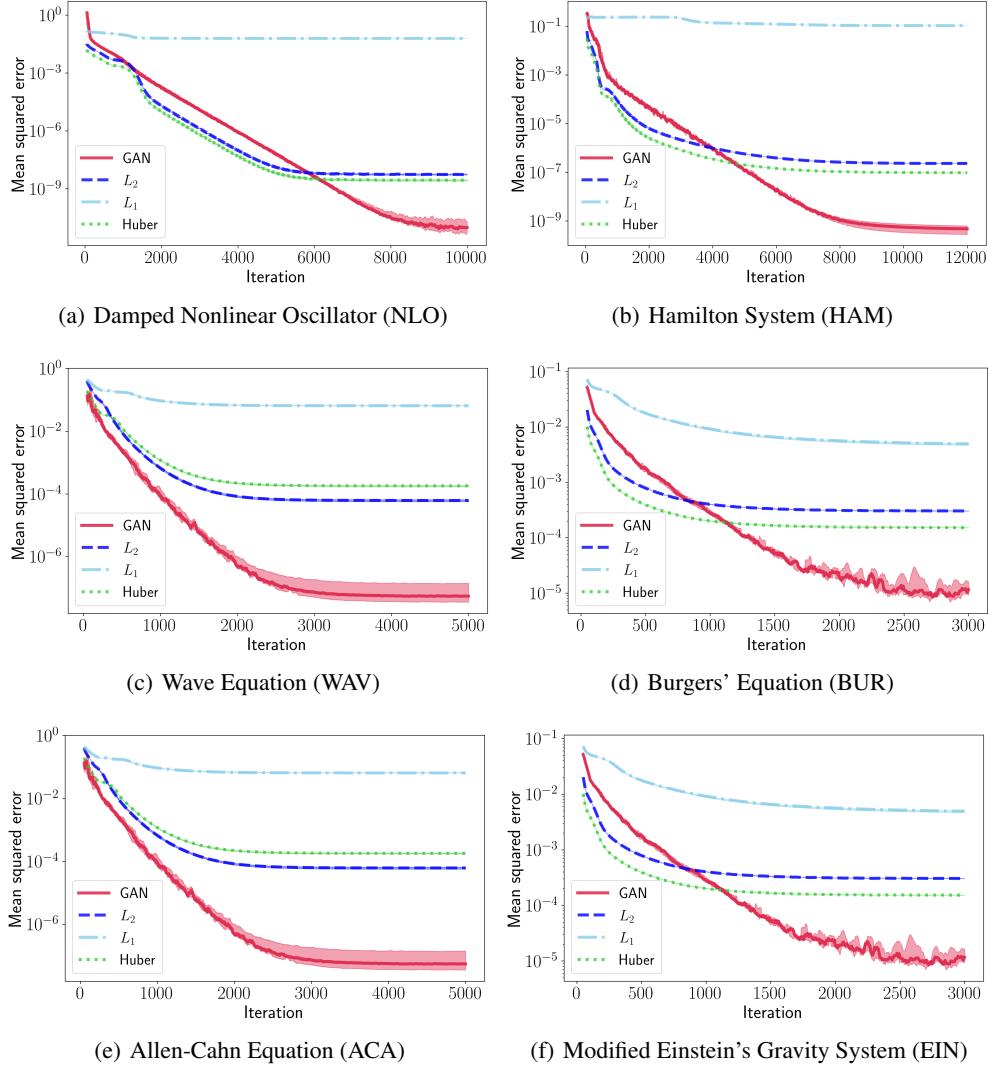


Figure 2. Mean squared errors vs. iteration for DEQGAN,  $L_2$ ,  $L_1$ , and Huber loss for six equations. We perform ten randomized trials and plot the median (bold) and (25, 75) percentile range (shaded). We smooth the values using a simple moving average with window size 50.

solution to the modified Einstein's gravity equations. Although this system has previously been solved using PINNs, the networks relied on a custom loss function that incorporated equation-specific parameters (Chantada et al., 2022). DEQGAN, however, is able to *automatically* learn a loss function that optimizes the generator to produce accurate solutions. DEQGAN solutions to four example equations are visualized in Figure 3, and similar plots for the other experiments are provided in Appendix A.2.

## 5.2. DEQGAN Training Stability: Ablation Study

In our experiments, we used instance noise to adaptively improve the training convergence of DEQGAN and employed residual monitoring to terminate poor-performing runs early. To quantify the increased robustness offered by

these techniques, we performed an ablation study comparing the percentage of high MSE ( $\geq 10^{-5}$ ) runs obtained by 500 randomized DEQGAN runs on the exponential decay equation. This experimental setup is detailed further in Appendix A.7.

Table 3 compares the percentage of high MSE runs with and without instance noise and residual monitoring. We see that adding instance noise decreased the percentage of runs with high MSE and that residual monitoring is highly effective at filtering out poor performing runs. When used together, these techniques eliminated all runs with MSE  $\geq 10^{-5}$ . These results agree with previous works, which have found that instance noise can improve the convergence of other GAN training algorithms (Sønderby et al., 2016; Arjovsky & Bottou, 2017). Further, they suggest that residual

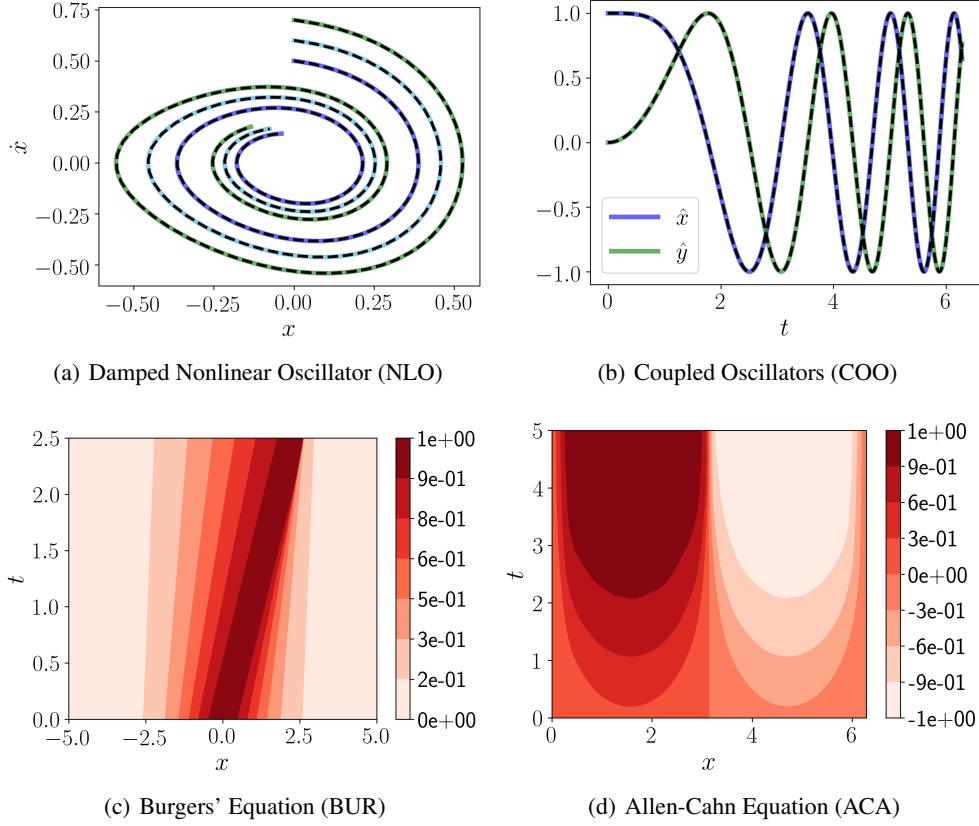


Figure 3. Visualization of DEQGAN solutions to four equations. The top left figure plots the phase space of the DEQGAN solutions (solid color lines) obtained for three initial conditions on the NLO problem, which is solved as a second-order ODE, and known solutions computed by a numerical integrator (dashed black lines). The figure to the right plots the DEQGAN solution to the COO problem, which is solved as a system of two first-order ODEs. The second row shows contour plots of the solutions obtained by DEQGAN on the BUR and ACA problems, both nonlinear PDEs.

monitoring provides a useful performance metric that could be applied to other PINN methods for solving differential equations.

## 6. Conclusion

PINNs offer a promising approach to solving differential equations and to applying deep learning methods to challenging problems in science and engineering. Classical PINNs, however, lack a theoretical justification for the use of any particular loss function. In this work, we presented Differential Equation GAN (DEQGAN), a novel method that leverages GAN-based adversarial training to “learn” the loss function for solving differential equations with PINNs. We demonstrated the advantage of this approach in comparison to using classical PINNs with pre-specified loss functions, which showed varied performance and failed to converge to an accurate solution to the modified Einstein’s gravity equations. In general, we demonstrated that our method can obtain multiple orders of magnitude lower mean

squared errors than PINNs that use  $L_2$ ,  $L_1$  and Huber loss functions, including on highly nonlinear PDEs and systems of ODEs. Further, we showed that DEQGAN achieves solution accuracies that are competitive with the fourth-order Runge Kutta and second-order finite difference numerical methods. Finally, we found that instance noise improved training stability and that residual monitoring provides a useful performance metric for PINNs. While the equation residuals are a good measure of solution quality, PINNs lack the error bounds enjoyed by numerical methods. Formalizing these bounds is an interesting avenue for future work and would enable PINNs to be more safely deployed in real-world applications. Further, while our results evidence the advantage of “learning the loss function” with a GAN, understanding exactly what the discriminator learns is an open problem. Post-hoc explainability methods, for example, might provide useful tools for characterizing the differences between classical losses and the loss functions learned by DEQGAN, which could deepen our understanding of PINN optimization more generally.

## References

- Arjovsky, M. and Bottou, L. Towards principled methods for training generative adversarial networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL [https://openreview.net/forum?id=Hk4\\_qw5xe](https://openreview.net/forum?id=Hk4_qw5xe).
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein gan, 2017.
- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1409.0473>.
- Bertalan, T., Dietrich, F., Mezić, I., and Kevrekidis, I. G. On learning hamiltonian systems from data. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 29(12):121107, dec 2019. doi: 10.1063/1.5128231. URL <https://doi.org/10.1063%2F1.5128231>.
- Berthelot, D., Schumm, T., and Metz, L. BEGAN: boundary equilibrium generative adversarial networks. *CoRR*, abs/1703.10717, 2017. URL <http://arxiv.org/abs/1703.10717>.
- Brunton, S. L. and Kutz, J. N. *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*. Cambridge University Press, 2019. doi: 10.1017/9781108380690.
- Chantada, A. T., Landau, S. J., Protopapas, P., Scóccola, C. G., and Garraffo, C. Cosmological informed neural networks to solve the background dynamics of the universe, 2022. URL <https://arxiv.org/abs/2205.02945>.
- Chen, F., Sondak, D., Protopapas, P., Mattheakis, M., Liu, S., Agarwal, D., and Di Giovanni, M. Neurodiffeq: A python package for solving differential equations with neural networks. *Journal of Open Source Software*, 5(46):1931, 2020.
- Choudhary, A., Lindner, J., Holliday, E., Miller, S., Sinha, S., and Ditto, W. Physics-enhanced neural networks learn order and chaos. *Physical Review E*, 101, 06 2020. doi: 10.1103/PhysRevE.101.062207.
- Dabney, W., Rowland, M., Bellemare, M. G., and Munos, R. Distributional reinforcement learning with quantile regression. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Desai, S., Mattheakis, M., Joy, H., Protopapas, P., and Roberts, S. One-shot transfer learning of physics-informed neural networks, 2021. URL <https://arxiv.org/abs/2110.11286>.
- Dissanayake, M. and Phan-Thien, N. Neural-network-based approximations for solving partial differential equations. *Communications in Numerical Methods in Engineering*, 10(3):195–201, 1994.
- Flamant, C., Protopapas, P., and Sondak, D. Solving differential equations using neural network solution bundles, 2020. URL <https://arxiv.org/abs/2006.14372>.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial networks, 2014.
- Greydanus, S., Dzamba, M., and Yosinski, J. Hamiltonian neural networks, 2019. URL <https://arxiv.org/abs/1906.01563>.
- Grohs, P., Hornung, F., Jentzen, A., and von Wurstemberger, P. A proof that artificial neural networks overcome the curse of dimensionality in the numerical approximation of black-scholes partial differential equations, 2018. URL <https://arxiv.org/abs/1809.02362>.
- Gu, S., Holly, E., Lillicrap, T., and Levine, S. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *2017 IEEE international conference on robotics and automation (ICRA)*, pp. 3389–3396. IEEE, 2017.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. Improved training of wasserstein gans, 2017.
- Hagge, T., Stinis, P., Yeung, E., and Tartakovsky, A. M. Solving differential equations with unknown constitutive relations as recurrent neural networks, 2017.
- Han, J., Jentzen, A., and E, W. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018. ISSN 0027-8424. doi: 10.1073/pnas.1718942115. URL <https://www.pnas.org/content/115/34/8505>.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL <https://arxiv.org/abs/1512.03385>.
- Hecht-Nielsen, R. Theory of the backpropagation neural network. In *Neural networks for perception*, pp. 65–93. Elsevier, 1992.

- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Klambauer, G., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a nash equilibrium. *CoRR*, abs/1706.08500, 2017. URL <http://arxiv.org/abs/1706.08500>.
- Huber, P. J. Robust estimation of a location parameter. *Ann. Math. Statist.*, 35(1):73–101, 03 1964. doi: 10.1214/aoms/1177703732. URL <https://doi.org/10.1214/aoms/1177703732>.
- Karnewar, A., Wang, O., and Iyengar, R. S. MSG-GAN: multi-scale gradient GAN for stable image synthesis. *CoRR*, abs/1903.06048, 2019. URL <http://arxiv.org/abs/1903.06048>.
- Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. *CoRR*, abs/1812.04948, 2018. URL <http://arxiv.org/abs/1812.04948>.
- Kodali, N., Abernethy, J. D., Hays, J., and Kira, Z. How to train your DRAGAN. *CoRR*, abs/1705.07215, 2017. URL <http://arxiv.org/abs/1705.07215>.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 25*, pp. 1097–1105. Curran Associates, Inc., 2012.
- Lagaris, I., Likas, A., and Fotiadis, D. Artificial neural networks for solving ordinary and partial differential equations. *IEEE Transactions on Neural Networks*, 9(5):987–1000, 1998. ISSN 1045-9227. doi: 10.1109/72.712178. URL <http://dx.doi.org/10.1109/72.712178>.
- Larsen, A. B. L., Sønderby, S. K., and Winther, O. Autoencoding beyond pixels using a learned similarity metric. *CoRR*, abs/1512.09300, 2015. URL <http://arxiv.org/abs/1512.09300>.
- Ledig, C., Theis, L., Huszar, F., Caballero, J., Aitken, A. P., Tejani, A., Totz, J., Wang, Z., and Shi, W. Photo-realistic single image super-resolution using a generative adversarial network. *CoRR*, abs/1609.04802, 2016. URL <http://arxiv.org/abs/1609.04802>.
- Liaw, R., Liang, E., Nishihara, R., Moritz, P., Gonzalez, J. E., and Stoica, I. Tune: A research platform for distributed model selection and training. *CoRR*, abs/1807.05118, 2018. URL <http://arxiv.org/abs/1807.05118>.
- Mattheakis, M., Protopapas, P., Sondak, D., Giovanni, M. D., and Kaxiras, E. Physical symmetries embedded in neural networks, 2019.
- Mattheakis, M., Sondak, D., Dogra, A. S., and Protopapas, P. Hamiltonian neural networks for solving differential equations, 2020.
- Mattheakis, M., Joy, H., and Protopapas, P. Unsupervised reservoir computing for solving ordinary differential equations, 2021. URL <https://arxiv.org/abs/2108.11417>.
- Mirza, M. and Osindero, S. Conditional generative adversarial nets, 2014.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. *CoRR*, abs/1802.05957, 2018. URL <http://arxiv.org/abs/1802.05957>.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Paticchio, A., Scarlatti, T., Mattheakis, M., Protopapas, P., and Brambilla, M. Semi-supervised neural networks solve an inverse problem for modeling covid-19 spread. 2020. doi: 10.48550/ARXIV.2010.05074. URL <https://arxiv.org/abs/2010.05074>.
- Piscopo, M. L., Spannowsky, M., and Waite, P. Solving differential equations with neural networks: Applications to the calculation of cosmological phase transitions. *Phys. Rev. D*, 100:016002, Jul 2019. doi: 10.1103/PhysRevD.100.016002. URL <https://link.aps.org/doi/10.1103/PhysRevD.100.016002>.
- Raissi, M. Forward-backward stochastic neural networks: Deep learning of high-dimensional partial differential equations. *arXiv preprint arXiv:1804.07010*, 2018.
- Raissi, M., Perdikaris, P., and Karniadakis, G. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686 – 707, 2019. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2018.10.045>. URL <http://www.sciencedirect.com/science/article/pii/S0021999118307125>.
- Riess, A. G., Filippenko, A. V., Challis, P., Clocchiatti, A., Diercks, A., Garnavich, P. M., Gilliland, R. L., Hogan, C. J., Jha, S., Kirshner, R. P., Leibundgut, B., Phillips, M. M., Reiss, D., Schmidt, B. P., Schommer, R. A., Smith, R. C., Spyromilio, J., Stubbs, C., Suntzeff, N. B., and Tonry, J. Observational evidence from supernovae for an accelerating universe and a cosmological constant. *The Astronomical Journal*, 116(3):1009–1038, sep 1998. doi: 10.1086/300499. URL <https://doi.org/10.1086%2F300499>.

- Salimans, T., Goodfellow, I. J., Zaremba, W., Cheung, V., Radford, A., and Chen, X. Improved techniques for training gans. *CoRR*, abs/1606.03498, 2016. URL <http://arxiv.org/abs/1606.03498>.
- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- Sirignano, J. and Spiliopoulos, K. Dgm: A deep learning algorithm for solving partial differential equations. *Journal of Computational Physics*, 375:1339–1364, 2018.
- Sønderby, C. K., Caballero, J., Theis, L., Shi, W., and Huszár, F. Amortised MAP inference for image super-resolution. *CoRR*, abs/1610.04490, 2016. URL <http://arxiv.org/abs/1610.04490>.
- Stevens, B. and Colonius, T. Finitenet: A fully convolutional lstm network architecture for time-dependent partial differential equations, 2020.
- Subramanian, A., Wong, M.-L., Borker, R., and Nimmagadda, S. Turbulence enrichment using generative adversarial networks, 2018. URL [http://cs230.stanford.edu/files\\_winter\\_2018/projects/6939636.pdf](http://cs230.stanford.edu/files_winter_2018/projects/6939636.pdf).
- Sutskever, I., Vinyals, O., and Le, Q. V. Sequence to sequence learning with neural networks. *CoRR*, abs/1409.3215, 2014. URL <http://arxiv.org/abs/1409.3215>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *CoRR*, abs/1706.03762, 2017. URL <http://arxiv.org/abs/1706.03762>.
- Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Jarrod Millman, K., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C., Polat, İ., Feng, Y., Moore, E. W., Vand erPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., and Contributors, S. . SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: <https://doi.org/10.1038/s41592-019-0686-2>.
- Wang, C., Horby, P. W., Hayden, F. G., and Gao, G. F. A novel coronavirus outbreak of global health concern. *The Lancet*, 395(10223):470–473, 2020.
- Yang, L., Zhang, D., and Karniadakis, G. E. Physics-informed generative adversarial networks for stochastic differential equations, 2018.
- Zeng, S., Zhang, Z., and Zou, Q. Adaptive deep neural networks methods for high-dimensional partial differential equations. *Journal of Computational Physics*, pp. 111232, 2022. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2022.111232>. URL <https://www.sciencedirect.com/science/article/pii/S0021999122002947>.

## A. Appendix

### A.1. Classical Loss Functions

A plot of the various classical loss functions is provided in Figure 4.

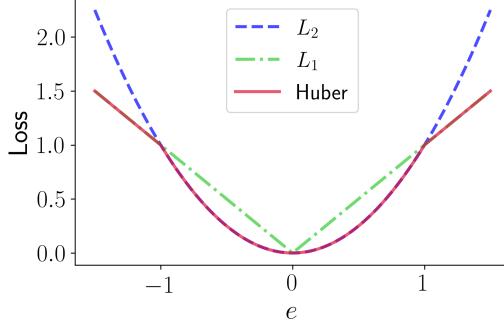


Figure 4. Comparison of  $L_2$ ,  $L_1$ , and Huber loss functions. The Huber loss is equal to  $L_2$  for  $e \leq 1$  and to  $L_1$  for  $e > 1$ .

### A.2. Description of Experiments

#### A.2.1. EXPONENTIAL DECAY (EXP)

Consider a model for population decay  $x(t)$  given by the exponential differential equation

$$\dot{x}(t) + x(t) = 0, \quad (13)$$

with  $x(0) = 1$  and  $t \in [0, 10]$ . The ground truth solution  $x(t) = e^{-t}$  can be obtained analytically, which we use to calculate the mean squared error of the predicted solution.

To set up the problem for DEQGAN, we define  $LHS = \dot{x} + x$  and  $RHS = 0$ . Figure 5 presents the results from training DEQGAN on this equation.

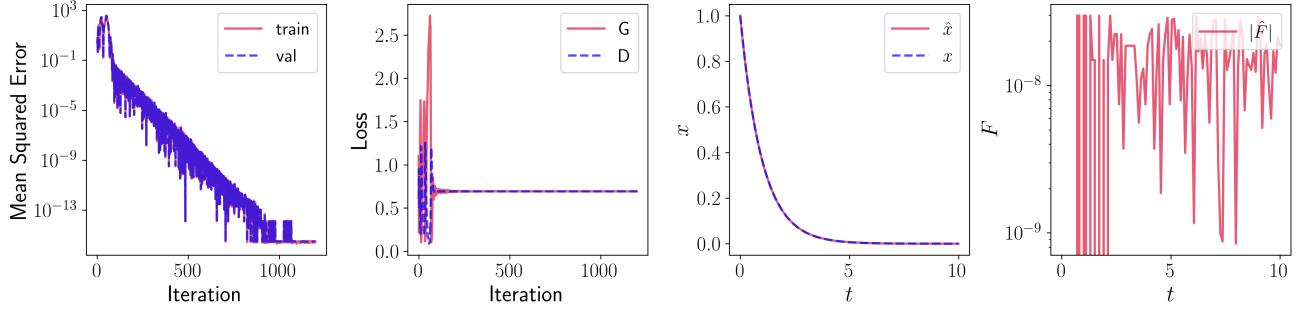


Figure 5. Visualization of DEQGAN training for the exponential decay problem. The left-most figure plots the mean squared error vs. iteration. To the right, we plot the value of the generator (G) and discriminator (D) losses at each iteration. Right of this we plot the prediction of the generator  $\hat{x}$  and the true analytic solution  $x$  as functions of time  $t$ . The right-most figure plots the absolute value of the residual of the predicted solution  $\hat{F}$ .

#### A.2.2. SIMPLE HARMONIC OSCILLATOR (SHO)

Consider the motion of an oscillating body  $x(t)$ , which can be modeled by the simple harmonic oscillator differential equation

$$\ddot{x}(t) + x(t) = 0, \quad (14)$$

with  $x(0) = 0$ ,  $\dot{x}(0) = 1$ , and  $t \in [0, 2\pi]$ . This differential equation can be solved analytically and has an exact solution  $x(t) = \sin t$ .

Here we set  $LHS = \ddot{x} + x$  and  $RHS = 0$ . Figure 6 plots the results of training DEQGAN on this problem.

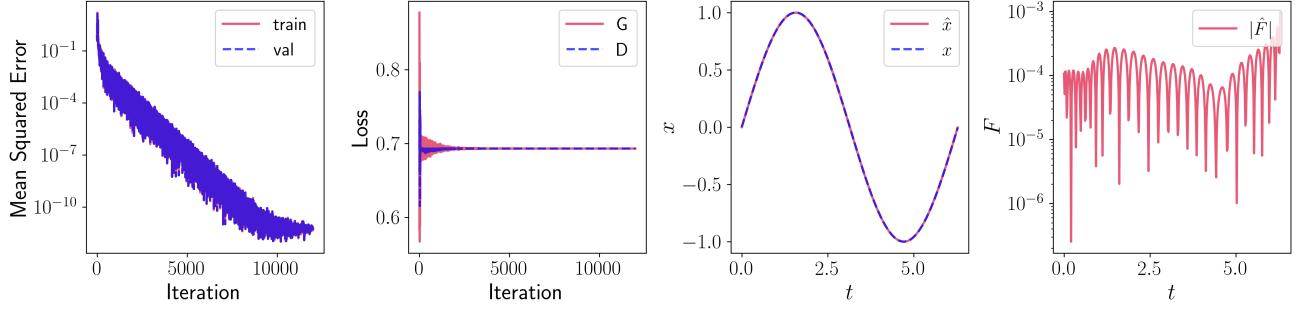


Figure 6. Visualization of DEQGAN training for the simple harmonic oscillator problem.

#### A.2.3. DAMPED NONLINEAR OSCILLATOR (NLO)

Further increasing the complexity of the differential equations being considered, consider a less idealized oscillating body subject to additional forces, whose motion  $x(t)$  we can described by the nonlinear oscillator differential equation

$$\ddot{x}(t) + 2\beta\dot{x}(t) + \omega^2 x(t) + \phi x(t)^2 + \epsilon x(t)^3 = 0, \quad (15)$$

with  $\beta = 0.1$ ,  $\omega = 1$ ,  $\phi = 1$ ,  $\epsilon = 0.1$ ,  $x(0) = 0$ ,  $\dot{x}(0) = 0.5$ , and  $t \in [0, 4\pi]$ . This equation does not admit an analytical solution. Instead, we use the high-quality solver provided by SciPy's `solve_ivp` (Virtanen et al., 2020).

We set  $LHS = \ddot{x} + 2\beta\dot{x} + \omega^2 x + \phi x^2 + \epsilon x^3 = 0$  and  $RHS = 0$ . Figure 7 plots the results obtained from training DEQGAN on this equation.

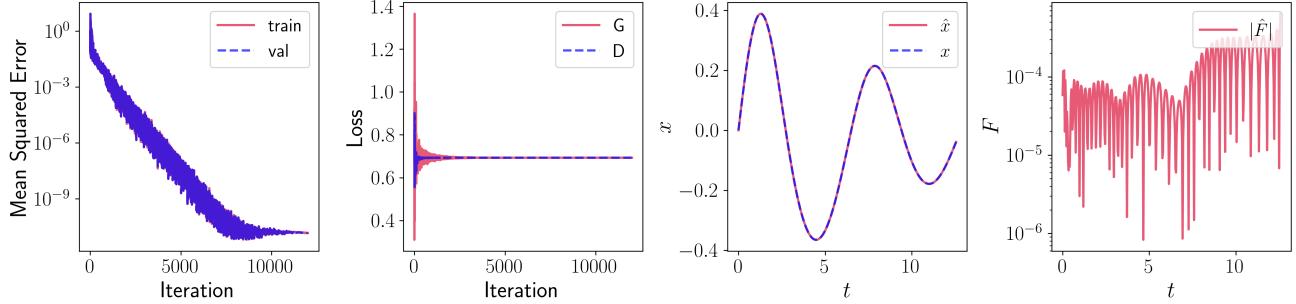


Figure 7. Visualization of DEQGAN training for the nonlinear oscillator problem.

#### A.2.4. COUPLED OSCILLATORS (COO)

Consider the system of ordinary differential equations given by

$$\begin{cases} \dot{x}(t) = -ty \\ \dot{y}(t) = tx \end{cases} \quad (16)$$

with  $x(0) = 1$ ,  $y(0) = 0$ , and  $t \in [0, 2\pi]$ . This equation has an exact analytical solution given by

$$\begin{cases} x = \cos\left(\frac{t^2}{2}\right) \\ y = \sin\left(\frac{t^2}{2}\right) \end{cases} \quad (17)$$

Here we set

$$LHS = \left[ \frac{dx}{dt} + ty, \frac{dy}{dt} - xy \right]^T \quad (18)$$

and  $RHS = [0, 0]^T$ . Figure 8 plots the result of training DEQGAN on this problem.

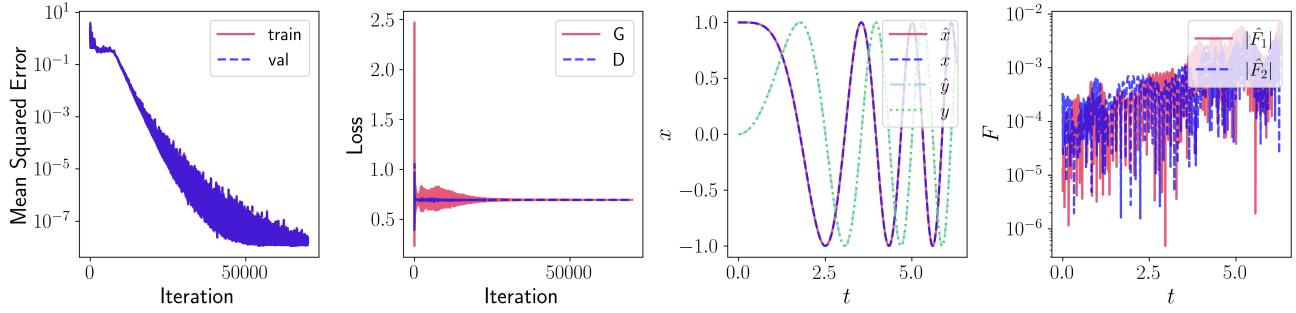


Figure 8. Visualization of DEQGAN training for the coupled oscillators system of equations. In the third figure, we plot the predictions of the generator  $\hat{x}, \hat{y}$  and the true analytic solutions  $x, y$  as functions of time  $t$ . The right-most figure plots the absolute value of the residuals of the predicted solution  $\hat{F}_j$  for each equation  $j$ .

#### A.2.5. SIR EPIDEMIOLOGICAL MODEL (SIR)

Given the ongoing pandemic of novel coronavirus (COVID-19) (Wang et al., 2020), we consider an epidemiological model of infectious disease spread given by a system of ordinary differential equations. Specifically, consider the Susceptible  $S(t)$ , Infected  $I(t)$ , Recovered  $R(t)$  model for the spread of an infectious disease over time  $t$ . The model is defined by a system of three ordinary differential equations

$$\begin{cases} \dot{S}(t) = -\beta \frac{IS}{N} \\ \dot{I}(t) = \beta \frac{IS}{N} - \gamma I \\ \dot{R}(t) = \gamma I \end{cases} \quad (19)$$

where  $\beta = 3, \gamma = 1$  are given constants related to the infectiousness of the disease,  $N = S + I + R$  is the (constant) total population,  $S(0) = 0.99, I(0) = 0.01, R(0) = 0$ , and  $t \in [0, 10]$ . As this system has no analytical solution, we use SciPy's `solve_ivp` solver (Virtanen et al., 2020) to obtain ground truth solutions.

We set  $LHS$  to be the vector

$$LHS = \left[ \frac{dS}{dt} + \beta \frac{IS}{N}, \frac{dI}{dt} - \beta \frac{IS}{N} + \gamma I, \frac{dR}{dt} - \gamma I \right]^T \quad (20)$$

and  $RHS = [0, 0, 0]^T$ . We present the results of training DEQGAN to solve this system of differential equations in Figure 9.

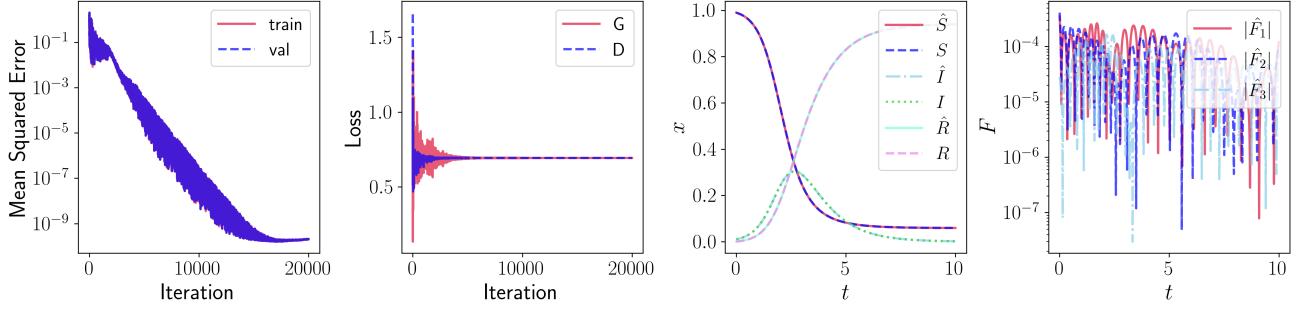


Figure 9. Visualization of DEQGAN training for the SIR system of equations.

#### A.2.6. HAMILTONIAN SYSTEM (HAM)

Consider a particle moving through a potential  $V$ , the trajectory of which is described by the system of ordinary differential equations

$$\begin{cases} \dot{x}(t) = p_x \\ \dot{y}(t) = p_y \\ \dot{p}_x(t) = -V_x \\ \dot{p}_y(t) = -V_y \end{cases} \quad (21)$$

with  $x(0) = 0, y(0) = 0.3, p_x(0) = 1, p_y(0) = 0$ , and  $t \in [0, 1]$ .  $V_x$  and  $V_y$  are the  $x$  and  $y$  derivatives of the potential  $V$ , which we construct by summing ten random bivariate Gaussians

$$V = -\frac{A}{2\pi\sigma^2} \sum_{i=1}^{10} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}(t) - \mu_i\|_2^2\right) \quad (22)$$

where  $\mathbf{x}(t) = [x(t), y(t)]^T$ ,  $A = 0.1, \sigma = 0.1$ , and each  $\mu_i$  is sampled from  $[0, 1] \times [0, 1]$  uniformly at random. As before, we use SciPy to obtain ground-truth solutions.

We set  $LHS$  to be the vector

$$LHS = \left[ \frac{dx}{dt} - p_x, \frac{dy}{dt} - p_y, \frac{dp_x}{dt} + V_x, \frac{dp_y}{dt} + V_y \right]^T \quad (23)$$

and  $RHS = [0, 0, 0, 0]^T$ . We present the results of training DEQGAN to solve this system of differential equations in Figure 10.

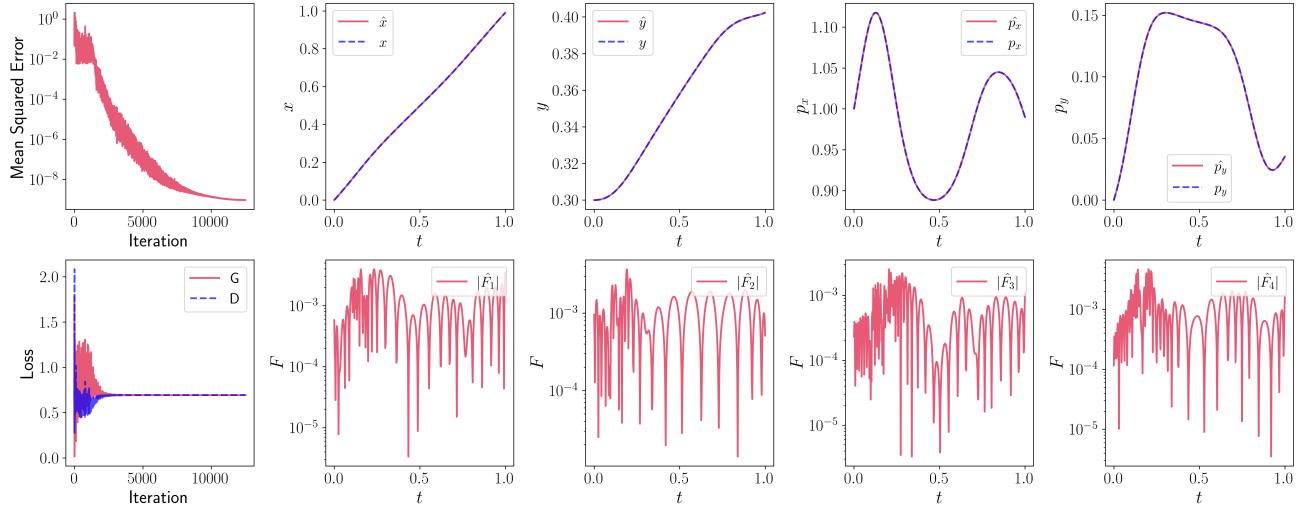


Figure 10. Visualization of DEQGAN training for the Hamiltonian system of equations. For ease of visualization, we plot the predictions and residuals for each equation separately.

#### A.2.7. MODIFIED EINSTEIN'S GRAVITY SYSTEM (EIN)

The most challenging system of ODEs we consider comes from Einstein's theory of general relativity. Following observations from type Ia supernovae in 1998 (Riess et al., 1998), several cosmological models have been proposed to explain the accelerated expansion of the universe. Some of these rely on the existence of unobserved forms such as dark energy and dark matter, while others directly modify Einstein's theory.

Hu-Sawicky  $f(R)$  gravity is one model that falls under this category. Chantada et al. (2022) show how the following system of five ODEs can be derived from the modified field equations implied by this model.

$$\begin{cases} \dot{x}(z) = \frac{1}{z+1}(-\Omega - 2v + x + 4y + xv + x^2) \\ \dot{y}(z) = \frac{-1}{z+1}(vx\Gamma(r) - xy + 4y - 2yv) \\ \dot{v}(z) = \frac{-v}{z+1}(x\Gamma(r) + 4 - 2v) \\ \dot{\Omega}(z) = \frac{\Omega}{z+1}(-1 + 2v + x) \\ \dot{r}(z) = \frac{-r\Gamma(r)x}{z+1} \end{cases} \quad (24)$$

where

$$\Gamma(r) = \frac{(r+b)[(r+b)^2 - 2b]}{4br}. \quad (25)$$

The initial conditions are given by

$$\begin{cases} x_0 = 0 \\ y_0 = \frac{\Omega_{m,0}(1+z_0)^3 + 2(1-\Omega_{m,0})}{2[\Omega_{m,0}(1+z_0)^3 + (1-\Omega_{m,0})]} \\ v_0 = \frac{\Omega_{m,0}(1+z_0)^3 + 4(1-\Omega_{m,0})}{2[\Omega_{m,0}(1+z_0)^3 + (1-\Omega_{m,0})]} \\ \Omega_0 = \frac{\Omega_{m,0}(1+z_0)^3}{\Omega_{m,0}(1+z_0)^3 + (1-\Omega_{m,0})} \\ r_0 = \frac{\Omega_{m,0}(1+z_0)^3 + 4(1-\Omega_{m,0})}{(1-\Omega_{m,0})} \end{cases} \quad (26)$$

where  $z_0 = 10$ ,  $\Omega_{m,0} = 0.15$ ,  $b = 5$  and we solve the system for  $z \in [0, z_0]$ . While the physical interpretation of the various parameters is beyond the scope of this paper, we note that Equations 24 and 25 exhibit a high degree of non-linearity. Ground truth solutions are again obtained using SciPy, and the results obtained by DEQGAN are shown in Figure 11.

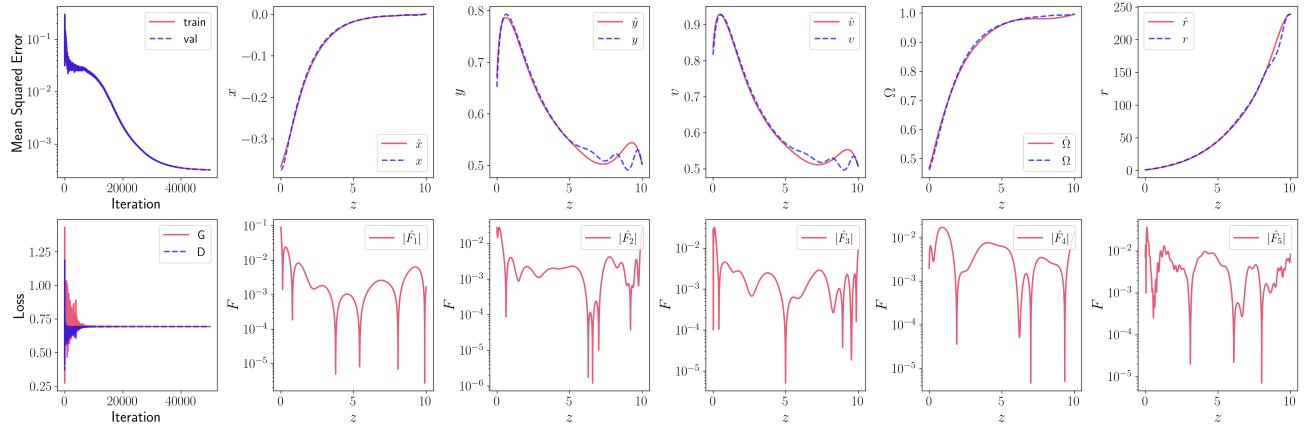


Figure 11. Visualization of DEQGAN training for the modified Einstein's gravity system of equations. For ease of visualization, we plot the predictions and residuals for each equation separately.

#### A.2.8. POISSON EQUATION (POS)

Consider the Poisson partial differential equation (PDE) given by

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 2x(y-1)(y-2x+xy+2)e^{x-y} \quad (27)$$

where  $(x, y) \in [0, 1] \times [0, 1]$ . The equation is subject to Dirichlet boundary conditions on the edges of the unit square

$$\begin{aligned} u(x, y) \Big|_{x=0} &= 0 \\ u(x, y) \Big|_{x=1} &= 0 \\ u(x, y) \Big|_{y=0} &= 0 \\ u(x, y) \Big|_{y=1} &= 0. \end{aligned} \quad (28)$$

The analytical solution is

$$u(x, y) = x(1-x)y(1-y)e^{x-y}. \quad (29)$$

We use the two-dimensional Dirichlet boundary adjustment formulae provided in Chen et al. (2020). To set up the problem for DEQGAN we let

$$LHS = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} - 2x(y-1)(y-2x+xy+2)e^{x-y} \quad (30)$$

and  $RHS = 0$ . We present the results of training DEQGAN on this problem in Figure 12.

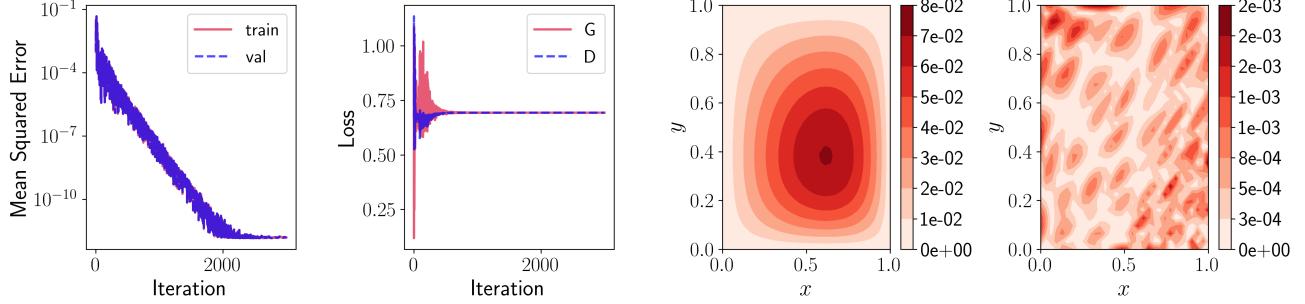


Figure 12. Visualization of DEQGAN training for the Poisson equation. In the third figure, we plot the prediction of the generator  $\hat{u}$  as a function of position  $(x, y)$ . The right-most figure plots the absolute value of the residual  $\hat{F}$ , as a function of  $(x, y)$ .

#### A.2.9. HEAT EQUATION (HEA)

We consider the time-dependent heat (diffusion) equation given by

$$\frac{\partial u}{\partial t} = \kappa \frac{\partial^2 u}{\partial x^2} \quad (31)$$

where  $\kappa = 1$  and  $(x, t) \in [0, 1] \times [0, 0.2]$ . The equation is subject to an initial condition and Dirichlet boundary conditions given by

$$\begin{aligned} u(x, y) \Big|_{t=0} &= \sin(\pi x) \\ u(x, y) \Big|_{x=0} &= 0 \\ u(x, y) \Big|_{x=1} &= 0 \end{aligned} \quad (32)$$

and has an analytical solution

$$u(x, y) = e^{-\kappa\pi^2 t} \sin(\pi x). \quad (33)$$

The results obtained by DEQGAN on this problem are shown in Figure 13.

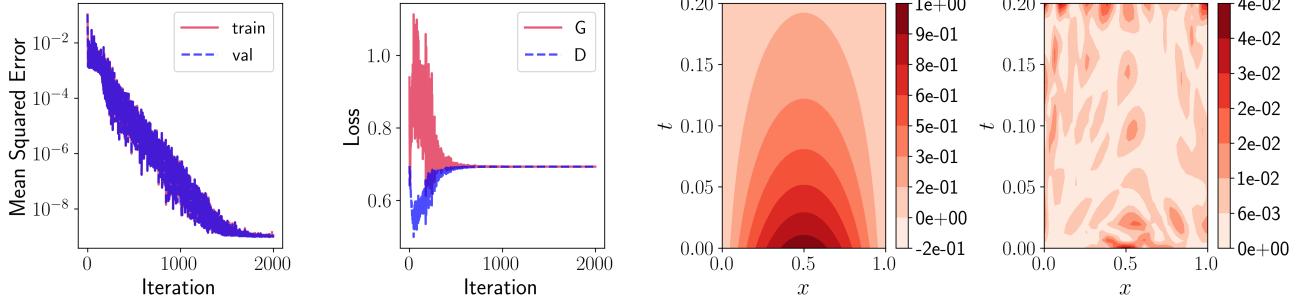


Figure 13. Visualization of DEQGAN training for the heat equation. In the third figure, we plot the prediction of the generator  $\hat{u}$  as a function of position  $(x, t)$ . The right-most figure plots the absolute value of the residual  $\hat{F}$ , as a function of  $(x, t)$ .

#### A.2.10. WAVE EQUATION (WAV)

Consider the time-dependent wave equation given by

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2} \quad (34)$$

where  $c = 1$  and  $(x, t) \in [0, 1] \times [0, 1]$ . This formulation is very similar to the heat equation but involves a second order derivative with respect to time. We subject the equation to the same initial condition and boundary conditions as 32 but require an added Neumann condition due to the equation's second time derivative.

$$\begin{aligned} u(x, y) \Big|_{t=0} &= \sin(\pi x) \\ u_t(x, y) \Big|_{t=0} &= 0 \\ u(x, y) \Big|_{x=0} &= 0 \\ u(x, y) \Big|_{x=1} &= 0 \end{aligned} \quad (35)$$

This yields the analytical solution

$$u(x, y) = \cos(c\pi t) \sin(\pi x). \quad (36)$$

The results of training DEQGAN on this problem are shown in Figure 13.

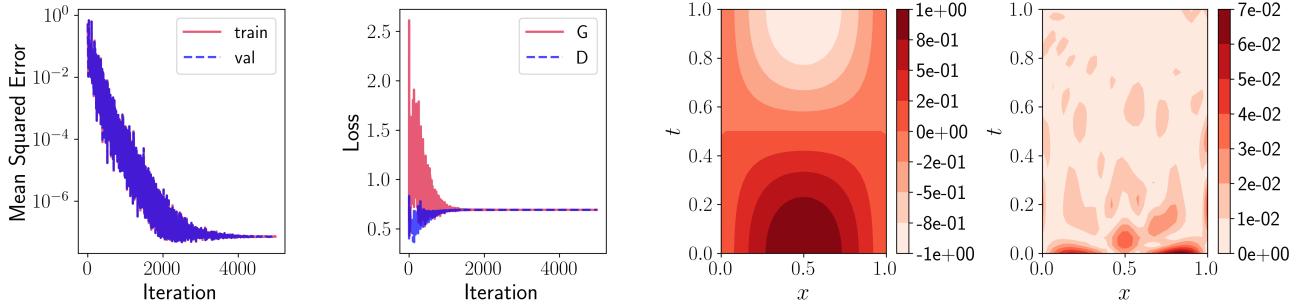


Figure 14. Visualization of DEQGAN training for the wave equation.

### A.2.11. BUGERS' EQUATION (BUR)

Moving to non-linear PDEs, we consider the viscous Burgers' equation given by

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = \nu \frac{\partial^2 u}{\partial x^2} \quad (37)$$

where  $\nu = 0.001$  and  $(x, t) \in [-5, 5] \times [0, 2.5]$ . To specify the equation, we use the following initial condition and Dirichlet boundary conditions:

$$\begin{aligned} u(x, y) \Big|_{t=0} &= \frac{1}{\cosh(x)} \\ u(x, y) \Big|_{x=-5} &= 0 \\ u(x, y) \Big|_{x=5} &= 0 \end{aligned} \quad (38)$$

As this equation has no analytical solution, we use the fast Fourier transform (FFT) method (Brunton & Kutz, 2019) to obtain ground truth solutions. The results obtained by DEQGAN are summarized by Figure 15. As time progresses, we see the formation of a “shock wave” that becomes increasingly steep but remains smooth due to the regularizing diffusive term  $\nu u_{xx}$ .

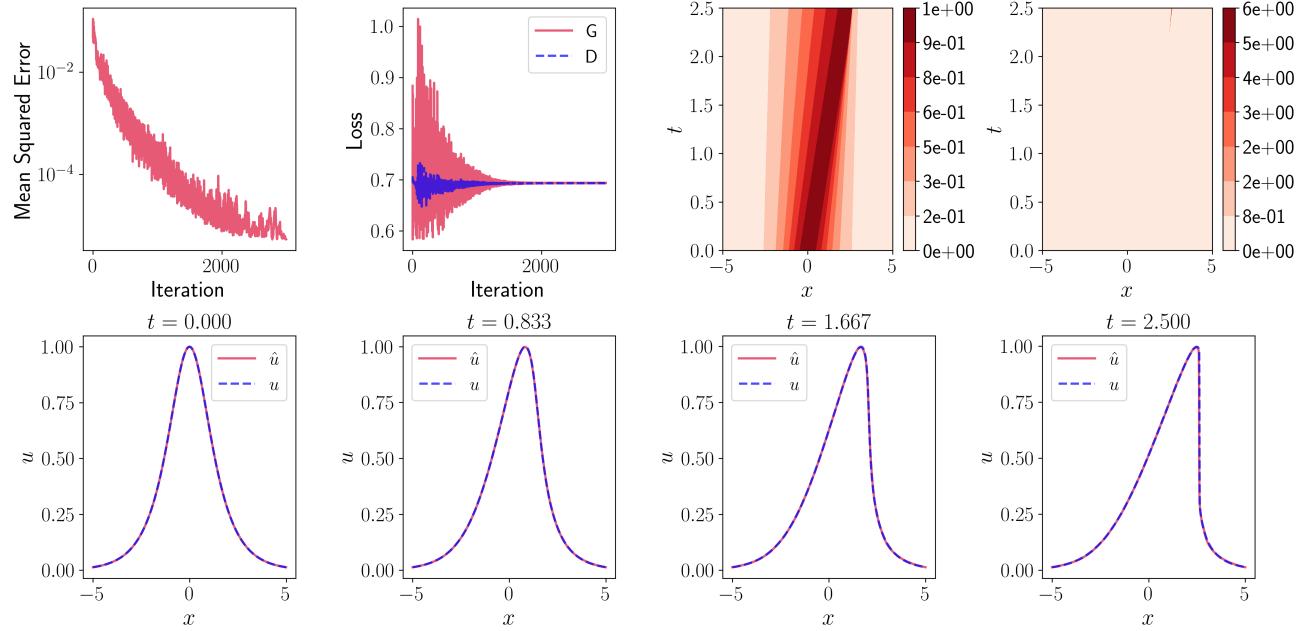


Figure 15. Visualization of DEQGAN training for Burgers' equation. The plots in the second row show “snapshots” of the 1D wave at different points along the time domain.

### A.2.12. ALLEN-CAHN EQUATION (ACA)

Finally, we consider the Allen-Cahn PDE, a well-known reaction-diffusion equation given by

$$\frac{\partial u}{\partial t} - \epsilon \frac{\partial^2 u}{\partial x^2} - u + u^3 = 0 \quad (39)$$

where  $\epsilon = 0.001$  and  $(x, t) \in [0, 2\pi] \times [0, 5]$ . We subject the equation to an initial condition and Dirichlet boundary conditions given by

$$\begin{aligned} u(x, y) \Big|_{t=0} &= \frac{1}{4} \sin(x) \\ u(x, y) \Big|_{x=0} &= 0 \\ u(x, y) \Big|_{x=2\pi} &= 0 \end{aligned} \quad (40)$$

The results are shown in Figure 16. We see that as time progresses, the sinusoidal initial condition transforms into a square wave, becoming very steep at the turning points of the solution.

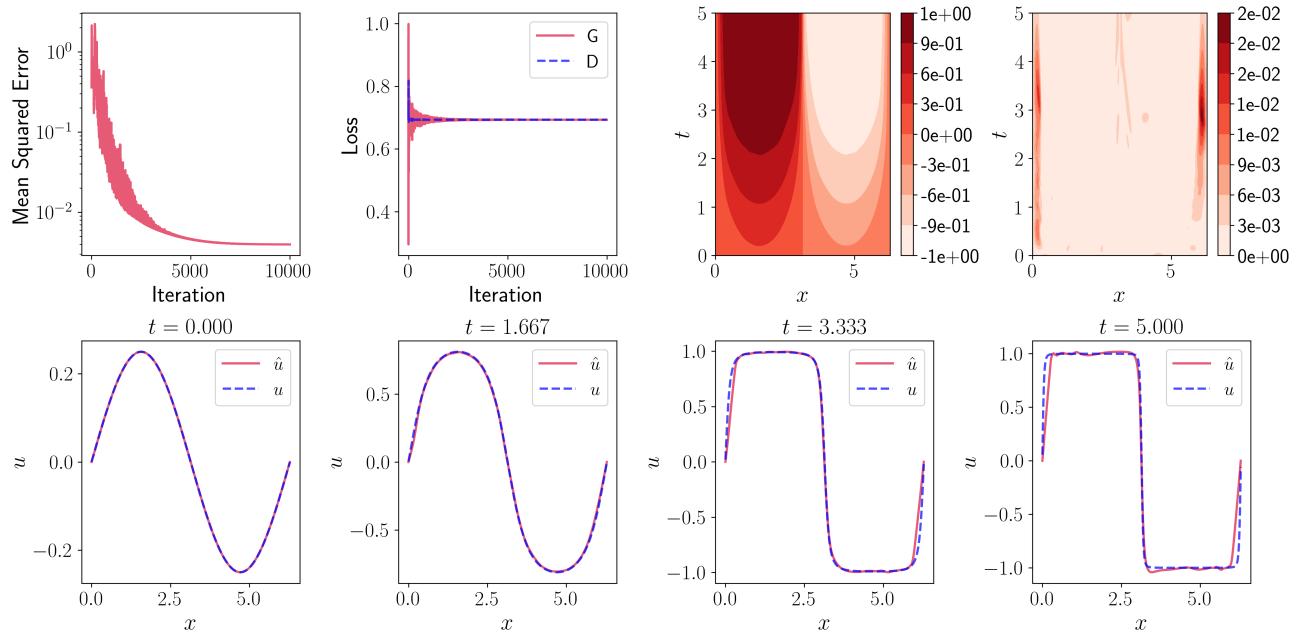


Figure 16. Visualization of DEQGAN training for the Allen-Cahn equation. The plots in the second row show “snapshots” of the 1D wave at different points along the time domain.

### A.3. Method Comparison for Other Experiments

Figure 17 visualizes the training results achieved by DEQGAN and the alternative unsupervised neural networks that use  $L_2$ ,  $L_1$  and Huber loss functions for the remaining six problems.

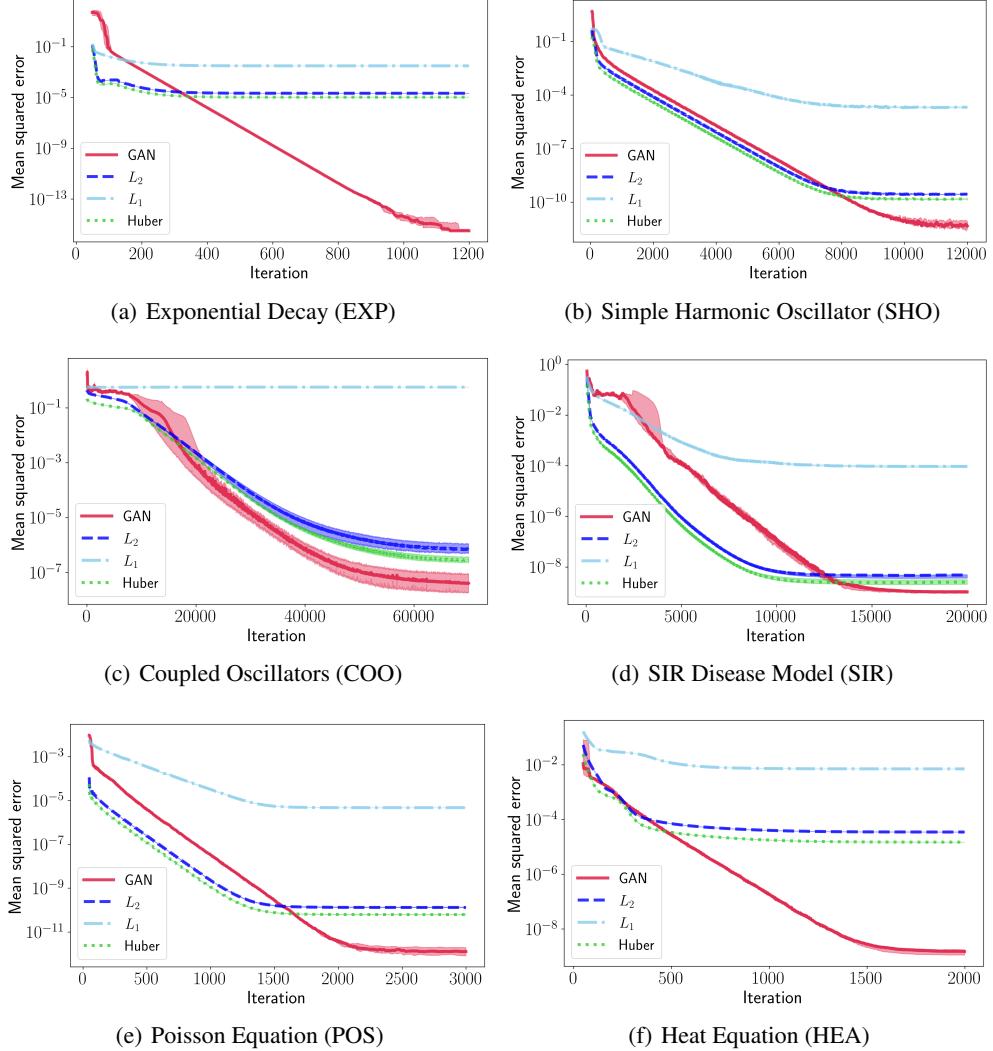


Figure 17. Mean squared errors vs. iteration for DEQGAN,  $L_2$ ,  $L_1$ , and Huber loss for various equations. We perform ten randomized trials and plot the median (bold) and (25, 75) percentile range (shaded). We smooth the values using a simple moving average with window size 50.

#### A.4. DEQGAN Hyperparameters

We used Ray Tune (Liaw et al., 2018) to tune DEQGAN hyperparameters for each differential equation. Tables 4 and 5 summarize these hyperparameter values for the ODE and PDE problems, respectively. The experiments and hyperparameter tuning conducted for this research totaled 13,272 hours of compute performed on Intel Cascade Lake CPU cores belonging to an internal cluster.

Table 4. Hyperparameter Settings for DEQGAN (ODEs)

HYPERPARAMETER	EXP	SHO	NLO	COO	SIR	HAM	EIN
NUM. ITERATIONS	1200	12000	12000	70000	20000	12500	50000
NUM. GRID POINTS	100	400	400	800	800	400	1000
$G$ UNITS/LAYER	40	40	40	40	50	40	40
$G$ NUM. LAYERS	2	3	4	5	4	5	4
$D$ UNITS/LAYER	20	50	20	40	50	50	30
$D$ NUM. LAYERS	4	3	2	2	4	2	2
ACTIVATIONS	tanh	tanh	tanh	tanh	tanh	tanh	tanh
$G$ LEARNING RATE	0.094	0.005	0.010	0.004	0.006	0.017	0.011
$D$ LEARNING RATE	0.012	0.0004	0.021	0.082	0.012	0.019	0.006
$G \beta_1$ (ADAM)	0.491	0.363	0.225	0.603	0.278	0.252	0.202
$G \beta_2$ (ADAM)	0.319	0.752	0.331	0.614	0.777	0.931	0.975
$D \beta_1$ (ADAM)	0.542	0.584	0.362	0.412	0.018	0.105	0.154
$D \beta_2$ (ADAM)	0.264	0.453	0.551	0.110	0.908	0.869	0.797
EXPONENTIAL LR DECAY ( $\gamma$ )	0.978	0.980	0.999	0.992	0.9996	0.985	0.996
DECAY STEP SIZE	3	19	15	16	11	13	17

Table 5. Hyperparameter Settings for DEQGAN (PDEs)

HYPERPARAMETER	POS	HEA	WAV	BUR	ACA
NUM. ITERATIONS	3000	2000	5000	3000	10000
NUM. GRID POINTS	$32 \times 32$	$32 \times 32$	$32 \times 32$	$64 \times 64$	$64 \times 64$
$G$ UNITS/LAYER	50	40	50	50	50
$G$ NUM. LAYERS	4	4	4	3	2
$D$ UNITS/LAYER	30	30	50	20	30
$D$ NUM. LAYERS	2	2	2	5	2
ACTIVATIONS	tanh	tanh	tanh	tanh	tanh
$G$ LEARNING RATE	0.019	0.010	0.012	0.012	0.020
$D$ LEARNING RATE	0.021	0.001	0.088	0.005	0.013
$G \beta_1$ (ADAM)	0.139	0.230	0.295	0.185	0.436
$G \beta_2$ (ADAM)	0.369	0.657	0.358	0.594	0.910
$D \beta_1$ (ADAM)	0.745	0.120	0.575	0.093	0.484
$D \beta_2$ (ADAM)	0.759	0.251	0.133	0.184	0.297
EXPONENTIAL LR DECAY ( $\gamma$ )	0.957	0.950	0.953	0.954	0.983
DECAY STEP SIZE	3	10	18	20	15

### A.5. Non-GAN Hyperparameter Tuning

Table 6 presents the minimum mean squared errors obtained after tuning hyperparameters for the alternative unsupervised neural network methods that use  $L_1$ ,  $L_2$  and Huber loss functions.

Table 6. Experimental Results With Non-GAN Hyperparameter Tuning

Key	Mean Squared Error				
	$L_1$	$L_2$	Huber	DEQGAN	Traditional
EXP	$1 \cdot 10^{-4}$	$4 \cdot 10^{-8}$	$2 \cdot 10^{-8}$	$3 \cdot 10^{-16}$	$2 \cdot 10^{-14}$ (RK4)
SHO	$1 \cdot 10^{-5}$	$1 \cdot 10^{-9}$	$5 \cdot 10^{-10}$	$4 \cdot 10^{-13}$	$1 \cdot 10^{-11}$ (RK4)
NLO	$1 \cdot 10^{-4}$	$3 \cdot 10^{-10}$	$1 \cdot 10^{-10}$	$1 \cdot 10^{-12}$	$4 \cdot 10^{-11}$ (RK4)
COO	$5 \cdot 10^{-1}$	$2 \cdot 10^{-7}$	$3 \cdot 10^{-7}$	$1 \cdot 10^{-8}$	$2 \cdot 10^{-9}$ (RK4)
SIR	$9 \cdot 10^{-6}$	$1 \cdot 10^{-10}$	$1 \cdot 10^{-10}$	$1 \cdot 10^{-10}$	$5 \cdot 10^{-13}$ (RK4)
HAM	$4 \cdot 10^{-5}$	$1 \cdot 10^{-8}$	$6 \cdot 10^{-9}$	$1 \cdot 10^{-10}$	$7 \cdot 10^{-14}$ (RK4)
EIN	$5 \cdot 10^{-2}$	$2 \cdot 10^{-2}$	$1 \cdot 10^{-2}$	$4 \cdot 10^{-4}$	$4 \cdot 10^{-7}$ (RK4)
POS	$9 \cdot 10^{-6}$	$1 \cdot 10^{-10}$	$1 \cdot 10^{-10}$	$4 \cdot 10^{-13}$	$3 \cdot 10^{-10}$ (FD)
HEA	$1 \cdot 10^{-4}$	$4 \cdot 10^{-8}$	$2 \cdot 10^{-8}$	$6 \cdot 10^{-10}$	$4 \cdot 10^{-7}$ (FD)
WAV	$4 \cdot 10^{-4}$	$6 \cdot 10^{-7}$	$2 \cdot 10^{-7}$	$1 \cdot 10^{-8}$	$7 \cdot 10^{-5}$ (FD)
BUR	$1 \cdot 10^{-3}$	$1 \cdot 10^{-4}$	$9 \cdot 10^{-5}$	$4 \cdot 10^{-6}$	$1 \cdot 10^{-3}$ (FD)
ACA	$5 \cdot 10^{-2}$	$1 \cdot 10^{-2}$	$3 \cdot 10^{-3}$	$5 \cdot 10^{-3}$	$2 \cdot 10^{-4}$ (FD)

### A.6. Residual Monitoring

Figure 18 shows several examples of how we detect bad training runs by monitoring the variance of the  $L_1$  norm of the *LHS* (vector of equation residuals) in the first 25% of training iterations. Because the *LHS* may oscillate initially even for successful runs, we use a patience window in the first 15% of iterations. In all three equations below, we terminate runs if the variance of the residual  $L_1$  norm over 20 iterations exceeds 0.01.

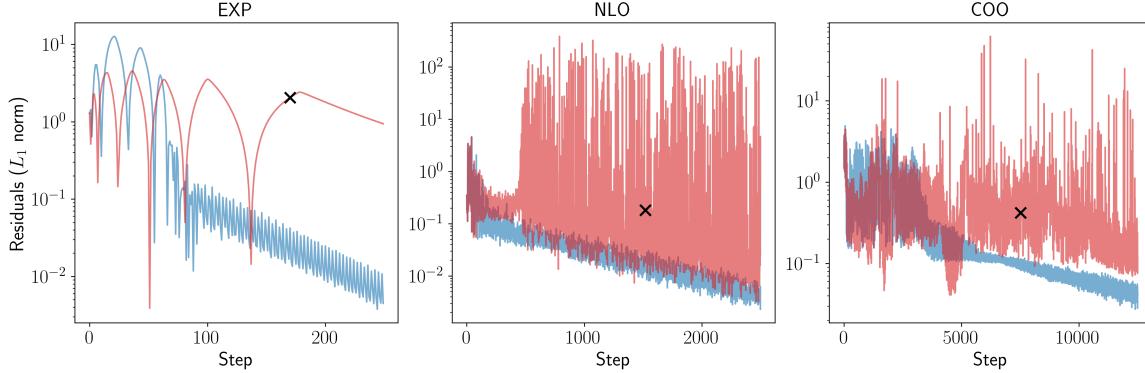


Figure 18. Equation residuals in the first 25% of training runs that ended with high (red) and low (blue) mean squared error for the exponential decay (EXP), non-linear oscillator (NLO) and coupled oscillators (COO) problems. The black crosses show the point at which the high MSE runs were terminated early.

### A.7. Ablation Study

To quantify the increased robustness offered by instance noise and residual monitoring, we performed an ablation study comparing the percentage of high MSE ( $\geq 10^{-5}$ ) runs obtained by 500 randomized DEQGAN runs for the exponential decay equation with and without using these techniques.

Figure 19 plots the results of these 500 DEQGAN experiments with instance noise added. For each experiment, we uniformly selected a random seed controlling model weight initialization as an integer from the range [0, 9], as well as

separate learning rates for the discriminator and generator in the range  $[0.01, 0.1]$ . We then recorded the final mean squared error after running DEQGAN training for 1000 iterations. The red lines represent runs which would be terminated early by our residual monitoring method, while the blue lines represent those which would be run to completion.

Figure 19 shows that the large majority of hyperparameter settings tested with the addition of instance noise resulted in low mean squared errors. Further, residual monitoring was able to detect all runs with  $\text{MSE} \geq 10^{-5}$  within the first 25% of training iterations. Approximately half of the MSE runs in  $[10^{-8}, 10^{-5}]$  would be terminated, while 96% of runs with  $\text{MSE} \leq 10^{-8}$  would be run to completion.

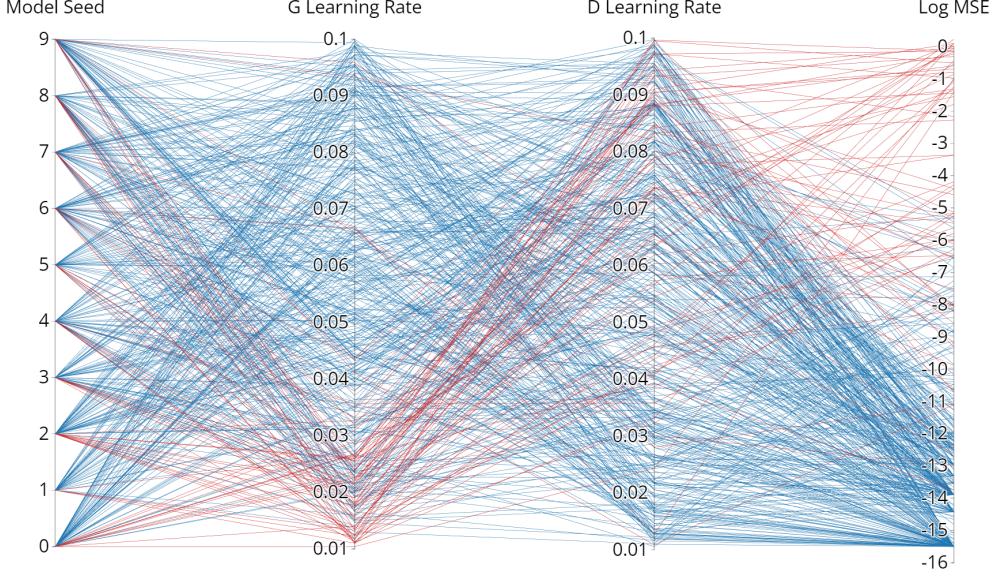


Figure 19. Parallel plot showing the results of 500 DEQGAN experiments on the exponential decay equation with instance noise. The red lines represent runs which would be terminated early by monitoring the variance of the equation residuals in the first 25% of training iterations. The mean squared error is plotted on a  $\log_{10}$  scale.