
Can Neural Networks Improve Classical Optimization of Inverse Problems?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Finding the values of model parameters from data is an essential task in science.
2 While iterative optimization algorithms like BFGS can find solutions to inverse
3 problems with machine precision for simple problems, their reliance on local in-
4 formation limits their effectiveness for complex problems involving local minima,
5 chaos, or zero-gradient regions. This study explores the potential for overcoming
6 these limitations by jointly optimizing multiple examples. To achieve this, we
7 employ neural networks to reparameterize the solution space and leverage the
8 training procedure as an alternative to classical optimization. This approach is as
9 versatile as traditional optimizers and does not require additional information about
10 the inverse problems, meaning it can be added to existing general-purpose opti-
11 mization libraries. We evaluate the effectiveness of this approach by comparing it
12 to traditional optimization on various inverse problems involving complex physical
13 systems, such as the incompressible Navier-Stokes equations. Our findings reveal
14 significant improvements in the accuracy of the obtained solutions.

15 1 Introduction

16 Estimating model parameters by solving inverse problems [Tar05] is a central task in scientific
17 research, from detecting gravitational waves [GH18] to controlling plasma flows [MLA⁺19] to
18 searching for neutrinoless double-beta decay [AAA⁺13, AAA⁺18]. Iterative optimization algorithms,
19 such as limited-memory BFGS [LN89] or Gauss-Newton [GM78], are often employed for solving
20 unconstrained parameter estimation problems [PTVF07]. These algorithms offer advantages such
21 as ease of use, broad applicability, quick convergence, and high accuracy, typically limited only
22 by noise in the observations and floating point precision. However, they face several fundamental
23 problems that are rooted in the fact that these algorithms rely on local information, i.e., objective
24 values $L(x_k)$ and derivatives close to the current solution estimate x_k , such as the gradient $\partial L/\partial x|_{x_k}$
25 and the Hessian matrix $\partial^2 L/\partial x^2|_{x_k}$. Acquiring non-local information can be done in low-dimensional
26 solution spaces, but the curse of dimensionality prevents this approach for high-dimensional problems.
27 These limitations lead to poor performance or failure in various problem settings:

- 28 • *Local optima* attract the optimizer in the absence of a counter-acting force. Although using a
29 large step size or adding momentum to the optimizer can help to traverse small local minima,
30 local optimizers are fundamentally unable to avoid this issue.
- 31 • *Flat regions* can cause optimizers to become trapped along one or multiple directions.
32 Higher-order solvers can overcome this issue when the Hessian only vanishes proportionally
33 with the gradient, but all local optimizers struggle in zero-gradient regions.
- 34 • *Chaotic regions*, characterized by rapidly changing gradients, are extremely hard to optimize.
35 Iterative optimizers typically decrease their step size to compensate, which prevents the
36 optimization from progressing on larger scales.

37 In many practical cases, a *set* of observations is available, comprising many individual parameter
38 estimation problems, e.g., when repeating experiments multiple times or collecting data over a time
39 frame [CCC⁺19, DJO⁺18, GH18, AAA⁺13, MAL13] and, even in the absence of many recorded
40 samples, synthetic data can be generated to supplement the data set. Given such a set of inverse
41 problems, we pose the question: *Can we find better solutions x_i to general inverse problems by
42 optimizing them jointly instead of individually, without requiring additional information about the
43 problems?*

44 To answer this question, we employ neural networks to formulate a joint optimization problem.
45 Neural networks as general function approximators are a natural and straightforward way to enable
46 joint optimization of multiple a priori independent examples. They have been extensively used in
47 the field of machine learning [GBCB16], and a large number of network architectures have been
48 developed, from multilayer perceptrons (MLPs) [Hay94] to convolutional networks (CNNs) [KSH12]
49 to transformers [VSP⁺17]. Overparameterized neural network architectures typically smoothly
50 interpolate the training data [BHM18, BPL21], allowing them to generalize, i.e., make predictions
51 about data the network was not trained on.

52 It has recently been shown that this generalization capability or *inductive bias* benefits the optimization
53 of individual problems with grid-like solution spaces by implicitly adding a prior to the optimization
54 based on the network architecture [UVL18, HSDG19]. However, these effects have yet to be
55 investigated for general inverse problems or in the context of joint optimization. We propose using
56 the training process of a neural network as a drop-in component for traditional optimizers like BFGS
57 without requiring additional data, configuration, or tuning. Instead of making predictions about new
58 data after training, our objective is to solve only the problems that are part of the training set, i.e.,
59 the training itself produces the solutions to the inverse problems, and the network is never used for
60 inference. These solutions can also be combined with an iterative optimizer to improve accuracy.
61 Unlike related machine learning applications [KAT⁺19, SGGP⁺20, SFK⁺21, RT21, SHT22, HKT21,
62 SF18, RPM20, ALGS⁺22], where a significant goal is accelerating time-intensive computations, we
63 accept a higher computational demand if the resulting solutions are more accurate.

64 To quantify the gains in accuracy that can be obtained, we compare this approach to classical
65 optimization as well as related techniques on four experiments involving difficult inverse problems:
66 (i) a curve fit with many local minima, (ii) a billiards-inspired rigid body simulation featuring zero-
67 gradient areas, (iii) a chaotic system governed by the Kuramoto–Sivashinsky equation and (iv) an
68 incompressible fluid system that is only partially observable. We compare joint optimization to direct
69 iterative methods and related techniques in each experiment.

70 2 Related work

71 Neural networks have become popular tools to model physical processes, either completely replac-
72 ing physics solvers [KAT⁺19, SGGP⁺20, SFK⁺21, RT21] or improving them [TSSP17, UBF⁺20,
73 KSA⁺21]. This can improve performance since network evaluations and solvers may be run at lower
74 resolution while maintaining stability and accuracy. Additionally, it automatically yields a different-
75 iable forward process which can then be used to solve inverse problems [SF18, RPM20, ALGS⁺22],
76 similar to how style transfer optimizes images [GEB16].

77 Alternatively, neural networks can be used as regularizers to solve inverse problems on sparse tomog-
78 raphy data [LSAH20] or employed recurrently for image denoising and super-resolution [PW17].
79 Recent works have also explored them for predicting solutions to inverse problems [HKT21, SHT22].
80 In these settings, neural networks are trained offline and then used to infer solutions to new inverse
81 problems, eliminating the iterative optimization process at test time.

82 Underlying many of these approaches are differentiable simulations required to obtain gradients of
83 the inverse problem. These can be used in iterative optimization or to train neural networks. Many
84 recent software packages have demonstrated this use of differentiable simulations, with general
85 frameworks [HAL⁺20, SC19, HKT20] and specialized simulators [TLQL21, LLK19].

86 Physics-informed neural networks [RPK19] encode solutions to optimization problems in the network
87 weights themselves. They model a continuous solution to an ODE or PDE and are trained by
88 formulating a loss function based on the differential equation, and have been explored for a variety

89 of directions [YZWX19, LPY+21, KGZ+21]. However, as these approaches rely on loss terms
 90 formulated with neural network derivatives, they do not apply to general inverse problems.

91 3 Reparameterizing inverse problems with neural networks

92 We consider a set of n similar inverse problems where we take *similar* to mean we can express all of
 93 them using a function $F(\xi_i | x_i)$ conditioned on a problem-specific vector x_i with $i = 1, \dots, n$. Each
 94 inverse problem then consists of finding optimal parameters ξ_i^* such that a desired or observed output
 95 y_i is reproduced, i.e.

$$\xi_i^* = \arg \min_{\xi_i} \mathcal{L}(F(\xi_i | x_i), y_i), \quad (1)$$

96 where \mathcal{L} denotes an error measure, such as the squared L^2 norm $\|\cdot\|_2^2$. We assume that F is
 97 differentiable and can be approximately simulated, i.e., the observed output y_i may not be reproducible
 98 exactly using F due to hidden information or stochasticity.

99 A common approach to finding ξ_i^* is performing a nonlinear optimization, minimizing \mathcal{L} using
 100 the gradients $\frac{\partial \mathcal{L}}{\partial F} \frac{\partial F}{\partial \xi_i}$. In strictly convex optimization, many optimizers guarantee convergence to
 101 the global optimum in these circumstances. However, when considering more complex problems,
 102 generic optimizers often fail to find the global optimum due to local optima, flat regions, or chaotic
 103 regions. Trust region methods [Yua00] can be used on low-dimensional problems but scale poorly to
 104 higher-dimensional problems. Without further domain-specific knowledge, these methods are limited
 105 to individually optimizing all n inverse problems.

106 Instead of improving the optimizer itself, we want to investigate whether better solutions can be found
 107 by jointly optimizing all problems. However, without domain-specific knowledge, it is unknown
 108 which parameters of ξ_i are shared among multiple problems. We therefore first reparameterize the
 109 full solution vectors ξ_i using a set of functions $\hat{\xi}_i$, setting $\xi_i \equiv \hat{\xi}_i(\theta)$ where θ represents a set of
 110 shared parameters. With this change, the original parameters ξ_i become functions of θ , allowing θ
 111 to be jointly optimized over all problems. Here, the different $\hat{\xi}_i$ can be considered transformation
 112 functions mapping θ to the actual solutions ξ_i , similar to transforming Cartesian to polar coordinates.
 113 Second, we sum the errors of all examples to define the overall objective function $L = \sum_{i=1}^n \mathcal{L}_i$.

114 For generality, all $\hat{\xi}_i(\theta)$ should be able to approximate arbitrary functions. We implement them as an
 115 artificial neural network \mathcal{N} with weights θ : $\hat{\xi}_i(\theta) \equiv \mathcal{N}(x_i, y_i | \theta)$. Inserting these changes into Eq. 1
 116 yields the reparameterized optimization problem

$$\xi_i^* = \hat{\xi}_i(\theta^*), \quad \theta^* = \operatorname{argmin}_{\theta} \sum_{i=1}^n \mathcal{L}(F(\mathcal{N}(x_i, y_i | \theta) | x_i), y_i). \quad (2)$$

117 We see that the joint optimization with reparameterization
 118 strongly resembles standard formulations of neural network
 119 training where (x_i, y_i) is the input to the network and $F \circ$
 120 L represents the effective loss function. However, from
 121 the viewpoint of optimizing inverse problems, the network
 122 is not primarily a function of (x_i, y_i) but rather a set of
 123 transformation functions of θ , each corresponding to a fixed
 124 and discrete (x_i, y_i) . Figure 1 shows the computational graph
 125 corresponding to Eq. 2.

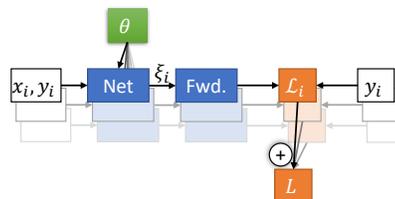


Figure 1: Reparameterized optimization

126 While the tasks of optimizing inverse problems and learning
 127 patterns from data may seem unrelated at first, there is a
 128 strong connection between the two. The inductive bias of a chosen network architecture, which
 129 enables generalization, also affects the update direction of classical optimizers under reparam-
 130 eterization. This can be seen most clearly if we consider gradient descent steps. There, individual
 131 optimization yields the updates $\Delta \xi_i = -\eta \frac{\partial \mathcal{L}_i}{\partial \xi_i}$ with η denoting the step size. After reparameterization,
 132 the updates are $\Delta \theta = -\eta \sum_i \frac{\partial \mathcal{L}_i}{\partial \xi_i} \frac{\partial \mathcal{N}}{\partial \theta}$. As we can see, $\frac{\partial \mathcal{N}}{\partial \theta}$, which is independent of the specific
 133 example, now contributes a large part to the update direction, allowing for cross-talk between the
 134 different optimization problems.

135 Despite the similarities to machine learning, the different use case of this setup leads to differences
 136 in the training procedure. For example, while overfitting is usually seen as undesirable in machine
 137 learning, we want the solutions to our inverse problems to be as accurate as possible, i.e. we want to
 138 "overfit" to the data. Consequently, we do not have to worry about the curvature at θ^* and will not use
 139 mini-batches for training the reparameterization network.

140 **Supervised training.** Our main goal is obtaining an optimization scheme that works exactly like
 141 classical optimizers, only requiring the forward process F , x_i in the form of a numerical simulator,
 142 and desired outputs y_i . However, if we additionally have a prior on the solution space $P(\xi)$, we can
 143 generate synthetic training data $\{(x_j, y_i), \xi_j\}$ with $y_j = F(x_j, \xi_j)$ by sampling $\xi_i \sim P(\xi)$. Using
 144 this data set, we can alternatively train \mathcal{N} with the supervised objective

$$\tilde{L} = \sum_j \|\mathcal{N}(x_j, y_j) - \xi_j\|_2^2. \quad (3)$$

145 Since \mathcal{N} has the same inputs and outputs, we can use the same network architecture as above
 146 and the solutions to the original inverse problems can be obtained as $\xi_i = \mathcal{N}(x_i, y_i)$. While this
 147 method requires domain knowledge in the form of the distributions $P(x)$ and $P(\xi)$, it has the distinct
 148 advantage of being independent of the characteristics of F . For example, if F is chaotic, directly
 149 optimizing through F can yield very large and unstable gradients, while the loss landscape of \tilde{L} can
 150 still be smooth. However, we cannot expect the inferred solutions to be highly accurate as the network
 151 is not trained on the inverse problems we want to solve and, thus, has to interpolate. Additionally, this
 152 method is only suited to unimodal problems, i.e. inverse problems with a unique global minimum.
 153 On multimodal problems, the network cannot be prevented from learning an interpolation of possible
 154 solutions, which may result in poor accuracy.

155 **Refinement** Obtaining a high accuracy on the inverse problems of interest is generally difficult
 156 when the training set size is limited, which can result in suboptimal solutions. This is especially
 157 problematic when the global minima are narrow and no direct feedback from F is available, as in the
 158 case of supervised training. To ensure that all learned methods have the potential to compete with
 159 gradient-based optimizers like BFGS, we pass the solution estimates for ξ to a secondary refinement
 160 stage where they are used as an initial guess for BFGS. The refinement uses the true gradients of F
 161 to find a nearby minimum of \mathcal{L} .

162 4 Experiments

163 We perform a series of numerical experiments to test the convergence properties of the reparameterized
 164 joint optimization. An overview of the experiments is given in Tab. 1 and additional details of all
 165 performed experiments can be found in Appendix B. We run each experiment and method multiple
 166 times, varying the neural network initializations and data sets to obtain statistically significant results.

167 To test the capabilities of the algorithms as a black-box extension of generic optimizers, all experi-
 168 ments use off-the-shelf neural network architectures and only require hyperparameter tuning in terms
 169 of decreasing the Adam [KB15] learning rate until stable convergence is reached. We then compare
 170 the reparameterized optimization to BFGS [LN89], a popular classical solver for unconstrained
 171 optimization problems, and to the neural adjoint method, which has been shown to outperform
 172 various other neural-network-based approaches for solving inverse problems [RPM20].

173 **Neural adjoint** The neural adjoint method relies on an approximation of the forward process
 174 by a surrogate neural network $S(x_i, \xi_i | \theta)$. We first train the surrogate on an independent data set

Table 1: Overview of numerical experiments.

Experiment	$\nabla = 0$ areas	Chaotic	x_i known	$P(\xi)$ known
Wave packet localization	No	No	No	Yes
Billiards	Yes	No	Yes	No
Kuramoto–Sivashinsky	No	Yes	Yes	Yes
Incompr. Navier-Stokes	No	Yes	No	Yes

175 generated from the same distribution as the inverse problems and contains many examples. We use
 176 the same examples as for the supervised approach outlined above but switch the labels to match the
 177 network design, $\{(x_i, \xi_i), y_i\}$. After training, the weights θ are frozen and BFGS is used to optimize
 178 ξ_i on the proxy process $\hat{F}(\xi_i | x_i) = S(\xi_i, x_i) + B(\xi_i)$ where B denotes a boundary loss term (see
 179 Appendix A). With the loss function \mathcal{L} from Eq. 1, this yields the effective objective $\mathcal{L}(F(\xi_i | x_i), y_i)$
 180 for solving the inverse problems. Like with the other methods, the result of the surrogate optimization
 181 is then used as a starting point for the refinement stage described above.

182 4.1 Wave packet localization

183 First, we consider a 1D curve fit. A noisy signal $u(t)$ containing a wave packet centered at t_0
 184 is measured, resulting in the observed data $u(t) = A \cdot \sin(t - t_0) \cdot \exp(-(t - t_0)^2 / \sigma^2) + \epsilon(t)$ where
 185 $\epsilon(t)$ denotes random noise and $t = 1, \dots, 256$. An example waveform is shown in Fig. 2a. For fixed
 186 A and σ , the task is to locate the wave packet, i.e. retrieve t_0 . This task is difficult for optimization
 187 algorithms because the loss landscape (Fig. 2b) contains many local optima that must be traversed.
 188 This results in alternating gradient directions when traversing the parameter space, with maximum
 189 magnitude near the correct solution.

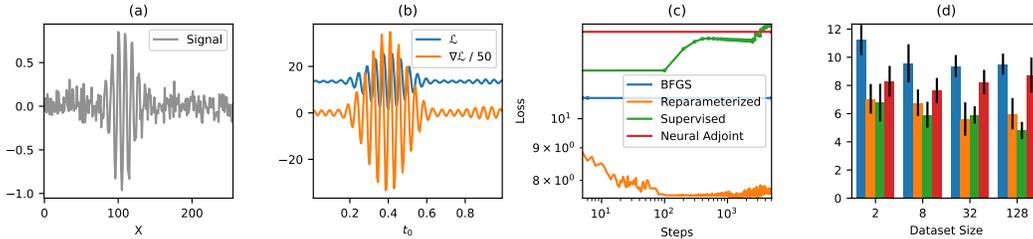


Figure 2: Wave packet localization. (a) Example waveform $u(t)$, (b) corresponding loss and gradient landscape for t_0 , (c) optimization curves without refinement, (d) refined loss L/n by the number of examples n , mean and standard deviation over multiple network initializations and data sets.

190 We generate the inverse problems by sampling random t_0 and $\epsilon(t)$ from ground truth prior distributions
 191 and simulating the corresponding outputs $u(t) = F\epsilon(t | t_0)$. Because the noise distribution $\epsilon(t)$ is
 192 not available to any of the optimization methods, a perfect solution with $\mathcal{L} = 0$ is impossible.

193 Fig. 2c shows the optimization process. Iterative optimizers like BFGS get stuck in local minima
 194 quickly on this task. In most examples, BFGS moves a considerable distance in the first iteration and
 195 then quickly halts. However, due to the oscillating gradient directions, this initial step is likely to
 196 propel the estimate away from the global optimum, leading many solutions to lie further from the
 197 actual optimum than the initial guess.

198 The neural adjoint method finds better solutions than BFGS for about a third of examples for $n = 256$
 199 (see Tab. 2). In many cases, the optimization progresses towards the boundary and gets stuck once
 200 the boundary loss B balances the gradients from the surrogate network.

201 To reparameterize the problem, we create a neural network \mathcal{N} that maps the 256 values of the observed
 202 signal $u(t)$ to the unknown value t_0 . We chose a standard architecture inspired by image classification
 203 networks [SZ14] and train it according to Eq. 2. The network consists of five convolutional layers
 204 with ReLU activation functions, batch normalization, and max-pooling layers, followed by two
 205 fully-connected layers. During the optimization, the estimate of t_0 repeatedly moves from minimum
 206 to minimum until settling after around 500 iterations. Like BFGS, most examples do not converge to
 207 the global optimum and stop at a local minimum instead. However, the cross-talk between different
 208 examples, induced by the shared parameters θ and the summation of the individual loss functions,
 209 regularizes the movement in t_0 space, preventing solutions from moving far away from the global
 210 optimum. Meanwhile, the feedback from the analytic gradients of F ensures that each example finds
 211 a locally optimal solution. Overall, this results in around 80% of examples finding a better solution
 212 than BFGS.

213 For supervised training of \mathcal{N} , we use the same training data set as for the neural adjoint method. This
 214 approach’s much smoother loss landscape lets all solution estimates progress close to the ground

215 truth. However, lacking the gradient feedback from the forward process \mathcal{F} , the inferred solutions
 216 are slightly off from the actual solution and, since the highest loss values are close to the global
 217 optimum, this raises the overall loss during training even though the solutions are approaching the
 218 global optima. This phenomenon gets resolved with solution refinement using BFGS.

219 Fig. 2d shows the results for different numbers of inverse problems and training set sizes n . Since
 220 BFGS optimizes each example independently, the data set size has no influence on its performance.
 221 Variances in the mean final loss indicate that the specific selection of inverse problems may be slightly
 222 easier or harder to solve than the average. The neural adjoint method and reparameterized optimization
 223 both perform better than BFGS with the reparameterized optimization producing lower loss values.
 224 However, both do not scale with n in this example. This feature can only be observed with supervised
 225 training whose solution quality noticeably increases with n . This is due to the corresponding increase
 226 in training set size, which allows the model to improve generalization and does not depend on
 227 the number of tested inverse problems. For $n \geq 32$, supervised training in combination with the
 228 above-mentioned solution refinement consistently outperforms all other methods.

229 A detailed description of the network architecture along with additional learning curves, parameter
 230 evolution plots as well as the performance on further data set sizes n can be found in Appendix B.1.

231 4.2 Billiards

232 Next, we consider a rigid-body setup inspired by differentiable billiards simulations of previous work
 233 [HAL⁺20]. The task consists of finding the optimal initial velocity \vec{v}_0 of a cue ball so it hits another
 234 ball, imparting momentum in a non-elastic collision to make the second ball come to rest at a fixed
 235 target location. This setup is portrayed in Fig. 3a and the corresponding loss landscape for a fixed x
 236 velocity in Fig. 3b. A collision only occurs if \vec{v}_0 is large enough and pointed towards the other ball.
 237 Otherwise, the second ball stays motionless, resulting in a constant loss value and $\frac{\partial \mathcal{L}}{\partial \vec{v}_0} = 0$.

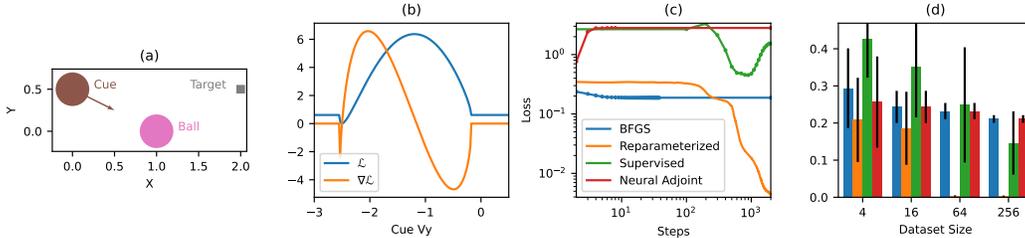


Figure 3: Billiards experiment. **(a)** Task: the cue ball must hit the other ball so that it comes to rest at the target, **(b)** corresponding loss and gradient landscape for v_y , **(c)** optimization curves without refinement, **(d)** refined loss L/n by number of examples n , mean and standard deviation over multiple network initializations and data sets.

238 This property prevents classical optimizers from converging if they hit such a region in the solution
 239 space. The optimization curves are shown in Fig. 3c. BFGS only converges for those examples where
 240 the cue ball already hits the correct side of the other ball.

241 For reparameterization, we employ a fully-connected neural network \mathcal{N} with three hidden layers
 242 using Sigmoid activation functions and positional encoding. The joint optimization with \mathcal{N} drastically
 243 improves the solutions. While for $n \leq 32$ only small differences to BFGS can be observed, access to
 244 more inverse problems lets gradients from some problems steer the optimization of others that get
 245 no useful feedback. This results in almost all problems converging to the solution for $n \geq 64$ (see
 246 Fig. 3d).

247 In this experiment, the distribution of the solutions $P(\vec{v}_0)$ is not available as hitting the target precisely
 248 requires a specific velocity \vec{v}_0 that is unknown a-priori. We can, however, generate training data with
 249 varying \vec{v}_0 and observe the final positions of the balls, then train a supervised \mathcal{N} as well as a surrogate
 250 network for the neural adjoint method on this data set. However, this is less efficient as most of the
 251 examples in the data set do not result in an optimal collision.

252 The neural adjoint method fails to approach the true solutions and instead gets stuck on the training
 253 data boundary in solution space. Likewise, the supervised model cannot accurately extrapolate the
 254 true solution distribution from the sub-par training set.

255 4.3 Kuramoto–Sivashinsky equation

256 The Kuramoto–Sivashinsky (KS) equation, originally developed to model the unstable behavior of
 257 flame fronts [Kur78], models a chaotic one-dimensional system, $\dot{u}(t) = -\frac{\partial^2 u}{\partial x^2} - \frac{\partial^4 u}{\partial x^4} - u \cdot \nabla u$. We
 258 consider a two-parameter inverse problem involving the forced KS equation with altered advection
 259 strength,

$$\dot{u}(t) = \alpha \cdot G(x) - \frac{\partial^2 u}{\partial x^2} - \frac{\partial^4 u}{\partial x^4} - \beta \cdot u \cdot \nabla u,$$

260 where $G(x)$ is a fixed time-independent forcing term and $\alpha, \beta \in \mathbb{R}$ denote the unknown parameters
 261 governing the evolution. Each inverse problem starts from a randomly generated initial state $u(t=0)$
 262 and is simulated until $t=25$, by which point the system becomes chaotic but is still smooth enough
 263 to allow for gradient-based optimization. We constrain $\alpha \in [-1, 1]$, $\beta \in [\frac{1}{2}, \frac{3}{2}]$ to keep the system
 264 numerically stable. Fig. 4a shows example trajectories of this setup and the corresponding gradient
 265 landscape of $\frac{\partial \mathcal{L}}{\partial \beta} \Big|_{\alpha=\alpha^*}$ for the true value of α is shown in Fig. 4b.

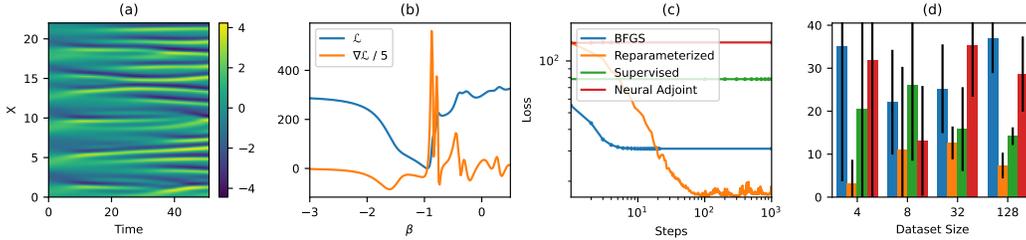


Figure 4: Kuramoto–Sivashinsky experiment. **(a)** Example trajectory, **(b)** corresponding loss and gradient landscape for β , **(c)** optimization curves without refinement, **(d)** refined loss L/n by number of examples n , mean and standard deviation over multiple network initializations and data sets.

266 Fig. 4c shows the optimization curves for finding α, β . Despite the complex nature of the loss
 267 landscape, BFGS manages to find the correct solution in about 60% of cases. The reparameterized
 268 optimization, based on a similar network architecture as for the wavepacket experiment but utilizing
 269 2D convolutions, finds the correct solutions in over 80% of cases but, without refinement, the accuracy
 270 stagnates far from machine precision. Refining these solutions with BFGS, as described above, sees
 271 the accuracy of these cases decrease to machine precision in 4 to 17 iterations, less than the 12 to 22
 272 that BFGS requires when initialized from the distribution mean $\mathbb{E}[P(\xi)]$.

273 Supervised training with refinement produces better solutions in 58% of examples, averaged over
 274 the shown n . The unrefined solutions benefit from larger n on this example because of the large
 275 number of possible observed outputs that the KS equation can produce for varying α, β . At $n=2$,
 276 all unrefined solutions are worse than BFGS while for $n \geq 64$ around 20% of problems find better
 277 solutions. With refinement, these number jump to 50% and 62%.

278 This property also makes it hard for a surrogate network, required by the neural adjoint method,
 279 to accurately approximate the KS equation, causing the following adjoint optimization to yield
 280 inaccurate results that fail to match BFGS even after refinement.

281 4.4 Incompressible Navier-Stokes

282 Incompressible Newtonian fluids are described by the Navier-Stokes equations,

$$\dot{u}(\vec{x}, t) = \nu \nabla^2 u - u \cdot \nabla u - \nabla p \quad \text{s.t.} \quad \nabla^2 p = \nabla \cdot v$$

283 with $\nu \geq 0$. As they can result in highly complex dynamics [BB67], they represent a particularly
 284 challenging test case, which is relevant for a variety of real-world problems [Pop00]. We consider

285 a setup similar to particle imaging velocimetry [Gra97] in which the velocity in the upper half of
 286 a two-dimensional domain with obstacles can be observed. The velocity is randomly initialized in
 287 the whole domain and a localized force is applied near the bottom of the domain at $t = 0$. The task
 288 is to reconstruct the position x_0 and initial velocity \vec{v}_0 of this force region by observing the initial
 289 and final velocity field only in the top half of the domain. The initial velocity in the bottom half is
 290 unknown and cannot be recovered, making a perfect fit impossible. Fig. 5a,b show an example initial
 291 and final state of the system. The final velocity field is measured at $t = 56$ by which time fast eddies
 292 have dissipated significantly.

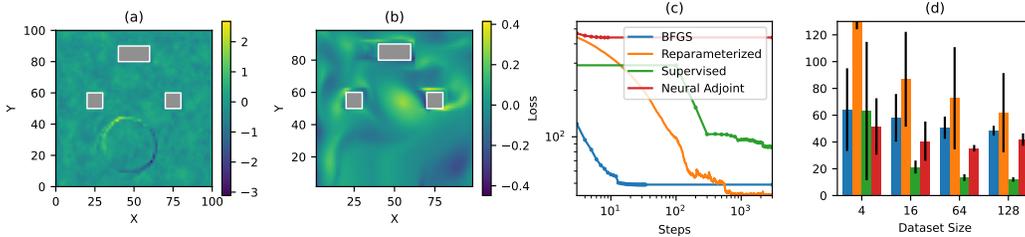


Figure 5: Fluid experiment. **(a,b)** Example initial and final velocity fields, obstacles in gray. Only the upper half, $y \geq 50$, is observed. **(c)** Optimization curves without refinement, **(d)** refined loss L/n by the number of examples n , mean and standard deviation over multiple network initializations and data sets.

293 Fig. 5c shows the optimization curves. On this problem, BFGS converges to some optimum in all
 294 cases, usually within 10 iterations, sometimes requiring up to 40 iterations. However, many examples
 295 get stuck in local optima.

296 For joint optimization, we reparameterize the solution space using a network architecture similar
 297 to the previous experiment, featuring four 2D convolutional layers and two fully-connected layers.
 298 For all tested n , the reparameterized optimization produces larger mean loss values than BFGS,
 299 especially for small n . This results from about 10% of examples seeing higher than average loss
 300 values. Nonetheless, 66.7% of the inverse problems are solved more accurately than BFGS on average
 301 for $n > 4$.

302 The neural adjoint method nearly always converges to solutions within the training set parameter
 303 space, not relying on the boundary loss. With solution refinement, this results in a mean loss that
 304 seems largely independent of n and is slightly lower than the results from direct BFGS optimization.
 305 However, most of this improvement comes from the secondary refinement stage which runs BFGS on
 306 the true F . Without solutions refinement, the neural adjoint method yields inaccurate results, losing
 307 to BFGS in 98.2% of cases.

308 Supervised training does not suffer from examples getting stuck in a local minimum early on. The
 309 highest-loss solutions, which contribute the most to L , are about an order of magnitude better than
 310 the worst BFGS solutions, leading to a much smaller total loss for $n \geq 16$. With solution refinement,
 311 64%, 73% and 72% of examples yield a better solution than BFGS for $n = 16, 64, 128$, respectively.

312 5 Discussion

313 In our experiments, we have focused on relatively small data sets of between 2 and 256 examples to
 314 quantify the worst-case for machine learning methods and observe trends. Using off-the-shelf neural
 315 network architectures and optimizers with no tuning to the specific problem, joint optimization finds
 316 better solutions than BFGS in an average of 69% of tested problems. However, to achieve the best
 317 accuracy, the solution estimates must be passed to a classical optimizer for refinement as training
 318 the network to this level of accuracy would take an inordinate amount of time and large data sets.
 319 Tuning the architectures to the specific examples could lead to further improvements in performance
 320 but would make the approach domain-dependent.

321 When training data including ground truth solutions are available or can be generated, supervised
 322 learning can sidestep many difficulties that complex loss landscapes pose, such as local minima,

Table 2: Fraction of inverse problems for which neural-network-based methods with refinement find better or equal solutions than BFGS. Mean over multiple seeds and all n shown in subfigures (d).

Experiment	Reparameterized		Supervised		Neural Adjoint	
	Better	Equal	Better	Equal	Better	Equal
Wave packet fit	86.0%	1.8%	65.1%	14.4%	40.2%	47.4%
Billiards	61.7%	9.0%	27.0%	27.2%	1.6%	98.4%
Kuramoto–Sivashinsky	62.3%	0.0%	57.7%	0.0%	23.9%	62.2%
Incompr. Navier-Stokes	64.1%	0.0%	66.2%	0.1%	56.9%	0.1%

323 alternating gradient directions, or zero-gradient areas. This makes supervised learning another
 324 promising alternative to direct optimization, albeit a more involved one.

325 The neural adjoint method, on the other hand, yields only very minor improvements over BFGS
 326 optimization in our experiments, despite the surrogate network successfully learning to reproduce the
 327 training data. This is not surprising as the neural adjoint method tries to approximate the original
 328 loss landscape which is often difficult to optimize. Improvements over BFGS must therefore come
 329 from regularization effects and exposure to a larger part of the solution space. The fact that the neural
 330 adjoint method with solution refinement produces similar results almost independent of the number
 331 of data points n shows that the joint optimization has little benefit here. Instead, the refinement stage,
 332 which treats all examples independently, dominates the final solution quality. Note that the neural
 333 adjoint method is purely data-driven and does not require an explicit form for the forward process F ,
 334 making it more widely applicable than the setting considered here.

335 Tab. 2 summarizes the improvements over classical optimizations for all methods. A corresponding
 336 table without solution refinement can be found in Appendix B. Considering that reparameterized
 337 optimization is the only network-based method that does not require domain-specific information
 338 and nevertheless shows the biggest improvement overall, we believe it is the most attractive variant
 339 among the three learned versions. Inverse problems for which reparameterized training does not
 340 find good solutions are easy to identify by their outlier loss values. In these cases, one could simply
 341 compare the solution to a reference solution obtained via direct optimization, and choose the best
 342 result.

343 **Limitations** We have only considered unconstrained optimization problems in this work, enforcing
 344 hard constraints by running bounded parameters through a scaled tanh function which naturally
 345 clamps out-of-bounds values in a differentiable manner.

346 The improved solutions found by joint optimization come with an increased computational cost
 347 compared to direct optimization. The time it took to train the reparameterization networks was 3x to
 348 6x longer for the first three experiments and 22x for the fluids experiment.

349 6 Conclusions and outlook

350 We have investigated the effects of joint optimization of multiple inverse problems by reparameterizing
 351 the solution space using a neural network, showing that joint optimization can often find better
 352 solutions than classical optimization techniques. Since our reparameterization approach does not
 353 require any more information than classical optimizers, it can be used as a drop-in replacement. This
 354 could be achieved by adding a function or option to existing optimization libraries that internally
 355 sets up a standard neural network with the required number of inputs and outputs and runs the
 356 optimization, hiding details of the training process, network architecture, and hyperparameters from
 357 the user while making the gains in optimization accuracy conveniently accessible. To facilitate this,
 358 we will make the full source code publicly available.

359 From accelerating matrix multiplications [FBH⁺22] to solving systems of linear equations [CHL⁺23,
 360 SSHR19], it is becoming increasingly clear that machine learning methods can be applied to purely
 361 numerical problems outside of typical big data settings, and our results show that this also extends to
 362 solving nonlinear inverse problems.

363 **References**

- 364 [AAA⁺13] M Agostini, M Allardt, E Andreotti, AM Bakalyarov, M Balata, I Barabanov, M Barnabé
365 Heider, N Barros, L Baudis, C Bauer, et al. Pulse shape discrimination for gerda phase
366 i data. *The European Physical Journal C*, 73(10):2583, 2013.
- 367 [AAA⁺18] Craig E Aalseth, N Abgrall, Estanislao Aguayo, SI Alvis, M Amman, Isaac J Arnquist,
368 FT Avignone III, Henning O Back, Alexander S Barabash, PS Barbeau, et al. Search
369 for neutrinoless double- β decay in ge 76 with the majorana demonstrator. *Physical*
370 *review letters*, 120(13):132502, 2018.
- 371 [ALGS⁺22] Kelsey R Allen, Tatiana Lopez-Guevara, Kimberly Stachenfeld, Alvaro Sanchez-
372 Gonzalez, Peter Battaglia, Jessica Hamrick, and Tobias Pfaff. Physical design using
373 differentiable learned simulators. *arXiv preprint arXiv:2202.00728*, 2022.
- 374 [BB67] Cx K Batchelor and George Keith Batchelor. *An introduction to fluid dynamics*. Cam-
375 bridge university press, 1967.
- 376 [BHM18] Mikhail Belkin, Daniel J Hsu, and Partha Mitra. Overfitting or perfect fitting? risk
377 bounds for classification and regression rules that interpolate. *Advances in neural*
378 *information processing systems*, 31, 2018.
- 379 [BPL21] Randall Balestriero, Jerome Pesenti, and Yann LeCun. Learning in high dimension
380 always amounts to extrapolation. *arXiv preprint arXiv:2110.09485*, 2021.
- 381 [CCC⁺19] Giuseppe Carleo, Ignacio Cirac, Kyle Cranmer, Laurent Daudet, Maria Schuld, Naftali
382 Tishby, Leslie Vogt-Maranto, and Lenka Zdeborová. Machine learning and the physical
383 sciences. *Reviews of Modern Physics*, 91(4):045002, 2019.
- 384 [CHL⁺23] Salvatore Cali, Daniel C Hackett, Yin Lin, Phiala E Shanahan, and Brian Xiao. Neural-
385 network preconditioners for solving the dirac equation in lattice gauge theory. *Physical*
386 *Review D*, 107(3):034508, 2023.
- 387 [DJO⁺18] S Delaquis, MJ Jewell, I Ostrovskiy, M Weber, T Ziegler, J Dalmasson, LJ Kaufman,
388 T Richards, JB Albert, G Anton, et al. Deep neural networks for energy and position
389 reconstruction in exo-200. *Journal of Instrumentation*, 13(08):P08023, 2018.
- 390 [FBH⁺22] Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-
391 Paredes, Mohammadamin Barekatin, Alexander Novikov, Francisco J R Ruiz, Julian
392 Schrittwieser, Grzegorz Swirszcz, et al. Discovering faster matrix multiplication algo-
393 rithms with reinforcement learning. *Nature*, 610(7930):47–53, 2022.
- 394 [GBCB16] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*.
395 MIT press Cambridge, 2016.
- 396 [GEB16] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using
397 convolutional neural networks. In *Proceedings of the IEEE conference on computer*
398 *vision and pattern recognition*, pages 2414–2423, 2016.
- 399 [GH18] Daniel George and EA Huerta. Deep learning for real-time gravitational wave detection
400 and parameter estimation: Results with advanced ligo data. *Physics Letters B*, 778:64–
401 70, 2018.
- 402 [GM78] Philip E Gill and Walter Murray. Algorithms for the solution of the nonlinear least-
403 squares problem. *SIAM Journal on Numerical Analysis*, 15(5):977–992, 1978.
- 404 [Gra97] Ian Grant. Particle image velocimetry: a review. *Proceedings of the Institution of*
405 *Mechanical Engineers, Part C: Journal of Mechanical Engineering Science*, 211(1):55–
406 76, 1997.
- 407 [HAL⁺20] Yuanming Hu, Luke Anderson, Tzu-Mao Li, Qi Sun, Nathan Carr, Jonathan Ragan-
408 Kelley, and Frédo Durand. DiffTaichi: Differentiable programming for physical simula-
409 tion. *International Conference on Learning Representations (ICLR)*, 2020.

- 410 [Hay94] Simon Haykin. *Neural networks: a comprehensive foundation*. Prentice Hall PTR,
411 1994.
- 412 [HKT20] Philipp Holl, Vladlen Koltun, and Nils Thuerey. Learning to control pdes with differ-
413 entiable physics. In *International Conference on Learning Representations (ICLR)*,
414 2020.
- 415 [HKT21] Philipp Holl, Vladlen Koltun, and Nils Thuerey. Scale-invariant learning by physics
416 inversion. *arXiv preprint arXiv:2109.15048*, 2021.
- 417 [HSDG19] Stephan Hoyer, Jascha Sohl-Dickstein, and Sam Greydanus. Neural reparameterization
418 improves structural optimization. *arXiv preprint arXiv:1909.04240*, 2019.
- 419 [KAT⁺19] Byungsoo Kim, Vinicius C. Azevedo, Nils Thuerey, Theodore Kim, Markus Gross,
420 and Barbara Solenthaler. Deep fluids: A generative network for parameterized fluid
421 simulations. *Computer Graphics Forum*, 2019.
- 422 [KB15] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In
423 *International Conference on Learning Representations (ICLR)*, 2015.
- 424 [KGZ⁺21] Aditi Krishnapriyan, Amir Gholami, Shandian Zhe, Robert Kirby, and Michael W
425 Mahoney. Characterizing possible failure modes in physics-informed neural networks.
426 *Advances in Neural Information Processing Systems*, 34:26548–26560, 2021.
- 427 [KSA⁺21] Dmitrii Kochkov, Jamie A. Smith, Ayya Alieva, Qing Wang, Michael P. Brenner,
428 and Stephan Hoyer. Machine learning accelerated computational fluid dynamics.
429 *arXiv:2102.01010 [physics]*, 2021.
- 430 [KSH12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with
431 deep convolutional neural networks. In *Advances in Neural Information Processing*
432 *Systems*, 2012.
- 433 [Kur78] Yoshiki Kuramoto. Diffusion-induced chaos in reaction systems. *Progress of Theoretical*
434 *Physics Supplement*, 64:346–367, 1978.
- 435 [LLK19] Junbang Liang, Ming Lin, and Vladlen Koltun. Differentiable cloth simulation for
436 inverse problems. In *Advances in Neural Information Processing Systems*, pages
437 771–780, 2019.
- 438 [LN89] Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale
439 optimization. *Mathematical programming*, 45(1-3):503–528, 1989.
- 440 [LPY⁺21] Lu Lu, Raphael Pestourie, Wenjie Yao, Zhicheng Wang, Francesc Verdugo, and
441 Steven G Johnson. Physics-informed neural networks with hard constraints for in-
442 verse design. *SIAM Journal on Scientific Computing*, 43(6):B1105–B1132, 2021.
- 443 [LSAH20] Housen Li, Johannes Schwab, Stephan Antholzer, and Markus Haltmeier. Nett: Solving
444 inverse problems with deep neural networks. *Inverse Problems*, 36(6):065005, 2020.
- 445 [MAL13] Kohta Murase, Markus Ahlers, and Brian C Lacki. Testing the hadronuclear origin of
446 pev neutrinos observed with iccube. *Physical Review D*, 88(12):121301, 2013.
- 447 [MLA⁺19] Rajesh Maingi, Arnold Lumsdaine, Jean Paul Allain, Luis Chacon, SA Gourlay,
448 CM Greenfield, JW Hughes, D Humphreys, V Izzo, H McLean, et al. Summary
449 of the fesac transformative enabling capabilities panel report. *Fusion Science and*
450 *Technology*, 75(3):167–177, 2019.
- 451 [Pop00] Stephen Pope. *Turbulent Flows*. Cambridge University Press, 2000.
- 452 [PTVF07] William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery.
453 *Numerical Recipes*. Cambridge University Press, 3 edition, 2007.
- 454 [PW17] Patrick Putzky and Max Welling. Recurrent inference machines for solving inverse
455 problems. *arXiv preprint arXiv:1706.04008*, 2017.

- 456 [RPK19] Maziar Raissi, Paris Perdikaris, and George Karniadakis. Physics-informed neural
457 networks: A deep learning framework for solving forward and inverse problems in-
458 volving nonlinear partial differential equations. *Journal of Computational Physics*,
459 378:686–707, 2019.
- 460 [RPM20] Simiao Ren, Willie Padilla, and Jordan Malof. Benchmarking deep inverse models
461 over time, and the neural-adjoint method. *Advances in Neural Information Processing*
462 *Systems*, 33:38–48, 2020.
- 463 [RT21] Stephan Rasp and Nils Thuerey. Data-driven medium-range weather prediction with a
464 resnet pretrained on climate simulations: A new model for weatherbench. *Journal of*
465 *Advances in Modeling Earth Systems*, 13(2):e2020MS002405, 2021.
- 466 [SC19] Samuel S Schoenholz and Ekin D Cubuk. Jax, md: End-to-end differentiable, hardware
467 accelerated, molecular dynamics in pure python. *arXiv:1912.04232*, 2019.
- 468 [SF18] Connor Schenck and Dieter Fox. Spnets: Differentiable fluid dynamics for deep neural
469 networks. In *Conference on Robot Learning*, pages 317–335, 2018.
- 470 [SFK⁺21] Kimberly Stachenfeld, Drummond B Fielding, Dmitrii Kochkov, Miles Cranmer, Tobias
471 Pfaff, Jonathan Godwin, Can Cui, Shirley Ho, Peter Battaglia, and Alvaro Sanchez-
472 Gonzalez. Learned coarse models for efficient turbulence simulation. *arXiv preprint*
473 *arXiv:2112.15275*, 2021.
- 474 [SGGP⁺20] Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec,
475 and Peter Battaglia. Learning to simulate complex physics with graph networks. In
476 *International conference on machine learning*, pages 8459–8468. PMLR, 2020.
- 477 [SHT22] Patrick Schnell, Philipp Holl, and Nils Thuerey. Half-inverse gradients for physical
478 deep learning. *arXiv preprint arXiv:2203.10131*, 2022.
- 479 [SSHR19] Johannes Sappl, Laurent Seiler, Matthias Harders, and Wolfgang Rauch. Deep learning
480 of preconditioners for conjugate gradient solvers in urban water related problems. *arXiv*
481 *preprint arXiv:1906.06925*, 2019.
- 482 [SZ14] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-
483 scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 484 [Tar05] Albert Tarantola. *Inverse problem theory and methods for model parameter estimation*.
485 SIAM, 2005.
- 486 [TLQL21] Tetsuya Takahashi, Junbang Liang, Yi-Ling Qiao, and Ming C Lin. Differentiable fluids
487 with solid coupling for learning and control. In *Proceedings of the AAAI Conference on*
488 *Artificial Intelligence*, volume 35(7), pages 6138–6146, 2021.
- 489 [TSSP17] Jonathan Tompson, Kristofer Schlachter, Pablo Sprechmann, and Ken Perlin. Acceler-
490 ating eulerian fluid simulation with convolutional networks. In *Proceedings of Machine*
491 *Learning Research*, pages 3424–3433, 2017.
- 492 [UBF⁺20] Kiwon Um, Robert Brand, Yun Raymond Fei, Philipp Holl, and Nils Thuerey. Solver-
493 in-the-loop: Learning from differentiable physics to interact with iterative pde-solvers.
494 *Advances in Neural Information Processing Systems*, 2020.
- 495 [UVL18] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *Pro-*
496 *ceedings of the IEEE conference on computer vision and pattern recognition*, pages
497 9446–9454, 2018.
- 498 [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
499 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in*
500 *Neural Information Processing Systems*, pages 5998–6008, 2017.
- 501 [Yua00] Ya-xiang Yuan. A review of trust region algorithms for optimization. In *Iciam*, volume
502 99(1), pages 271–282, 2000.
- 503 [YZWX19] XIA Yang, Suhaib Zafar, J-X Wang, and Heng Xiao. Predictive large-eddy-simulation
504 wall modeling via physics-informed neural networks. *Physical Review Fluids*,
505 4(3):034602, 2019.