

Universal Differential Equations for Scientific Machine Learning: Supplemental Information

Christopher Rackauckas^{a,b}, Yingbo Ma^c, Julius Martensen^d, Collin Warner^a, Kirill Zubov^e, Rohit Supekar^a, Dominic Skinner^a, Ali Ramadhan^a, and Alan Edelman^a

^aMassachusetts Institute of Technology

^bUniversity of Maryland, Baltimore

^cJulia Computing

^dUniversity of Bremen

^eSaint Petersburg State University

August 26, 2020

1 DiffEqFlux.jl Pullback Construction

The DiffEqFlux.jl pullback construction is not based on just one method but instead has a dispatch-based mechanism for choosing between different adjoint implementations. At a high level, the library defines the pullback on the differential equation solve function, and thus using a differential equation inside of a larger program leads to this chunk as being a single differentiable primitive that is inserted into the back pass of Flux.jl when encountered by overloading the Zygote.jl [1] and ChainRules.jl rule sets. For any ChainRules.jl-compliant reverse-mode AD package in the Julia language, when a differential equation solve is encountered in any Julia library during the backwards pass, the adjoint method is automatically swapped in to be used for the backpropagation of the solver. The choice of the adjoint is chosen by the type of the sensealg keyword argument which are fully described below.

1.1 Backpropagation-Accelerated DAE Adjoints for Index-1 DAEs with Constraint Equations

Before describing the modes, we first describe the adjoint of the differential equation with constraints. The constrained ordinary differential equation:

$$u' = \tilde{f}(u, p, t), \tag{1}$$

$$0 = c(u, p, t), \tag{2}$$

16 can be rewritten in mass matrix form:

$$Mu' = f(u, p, t). \quad (3)$$

17 We wish to solve for some cost function $G(u, p)$ evaluated throughout the dif-
18 ferential equation, i.e.:

$$G(u, p) = \int_{t_0}^T g(u, p, t) dt, \quad (4)$$

19 To derive this adjoint, introduce the Lagrange multiplier λ to form:

$$I(p) = G(p) - \int_{t_0}^T \lambda^*(Mu' - f(u, p, t)) dt, \quad (5)$$

20 Since $u' = f(u, p, t)$, we have that:

$$\frac{dG}{dp} = \frac{dI}{dp} = \int_{t_0}^T (g_p + g_u s) dt - \int_{t_0}^T \lambda^*(Ms' - f_u s - f_p) dt, \quad (6)$$

21 for s_i being the sensitivity of the i th variable. After applying integration by
22 parts to $\lambda^* Ms'$, we require that:

$$M^* \lambda' = -\frac{df}{du}^* \lambda - \left(\frac{dg}{du}\right)^*, \quad (7)$$

$$\lambda(T) = 0, \quad (8)$$

23 to zero out a term and get:

$$\frac{dG}{dp} = \lambda^*(t_0) M \frac{du}{dp}(t_0) + \int_{t_0}^T (g_p + \lambda^* f_p) dt. \quad (9)$$

24 If G is discrete, then it can be represented via the Dirac delta:

$$G(u, p) = \int_{t_0}^T \sum_{i=1}^N \|d_i - u(t_i, p)\|^2 \delta(t_i - t) dt, \quad (10)$$

25 in which case

$$g_u(t_i) = 2(d_i - u(t_i, p)), \quad (11)$$

26 at the data points (t_i, d_i) . Therefore, the derivative of an ODE solution with
27 respect to a cost function is given by solving for λ^* using an ODE for λ^T in
28 reverse time, and then using that to calculate $\frac{dG}{dp}$. At each time point where
29 discrete data is found, λ is then changed using a callback (discrete event han-
30 dling) by the amount g_u to represent the Dirac delta portion of the integral.
31 Lastly, we note that $\frac{dG}{du_0} = -\lambda(0)$ in this formulation.

We have to take care of consistent initialization in the case of semi-explicit index-1 DAEs. We need to satisfy the system of equations

$$M^* \Delta \lambda^d = h_{u^d}^* \lambda^a + g_{u^d}^* \quad (12)$$

$$0 = h_{u^a}^* \lambda^a + g_{u^a}^*, \quad (13)$$

where d and a denote differential and algebraic variables, and f and g denote differential and algebraic equations respectively. Combining the above two equations, we know that we need to increment the differential part of λ by

$$-h_{u^d}^* (h_{u^a}^*)^{-1} g_{u^a}^* + g_{u^d}^* \quad (14)$$

at each callback. Additionally, the ODEs

$$\mu' = -\lambda^* \frac{\partial f}{\partial p} \quad (15)$$

with $\mu(T) = 0$ can be appended to the system of equations to perform the quadrature for $\frac{dG}{dp}$.

1.2 Current Adjoint Calculation Methods

From this setup we have the following 8 possible modes for calculating the adjoint, with their pros and cons.

1. **QuadratueAdjoint**: a quadrature-based approach. This utilizes interpolation of the forward solution provided by `DifferentialEquations.jl` to calculate $u(t)$ at arbitrary time points for doing the calculations with respect to $\frac{df}{du}$ in the reverse ODE of λ . From this a continuous interpolatable $\lambda(t)$ is generated, and the integral formula for $\frac{dG}{dp}$ is calculated using the `QuadGK.jl` implementation of Gauss-Kronrod quadrature. While this approach is memory heavy due to requiring the interpolation of the forward and reverse passes, it can be the fastest version for cases where the number of ODE/DAE states is small and the number of parameters is large since the `QuadGK` quadrature can converge faster than ODE/DAE-based versions of quadrature. This method requires an ODE or a DAE.
2. **InterpolatingAdjoint**: a checkpointed interpolation approach. This approach solves the $\lambda(t)$ ODE in reverse using an interpolation of $u(t)$, but appends the equations for $\mu(t)$ and thus does not require saving the time-series trajectory of $\lambda(t)$. For checkpointing, a scheme similar to that found in `SUNDIALS` [2] is used. Points (u_k, t_k) from the forward solution are chosen as the interval points. Whenever the backwards pass enters a new interval, the ODE is re-solved on $t \in [t_{k-1}, t_k]$ with a continuous interpolation provided by `DifferentialEquations.jl`. For the reverse pass, the `tstops` argument is set for each t_k , ensuring that no backwards integration step lies in two checkpointing intervals. This requires at most a total of

62 two forward solutions of the ODE and the memory required to hold the
 63 interpolation of the solution between two consecutive checkpoints. Note
 64 that making the checkpoints at the start and end of the integration in-
 65 terval makes this equivalent to a non-checkpointed interpolated approach
 66 which replaces the quadrature with an ODE/SDE/DAE solve for memory
 67 efficiency. This method tends to be both stable and require a minimal
 68 amount of memory, and is thus the default. This method requires an
 69 ODE, SDE, or a DAE.

70 3. BacksolveAdjoint: a checkpointed backwards solution approach. Follow-
 71 ing [3], after a forward solution, this approach solves the $u(t)$ equation in
 72 reverse along with the $\lambda(t)$ and $\mu(t)$ ODEs. Thus, since no interpolations
 73 are required, it requires $\mathcal{O}(1)$ memory. Unfortunately, many theoretical
 74 results show that backwards solution of ODEs is not guaranteed to be stable,
 75 and testing this adjoint on the universal partial differential equations
 76 like the diffusion-advection example of this paper showcases that it can be
 77 divergent and is thus not universally applicable, especially in cases of stiff-
 78 ness. Thus for stability we modify this approach by allowing checkpoints
 79 (u_k, t_k) at which the reverse integration is reset, i.e. $u(t_k) = u_k$, and the
 80 backsolve is then continued. The `tstops` argument is set in the integrator
 81 to require that each checkpoint is hit exactly for this resetting to occur.
 82 By doing so, the resulting method utilizes $\mathcal{O}(1)$ memory + the number
 83 of checkpoints required for stability, making it take the least memory ap-
 84 proach. However, the potential divergence does lead to small errors in the
 85 gradient, and thus for highly stiff equations we have found that this is
 86 only applicable to a certain degree of tolerance like 10^{-6} given reasonable
 87 numbers of checkpoints. When applicable this can be the most efficient
 88 method for large memory problems. This method requires an ODE, SDE,
 89 or a DAE.

90 4. ForwardSensitivity: a forward sensitivity approach. From $u' = f(u, p, t)$,
 91 the chain rule gives $\frac{d}{dt} \frac{du}{dp} = \frac{df}{du} \frac{du}{dp} + \frac{df}{dp}$ which can be appended to the
 92 original equations to give $\frac{du}{dp}$ as a time series, which can then be used
 93 to compute $\frac{dG}{dp}$. While the computational cost of the adjoint methods
 94 scales like $\mathcal{O}(N + P)$ for N differential equations and P parameters, this
 95 approach scales like $\mathcal{O}(NP)$ and is thus only applicable to models with
 96 small numbers of parameters (thus excluding neural networks). However,
 97 when the universal approximator has small numbers of parameters, this
 98 can be the most efficient approach. This method requires an ODE or a
 99 DAE.

100 5. ForwardDiffSensitivity: a forward-mode automatic differentiation approach,
 101 using ForwardDiff.jl [4] to calculate the forward sensitivity equations,
 102 i.e. an AD-generated implementation of forward-mode “discretize-then-
 103 optimize”. Because it utilizes a forward-mode approach, the scaling matches
 104 that of the forward sensitivity approach and it tends to have similar perfor-

105 mance characteristics. This method applies to any Julia-based differential
106 equation solver.

107 6. TrackerAdjoint: a Tracker-driven taped-based reverse-mode discrete ad-
108 joint sensitivity, i.e. an AD-generated implementation of reverse-mode
109 “discretize-then-optimize”. This is done by using the TrackedArray con-
110 structs of Tracker.jl [5] to build a Wengert list (or tape) of the forward
111 execution of the ODE solver which is then reversed. This method applies
112 to any Julia-based differential equation solver.

113 7. ZygoteAdjoint: a Zygote-driven source-to-source reverse-mode discrete
114 adjoint sensitivity, i.e. an AD-generated implementation of reverse-mode
115 “discretize-then-optimize”. This utilizes the Zygote.jl [1] system directly
116 on the differential equation solvers to generate a source code for the re-
117 verse pass of the solver itself. Currently this is only directly applicable to
118 a few differential equation solvers, but is under heavy development.

119 8. ReverseDiffAdjoint: A ReverseDiff.jl taped-based reverse-mode discrete
120 adjoint sensitivity, i.e. an AD-generated implementation of reverse-mode
121 “discretize-then-optimize”. In contrast to TrackerAdjoint, this methodol-
122 ogy can be substantially faster due to its ability to precompile the tape
123 but only supports calculations on the CPU.

124 For each of the non-AD approaches, there are the following choices for how
125 the Jacobian-vector products Jv (jvp) of the forward sensitivity equations and
126 the vector-Jacobian products $v'J$ (vjp) of the adjoint sensitivity equations are
127 computed:

128 1. Automatic differentiation for the jvp and vjp. In this approach, automatic
129 differentiation is utilized for directly calculating the jvps and vjps. For-
130 wardDiff.jl with a single dual dimension is applied at $f(u + \lambda\epsilon)$ to calculate
131 $\frac{df}{du}\lambda$ where ϵ is a dual dimensional signifier. For the vector-Jacobian prod-
132 ucts, a forward pass at $f(u)$ is utilized and the backwards pass is seeded
133 at λ to compute the $\lambda' \frac{df}{du}$ (and similarly for $\frac{df}{dp}$). Note that if f is a neu-
134 ral network, this implies that this product is computed by starting the
135 backpropogation of the neural network with λ and the vjp is the result-
136 ing return. Three methods are allowed to be chosen for performing the
137 internal vjp calculations:

138 (a) Zygote.jl source-to-source transformation based vjps. Note that only
139 non-mutating differential equation function definitions are supported
140 in this mode. This mode is the most efficient in the presence of neural
141 networks.

142 (b) ReverseDiff.jl tape-based vjps. This allows for JIT-compilation of
143 the tape for accelerated computation. This is the fastest vjp choice
144 in the presence of heavy scalar operations like in chemical reaction
145 networks, but is not compatible with GPU acceleration.

146 (c) Tracker.jl with arrays of tracked real values is utilized on mutating
147 functions.

148 The internal calculation of the vjp on a general UDE recurses down to
149 primitives and embeds optimized backpropogations of the internal neural
150 networks (and other universal approximators) for the calculation of this
151 product when this option is used.

152 2. Numerical differentiation for the jvp and vjp. In this approach, finite
153 differences is utilized for directly calculating the jvps and vjps. For a
154 small but finite ϵ , $(f(u + \lambda\epsilon) - f(u)) / \epsilon$ is used to approximate $\frac{df}{du}\lambda$. For
155 vjps, a finite difference gradient of $\lambda'f(u)$ is used.

156 3. Automatic differentiation for Jacobian construction. In this approach,
157 (sparse) forward-mode automatic differentiation is utilized by a combi-
158 nation of ForwardDiff.jl [4] with SparseDiffTools.jl for color-vector based
159 sparse Jacobian construction. After forming the Jacobian, the jvp or vjp
160 is calculated.

161 4. Numerical differentiation for Jacobian construction. In this approach,
162 (sparse) numerical differentiation is utilized by a combination of DiffEqD-
163 iffTools.jl with SparseDiffTools.jl for color-vector based sparse Jacobian
164 construction. After forming the Jacobian, the jvp or vjp is calculated.

165 In total this gives 48 different adjoint method approaches, each with different
166 performance characteristics and limitations. A full performance analysis which
167 measures the optimal adjoint approach for various UDEs has been omitted from
168 this paper, since the combinatorial nature of the options requires a considerable
169 amount of space to showcase the performance advantages and generality disad-
170 vantages between each of the approaches. A follow-up study focusing on accu-
171 rate performance measurements of the adjoint choice combinations on families
172 of UDEs is planned.

173 2 Benchmarks

174 2.1 ODE Solve Benchmarks

The three ODE benchmarks utilized the Lorenz equations (LRNZ) weather prediction model from [6] and the standard ODE IVP Testset [7, 8]:

$$\frac{dx}{dt} = \sigma(y - x) \tag{16}$$

$$\frac{dy}{dt} = x(\rho - z) - y \tag{17}$$

$$\frac{dz}{dt} = xy - \beta z \tag{18}$$

The 28 ODE benchmarks utilized the Pleiades equation (PLEI) celestial mechanics simulation from [6] and the standard ODE IVP Testset [7, 8]:

$$x_i'' = \sum_{j \neq i} m_j (x_j - x_i) / r_{ij} \quad (19)$$

$$y_i'' = \sum_{j \neq i} m_j (y_j - y_i) / r_{ij} \quad (20)$$

175 where

$$r_{ij} = ((x_i - x_j)^2 + (y_i - y_j)^2)^{3/2} \quad (21)$$

on $t \in [0, 3]$ with initial conditions:

$$u(0) = [3.0, 3.0, -1.0, -3.0, 2.0, -2.0, 2.0, 3.0, -3.0, 2.0, 0, 0, \quad (22)$$

$$-4.0, 4.0, 0, 0, 0, 0, 0, 1.75, -1.5, 0, 0, 0, -1.25, 1, 0, 0] \quad (23)$$

176 written in the form $u = [x_i, y_i, x_i', y_i']$.

177 The rest of the benchmarks were derived from a discretization a two-dimensional
 178 reaction diffusion equation, representing systems biology, combustion mechan-
 179 ics, spatial ecology, spatial epidemiology, and more generally physical PDE equa-
 180 tions:

$$A_t = D\Delta A + \alpha_A(x) - \beta_A A - r_1 AB + r_2 C \quad (24)$$

$$B_t = \alpha_B - \beta_B B - r_1 AB + r_2 C \quad (25)$$

$$C_t = \alpha_C - \beta_C C + r_1 AB - r_2 C \quad (26)$$

181 where $\alpha_A(x) = 1$ if $x > 80$ and 0 otherwise on the domain $x \in [0, 100]$,
 182 $y \in [0, 100]$, and $t \in [0, 10]$ with zero-flux boundary conditions. The diffu-
 183 sion constant D was chosen as 100 and all other parameters were left at 1.0.
 184 In this parameter regime the ODE was non-stiff as indicated by solves with
 185 implicit methods not yielding performance advantages. The diffusion operator
 186 was discretized using the second order finite difference stencil on an $N \times N$ grid,
 187 where N was chosen to be 16, 32, 64, 128, 256, and 512. To ensure fairness, the
 188 torchdiffeq functions were compiled using torchscript. The code for reproducing
 189 the benchmark can be found at:

190 <https://gist.github.com/ChrisRackauckas/cc6ac746e2dfd285c28e0584a2bfd320>

191 2.2 Neural ODE Training Benchmark

192 The spiral neural ODE from [3] was used as the benchmark for the training of
 193 neural ODEs. The data was generated according from the form:

$$u' = Au^3 \quad (27)$$

194 where $A = [-0.1, 2.0; -2.0, -0.1]$ on $t \in [0, 1.5]$ where data was taken at 30
 195 evenly spaced points. Each of the software packages trained the neural ODE

196 for 500 iterations using ADAM with a learning rate of 0.05. The defaults using
 197 the SciML software resulted in a final loss of 4.895287e-02 in 7.4 seconds, the
 198 optimized version (choosing BacksolveAdjoint with compiled ReverseDiff vector-
 199 jacobian products) resulted in a final loss of 2.761669e-02 in 2.7 seconds, while
 200 torchdiffeq achieved a final loss of 0.0596 in 289 seconds. To ensure fairness, the
 201 torchdiffeq functions were compiled using torchscript. Code to reproduce the
 202 benchmark can be found at:

203 <https://gist.github.com/ChrisRackauckas/4a4d526c15cc4170ce37da837bfc32c4>

204 2.3 SDE Solve Benchmark

205 The torchsde benchmarks were created using the geometric Brownian motion
 206 example from the torchsde README. The SDE was a 4 independent geometric
 207 Brownian motions:

$$dX_t = \mu X_t dt + \sigma X_t dW_t \quad (28)$$

208 where $\mu = 0.5$ and $\sigma = 1.0$. Both software solved the SDE 100 times using the
 209 SRI method [9] with fixed time step chosen to give 20 evenly spaced steps for
 210 $t \in [0, 1]$. The SciML ecosystem solvers solved the equation 100 times in 0.00115
 211 seconds, while torchsde v0.1 took 1.86 seconds. We contacted the author who
 212 rewrote the Brownian motion portions into C++ and linked it to torchsde as
 213 v0.1.1 and this improved the timing to roughly 5 seconds, resulting in a final
 214 performance difference of approximately 1,600x. The code to reproduce the
 215 benchmarks and the torchsde author’s optimization notes can be found at:

216 <https://gist.github.com/ChrisRackauckas/6a03e7b151c86b32d74b41af54d495c6>

217 3 Sparse Identification of Missing Model Terms 218 via Universal Differential Equations

219 The SINDy algorithm [10, 11, 12] enables data-driven discovery of governing
 220 equations from data. Notice that to use this method, derivative data \dot{X} is
 221 required. While in most publications on the subject this [10, 11, 12] information
 222 is assumed. However, for our studies we assume that only the time series infor-
 223 mation is available. Here we modify the algorithm to apply to only subsets of
 224 the equation in order to perform equation discovery specifically on the trained
 225 neural network, and in our modification the \dot{X} term is replaced with $U_\theta(t)$, the
 226 output of the universal approximator, and thus is directly computable from any
 227 trained UDE.

228 After training the UDE, choose a set of state variables:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^T(t_1) \\ \mathbf{x}^T(t_2) \\ \vdots \\ \mathbf{x}^T(t_m) \end{bmatrix} = \begin{bmatrix} x_1(t_1) & x_2(t_1) & \cdots & x_n(t_1) \\ x_1(t_2) & x_2(t_2) & \cdots & x_n(t_2) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(t_m) & x_2(t_m) & \cdots & x_n(t_m) \end{bmatrix} \quad (29)$$

229 and compute a the action of the universal approximator on the chosen states:

$$\dot{\mathbf{X}} = \begin{bmatrix} \mathbf{x}^T(t_1) \\ \dot{\mathbf{x}}^T(t_2) \\ \vdots \\ \mathbf{x}^T(t_m) \end{bmatrix} = \begin{bmatrix} \dot{x}_1(t_1) & \dot{x}_2(t_1) & \cdots & \dot{x}_n(t_1) \\ \dot{x}_1(t_2) & \dot{x}_2(t_2) & \cdots & \dot{x}_n(t_2) \\ \vdots & \vdots & \ddots & \vdots \\ \dot{x}_1(t_m) & \dot{x}_2(t_m) & \cdots & \dot{x}_n(t_m) \end{bmatrix} \quad (30)$$

230 Then evaluate the observations in a basis $\Theta(X)$. For example:

$$\Theta(\mathbf{X}) = [1 \quad \mathbf{X} \quad \mathbf{X}^{P_2} \quad \mathbf{X}^{P_3} \quad \cdots \quad \sin(\mathbf{X}) \quad \cos(\mathbf{X}) \quad \cdots] \quad (31)$$

231 where X^{P_i} stands for all P_i th order polynomial terms such as

$$\mathbf{X}^{P_2} = \begin{bmatrix} x_1^2(t_1) & x_1(t_1)x_2(t_1) & \cdots & x_2^2(t_1) & \cdots & x_n^2(t_1) \\ x_1^2(t_2) & x_1(t_2)x_2(t_2) & \cdots & x_2^2(t_2) & \cdots & x_n^2(t_2) \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_1^2(t_m) & x_1(t_m)x_2(t_m) & \cdots & x_2^2(t_m) & \cdots & x_n^2(t_m) \end{bmatrix} \quad (32)$$

232 Using these matrices, find this sparse basis Ξ over a given candidate library
 233 Θ by solving the sparse regression problem $\dot{X} = \Theta\Xi$ with L_1 regularization,
 234 i.e. minimizing the objective function $\|\dot{\mathbf{X}} - \Theta\Xi\|_2 + \lambda\|\Xi\|_1$. This method and
 235 other variants of SInDy applied to UDEs, along with specialized optimizers for
 236 the LASSO L_1 optimization problem, have been implemented by the authors
 237 and collaborators DataDrivenDiffEq.jl library for use with the DifferentialEqua-
 238 tions.jl training framework.

239 3.1 Application to the Partial Reconstruction of the Lotka- 240 Volterra Equations

241 On the Lotka-Volterra equations, we trained a UDE model against a trajec-
 242 tory of 31 points with a constant step size $\Delta t = 0.1$ starting from $x_0 =$
 243 0.44249296 , $y_0 = 4.6280594$ to recover the function $U_\theta(x, y)$. The parame-
 244 ters are chosen to be $\alpha = 1.3$, $\beta = 0.9$, $\gamma = 0.8$, $\delta = 1.8$. The trajectory has
 245 been perturbed with additive noise drawn from a normal distribution scaled by
 246 10^{-3} .

247 The neural network consists of an input layer, one hidden layers with 32
 248 neurons and a linear output layer. The input and hidden layer have hyperbolic
 249 tangent activation functions. We trained for 200 iterations with ADAM with a
 250 learning rate $\gamma = 10^{-2}$. We then switched to BFGS with an initial stepnorm
 251 of $\gamma = 10^{-2}$ setting the maximum iterations to 10000. Typically the training
 252 converged after 400-600 iterations in total. The loss was chosen was the L2 loss
 253 $\mathcal{L} = \sum_i (u_\theta(t_i) - d_i)^2$. The training converged to a final loss between 10^{-2} and
 254 10^{-5} for the knowledge-enhanced neural network.

255 From the trained neural network, data was sampled over the original tra-
 256 jectory and fitted using the SR3 method with varying threshold between 10^{-7}
 257 and 10^3 logarithmically spaced to ensure a sparse solution. A pareto optimal

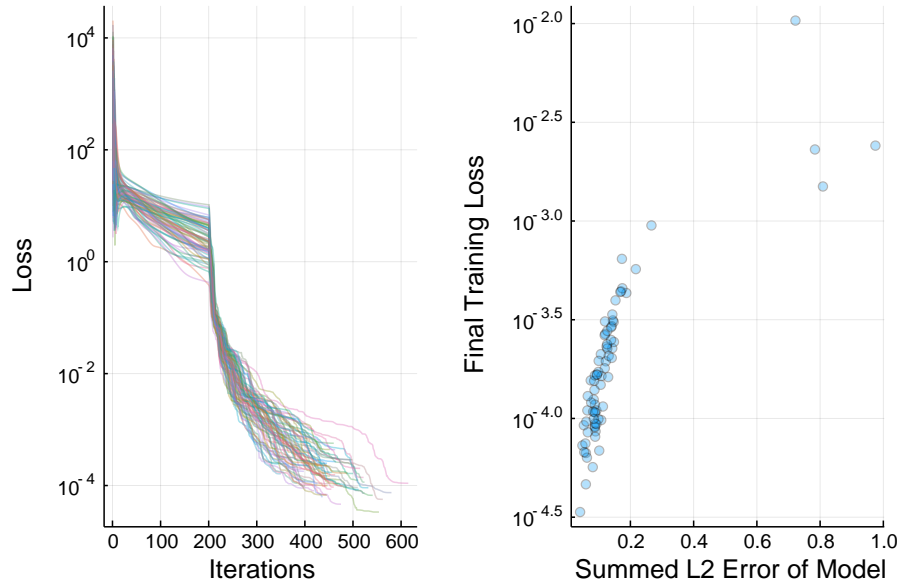


Figure 1: Loss trajectories and errors for 74 runs of the partially reconstruction of the Lotka Volterra equations under noisy data.

258 solution was selected via the L1 norm of the coefficients and the L2 norm of the
 259 corresponding error between the differential data and estimate. The knowledge-
 260 enhanced neural network returned zero for all terms except non-zeros on the
 261 quadratic terms with nearly the correct coefficients that were then fixed using
 262 an additional sparse regression over the quadratic terms. The resulting param-
 263 eters extracted are $\beta \approx 0.9239$ and $\gamma \approx 0.8145$. Performing a sparse identification
 264 on the ideal, full solution and numerical derivatives computed via an interpo-
 265 lating spline resulted in a failure.

266 To ensure reproducible results, the training procedure has been evaluated
 267 74 times given the initial trajectory. The training failed once (1.35 %) to
 268 exactly identify the original equation in symbolic terms. The final loss has a
 269 mean of 0.0004 with a standard error of 0.0012, an upper and lower bound of
 270 0.0103 and $3.3573 \cdot 10^{-5}$ respectively. The training converged between 241 and
 271 615 iterations. The sum of the L2 error of the recovered model on the training
 272 data is centered around 0.1472 with a standard deviation of 0.1688. The results
 273 are shown in Figure 1.

274 **3.2 Discovery of Robertson’s Equations with Prior Con-** 275 **servations Laws**

276 On Robertson’s equations, we trained a UDAE model against a trajectory of
 277 10 points on the timespan $t \in [0.0, 1.0]$ starting from $y_1 = 1.0$, $y_2 = 0.0$, and

278 $y_3 = 0.0$. The parameters of the generating equation were $k_1 = 0.04$, $k_2 = 3e7$,
 279 and $k_3 = 1e4$. The universal approximator was a neural network with one
 280 hidden layers of size 64. The equation was trained using the BFGS optimizer
 281 to a loss of $9e - 6$.

282 4 Model-based Learning for the Fisher-KPP Equa- 283 tions

284 To generate training data for the 1D Fisher-KPP equation 15, we take the
 285 growth rate and the diffusion coefficient to be $r = 1$ and $D = 0.01$ respectively.
 286 The equation is numerically solved in the domain $x \in [0, 1]$ and $t \in [0, T]$
 287 using a 2nd order central difference scheme for the spatial derivatives and the
 288 time-integration is done using the Tsitouras 5/4 Runge-Kutta method. We
 289 implement periodic boundary condition $\rho(x = 0, t) = \rho(x = 1, t)$ and initial
 290 condition $\rho(x, t = 0) = \rho_0(x)$ is taken to be a localized function given by

$$\rho_0(x) = \frac{1}{2} \left(\tanh \left(\frac{x - (0.5 - \Delta/2)}{\Delta/10} \right) - \tanh \left(\frac{x - (0.5 + \Delta/2)}{\Delta/10} \right) \right), \quad (33)$$

291 with $\Delta = 0.2$ which represents the width of the region where $\rho \simeq 1$. The data
 292 are saved at evenly spaced points with $\Delta x = 0.04$ and $\Delta t = 0.5$.

293 In the UPDE 16, the growth neural network $\text{NN}_\theta(\rho)$ has 4 densely connected
 294 layers with 10, 20, 20 and 10 neurons each and tanh activation functions. The
 295 diffusion operator is represented by a CNN that operates on an input vector of
 296 arbitrary size. It has 1 hidden layer with a 3×1 filter $[w_1, w_2, w_3]$ without any
 297 bias. To implement periodic boundary conditions for the UPDE at each time
 298 step, the vector of values at different spatial locations $[\rho_1, \rho_2, \dots, \rho_{N_x}]$ is padded
 299 with ρ_{N_x} to the left and ρ_1 to the right. This also ensures that the output of
 300 the CNN is of the same size as the input. The weights of both the neural
 301 networks and the diffusion coefficient are simultaneously trained to minimize
 302 the loss function

$$\mathcal{L} = \sum_i (\rho(x_i, t_i) - \rho_{\text{data}}(x_i, t_i))^2 + \lambda |w_1 + w_2 + w_3|, \quad (34)$$

303 where λ is taken to be 10^2 (note that one could also structurally enforce $w_3 =$
 304 $-(w_1 + w_2)$). The second term in the loss function enforces that the differential
 305 operator that is learned is conservative—that is, the weights sum to zero. The
 306 training is done using the ADAM optimizer with learning rate 10^{-3} .

307 **5 Adaptive Solving for the 100 Dimensional Hamilton-**
 308 **Jacobi-Bellman Equation**

309 **5.1 Forward-Backwards SDE Formulation**

310 Consider the class of semilinear parabolic PDEs, in finite time $t \in [0, T]$ and
 311 d -dimensional space $x \in \mathbb{R}^d$, that have the form:

$$\begin{aligned} \frac{\partial u}{\partial t}(t, x) + \frac{1}{2} \text{tr}(\sigma \sigma^T(t, x) (\text{Hess}_x u)(t, x)) \\ + \nabla u(t, x) \cdot \mu(t, x) \\ + f(t, x, u(t, x), \sigma^T(t, x) \nabla u(t, x)) = 0, \end{aligned} \quad (35)$$

312 with a terminal condition $u(T, x) = g(x)$. In this equation, tr is the trace of
 313 a matrix, σ^T is the transpose of σ , ∇u is the gradient of u , and $\text{Hess}_x u$ is the
 314 Hessian of u with respect to x . Furthermore, μ is a vector-valued function, σ
 315 is a $d \times d$ matrix-valued function and f is a nonlinear function. We assume that
 316 μ , σ , and f are known. We wish to find the solution at initial time, $t = 0$, at
 317 some starting point, $x = \zeta$.

318 Let W_t be a Brownian motion and take X_t to be the solution to the stochastic
 319 differential equation

$$dX_t = \mu(t, X_t)dt + \sigma(t, X_t)dW_t \quad (36)$$

320 with a terminal condition $u(T, x) = g(x)$. With initial condition $X(0) = \zeta$
 321 has shown that the solution to 17 satisfies the following forward-backward SDE
 322 (FBSDE) [13]:

$$\begin{aligned} u(t, X_t) - u(0, \zeta) = \\ - \int_0^t f(s, X_s, u(s, X_s), \sigma^T(s, X_s) \nabla u(s, X_s)) ds \\ + \int_0^t [\nabla u(s, X_s)]^T \sigma(s, X_s) dW_s, \end{aligned} \quad (37)$$

323 with terminating condition $g(X_T) = u(X_T, W_T)$. Notice that we can combine
 324 36 and 37 into a system of $d + 1$ SDEs:

$$\begin{aligned} dX_t &= \mu(t, X_t)dt + \sigma(t, X_t)dW_t, \\ dU_t &= f(t, X_t, U_t, \sigma^T(t, X_t) \nabla u(t, X_t))dt \\ &\quad + [\sigma^T(t, X_t) \nabla u(t, X_t)]^T dW_t, \end{aligned} \quad (38)$$

325 where $U_t = u(t, X_t)$. Since X_0 , μ , σ , and f are known from the choice of
 326 model, the remaining unknown portions are the functional $\sigma^T(t, X_t) \nabla u(t, X_t)$
 327 and initial condition $U(0) = u(0, \zeta)$, the latter being the point estimate solution
 328 to the PDE.

329 To solve this problem, we approximate both unknown quantities by universal
 330 approximators:

$$\begin{aligned}\sigma^T(t, X_t)\nabla u(t, X_t) &\approx U_{\theta_1}^1(t, X_t), \\ u(0, X_0) &\approx U_{\theta_2}^2(X_0),\end{aligned}\tag{39}$$

331 Therefore we can rewrite 38 as a stochastic UDE of the form:

$$\begin{aligned}dX_t &= \mu(t, X_t)dt + \sigma(t, X_t)dW_t, \\ dU_t &= f(t, X_t, U_t, U_{\theta_1}^1(t, X_t))dt + [U_{\theta_1}^1(t, X_t)]^T dW_t,\end{aligned}\tag{40}$$

332 with initial condition $(X_0, U_0) = (X_0, U_{\theta_2}^2(X_0))$.

333 To be a solution of the PDE, the approximation must satisfy the terminating
 334 condition, and thus we define our loss to be the expected difference between the
 335 approximating solution and the required terminating condition:

$$l(\theta_1, \theta_2 | X_T, U_T) = \mathbb{E}[\|g(X_T) - U_T\|].\tag{41}$$

336 Finding the parameters (θ_1, θ_2) which minimize this loss function thus give
 337 rise to a BSDE which solves the PDE, and thus $U_{\theta_2}^2(X_0)$ is the solution to the
 338 PDE once trained.

339 5.2 The LQG Control Problem

340 This PDE can be rewritten into the canonical form by setting:

$$\begin{aligned}\mu &= 0, \\ \sigma &= \bar{\sigma}I, \\ f &= -\alpha \|\sigma^T(s, X_s)\nabla u(s, X_s)\|^2,\end{aligned}\tag{42}$$

341 where $\bar{\sigma} = \sqrt{2}$, $T = 1$ and $X_0 = (0, \dots, 0) \in R^{100}$. The universal stochastic
 342 differential equation was then supplemented with a neural network as the ap-
 343 proximator. The initial condition neural network was had 1 hidden layer of size
 344 110, and the $\sigma^T(t, X_t)\nabla u(t, X_t)$ neural network had two layers both of size 110.
 345 For the example we chose $\lambda = 1$. This was trained with the LambaEM method
 346 of DifferentialEquations.jl [14] with relative and absolute tolerances set at $1e-4$
 347 using 500 training iterations and using a loss of 100 trajectories per epoch.

348 On this problem, for an arbitrary g , one can show with Itô's formula that:

$$u(t, x) = -\frac{1}{\lambda} \ln \left(\mathbb{E} \left[\exp \left(-\lambda g(x + \sqrt{2}W_{T-t}) \right) \right] \right),\tag{43}$$

349 which was used to calculate the error from the true solution.

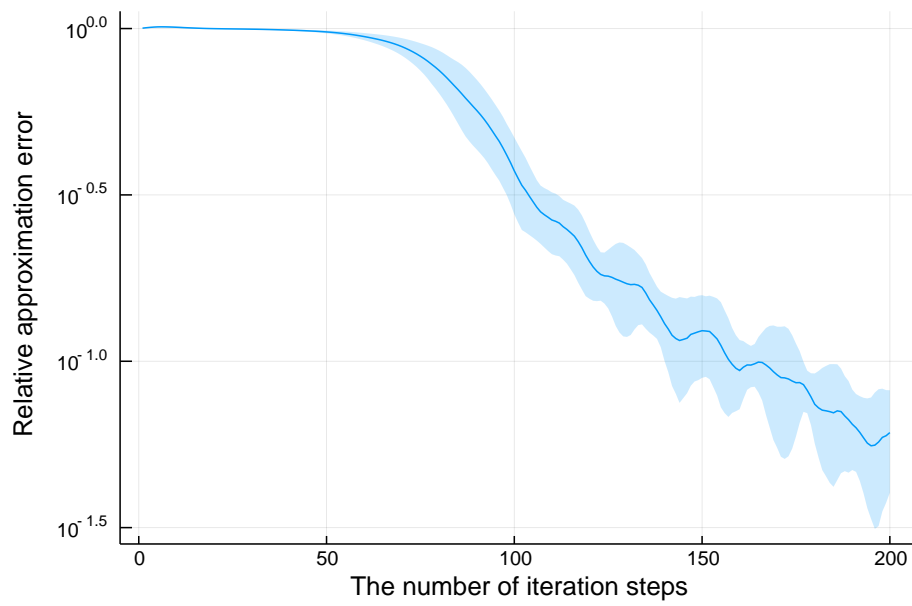


Figure 2: Adaptive solution of the 100-dimensional Hamilton-Jacobi-Bellman equation. This demonstrates that as the universal approximators $U_{\theta_1}^1$ and $U_{\theta_2}^2$ converge to satisfy the terminating condition, $U_{\theta_2}^2$ network converges to the solution of Equation 20.

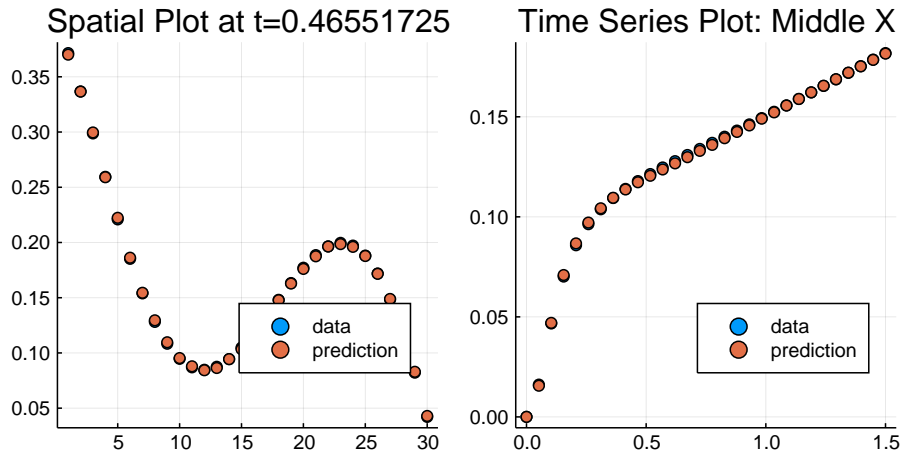


Figure 3: Reduction of the Boussinesq equations. On the left is the comparison between the training data (blue) and the trained UPDE (orange) over space at the 10th fitting time point, and on the right is the same comparison shown over time at spatial midpoint.

350 6 Reduction of the Boussinesq Equations

351 As a test for the diffusion-advection equation parameterization approach, data
 352 was generated from the diffusion-advection equations using the missing function
 353 $\overline{wT} = \cos(\sin(T^3)) + \sin(\cos(T^2))$ with N spatial points discretized by a finite
 354 difference method with $t \in [0, 1.5]$ with Neumann zero-flux boundary conditions.
 355 A neural network with two hidden layers of size 8 and tanh activation functions
 356 was trained against 30 data points sampled from the true PDE. The UPDE was
 357 fit by using the ADAM optimizer with learning rate 10^{-2} for 200 iterations and
 358 then ADAM with a learning rate of 10^{-3} for 1000 iterations. The resulting fit
 359 is shown in 3 which resulted in a final loss of approximately 0.007. We note
 360 that the stabilized adjoints were required for this equation, i.e. the backsolve
 361 adjoint method was unstable and results in divergence and thus cannot be used
 362 on this type of equation. The trained neural network had a forward pass that
 363 took around 0.9 seconds.

364 For the benchmark against the full Boussinesq equations, we utilized Oceanani-
 365 gans.jl [15]. It was set to utilize adaptive time stepping to maximize the time
 366 step according to the CFL condition number (capped at $\text{CFL} \leq 0.3$) and
 367 matched the boundary conditions, along with setting periodic boundary condi-
 368 tions in the horizontal dimension. The Boussinesq simulation used $128 \times 128 \times 128$
 369 spatial points, a larger number than the parameterization, in order to accurately
 370 resolve the mean statistics of the 3-dimensional dynamics as is commonly re-
 371 quired in practice [16]. The resulting simulation took 13,737 seconds on the same
 372 computer used for the neural diffusion-advection approach, demonstrating the
 373 approximate 15,000x acceleration.

374 **7 Automated Derivation of Closure Relations**
 375 **for Viscoelastic Fluids**

The full FENE-P model is:

$$\boldsymbol{\sigma} + g\left(\frac{\lambda}{f(\boldsymbol{\sigma})}\boldsymbol{\sigma}\right) = \frac{\eta}{f(\boldsymbol{\sigma})}\dot{\boldsymbol{\gamma}}, \quad (44)$$

$$f(\boldsymbol{\sigma}) = \frac{L^2 + \frac{\lambda(L^2-3)}{L^2\eta}\text{Tr}(\boldsymbol{\sigma})}{L^2 - 3}, \quad (45)$$

where

$$g(\mathbf{A}) = \frac{D\mathbf{A}}{Dt} - (\nabla\mathbf{u}^T)\mathbf{A} - \mathbf{A}(\nabla\mathbf{u}^T),$$

376 is the upper convected derivative, and L , η , λ are parameters [17]. For a one
 377 dimensional strain rate, $\dot{\boldsymbol{\gamma}} = \dot{\boldsymbol{\gamma}}_{12} = \dot{\boldsymbol{\gamma}}_{21} \neq 0$, $\dot{\boldsymbol{\gamma}}_{ij} = 0$ else, the one dimensional
 378 stress required is $\boldsymbol{\sigma} = \sigma_{12}$. However, σ_{11} and σ_{22} are both non-zero and store
 379 memory of the deformation (normal stresses). The Oldroyd-B model is the
 380 approximation:

$$G(t) = 2\eta\delta(t) + G_0e^{-t/\tau}, \quad (46)$$

with the exact closure relation:

$$\boldsymbol{\sigma}(t) = \eta\dot{\boldsymbol{\gamma}}(t) + \boldsymbol{\phi}, \quad (47)$$

$$\frac{d\boldsymbol{\phi}}{dt} = G_0\dot{\boldsymbol{\gamma}} - \boldsymbol{\phi}/\tau. \quad (48)$$

As an arbitrary nonlinear extension, train a UDE model using a single ad-
 ditional memory field against simulated FENE-P data with parameters $\lambda = 2$,
 $L = 2$, $\eta = 4$. The UDE model is of the form,

$$\boldsymbol{\sigma} = U_0(\boldsymbol{\phi}, \dot{\boldsymbol{\gamma}}) \quad (49)$$

$$\frac{d\boldsymbol{\phi}}{dt} = U_1(\boldsymbol{\phi}, \dot{\boldsymbol{\gamma}}) \quad (50)$$

381 where U_0, U_1 are neural networks each with a single hidden layer containing 4
 382 neurons. The hidden layer has a tanh activation function. The loss was taken as
 383 $\mathcal{L} = \sum_i (\boldsymbol{\sigma}(t_i) - \boldsymbol{\sigma}_{\text{FENE-P}}(t_i))^2$ for 100 evenly spaced time points in $t_i \in [0, 2\pi]$,
 384 and the system was trained using an ADAM iterator with learning rate 0.015.
 385 The fluid is assumed to be at rest before $t = 0$, making the initial stress also
 386 zero.

387 **References**

- 388 [1] Mike Innes, Alan Edelman, Keno Fischer, Chris Rackauckus, Elliot Saba,
 389 Viral B Shah, and Will Tebbutt. Zygote: A differentiable programming
 390 system to bridge machine learning and scientific computing. *arXiv preprint*
 391 *arXiv:1907.07587*, 2019.

- 392 [2] Alan C Hindmarsh, Peter N Brown, Keith E Grant, Steven L Lee, Radu
393 Serban, Dan E Shumaker, and Carol S Woodward. SUNDIALS: Suite of
394 nonlinear and differential/algebraic equation solvers. *ACM Transactions*
395 *on Mathematical Software (TOMS)*, 31(3):363–396, 2005.
- 396 [3] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud.
397 Neural ordinary differential equations. In *Advances in neural information*
398 *processing systems*, pages 6571–6583, 2018.
- 399 [4] J. Revels, M. Lubin, and T. Papamarkou. Forward-mode automatic differ-
400 entiation in julia. *arXiv:1607.07892 [cs.MS]*, 2016.
- 401 [5] Michael Innes, Elliot Saba, Keno Fischer, Dhairya Gandhi, Marco Concetto
402 Rudilosso, Neethu Mariya Joy, Tejan Karmali, Avik Pal, and Viral Shah.
403 Fashionable modelling with flux. *CoRR*, abs/1811.01457, 2018.
- 404 [6] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving Ordinary Differential*
405 *Equations I (2nd Revised. Ed.): Nonstiff Problems*. Springer-Verlag, Berlin,
406 Heidelberg, 1993.
- 407 [7] Francesca Mazzia and Cecilia Magherini. Test set for initial value prob-
408 lem solvers, release 2.4. Technical Report 4, Department of Mathematics,
409 University of Bari, Italy, February 2008.
- 410 [8] Francesca Mazzia and Cecilia Magherini. *Test Set for Initial Value Problem*
411 *Solvers, release 2.4*. Department of Mathematics, University of Bari and
412 INdAM, Research Unit of Bari, February 2008.
- 413 [9] Andreas Rößler. Runge–kutta methods for the strong approximation of
414 solutions of stochastic differential equations. *SIAM Journal on Numerical*
415 *Analysis*, 48(3):922–952, 2010.
- 416 [10] Steven L Brunton, Joshua L Proctor, and J Nathan Kutz. Discovering gov-
417 erning equations from data by sparse identification of nonlinear dynamical
418 systems. *Proceedings of the National Academy of Sciences*, 113(15):3932–
419 3937, 2016.
- 420 [11] Niall M Mangan, Steven L Brunton, Joshua L Proctor, and J Nathan
421 Kutz. Inferring biological networks by sparse identification of nonlinear
422 dynamics. *IEEE Transactions on Molecular, Biological and Multi-Scale*
423 *Communications*, 2(1):52–63, 2016.
- 424 [12] Niall M Mangan, J Nathan Kutz, Steven L Brunton, and Joshua L Proctor.
425 Model selection for dynamical systems via sparse regression and informa-
426 tion criteria. *Proceedings of the Royal Society A: Mathematical, Physical*
427 *and Engineering Sciences*, 473(2204):20170009, 2017.
- 428 [13] Jianfeng Zhang. Backward stochastic differential equations. In *Backward*
429 *Stochastic Differential Equations*, pages 79–99. Springer, 2017.

- 430 [14] Christopher Rackauckas and Qing Nie. Differentialequations.jl – a performant and feature-rich ecosystem for solving differential equations in julia. *The Journal of Open Research Software*, 5(1), 2017. Exported from <https://app.dimensions.ai> on 2019/05/05.
- 431
432
433
- 434 [15] Ali Ramadhan, Gregory LeClaire Wagner, Chris Hill, Jean-Michel Campin, Valentin Churavy, Tim Besard, Andre Souza, Alan Edelman, John Marshall, and Raffaele Ferrari. Oceananigans.jl: Fast and friendly geophysical fluid dynamics on GPUs. *The Journal of Open Source Software*, 4(44):1965, 2020.
- 435
436
437
438
- 439 [16] Benoit Cushman-Roisin and Jean-Marie Beckers. Chapter 4 - equations governing geophysical flows. In Benoit Cushman-Roisin and Jean-Marie Beckers, editors, *Introduction to Geophysical Fluid Dynamics*, volume 101 of *International Geophysics*, pages 99 – 129. Academic Press, 2011.
- 440
441
442
- 443 [17] P.J. Oliveira. Alternative derivation of differential constitutive equations of the oldroyd-b type. *Journal of Non-Newtonian Fluid Mechanics*, 160(1):40 – 46, 2009. Complex flows of complex fluids.
- 444
445