

---

# An Extensible Benchmark Suite for Learning to Simulate Physical Systems

---

**Karl Otness**

Courant Institute of Mathematical Sciences  
New York University  
karl.otness@nyu.edu

**Arvi Gjoka**

Courant Institute of Mathematical Sciences  
New York University  
arvi.gjoka@nyu.edu

**Joan Bruna**

Courant Institute of Mathematical Sciences  
New York University

**Daniele Panozzo**

Courant Institute of Mathematical Sciences  
New York University

**Benjamin Peherstorfer**

Courant Institute of Mathematical Sciences  
New York University

**Teseo Schneider**

University of Victoria

**Denis Zorin**

Courant Institute of Mathematical Sciences  
New York University

## Abstract

1 Simulating physical systems is a core component of scientific computing, encom-  
2 passing a wide range of physical domains and applications. Recently, there has been  
3 a surge in data-driven methods to complement traditional numerical simulations  
4 methods, motivated by the opportunity to reduce computational costs and/or learn  
5 new physical models leveraging access to large collections of data. However, the  
6 diversity of problem settings and applications has led to a plethora of approaches,  
7 each one evaluated on a different setup and with different evaluation metrics. We  
8 introduce a set of benchmark problems to take a step towards unified benchmarks  
9 and evaluation protocols. We propose four representative physical systems, as well  
10 as a collection of both widely used classical time integrators and representative  
11 data-driven methods (kernel-based, MLP, CNN, Nearest-Neighbors). Our frame-  
12 work allows to evaluate objectively and systematically the stability, accuracy, and  
13 computational efficiency of data-driven methods. Additionally, it is configurable to  
14 permit adjustments for accommodating other learning tasks and for establishing a  
15 foundation for future developments in machine learning for scientific computing.

## 16 1 Introduction

17 Computational modeling of physical systems is a core task of scientific computing. Standard methods  
18 rely on discretizations of explicit models typically given in the form of partial differential equations  
19 (PDEs). Machine learning techniques can extend these techniques in a number of ways. In some  
20 cases, a closed system of analytic equations relating all variables may not be available (e.g., a  
21 constitutive relation for a material may not be known). In other cases, while a full analytic description  
22 of a system is available, a traditional solution may be too costly (e.g., turbulence) or can be sped

23 up substantially using data-driven reduced-order models. However, despite promising results, a  
24 successful adoption of these data-driven approaches into scientific computing pipelines requires  
25 a solid and exhaustive assessment of their performance—a challenging task given the diversity  
26 of physical systems, corresponding data-driven approaches, and the lack of standardized sets of  
27 problems, comparison protocols, and metrics.

28 We focus on the setting where the physical model is unavailable during training, mimicking situations  
29 in computational science and engineering with ample data and a lack of models. One can generally  
30 distinguish two different flavors of physical simulation with different associated computational  
31 cost: those that map a high-dimensional state space into another high-dimensional space (as in  
32 temporal integration schemes, mapping the state of the system at one time step to the next), or from  
33 a high-dimensional input space to a lower-dimensional output (as in surrogate models, mapping  
34 the initial conditions to a functional of the solution). While this distinction also applies to data-  
35 driven approaches, another critical aspect emerges, given by the choice of input data distribution.  
36 We identify two extremes: the *narrow* data regime, where initial conditions are sampled from a  
37 low-dimensional manifold (even within a high-dimensional state space), and the *wide* regime, where  
38 initial conditions span a truly high-dimensional space. As could be expected, narrow data regimes  
39 define an easier prediction task where data-driven methods can potentially ‘bypass the physics’,  
40 whereas wide regimes require models with enough encoded physical priors in order to beat the  
41 curse of dimensionality. Therefore, such choice of data distribution is a critical component of any  
42 data-driven physical simulation benchmark.

43 In this work, we introduce an extensible benchmark suite, including: (1) an extensible set of simple,  
44 yet representative, physical models with a range of training and evaluation (test) setups, as well  
45 as reference, high-accuracy numerical solutions to benchmark data-driven methods, (2) reference  
46 implementations of traditional time integration schemes, which are used as baselines for evaluation,  
47 and (3) implementations of widely used data-driven methods, including physics-agnostic multi-layer  
48 perceptrons (MLPs), convolutional neural networks (CNNs), kernel machines and non-parametric  
49 Nearest Neighbors. Our benchmark suite is modular, permitting extensions with limited code  
50 changes, and captures both ‘narrow’ and ‘wide’ regimes by appropriately parametrizing the set of  
51 initial conditions.

52 Our analysis reveals two important conclusions: First, even in the simplest physical models, current  
53 data-driven pipelines, while providing qualitatively acceptable solutions, are quantitatively far from  
54 directly numerically integrating physical models, and this performance gap appears unfeasible to  
55 close by merely scaling up the models and/or the dataset size. In other words, the cost of ignoring  
56 the physics is high, even for the simplest physics, and cannot in general be compensated by data,  
57 matching insights that have been obtained in other scientific computing settings [5, 54]. Next,  
58 and more importantly, our simple  $L^2$ -based nearest neighbor regressor is used to calibrate how  
59 ‘narrow’ the learning task is. Our finding is that even for seemingly complex systems, such as  
60 the incompressible Navier-Stokes systems, such naive predictor outperforms most deep-learning-  
61 based models in the narrow regime—thus providing a simple calibration of the true difficulty of the  
62 simulation task, that we advocate should be present in every future evaluation.

## 63 2 Related Work

64 Machine learning is used in physical simulation in a number of interrelated ways. Some important uses  
65 include reduced-order/surrogate modeling, learning constitutive models or more generally compact  
66 analytic representations from data. A unifying theme of these applications of machine learning is  
67 automatic construction of parametric models capable of reproducing the behavior of physical systems  
68 for a sufficiently broad range of initial data, boundary conditions and other system parameters.  
69 The purpose of these representations varies from acceleration (e.g., surrogate machine learning  
70 models are used to accelerate optimization), to automatic construction of multiscale models (learning  
71 macroscopic constitutive laws from microscopic simulation), to inferring compact descriptions of  
72 unknown representations from experimental data.

73 The purpose of our proposed benchmarks is to enable comparisons of different learning-based  
74 methods in terms of their accuracy and efficiency. We briefly review two streams of learning  
75 methods for physical systems: **(1)** One line of work aims to understand how neural networks can  
76 be structured and trained to reproduce known physical system behavior, with the goal of designing

77 general methods applicable in a variety of settings [8, 41, 4, 40, 36, 37, 25, 10, 49, 50]. Our  
 78 benchmark cases fit primarily into this category. **(2)** Another line of research aims to develop  
 79 a variety of techniques to accelerate solving PDEs. Typically, these methods are developed for  
 80 specific PDEs and a specific restricted range of problems. For example, fluid dynamics problems  
 81 [38, 16, 55], with particular applications to cardiovascular modeling [24, 19] and aerodynamics  
 82 [52]; or solid mechanics simulation tasks, including stresses [28, 23, 26, 15, 21, 22]. In cases where  
 83 the governing equations are not given, the learning task becomes approximating them from data  
 84 [29, 7, 1, 9, 2, 27, 3, 42, 43, 45, 44, 51, 34].

### 85 3 Background and problem setup

86 **PDEs, dynamical systems, and time integration.** Consider a time-dependent PDE of the form  
 87  $\partial_t u = \mathcal{L}(u)$ , where  $u$  is the unknown function and  $\mathcal{L}$  is a possibly nonlinear operator that includes  
 88 spatial derivatives of  $u$ . By discretizing in space, one obtains a dynamical system

$$\dot{x}(t) = f(x(t)) \quad (1)$$

89 with an  $N$ -dimensional state  $x(t) \in \mathbb{R}^N$  at time  $t \in [0, T]$ . The function  $f$  is assumed to be Lipschitz  
 90 to ensure solution uniqueness and the initial condition is denoted as  $x_0 \in \mathbb{R}^N$ . A PDE of a higher  
 91 order in time can be reduced to the first-order form in the standard way, e.g., if we have a second-  
 92 order system  $\ddot{q}(t) = f(q(t))$ , then we consider its formulation via position  $q$  and momentum  $p$  as  
 93 a first-order system with  $x = [q; p]: [\dot{q}(t); \dot{p}(t)] = [p(t); f(q(t))]$ . To numerically integrate (1), we  
 94 choose time steps  $0 = t_0 < t_1 < \dots < t_K = T$ . Then, a time integration scheme (e.g., [48, 12, 11])  
 95 gives an approximation  $x_k \approx x(t_k)$  of the state  $x(t_k)$  at each time step  $k = 1, \dots, K$ . A list of the  
 96 schemes we use along with details is given in Appendix A.

97 **Problem setup and learning problems.** Given  $M$  initial conditions  $x_0^{(1)}, \dots, x_0^{(M)} \in \mathbb{R}^N$  and the  
 98 corresponding  $M$  trajectories  $X^{(i)} = [x_0^{(i)}, \dots, x_K^{(i)}] \in \mathbb{R}^{N \times (K+1)}$ ,  $i = 1, \dots, M$  obtained with a  
 99 time integration scheme from dynamical system (1), we consider the following two learning problems,  
 100 both of which aim to learn the physical model of the problem, viewed as unknown, from trajectory  
 101 samples: **(1)** Learning an approximation  $\tilde{f}$  of the right-hand side function  $f$  in Eq. (1). This gives  
 102 an approximate  $\dot{\tilde{x}}(t) = \tilde{f}(\tilde{x}(t))$  that is then numerically integrated to produce a trajectory  $\tilde{X}$  for an  
 103 initial condition  $\tilde{x}_0$ . The aim is that  $\tilde{X}$  approximates well the true trajectory  $X$  obtained with  $f$  from  
 104 (1) for the same initial condition. **(2)** Directly learning the next steps in the trajectory from the current  
 105 one, i.e. predict  $x_k^{(i)}$  given  $x_{k-1}^{(i)}$ .

106 To assess the learned models, we evaluate them on their ability to produce good approximate  
 107 trajectories from randomly sampled initial conditions, by either integration or direct step prediction.  
 108 During evaluation, we use initial conditions drawn independently from those used to produce training  
 109 data, both from the same distribution as the training samples, as well as from a distribution with  
 110 support outside the training range. We train networks on data sets of various sizes. For details, see  
 111 Appendix B.

### 112 4 Benchmark Systems

113 We consider four physical systems, illustrated in Figure 1: a single oscillating spring, a one-  
 114 dimensional linear wave equation, a Navier-Stokes flow problem and a mesh of damped springs.  
 115 These systems represent a progression of complexity: the spring system is a linear system with  
 116 low-dimensional space of initial conditions and low-dimensional state; the wave equation is a low-  
 117 dimensional linear system with a (relatively) high-dimensional state space after discretization; the  
 118 Navier-Stokes equations are nonlinear and we consider a setup with low-dimensional initial condi-  
 119 tions and high-dimensional state space; finally, the spring mesh system has both high-dimensional  
 120 initial conditions as well as high-dimensional states. Additionally, the proposed spring system and  
 121 Navier-Stokes problems represent diffusion-dominated and advection-dominated (for sufficiently  
 122 low viscosity) PDE behaviors, as well as variability in initial conditions with fixed domain (spring  
 123 system) and variable domain (Navier-Stokes). These varying complexities provide an opportunity to  
 124 test methods on simpler systems and the ability to examine changing performance as system size  
 125 increases, both in terms of the state dimension, and the initial condition distribution. The ground truth  
 126 models for the spring, wave, and spring mesh systems with classical time integrators are implemented

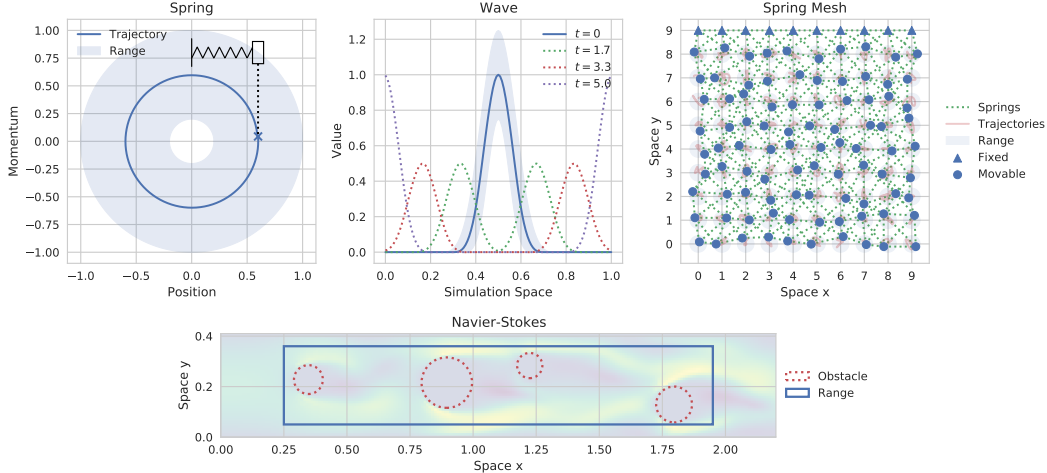


Figure 1: Representative visualizations of the three systems, depicting the results and ranges of initial condition sampling. Each two state components: for the Navier-Stokes system, a flow velocity and a pressure field, and for the other three a position  $q$  and momentum  $p$ .

127 using NumPy [13], SciPy [53], and accelerated, where possible, with Numba [20]. The Navier-Stokes  
 128 snapshots are generated using PolyFEM [46], a finite element library.

129 These systems were chosen in an effort to reflect the variety of systems used for testing in this  
 130 area, while unifying choices of particular formulations. Past works have chosen systems of the  
 131 types featured here: simple oscillators (both spring and pendulum [8]), particle systems with various  
 132 interaction laws (gravity, spring forces, charges, cloth simulations, etc. [6, 18, 41, 4, 33]), and fluid-  
 133 flow systems (with various sorts of obstacles, airfoils or cylinders [50, 33]). We make particular  
 134 selections here in an effort to unify systems of interest and facilitate comparisons across experiments  
 135 by providing a shared set of tasks which can be used for development and testing of machine learning  
 136 methods.

137 Some examples of initial condition selection for each system are illustrated in Figure 1. The ground  
 138 truth for the spring, wave, and spring mesh systems consists of the state variables  $(q, p)$  for position  
 139 and momentum, and their associated derivatives  $(\dot{q}, \dot{p})$ . For the Navier-Stokes system the state consists  
 140 of flow velocities, and a pressure field, along with approximated time derivatives for each.

141 Table 1 lists the parameters used to generate trajectories for training and evaluation. Training sets  
 142 of three sizes are generated, each containing the specified number of trajectories. The systems are  
 143 integrated at the listed time step sizes, but the ground truth data is subsampled further by the factor  
 144 shown after  $\div$  in the table: the integration schemes are run at a smaller time step and intermediate  
 145 computations are discarded. Each larger training set is a strict superset of its predecessor to ensure  
 146 that previous training samples are never removed.

#### 147 4.1 Spring

148 We simulate a simple one-dimensional oscillating spring. In this system, the spring has zero rest  
 149 length, and both the oscillating mass and spring constant are set to 1. The spring then exerts a force  
 150 inversely proportional to the position of the mass:  $q: \dot{p}(t) = -q$  and  $\dot{q}(t) = p$ .

151 The energy of the system is proportional to  $r = q^2 + p^2$  which is the radius of the circle in phase space.  
 152 To sample initial conditions, we first sample a radius uniformly, then choose an angle  $\theta$  uniformly.  
 153 This produces a uniform distribution over spring system energy levels and starts at an arbitrary point in  
 154 the cycle. The spring system has a closed-form solution:  $(q(t), p(t)) = (r \sin(t + \theta_0), r \cos(t + \theta_0))$   
 155 where  $r$  is the radius of the circle traced in phase space (the energy of the spring) and  $\theta_0$  is the  
 156 phase space angle at which the oscillation will start. While this closed form solution is useful,  
 157 for consistency with our other systems, we generate snapshots of the spring system by numerical  
 158 integration. Simulations of the spring system always run through one period. For “in-distribution”

159 training values, the radius is selected in the range  $(0.2, 1)$  and “out-of-distribution” radii are chosen  
 160 from  $(1, 1.2)$ .

## 161 4.2 Wave

162 This benchmark system is similar to the one used in Peng and Mohseni [32]. Consider the wave  
 163 equation with speed  $c = 0.1$

$$\partial_{tt}u = c^2 \partial_{xx}u, \quad (2)$$

164 on a one-dimensional spatial domain  $[0, 1]$  with periodic boundary conditions. We represent this  
 165 second-order system as a first-order system and discretize in space to obtain

$$\begin{bmatrix} \dot{q}(t) \\ \dot{p}(t) \end{bmatrix} = \begin{bmatrix} 0 & I \\ c^2 D_{xx} & 0 \end{bmatrix} \begin{bmatrix} q(t) \\ p(t) \end{bmatrix}, \quad (3)$$

166 where  $D_{xx} \in \mathbb{R}^{n \times n}$  corresponds to the three-point central difference approximation of the spatial  
 167 derivative  $\partial_{xx}$  and the matrices  $I$  and  $0$  are the identity and zero matrix, respectively, of appropriate  
 168 size. We discretize in space with  $n = 125$  evenly spaced grid points and evolve the system following  
 169 the dynamics described above.

170 Initial conditions are sampled with an initial pulse in the  $q$  component centered at 0.5. All initial  
 171 conditions have zero momentum. The initial pulse is produced by a spline kernel as described in [32]:

$$s(x) = \frac{10}{p_w} \cdot |x - 0.5|, \quad h(s) = p_h \cdot \begin{cases} 1 - \frac{3}{2}s^2 + \frac{3}{4}s^3 & \text{if } 0 \leq s \leq 1 \\ \frac{1}{4}(2-s)^3 & \text{if } 1 < s \leq 2 \\ 0 & \text{else} \end{cases}$$

172 where the width and height of the pulse are scaled by parameters  $p_w$  and  $p_h$ , respectively. The spline  
 173 kernel pulse is then  $h(s(x))$  for  $x \in [0, 1]$ , evaluated at the discretized grid points.

174 For “in-distribution” samples parameters  $p_w, p_h$  are both chosen uniformly in the range  $(0.75, 1.25)$   
 175 and “out-of-distribution” runs sample uniformly from  $(0.5, 0.75) \cup (1.25, 1.5)$ . All trajectories are  
 176 integrated until  $t = 5$  when the wave has traveled through half a period.

## 177 4.3 Spring Mesh

178 This system manipulates a square grid of particles connected by springs, in a two dimensional space,  
 179 and can be considered a simplified version of deformable surface and volume systems (cf. [33]). The  
 180 particles all have mass 1, and are arranged into a unit grid. Springs are added along the axis-aligned  
 181 edges and diagonally across each grid square, with rest lengths selected so that the regularly-spaced  
 182 particles are in a rest position.

183 In this work we use a  $10 \times 10$  grid where the top row of particles is fixed in place. Initial conditions  
 184 are sampled by choosing a perturbation for the position of each non-fixed spring. These perturbations  
 185 are chosen as uniform vectors inside a circle with radius 0.35. Out-of-distribution perturbations  
 186 are chosen uniformly in a ring with inner radius 0.35 and outer radius 0.45. The sampled initial  
 187 conditions all have zero momentum.

188 In this system, a spring between particles  $a$  and  $b$  exerts a force:

$$F_{ab} = -k \cdot (\|q_a - q_b\|_2 - \ell_{ab}) \frac{q_a - q_b}{\|q_a - q_b\|_2} - \gamma \dot{q}_a \quad (4)$$

189 where  $\ell_{ab}$  is the rest length of the spring,  $\gamma = 0.1$  is a parameter controlling the magnitude of an  
 190 underdamped velocity-based decay, and  $k = 1$  is the spring constant.

## 191 4.4 Navier-Stokes

192 We consider the standard Navier-Stokes equation over a domain  $\Omega$  (cf. [33, 50])

$$\left. \begin{aligned} \rho \frac{\partial u}{\partial t} + \rho(u \cdot \nabla)u - \nu \Delta u + \nabla p &= b \\ \nabla \cdot u &= 0 \\ u(0) &= u_0 \end{aligned} \right\} \text{ on } \Omega \times (0, T), \quad \left. \begin{aligned} u &= d \\ \nu \frac{\partial u}{\partial n} + pn &= g \end{aligned} \right\} \text{ on } \partial\Omega_D \times (0, T),$$

Table 1: Dataset sizes and simulation parameters

System	# Train Trajectories	# Eval Trajectories	Time Step Size	# Steps
Spring	10, 500, 1000	30	0.00781, $\div 128$	805
Wave	10, 25, 50	6	0.00049, $\div 8$	10204
Spring Mesh	25, 50, 100	15	0.00781, $\div 128$	805
Navier-Stokes	25, 50, 100	5	0.08, $\div 1$	65

193 where  $u: \Omega \times (0, T) \rightarrow R^2$  is the velocity at time  $t \in (0, T)$  of a fluid with kinematic viscosity  
194  $\nu$  and density  $\rho$ ,  $p: \Omega \times (0, T) \rightarrow R$  is the pressure and  $\partial\Omega_D$  and  $\partial\Omega_N$  are the Dirichlet and  
195 Neumann boundary conditions respectively. In our setup we use the finite element method (FEM)  
196 to solve the PDE using mixed discretization: quadratic polynomial for the velocity and linear for  
197 pressure. In our experiment the domain  $\Omega$  is a rectangle  $0.22 \times 0.41$  with a random generated  
198 set of circular obstacles. We start with  $u_0 = 0$  and specify a velocity on the left boundary of  
199  $u(0, y) = (6(1 - e^{-5t})(0.41 - y)y/0.1681, 0)$ , zero on the top and bottom, and zero Neumann on  
200 the right side ( $g = 0$ ). We solve the system using PolyFEM [46] using  $dt = 0.08$  and backward  
201 differentiation formula (BSF) of order 3 for the time integration.

202 We sample obstacles into two configurations: a single obstacle, or a set of four. In each case, we  
203 sample the obstacles leaving a margin of 0.05 between each circle, and a margin of 0.25 from the  
204 left and right sides, and 0.05 between the top and bottom. Otherwise, each obstacle is determined  
205 by first sampling a radius, then sampling a center from the valid space, respecting the margins. If  
206 the sampled obstacle is too close to an existing circle, it is discarded and a new sample is drawn.  
207 In-distribution obstacles have radii in the range  $(0.05, 0.1)$  and out-of-distribution radii are drawn  
208 from  $(0.025, 0.05)$ .

## 209 5 Numerical experiments

210 **Experimental Setup.** We apply several basic learning methods to the datasets developed in this  
211 work:  $k$ -nearest neighbor regressors, a neural network kernel method, several sizes of feed forward  
212 MLPs, and a variety of CNNs. Details of the architectures and the training protocol are provided in  
213 supplementary material, Appendix B. Each of the neural networks we consider is implemented using  
214 PyTorch [30].

215 The learning methods considered in this work are trained on one of two target task formulations  
216 described in Section 3. For derivative-based prediction, the training is conducted supervised on ground  
217 truth snapshots gathered from the underlying models. For each system we randomly sample initial  
218 conditions and each of these is then numerically integrated to produce a trajectory. Each trajectory  
219 includes state samples  $x$  as well as target derivatives  $\dot{x}$  used for training. For direct prediction, we  
220 no longer require numerical integration; instead we directly predict the trajectory in a sequential  
221 fashion. In this setting, we approximate  $f_\theta(x(t)) \approx x(t + \delta_t)$  for a discrete time step size  $\delta_t$ . For the  
222 derivative prediction we report results using the Leap-Frog integrator, the full results are available in  
223 the additional material.

224 We pick the same set of learning methods and apply it to both tasks independently to judge perfor-  
225 mance in each. For many systems the state is divided into position and momentum components:  
226  $x \equiv (q, p)$ . For the Navier-Stokes problem, the state  $x$  is made up of the flow velocity field, and  
227 the scalar field for pressure. After training, we produce rolled-out trajectories from held-out initial  
228 conditions either by combining with a numerical integrator in the case of derivative prediction, or  
229 in a directly recurrent fashion in the case of step prediction. Each neural network is instantiated in  
230 three independent copies, each of which is trained and evaluated across all sampled trajectories. We  
231 compute a per-step MSE against a ground truth value, we average these per-step MSEs to produce a  
232 per-trajectory error, and record these errors for analysis. Our experiments are designed to test several  
233 aspects of physical simulation. We highlight the most salient ones next, and report more extensive  
234 results in Appendix C.

235 **Training set size.** In general ML problems, one would expect additional training samples to  
236 systematically improve (in-distribution) evaluation performance. However, Figure 2 illustrates a

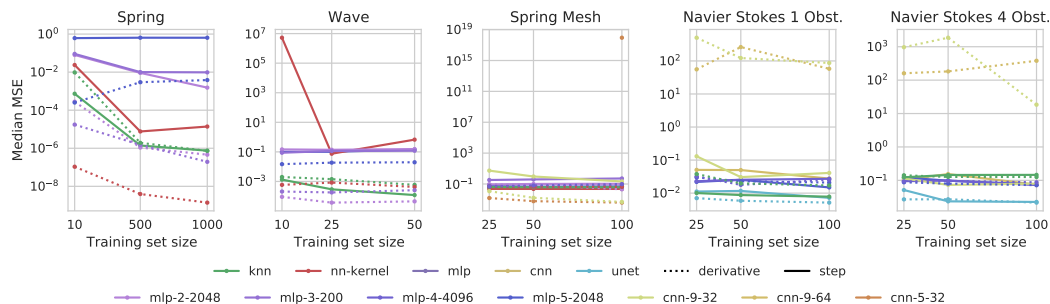


Figure 2: Median MSE error with respect to the training set size for our five different systems. For the same method we show different architecture choices.

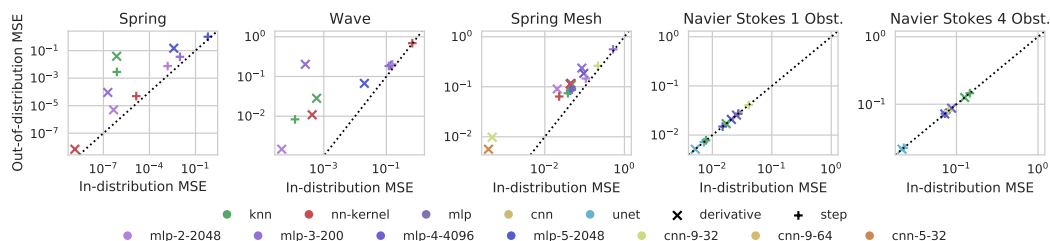


Figure 3: Median MSE for in-distribution evaluation sets vs. out-of-distribution evaluation sets for our five systems. Colors represent the same method and are repeated for different architecture choices. Marker shapes distinguish step and derivative prediction, and the dotted line is the identity line. Outliers for the spring-mesh and both Navier-Stokes systems were removed. Values on both axes were approximately  $10^{13}$  for the spring mesh and in the range  $10^2$ – $10^3$  for Navier-Stokes.

237 clear saturation of performance on the simplest systems when using neural networks as function  
 238 approximator, in contrast with non-parametric KNNs and the kernel method. We attribute this  
 239 saturation to an inherent gap between the training and evaluation objectives. While data-driven  
 240 methods are optimized to minimise next-step predictions, the final evaluation requires built-in  
 241 stability to prediction errors. Including regularisation strategies to incorporate stability, such as noise  
 242 injection [33], is shown to help, but not fully resolve this issue.

243 **Out-of-distribution evaluation.** For simplicity, we only examine the out-of-distribution error for  
 244 networks trained on the largest training set size. The added challenge of out-of-distribution samples  
 245 varies with the construction of each system. It is possible to get some idea of the difficulty increase  
 246 by examining the accuracy penalty for the KNNs, and comparing it to how well the more advanced  
 247 models are able to generalize.

248 Benefits of neural networks for generalization over KNN are visible across several systems in Figure 3,  
 249 particularly in the spring system for small MLPs for derivative prediction and nn-kernel in both  
 250 cases. The KNN suffers a significant increase in error while these methods produce only somewhat  
 251 worse predictions. Benefits are still present, though less pronounced, for the wave system derivative  
 252 prediction where neural networks increase in error, but the kernel method and small MLP maintain  
 253 a lower absolute error than the KNN. On the Navier-Stokes systems none of the methods suffer an  
 254 increase in error for out-of-distribution evaluation. The change in radius distribution for the obstacles  
 255 did not pose an additional challenge sufficient to produce a measurable change in error distribution.  
 256 We attribute this to low dimension of the initial condition space.

257 **Step and derivative prediction.** The step and derivative prediction instances of each learning  
 258 problem lead to different accuracy from the learning methods we test. While most physical systems  
 259 are naturally described in terms of their derivatives through corresponding ODE/PDEs, data-driven  
 260 simulations also offer the alternative of bypassing this differential formulation and predict the next  
 261 state directly. Such ‘cavalier’ approach avoids the compounding error amplification effects across

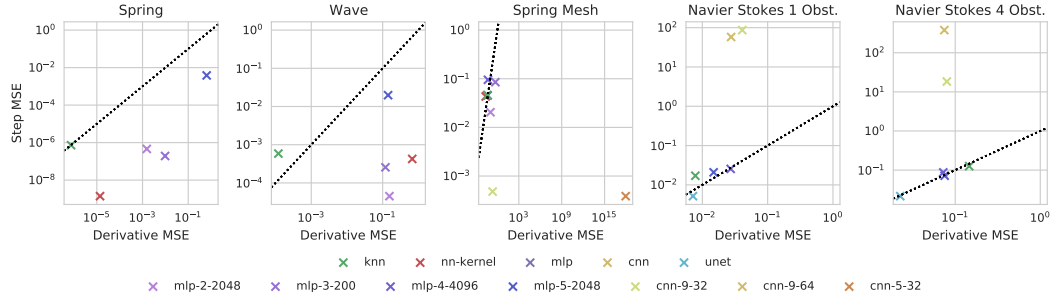


Figure 4: Median MSEs for derivative vs. step prediction on the same evaluation set. Results are displayed for each of our five systems. The dotted line is the identity line.

262 integration steps, at the expense of sample efficiency. Figure 4 illustrates these tradeoffs across our  
 263 systems.

264 An important example of this effect is the performance of CNNs on the spring mesh system (Figure 8  
 265 in the Appendix). When working through a numerical integrator and performing derivative prediction  
 266 they produce the lowest error of all methods tested, but following the same training protocol for step  
 267 prediction these architectures produce high errors, or are unstable. This case is likely an interaction  
 268 of the architecture with the specific learning task. For the spring mesh, step prediction requires  
 269 outputting the position of the particle which requires manipulating its global coordinates, while  
 270 derivative predictions allow the network to more easily act locally and compute only a relative motion  
 271 for the particle. The derivative prediction task better takes advantage of the spatial invariance of the  
 272 CNNs. This difference in performance reflects the importance of tailoring architectures to the specific  
 273 task, and some potential for neural network architectures to benefit from incorporating knowledge of  
 274 a system’s behavior.

275 **System and dataset complexity.** Several trends we observe correlate with the difficulty of learning  
 276 to simulate the system, and the variation in its behavior across the training and evaluation samples.  
 277 This is generally a combination of the system’s state dimension, and variation in its behavior,  
 278 approximated by the dimension of the distribution from which initial conditions are sampled.

279 This is particularly visible in Figure 5 in the performance of the KNN methods, and, in many cases,  
 280 the performance of simpler methods such as the small MLPs. On the simpler systems, such as the  
 281 spring and wave, the KNNs generally perform well because even though the wave system has a  
 282 relatively large state dimension of 125, like the spring its initial condition is sampled from only two  
 283 parameters and its behavior can be readily predicted from these. The Navier-Stokes system with  
 284 a single obstacle is another instance of this sort of behavior: the KNN is readily able to reproduce  
 285 flows it has not seen because a sampling of 100 obstacle positions is such that an evaluation sample is  
 286 close to a trajectory seen at training time. Therefore, small MLPs and the kernel method produce  
 287 similar performance. When the difficulty is increased by sampling four obstacles, the KNN and MLP  
 288 performances suffer, and larger networks such as the u-net are needed to maintain approximately the  
 289 same performance.

290 **Choice of numerical integrator.** For our derivative prediction tasks we combine our trained  
 291 methods with three explicit integrators with orders 1, 2, and 4. In most of our systems these produce  
 292 at most a small increase in accuracy, holding all other training and evaluation parameters equal.  
 293 However on the Navier-Stokes system the higher order integrators produce somewhat higher errors,  
 294 particularly for the u-net and the MLPs. This appears related to the approximated derivatives used for  
 295 training this system. The learned derivatives produce some small deviations which are compounded  
 296 when combining multiple derivative samples.

297 **Computational overheads.** Another important aspect to consider when applying learning methods  
 298 to physical simulation problems is the time required to compute each step, and the computational  
 299 overheads introduced by the lack of knowledge of the true system physics. With standard numerical  
 300 integration methods, it is generally possible to improve the quality of generated trajectories by  
 301 decreasing the size of the time step used during integration. We take advantage of this in order to



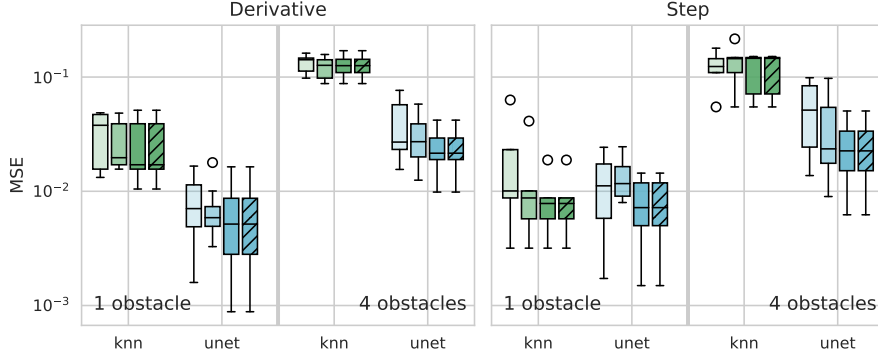


Figure 5: MSE distribution across trajectories in the evaluation set for KNNs and u-nets on our Navier-Stokes system, both the formulation with one obstacle, and the settings with four, as well as for derivative and step prediction. Box plot mid-lines show medians, box area is between first and third quantiles, whiskers extend  $1.5\times$  beyond the boxes, outliers are plotted past the whiskers. Darker colors denote increasing training set size. The final hatched box is the same network from the final un-hatched box, tested on an out-of-distribution evaluation set.

Table 2: Time comparison for derivative prediction against baseline numerical integrators

System	Architecture	Euler		Leapfrog		RK4	
		Time Ratio	Scaling	Time Ratio	Scaling	Time Ratio	Scaling
Spring	knn	367.0	$1\times$	405.8	$16\times$	311.3	$64\times$
	nn kernel	180.7	$1\times$	198.6	$1\times$	173.3	$1\times$
	mlp-2-2048	185.8	$1\times$	191.6	$16\times$	177.7	$64\times$
	mlp-3-200	237.6	$1\times$	237.5	$16\times$	227.2	$64\times$
	mlp-5-2048	473.8	$4\times$	369.6	$64\times$	360.9	$128\times$
Wave	knn	24,102.7	$8\times$	16,945.1	$256\times$	16,132.0	$256\times$
	nn kernel	35.3	$8\times$	22.3	$256\times$	19.9	$256\times$
	mlp-2-2048	25.4	$8\times$	16.5	$256\times$	14.2	$256\times$
	mlp-3-200	31.0	$16\times$	19.7	$256\times$	17.8	$256\times$
	mlp-5-2048	60.5	$16\times$	38.0	$256\times$	34.6	$256\times$
Spring Mesh	knn	708.6	$8\times$	626.1	$128\times$	690.4	$256\times$
	nn kernel	5.3	$8\times$	4.4	$128\times$	4.8	$256\times$
	mlp-2-2048	3.0	$8\times$	2.6	$128\times$	2.8	$256\times$
	mlp-3-200	3.4	$8\times$	3.2	$128\times$	3.4	$256\times$
	mlp-4-4096	7.8	$16\times$	7.1	$128\times$	7.7	$256\times$
	mlp-5-2048	7.1	$8\times$	6.2	$128\times$	5.3	$256\times$
	cnn-9-32	11.0	$2\times$	10.5	$32\times$	11.3	$64\times$
cnn-5-32	7.4	$1\times$	9.2	$32\times$	7.3	$64\times$	

302 estimate the time overheads of our learning methods relative to our baseline numerical integrators at  
 303 approximately corresponding error levels.

304 We numerically integrate each system at time step sizes scaled by powers of two. For each trajectory in  
 305 the derivative prediction setting, we find the smallest scaling factor at which the numerical integrator  
 306 exceeds the learning method’s error at their final shared time step, approximating the factor by which  
 307 numerical integration can be made faster until it begins to underperform the learned method.

308 Table 2 reports the results of these experiments. For each numerical integrator, the “scaling” column  
 309 reports the most common scaling factor found for each trajectory. The “time ratio” column represents  
 310 the learned method’s evaluation overhead (median times, counting only per-step network evaluation  
 311 costs, not numerical integration or data transfers). Note that the numerical integrator makes fewer  
 312 steps than the learned method so the overall trajectory time must be further adjusted by the scaling  
 313 factor.

314 In general, the neural networks are slower per-step by one or two orders of magnitude. KNNs are  
 315 slower by significantly larger factors, particularly for the wave system. This is likely partially due

316 to the default scikit-learn KNN implementation used, and due to the large size of the wave system  
317 training sets (large state dimension and large number of training snapshots). Scaling factors increase  
318 with the order of the integrator as higher-order integrators are more tolerant of large step sizes and  
319 maintain low error.

320 It is likely that these overheads could be reduced with more optimized implementations of both the  
321 numerical integrators and learned methods. The derivative prediction task is also constrained by its  
322 need to interact with the numerical integrator. In this setting the learned methods cannot be expected  
323 to outperform the quality of the solutions generated by the true system derivatives. This reflects a  
324 penalty resulting from a lack of knowledge of the true underlying system, and a penalty for learning  
325 from observations in this case. Step prediction without involving the numerical integrator potentially  
326 avoids some of these constraints, if learning is successful.

## 327 **6 Conclusions and Limitations**

328 The results in this work illustrate the performance achievable by applying common machine learning  
329 methods to the simulation problems in our proposed benchmark task. We envision three ways  
330 in which the results of this work might be used: (1) the datasets developed here can be used for  
331 training and evaluating new machine learning techniques in this area, (2) the simulation software  
332 can be used to generate new datasets from these systems of different sizes, different initial condition  
333 dimensionality and distribution; the training software could be used to assist in conducting further  
334 experiments, and (3) some of the trends seen in our results may help inform the design of future  
335 machine learning tasks for simulation.

336 For the first and second groups of downstream users, we have made available the pre-generated  
337 datasets used in this work, as well as the software used to produce them and carry out our experiments.  
338 These components allow carrying out the measurements we have made here, and permit further  
339 adjustments to be made. Documentation on using the datasets and the software is included in  
340 Appendix D.

341 For the third group, we highlighted a few trends that suggest useful steps to take in developing new  
342 problems and datasets in this area. First, we advise including several simple baseline methods when  
343 designing new tasks. In particular the inclusion of standard numerical integrators (for derivative-type  
344 problems) and KNNs are useful to evaluate the difficulty of the proposed task. Specifically, KNNs  
345 are useful for examining the performance achievable by memorizing the training set, and are thus  
346 *witnessing* an appropriate design of data distribution that captures the true high-dimensionality of  
347 the prediction task. As an example, in the Navier-Stokes examples some task formulations may  
348 inadvertently be simple to memorize, even if the complexity of the system itself may not immediately  
349 suggest it. The numerical integrators are likewise useful as baselines both to ensure that the derivative-  
350 learning is feasible even when achieving no error in predictions, and also to evaluate the penalty  
351 in accuracy which is incurred by operating without access to the true physics. We believe that in  
352 light of these observed trends, including baseline methods such as standard numerical integration  
353 schemes and simple learning methods such as the KNN is important in understanding tasks in this  
354 area. Including these assists in experiment design by helping to calibrate the difficulty of a target task.

355 **Limitations:** While our benchmark provides actionable conclusions on a wide array of simulation  
356 domains, it is currently focused at temporal integration, and as such it does not cover important settings  
357 in Scientific Computing. For instance, we do not currently include an instance of a surrogate model,  
358 which could provide different tradeoffs benefiting ML models. Additionally, we have focused on two  
359 setups for data-driven simulation (differential snapshot prediction and direct snapshot prediction), but  
360 other alternatives exist that might mitigate some of the shortcomings we observed; for instance by  
361 considering larger temporal contexts (as in [4]), as well as enforcing certain conservation laws into  
362 the model [8, 4]. Finally, we do not report a full timing analysis of data-driven models against the  
363 baseline integrators, although our preliminary analysis strongly suggests a substantial computational  
364 gap in favor of baseline integrators using default software implementations.

## 365 **References**

366 [1] A. C. Antoulas and B. D. O. Anderson. On the scalar rational interpolation problem. *IMA Journal of*  
367 *Mathematical Control & Information*, 3(2-3):61–88, 1986.

- 368 [2] A. C. Antoulas, C. Beattie, and S. Gugercin. *Interpolatory Methods for Model Reduction*. SIAM, 2020.
- 369 [3] S. L. Brunton, J. L. Proctor, and J. N. Kutz. Discovering governing equations from data by sparse  
370 identification of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 113(15):  
371 3932–3937, 2016.
- 372 [4] Z. Chen, J. Zhang, M. Arjovsky, and L. Bottou. Symplectic recurrent neural networks. In *International  
373 Conference on Learning Representations*, 2020.
- 374 [5] P. V. Coveney, E. R. Dougherty, and R. R. Highfield. Big data need big theory too. *Philosophical  
375 Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2080):  
376 20160153, 2016. doi: 10.1098/rsta.2016.0153. URL [https://royalsocietypublishing.org/doi/  
377 abs/10.1098/rsta.2016.0153](https://royalsocietypublishing.org/doi/abs/10.1098/rsta.2016.0153).
- 378 [6] M. Cranmer, A. Sanchez-Gonzalez, P. Battaglia, R. Xu, K. Cranmer, D. Spergel, and S. Ho. Discovering  
379 symbolic models from deep learning with inductive biases. *arXiv preprint*, 2020. URL [https://arxiv.  
380 org/abs/2006.11287](https://arxiv.org/abs/2006.11287).
- 381 [7] J. P. Crutchfield and B. S. Mcnamara. Equations of motion from a data series. *Complex Systems*, page 452,  
382 1987.
- 383 [8] S. Greydanus, M. Dzamba, and J. Yosinski. Hamiltonian neural networks. In H. Wal-  
384 lach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors,  
385 *Advances in Neural Information Processing Systems*, volume 32, pages 15379–15389. Cur-  
386 ran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/file/  
387 26cd8ecadce0d4efd6cc8a8725cbd1f8-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/26cd8ecadce0d4efd6cc8a8725cbd1f8-Paper.pdf).
- 388 [9] B. Gustavsen and A. Semlyen. Rational approximation of frequency domain responses by vector fitting.  
389 *Power Delivery, IEEE Transactions on*, 14(3):1052–1061, Jul 1999.
- 390 [10] E. Haghghat, M. Raissi, A. Moure, H. Gomez, and R. Juanes. A deep learning framework for solution and  
391 discovery in solid mechanics, 2020.
- 392 [11] E. Hairer and G. Wanner. *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic  
393 Problems*. Springer, 2009.
- 394 [12] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving Ordinary Differential Equations I: Nonstiff Problems*.  
395 Springer, 2009.
- 396 [13] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser,  
397 J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F.  
398 del R’io, M. Wiebe, P. Peterson, P. G’erard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi,  
399 C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, 2020. URL  
400 <https://doi.org/10.1038/s41586-020-2649-2>.
- 401 [14] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95,  
402 2007. doi: 10.1109/MCSE.2007.55.
- 403 [15] A. Khadilkar, J. Wang, and R. Rai. Deep learning–based stress prediction for bottom-up sla 3d print-  
404 ing process. *The International Journal of Advanced Manufacturing Technology*, 102(5):2555–2569,  
405 Jun 2019. ISSN 1433-3015. doi: 10.1007/s00170-019-03363-4. URL [https://doi.org/10.1007/  
s00170-019-03363-4](https://doi.org/10.1007/<br/>406 s00170-019-03363-4).
- 407 [16] B. Kim, V. C. Azevedo, N. Thuerey, T. Kim, M. Gross, and B. Solenthaler. Deep fluids: A generative  
408 network for parameterized fluid simulations. In *Computer Graphics Forum*, volume 38, pages 59–70.  
409 Wiley Online Library, 2019.
- 410 [17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
411 2014.
- 412 [18] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. Zemel. Neural relational inference for interacting  
413 systems. In *International Conference on Machine Learning*, pages 2688–2697. PMLR, 2018.
- 414 [19] G. Kissas, Y. Yang, E. Hwuang, W. R. Witschey, J. A. Detre, and P. Perdikaris. Machine learning in  
415 cardiovascular flows modeling: Predicting arterial blood pressure from non-invasive 4d flow mri data  
416 using physics-informed neural networks. *Computer Methods in Applied Mechanics and Engineering*, 358:  
417 112623, 2020.

- 418 [20] S. K. Lam, A. Pitrou, and S. Seibert. Numba: A llvm-based python jit compiler. In *Proceedings of the*  
419 *Second Workshop on the LLVM Compiler Infrastructure in HPC*, LLVM '15, New York, NY, USA, 2015.  
420 Association for Computing Machinery. ISBN 9781450340052. doi: 10.1145/2833157.2833162. URL  
421 <https://doi.org/10.1145/2833157.2833162>.
- 422 [21] Y. Li, H. Wang, K. Mo, and T. Zeng. Reconstruction of simulation-based physical field by reconstruction  
423 neural network method. *arXiv preprint arXiv:1805.00528*, 2018.
- 424 [22] Y. Lia, H. Wanga, W. Shuaia, H. Zhangb, and Y. Pengb. Image-based reconstruction for the impact  
425 problems by using dpnns.
- 426 [23] L. Liang, M. Liu, C. Martin, and W. Sun. A deep learning approach to estimate stress distribution: a  
427 fast and accurate surrogate of finite-element analysis. *Journal of The Royal Society Interface*, 15(138):  
428 20170844, 2018. doi: 10.1098/rsif.2017.0844. URL [https://royalsocietypublishing.org/doi/](https://royalsocietypublishing.org/doi/abs/10.1098/rsif.2017.0844)  
429 [abs/10.1098/rsif.2017.0844](https://royalsocietypublishing.org/doi/abs/10.1098/rsif.2017.0844).
- 430 [24] L. Liang, W. Mao, and W. Sun. A feasibility study of deep learning for predicting hemodynamics of human  
431 thoracic aorta. *Journal of Biomechanics*, 99:109544, 2020.
- 432 [25] L. Lu, X. Meng, Z. Mao, and G. E. Karniadakis. Deepxde: A deep learning library for solving differential  
433 equations. *arXiv preprint arXiv:1907.04502*, 2019.
- 434 [26] G. D. Maso Talou, T. P. Babarenda Gamage, M. Sagar, and M. P. Nash. Deep learning over reduced  
435 intrinsic domains for efficient mechanics of the left ventricle. *Frontiers in Physics*, 8:30, 2020.
- 436 [27] Y. Nakatsukasa, O. Sète, and L. N. Trefethen. The aaa algorithm for rational approximation. *SIAM Journal*  
437 *on Scientific Computing*, 40(3):A1494–A1522, 2018.
- 438 [28] Z. Nie, H. Jiang, and L. B. Kara. Stress field prediction in cantilevered structures using convolutional  
439 neural networks. *Journal of Computing and Information Science in Engineering*, 20(1), 2020.
- 440 [29] N. H. Packard, J. P. Crutchfield, J. D. Farmer, and R. S. Shaw. Geometry from a time series. *Physical*  
441 *Review Letters*, 45:712–716, 1980.
- 442 [30] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein,  
443 L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner,  
444 L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In  
445 H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances in*  
446 *Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- 447 [31] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,  
448 R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay.  
449 Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- 450 [32] L. Peng and K. Mohseni. Symplectic model reduction of hamiltonian systems. *SIAM Journal on Scientific*  
451 *Computing*, 38(1):A1–A27, 2016. URL <https://doi.org/10.1137/140978922>.
- 452 [33] T. Pfaff, M. Fortunato, A. Sanchez-Gonzalez, and P. W. Battaglia. Learning mesh-based simulation with  
453 graph networks. *arXiv preprint*, 2020. URL <https://arxiv.org/abs/2010.03409>.
- 454 [34] E. Qian, B. Kramer, B. Peherstorfer, and K. Willcox. Lift & learn: Physics-informed machine learning for  
455 large-scale nonlinear dynamical systems. *Physica D: Nonlinear Phenomena*, 2020.
- 456 [35] A. Rahimi, B. Recht, et al. Random features for large-scale kernel machines. In *NIPS*, volume 3, page 5.  
457 Citeseer, 2007.
- 458 [36] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics informed deep learning (part ii): Data-driven  
459 discovery of nonlinear partial differential equations. *arxiv. arXiv preprint arXiv:1711.10561*, 2017.
- 460 [37] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning  
461 framework for solving forward and inverse problems involving nonlinear partial differential equations.  
462 *Journal of Computational Physics*, 378:686–707, 2019.
- 463 [38] M. D. Ribeiro, A. Rehman, S. Ahmed, and A. Dengel. Deepcfd: Efficient steady-state laminar flow  
464 approximation with deep convolutional neural networks. *arXiv preprint arXiv:2004.08826*, 2020.
- 465 [39] A. Rudi, L. Carratino, and L. Rosasco. Falkon: An optimal large scale kernel method. *arXiv preprint*  
466 *arXiv:1705.10958*, 2017.

- 467 [40] A. Sanchez-Gonzalez, N. Heess, J. T. Springenberg, J. Merel, M. Riedmiller, R. Hadsell, and P. Battaglia.  
468 Graph networks as learnable physics engines for inference and control. In *International Conference on*  
469 *Machine Learning*, pages 4470–4479. PMLR, 2018.
- 470 [41] A. Sanchez-Gonzalez, V. Bapst, K. Cranmer, and P. Battaglia. Hamiltonian graph networks with ODE  
471 integrators. *arXiv preprint*, 2019. URL <https://arxiv.org/abs/1909.12790>.
- 472 [42] H. Schaeffer. Learning partial differential equations via data discovery and sparse optimization. *Proceedings*  
473 *of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 473(2197):20160446, 2017.
- 474 [43] H. Schaeffer, G. Tran, and R. Ward. Extracting sparse high-dimensional dynamics from limited data. *SIAM*  
475 *Journal on Applied Mathematics*, 78(6):3279–3295, 2018.
- 476 [44] P. Schmid. Dynamic mode decomposition of numerical and experimental data. *Journal of Fluid Mechanics*,  
477 656:5–28, 8 2010. ISSN 1469-7645.
- 478 [45] P. Schmid and S. J. Dynamic mode decomposition of numerical and experimental data. In *Bull. Amer.*  
479 *Phys. Soc., 61st APS meeting*, page 208. American Physical Society, 2008.
- 480 [46] T. Schneider, J. Dumas, X. Gao, D. Zorin, and D. Panozzo. Polyfem. <https://polyfem.github.io/>,  
481 2019.
- 482 [47] B. Schölkopf and A. J. Smola. *Learning with Kernels*. MIT Press, 1998.
- 483 [48] E. Süli and D. F. Mayers. *An Introduction to Numerical Analysis*, chapter Initial Value Problems for ODEs,  
484 pages 349–353. Cambridge University Press, 2003.
- 485 [49] A. M. Tartakovsky, C. O. Marrero, P. Perdikaris, G. D. Tartakovsky, and D. Barajas-Solano. Learning  
486 parameters and constitutive relationships with physics informed deep neural networks. *arXiv preprint*  
487 *arXiv:1808.03398*, 2018.
- 488 [50] N. Thuerey, K. Weißenow, L. Prantl, and X. Hu. Deep learning methods for reynolds-averaged navier-  
489 stokes simulations of airfoil flows. *AIAA Journal*, 58(1):25–36, 2020.
- 490 [51] J. H. Tu, C. W. Rowley, D. M. Luchtenburg, S. L. Brunton, and J. N. Kutz. On dynamic mode decomposition:  
491 Theory and applications. *Journal of Computational Dynamics*, 1(2):391–421, 2014.
- 492 [52] N. Umetani and B. Bickel. Learning three-dimensional flow for interactive aerodynamic design. *ACM*  
493 *Transactions on Graphics (TOG)*, 37(4):1–10, 2018.
- 494 [53] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson,  
495 W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J.  
496 Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde,  
497 J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro,  
498 F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for  
499 Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. URL <https://doi.org/10.1038/s41592-019-0686-2>.
- 501 [54] K. E. Willcox, O. Ghattas, and P. Heimbach. The imperative of physics-based modeling and inverse theory  
502 in computational science. *Nature Computational Science*, 1(3):166–168, Mar 2021. ISSN 2662-8457. doi:  
503 10.1038/s43588-021-00040-z. URL <https://doi.org/10.1038/s43588-021-00040-z>.
- 504 [55] Y. Xie, E. Franz, M. Chu, and N. Thuerey. TempoGAN: A temporally coherent, volumetric GAN for  
505 super-resolution fluid flow. *ACM Transactions on Graphics (TOG)*, 37(4):1–15, 2018.

## 506 Checklist

- 507 1. For all authors...
- 508 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
509 contributions and scope? [Yes]
- 510 (b) Did you describe the limitations of your work? [Yes] Yes, our conclusion in Section 6  
511 includes a discussion of the limitations of our work
- 512 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We do  
513 not expect any potential negative societal impacts.

- 514 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
515 them? [Yes] We have reviewed the ethics review guidelines
- 516 2. If you are including theoretical results...
- 517 (a) Did you state the full set of assumptions of all theoretical results? [N/A] No theoretical  
518 results
- 519 (b) Did you include complete proofs of all theoretical results? [N/A] No theoretical results
- 520 3. If you ran experiments (e.g. for benchmarks)...
- 521 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
522 mental results (either in the supplemental material or as a URL)? [Yes]
- 523 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
524 were chosen)? [Yes] We provide details of our experimental setup both in the main  
525 paper and in supplementary materials. We also make available the scripts used to run  
526 the tests reported here, which include settings for all parameters as well.
- 527 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
528 ments multiple times)? [Yes] We discuss trends in our experimental results along  
529 several different dimensions. We also provide plots illustrating the error distribution  
530 for our methods tested on several randomly-sampled initial conditions.
- 531 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
532 of GPUs, internal cluster, or cloud provider)? [Yes] We discuss our use of NYU's HPC  
533 system for training and evaluation and summarize the compute time required. These  
534 statistics are included in the supplementary materials appendix B.5.
- 535 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 536 (a) If your work uses existing assets, did you cite the creators? [Yes] We have included  
537 references to the core software libraries used
- 538 (b) Did you mention the license of the assets? [Yes] We document how to make use of our  
539 datasets and run our software and make clean their respective licenses
- 540 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
541 We develop a set of benchmark tasks, generated by configurable simulation software.  
542 We have included links to download the datasets and the software.
- 543 (d) Did you discuss whether and how consent was obtained from people whose data you're  
544 using/curating? [N/A] No human data used or curated
- 545 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
546 information or offensive content? [N/A] Our data is synthetic, the result of numerical  
547 simulation
- 548 5. If you used crowdsourcing or conducted research with human subjects...
- 549 (a) Did you include the full text of instructions given to participants and screenshots, if  
550 applicable? [N/A] No human participants
- 551 (b) Did you describe any potential participant risks, with links to Institutional Review  
552 Board (IRB) approvals, if applicable? [N/A] No human participants
- 553 (c) Did you include the estimated hourly wage paid to participants and the total amount  
554 spent on participant compensation? [N/A] No human participants