

Which Optimizer Works Best for Physics-Informed Neural Networks and Kolmogorov-Arnold Networks?

Elham Kiyani^a, Khemraj Shukla^a, Jorge F. Urbán^b, Jérôme Darbon^a, George Em Karniadakis^a

^a*Division of Applied Mathematics, Brown University, 182 George Street, Providence, RI 02912, USA*

^b*Departament de Física Aplicada, Universitat d'Alacant, Ap. Correus 99, Alacant, 03830, Comunitat Valenciana, Spain*

Abstract

Physics-Informed Neural Networks (PINNs) have revolutionized the computation of PDE solutions by integrating partial differential equations (PDEs) into the neural network's training process as soft constraints, becoming an important component of the scientific machine learning (SciML) ecosystem. More recently, physics-informed Kolmogorov-Arnold networks (PIKANs) have also shown to be effective and comparable in accuracy with PINNs. In their current implementation, both PINNs and PIKANs are mainly optimized using first-order methods like Adam, as well as quasi-Newton methods such as BFGS and its low-memory variant, L-BFGS. However, these optimizers often struggle with highly non-linear and non-convex loss landscapes, leading to challenges such as slow convergence, local minima entrapment, and (non)degenerate saddle points. In this study, we investigate the performance of Self-Scaled BFGS (SSBFGS), Self-Scaled Broyden (SSBroyden) methods and other advanced quasi-Newton schemes, including BFGS and L-BFGS with different line search strategies approaches. These methods dynamically rescale updates based on historical gradient information, thus enhancing training efficiency and accuracy. We systematically compare these optimizers -using both PINNs and PIKANs- on key challenging linear, stiff, multi-scale and non-linear PDEs, including the Burgers, Allen-Cahn, Kuramoto-Sivashinsky, and Ginzburg-Landau equations. Our findings provide state-of-the-art results with orders-of-magnitude accuracy improvements without the use of adaptive weights or any other enhancements typically employed in PINNs. More broadly, our results reveal insights into the effectiveness of second-order optimization strategies in significantly improving the convergence and accurate generalization of PINNs and PIKANs.

1. Introduction

This section begins by providing the necessary background and motivation for our study, highlighting the key challenges and objectives. Following this, we present an overview of the BFGS (Broyden–Fletcher–Goldfarb–Shanno), SSBFGS (Self-Scaled BFGS), and SSBroyden (Self-Scaled Broyden) optimization methods, discussing their theoretical foundations and main characteristics. Section 2 presents a comprehensive study comparing different optimization methods with various line search and trust-region strategies. Specifically, in Subsections 2.1.1 and 2.1.2, we discuss a wide range of scenarios for the Burgers equation. We extend our study to more complex PDEs, including the Allen-Cahn equation in Section 2.3, the Kuramoto-Sivashinsky equation in Subsection 2.4, and the Ginzburg-Landau equation in Section 2.5. The conclusions of our study are summarized in Section 3. Finally, in Appendix A [Appendix A](#), we review the performance of optimizers for the Lorenz system. Additionally, in Appendix B [Appendix B](#), we provide a quantitative analysis of BFGS and SSBroyden in minimizing the multi-dimensional Rosenbrock function, serving as a pedagogical example to illustrate the impact of dimensionality on optimization performance.

1.1. Background and Motivation

Physics-Informed Neural Networks (PINNs), introduced in 2017, are a groundbreaking development in scientific machine learning (SciML) [1, 2]. By seamlessly integrating the fundamental physical principles of a system with neural networks, PINNs offer a versatile and mesh-free framework for solving nonlinear partial differential equations (PDEs). Unlike traditional numerical methods, PINNs directly incorporate initial and boundary conditions as well as PDE residuals into their loss functions, enabling

them to address both forward problems (predicting solutions) and inverse problems (e.g., estimating unknown parameters or unknown functions). Their adaptability and scalability make PINNs particularly well-suited for tackling challenges in high-dimensional spaces and complex geometries that conventional methods struggle to handle. The network parameters are updated during training to minimize the loss function, resulting in a solution that meets the constraints applied in the loss function. More recently, physics-informed Kolmogorov-Anrold Networks (PIKANs) were introduced in [3, 4, 5] to solve PDEs, so in the present study, we aim to investigate both PINNs and PIKANs with respect to their optimization performance.

Since their introduction, PINNs have undergone numerous extensions and adaptations to enhance their applicability and address limitations in the original framework. These advancements include uncertainty quantification [6], domain decomposition techniques [7, 8, 9, 10], and the incorporation of alternative network architectures such as convolutional neural networks (CNNs) [11] and recurrent neural networks (RNNs) [12]. Innovative approaches like Generative Adversarial PINNs (GA-PINNs) [13], Physics-Informed Graph Convolutional Networks (GCNs) [14, 15], and Bayesian PINNs (B-PINNs) [16, 17, 18, 19] have further broadened the scope of PINNs. Physics-Informed Extreme Learning Machines (PIELM) [20] combine the computational efficiency of Extreme Learning Machines (ELMs) with the physics-informed principles of PINNs. The hp-VPINN method [21] integrates variational principles with neural networks, utilizing high-order polynomial test spaces for improved accuracy.

Extensive studies have also focused on error analysis and theoretical underpinnings of PINNs. Research on error estimates and convergence properties is presented in [22, 23, 24, 25]. Wang et al. [26] proposed a reformulation of the PINN loss function to explicitly incorporate physical causality during training, although its performance was limited to simple benchmark problems. Additionally, Wang et al. [27] established a theoretical foundation by deriving and analyzing the limiting Neural Tangent Kernel (NTK) of PINNs. This analysis has been extended in subsequent studies to justify and refine various PINN extensions [28, 29, 30, 31]. Furthermore, Anagnostopoulos et al. [32] investigated the learning dynamics of PINNs using the Information Bottleneck theory [33], identifying distinct training phases and proposing a loss weighting scheme to reduce generalization error.

One of the greatest challenges in neural network frameworks lies in the inherently non-convex nature of their optimization problems. As a result, a growing body of research has been dedicated to understanding their training dynamics (see e.g., [34, 35, 36]). Optimization methods are broadly categorized into first-order (e.g., stochastic gradient descent) [37, 38], high-order (e.g., Newton's method) [39, 40, 41], and heuristic derivative-free approaches [42, 43]. Compared to first-order methods, high-order optimization achieves faster convergence by leveraging curvature information but faces challenges in handling and storing the Hessian's inverse. To address this, Newton's method variants employ various techniques to approximate the Hessian matrix. For example, BFGS and L-BFGS (Limited-memory BFGS) approximate the Hessian or its inverse using rank-one or rank-two updates based on gradient information from previous iterations [44, 45, 46, 47, 48, 49, 50].

Currently, the standard optimization algorithms used in PINNs are Adam and quasi-Newton methods, such as BFGS or its low-memory variant, L-BFGS. BFGS and L-BFGS enhance significantly the performance of PINNs by incorporating second-order information of the trainable variable space to precondition it. Although these two optimizers are very commonly employed not only in PINNs but also in other optimization problems because of their theoretical and numerical properties, recent experience in PINNs has shown that other algorithms could outperform them for a great variety of problems [51, 52, 53]. Among all of them we can highlight the *Self-Scaled Broyden* algorithms [54], which update the approximation to the inverse Hessian matrix using a self-scaling technique. This advanced optimization technique enhances convergence by adaptively scaling gradient updates based on historical error information. Furthermore, [55] investigates self-scaling within quasi-Newton methods from the Broyden family, introducing innovative scaling schemes and updates. Urban et al. [53] further demonstrated that these advanced optimizers, when combined with rescaled loss functions, significantly enhance the efficiency and accuracy of PINNs, enabling smaller networks to solve problems with significantly greater accuracy. As Physics-Informed Neural Networks (PINNs) often suffer from conflicting gradients between multiple loss components (e.g., data and physics losses), resolving these conflicts is essential for stable and accurate training. Recent work by Wang et al. [56] rigorously proves that second-order optimizers, particularly Newton-type methods like SOAP, exhibit positive gradient alignment and can effectively mitigate such conflicts, leading to improved convergence in multi-task

optimization settings like PINNs.

In the seminal paper by [24], PINN (h_n) is formally described as a method that seeks to find a neural network minimizing a prescribed loss function (\mathcal{L}) within a class of neural networks \mathcal{H}_n , where n represents the number of parameters in the network. The resulting minimizer serves as an approximation to the solution of the PDE. A crucial question posed by [24] is: Does this sequence of neural networks converge to the solution of the PDE? The answer to this question depends on the regularity of the loss functions and the three types of errors observed during the network's training or optimization process.

For completeness, we define these three errors, which also serve as one of the motivations for this work. The errors, as outlined in [57, 58], are:

1. Approximation error.
2. Estimation error.
3. Optimization error.

For example, consider a function class \mathcal{H}_n and let u^* be the solution to the underlying PDE. Let m represent the number of training data points. The function h_m is the minimizer of the loss with m data points, while \hat{h} is the function in \mathcal{H}_n that minimizes the loss with infinitely many data points. The error between h_m and u^* represents the *approximation error*, while the error between \hat{h} and u^* is the *estimation error*. In practice, \tilde{h}_m denotes the approximation obtained after a finite number of training iterations, such as the result of one million iterations of gradient-based optimization. The error between \tilde{h}_m and \hat{h} is referred to as the *optimization error*. While the approximation error is well understood, the optimization error remains a challenging area due to the highly nonlinear and nonconvex nature of the objective function. These challenges are exacerbated by the inclusion of PDE terms in the loss function, leading to degenerate and non-degenerate saddle points. As a result, optimization often requires ad hoc tricks and tedious parameter fine-tuning through trial and error.

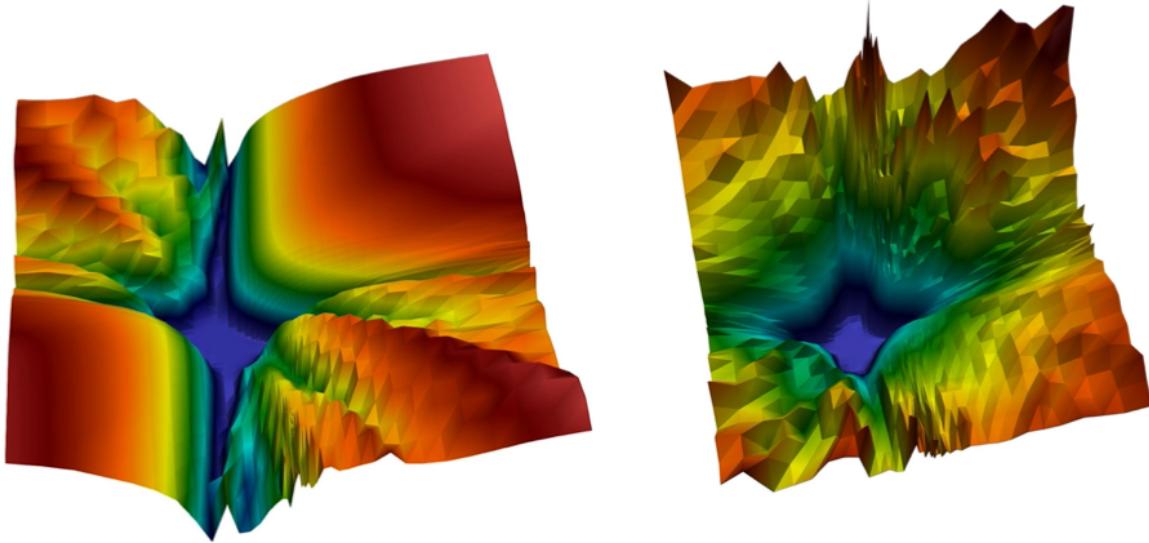


Figure 1: Loss landscape of $\mathcal{L} = \mathcal{L}_{\text{Data}} + \lambda \mathcal{L}_{\text{PDE}}$ for the viscous Burgers equation (16). In the left subfigure, the landscape is shown for $\lambda = 0$. This loss landscape is relatively smooth and exhibits well-defined convexity. However, it only fits the data (initial and boundary conditions) and does not converge to the solution within the domain. In contrast, the right subfigure illustrates the loss landscape for $\lambda = 1.0$. Unlike the smooth landscape in the left subfigure, this landscape contains numerous local minima and saddle points, potentially leading to (non)-degenerate saddle points. The loss landscape is computed by using the method proposed by [59].

To provide an intuitive explanation, Figure 1 illustrates the loss landscape $\mathcal{L} = \mathcal{L}_D + \lambda \mathcal{L}_{PDE}$ for the viscous Burgers Equation (16), using the method proposed by [59]. The left subfigure of Figure 1 corresponds to $\lambda = 0$, considering only the initial and boundary conditions. This loss landscape is smooth and exhibits well-defined convexity. However, it fits only the initial and boundary condition data and fails to converge to the solution within the domain. In contrast, the right subfigure depicts the loss landscape for $\lambda = 1.0$. Unlike the smooth landscape on the left, this landscape is crinkled with numerous local minima and saddle points, including potentially degenerate and non-degenerate saddle

points [60]. Consequently, the presence of multiple local minima can trap the optimizer, making it extremely challenging to minimize the loss function effectively. Motivated by this, in this work, we propose and analyze a suite of second-order optimizers and their applicability to linear, nonlinear, stiff, and chaotic PDEs whose solutions are computed by PINN and PIKAN types of approximation. Furthermore, we introduce a set of optimizers tailored to various classes of PDEs, achieving unprecedented convergence to machine precision with both single and double precision arithmetic. To the best of our knowledge, the current study is the most systematic and comprehensive investigation of how optimizers affect the performance of both PINNs and PIKAns.

1.2. Overview of BFGS and SSBroyden

In optimization, a fundamental challenge is determining the optimal direction and step size to transition from the current point \mathbf{x}_k to an improved solution. In general, two primary approaches address this problem: **line search methods** and **trust region methods**. Both rely on a quadratic model to approximate the objective function around the current iterate. The key distinction lies in how they utilize this approximation; line-search methods determine the step size along a chosen direction, while trust-region methods restrict the step to a pre-defined neighborhood where the model is considered reliable.

1.2.1. Line-search methods

A common approach in optimization is the **line-search strategy**, where the algorithm first selects a **search direction** \mathbf{p}_k that ideally points towards a region of lower function values. Once this direction is determined, the next crucial step is to decide the **step size** α_k , which dictates how far to move along the chosen direction to achieve sufficient improvement in the objective function f . This step involves solving approximately a one-dimensional minimization problem:

$$\min_{\alpha > 0} f(\mathbf{x}_k + \alpha \mathbf{p}_k). \quad (1)$$

The selection of both the search direction and step size plays a fundamental role in the convergence behavior and overall effectiveness of the optimization process. While an exact solution to (1) would maximize the benefit of the chosen direction \mathbf{p}_k , finding the exact minimum is often computationally prohibitive. Instead, line search methods typically rely on an approximate solution, evaluating a finite number of trial step lengths α_k until a suitable reduction in f is achieved. This search, which receives commonly the name of *inexact* line search, relies in a series of mathematical conditions to really ensure to obtain an satisfactory step length.

Once an acceptable step length α is identified, the algorithm updates the current iterate

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k. \quad (2)$$

At the new point \mathbf{x}_{k+1} , the process is repeated by selecting a new search direction \mathbf{p}_{k+1} and step length α_{k+1} . This iterative process continues until convergence criteria are satisfied.

In this paper, we explored two well-known methodologies that we briefly describe next.

Wolfe conditions

The Wolfe conditions are mathematical criteria used to ensure that a step size in iterative optimization methods satisfies specific properties of sufficiency. Given a point \mathbf{x}_k and a direction \mathbf{p}_k , the Wolfe conditions consist of the following two inequalities:

$$f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \leq f(\mathbf{x}_k) + c_1 \alpha_k \mathbf{p}_k^T \nabla f_k, \quad (3)$$

$$\nabla f(\mathbf{x}_k + \alpha_k \mathbf{p}_k)^T \mathbf{p}_k \geq c_2 \nabla f(\mathbf{x}_k)^T \mathbf{p}_k, \quad (4)$$

which commonly receive the names of *Armijo* and *sufficient decrease* conditions, respectively. The quantities c_1 and c_2 are constants that should follow $0 < c_1 < c_2 < 1$. While the first condition ensures a sufficient decrease in the function f , the second one rules out unacceptably short steps. The latter inequality is commonly replaced by taking absolute values at both sides

$$|\nabla f(\mathbf{x}_k + \alpha_k \mathbf{p}_k)^T \mathbf{p}_k| \leq c_2 |\nabla f(\mathbf{x}_k)^T \mathbf{p}_k|. \quad (5)$$

The conditions (3) and (5) together receive the name of Strong Wolfe conditions. For all tests we set $c_1 = 10^{-4}$ and $c_2 = 0.9$, which are the standard choices in the optimization literature.

Backtracking

By only using equation (3), sufficient progress between iterates is not guaranteed, as it does not exclude unacceptably low values of α_k . Because of that, the Wolfe or Strong Wolfe conditions introduce an additional condition given by (4) or (5), respectively. However, another line-search strategy that has proven numerically to be successful and does not need an additional condition apart from the Armijo one is backtracking [61]. In a backtracking strategy the step-length is chosen in a more systematic way; instead of evaluating the objective function multiple times for various step-lengths, we start with a reasonable (but not small) initial value of $\bar{\alpha}$. If the condition is satisfied, the step-size α is accepted (that is, we set $\alpha_k = \bar{\alpha}$), and the algorithm proceeds to the next iteration. Otherwise, the step size is reduced by multiplying it by the factor $\rho < 1$, i.e., $\bar{\alpha} \leftarrow \rho \bar{\alpha}$. The process is then repeated multiple times until the Armijo condition (3) is met. Additionally, we incorporate the following condition along with (3):

$$\mathbf{p}_k^\top \nabla f_k \leq 0. \quad (6)$$

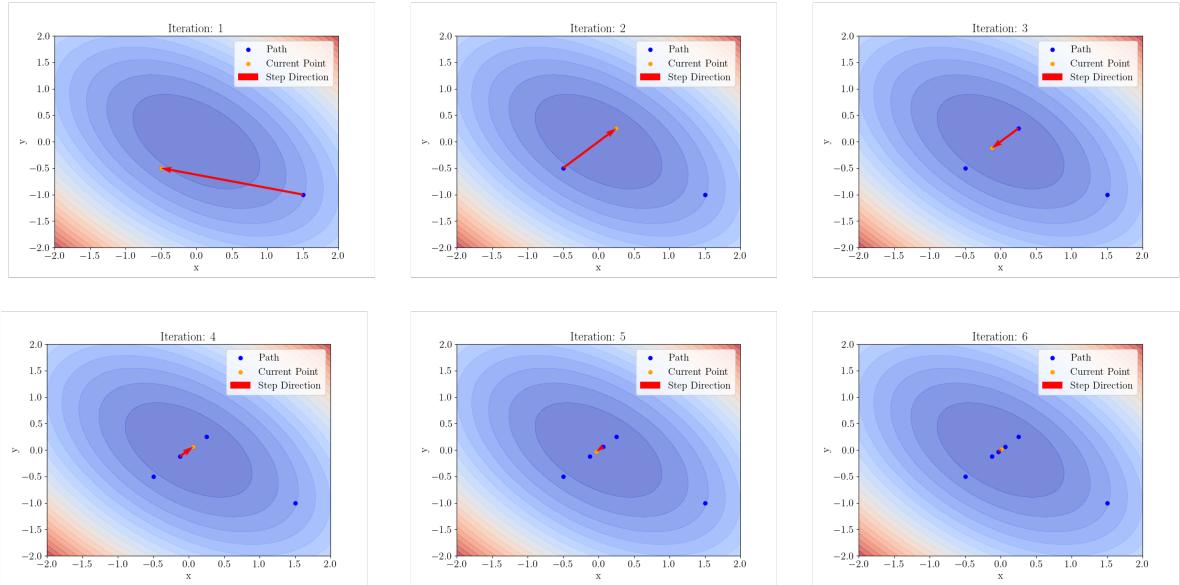


Figure 2: An example illustrating the minimization of a quadratic function, $f(x, y) = x^2 + y^2 + xy$, using the backtracking line-search algorithm is presented. The convergence trajectory is shown over six iterations. Notably, the search direction changes at each iteration, determined by evaluating the inequality condition. It is worth mentioning that backtracking achieves convergence to the minimum value of the level set faster than the trust-region approach.

1.2.2. Trust-region methods

Another class of optimization methods is the family of trust-region methods. The key distinction between trust-region and line-search methods is that the former defines a region (the trust region) around the current point \mathbf{x}_k and seeks a suitable point within this region that reduces the objective function. To achieve this, a trust-region algorithm first constructs a local quadratic approximation of the function, denoted as m_k , and then approximately solves the following subproblem:

$$\min_{\mathbf{p}} m_k(\mathbf{p}) = \nabla f_k^\top \mathbf{p} + \frac{1}{2} \mathbf{p}^\top \mathbf{B}_k \mathbf{p}$$

subject to:

$$\|\mathbf{p}\| \leq \Delta_k,$$

where Δ_k is the radius of the trust-region, and \mathbf{B}_k is some approximation of the Hessian matrix at \mathbf{x}_k . Then, the proposed step \mathbf{p}_k is evaluated. If it results in a significant improvement in the objective function, the step is accepted and the trust-region can be expanded. If the step yields poor results, it is rejected, and the trust-region is reduced. Finally, the iterate \mathbf{x}_k is updated simply as

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k,$$

that is, in a trust-region algorithm the direction and the step-length are calculated at the same time.

1.2.3. Quasi-Newton methods

A well-known example of a line-search direction is the steepest descent direction, which is given by the negative gradient at \mathbf{x}_k , that is,

$$\mathbf{p}_k = -\nabla f_k,$$

which ensures locally the steepest decrease in f at \mathbf{x}_k . Another well-known example is Newton's method, which stands out for its remarkable quadratic rate of convergence near the solution. By leveraging both first- and second-order derivatives of the function, Newton's method computes the following search direction \mathbf{p}_k by

$$\mathbf{p}_k = -(\nabla^2 f_k)^{-1} \nabla f_k,$$

where $\mathbf{H}_k \approx (\nabla^2 f_k)^{-1}$ is the Hessian matrix evaluated at the current iterate \mathbf{x}_k . In both cases, the search direction is then used to determine an appropriate step length α , ensuring a sufficient decrease in f .

Since Newton's search direction is derived from minimizing a local quadratic approximation of the objective function, it is highly efficient near the solution where the loss landscape is approximately convex and well-behaved. This quadratic model enables rapid convergence, often in just a few iterations, especially when the Hessian accurately captures local curvature. However, Newton's method requires computing and inverting the Hessian matrix, which becomes computationally prohibitive in high-dimensional settings, such as those encountered in deep learning or PINNs. Moreover, when the current iterate is far from the solution, the Hessian may be indefinite or poorly conditioned, leading to search directions that are not guaranteed to be descent directions and may even increase the loss. In contrast, the steepest descent (or gradient descent) method avoids the need for second-order information by using the negative gradient as the update direction. This makes each iteration relatively inexpensive and guarantees a descent direction at every step. However, its efficiency drastically deteriorates in the presence of *ill-conditioned*, which arises when the Hessian has eigenvalues that differ by several orders of magnitude. In such cases, the optimization trajectory may exhibit zigzagging behavior along the narrow valleys of the loss surface, causing slow convergence. This is particularly problematic in multi-objective or physics-informed settings, where different loss terms may induce conflicting gradients and anisotropic curvature. As a result, the steepest descent method often struggles in scenarios where curvature-aware updates, like those in quasi-Newton or second-order methods are essential for efficient and stable convergence [56].

To balance computational efficiency with convergence performance, quasi-Newton methods have emerged as an attractive alternative to full Newton's method. These approaches iteratively approximate the Hessian matrix using only first-order derivative information, thereby avoiding the explicit computation of second derivatives. Specifically, quasi-Newton methods achieve superlinear convergence by updating an approximation \mathbf{B}_k of the Hessian matrix, where $\mathbf{H}_k \equiv \mathbf{B}_k^{-1}$, using gradient information from previous iterations. These methods rely on update formulas, which ensure that the updated approximation satisfies the secant equation:

$$\mathbf{B}_{k+1} \mathbf{s}_k = \mathbf{y}_k, \tag{7}$$

where $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ represents the change in the iterates, and $\mathbf{y}_k = \nabla f_{k+1} - \nabla f_k$ represents the corresponding change in gradients. The specific form of the update rule defines the type of quasi-Newton method employed. Notable examples include the Symmetric Rank-One (SR1) and the widely used BFGS (Broyden–Fletcher–Goldfarb–Shanno) methods. Among these, the BFGS method stands out due to its ability to maintain symmetry and ensure positive definiteness of the Hessian approximation under mild conditions. This makes BFGS one of the most robust and commonly used quasi-Newton techniques. Once the search direction is determined—typically defined as

$$\mathbf{p}_k = -\mathbf{H}_k \nabla f_k, \tag{8}$$

a suitable step size α_k is selected (often via a line search strategy). The next iterate is then updated as $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$ followed by computing the \mathbf{s}_k and \mathbf{y}_k , and finally updating the Hessian approximation using the matrix \mathbf{B}_k as follows

$$\mathbf{B}_{k+1} = \mathbf{B}_k - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{s}_k^\top \mathbf{B}_k}{\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{y}_k^\top \mathbf{s}_k}. \tag{2.19}$$

To demonstrate the BFGS algorithm with the backtracking line search and trust-region, we present the convergence trajectory for a quadratic function defined as

$$f(x, y) = x^2 + y^2 + xy,$$

in Figures 2 and 3. We see that from iteration 1 to 6, the algorithm iteratively predicts the search direction. In this example, the trust-region radius is kept constant since $f(x)$ is convex. However, for PDE problems, we have implemented an adaptive trust-region size.

The choice between **backtracking line-search** and the **trust-region method** depends on the problem characteristics and specific implementation details. Backtracking line-search is often preferable for simpler or computationally inexpensive problems due to its straightforward nature and adaptability. In contrast, trust-region methods are more suitable for complex or constrained optimization tasks where careful control over the step size is essential. A key distinction is that backtracking line-search directly adjusts the step size without solving quadratic subproblems, making it computationally faster in cases where objective function and gradient evaluations are inexpensive. Ultimately, the choice between these methods depends on the specific problem and computational trade-offs.

The final piece of the algorithm is to choose the convergence criteria. The algorithm is considered to have converged if the norm of the gradient satisfies

$$\|\nabla f(\mathbf{x}_k)\| \leq \epsilon, \quad (9)$$

where $\|\nabla f(\mathbf{x}_k)\|$ is the Euclidean norm of the gradient at iteration k , $\epsilon > 0$ is a pre-specified tolerance level (a small positive value, e.g., 10^{-6}). The gradient of the objective function represents the direction and magnitude of the steepest ascent. At a local minimum, the gradient approaches zero. Thus, the norm of the gradient serves as a natural stopping criterion. A small tolerance (ϵ) improves precision but may increase computational cost. In some cases, the norm of the gradient stagnates near zero due to numerical precision issues, requiring additional criteria such as step size or objective function change thresholds.

In addition to using the gradient norm as a convergence criterion, the following supplementary criteria can be applied:

1. If $|f(\mathbf{x}_k) - f(\mathbf{x}_{k-1})| \leq \delta$, where δ is a small threshold.
2. If $\|\mathbf{x}_k - \mathbf{x}_{k-1}\| \leq \eta$, where η is a small tolerance.
3. Set a maximum number of iterations to prevent infinite loops in cases of slow convergence.

1.3. Broyden Family of Optimizers

Many quasi-Newton methods fall under the Broyden family of updates, characterized by the following general formula:

$$\mathbf{B}_{k+1} = \mathbf{B}_k - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{s}_k^\top \mathbf{B}_k}{\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{y}_k^\top \mathbf{s}_k} + \theta_k (\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k) \mathbf{w}_k \mathbf{w}_k^\top, \quad (10)$$

where θ_k is a scalar parameter that may vary at each iteration, and

$$\mathbf{w}_k = \frac{\mathbf{y}_k}{\mathbf{y}_k^\top \mathbf{s}_k} - \frac{\mathbf{B}_k \mathbf{s}_k}{\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k}.$$

The BFGS and DFP methods are special cases of the Broyden class. Specifically, setting $\theta_k = 0$ recovers the BFGS update, while $\theta_k = 1$ yields the DFP update. In practice, quasi-Newton methods typically work directly with \mathbf{H}_k , avoiding explicit matrix inversion. By applying the inverse to both sides of equation (10), and using the Sherman–Morrison formula, we obtain the following update for \mathbf{H}_k :

$$\mathbf{H}_{k+1} = \mathbf{H}_k - \frac{\mathbf{H}_k \mathbf{y}_k \mathbf{y}_k^\top \mathbf{H}_k}{\mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k} + \frac{\mathbf{s}_k \mathbf{s}_k^\top}{\mathbf{y}_k^\top \mathbf{s}_k} + \phi_k (\mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k) \mathbf{v}_k \mathbf{v}_k^\top, \quad (11)$$

where the intermediate quantities are defined as:

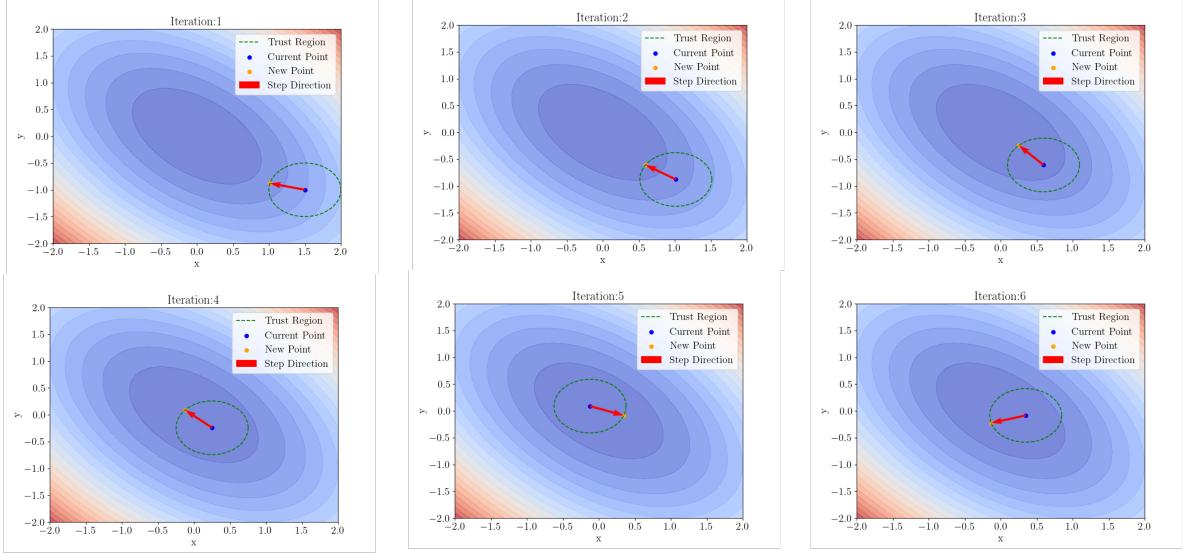


Figure 3: An example showing the minimization of a quadratic function $f(x, y) = x^2 + y^2 + xy$ using trust-region. The convergence is shown for six iteration. It is to be noted that at every iteration search direction changes computed by minimizing the line search criteria. In this example, we show that the radius of trust-region is constant as $f(x)$ is convex but for PDE problem we have implemented the daptive size of trust-region.

$$\begin{aligned} \mathbf{v}_k &= \frac{\mathbf{s}_k}{\mathbf{y}_k^\top \mathbf{s}_k} - \frac{\mathbf{H}_k \mathbf{y}_k}{\mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k}, \\ \phi_k &= \frac{1 - \theta_k}{1 + (h_k b_k - 1)\theta_k}, \\ b_k &= \frac{\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k}{\mathbf{y}_k^\top \mathbf{s}_k}, \\ h_k &= \frac{\mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k}{\mathbf{y}_k^\top \mathbf{s}_k}. \end{aligned}$$

Additional useful expressions include:

$$\begin{aligned} a_k &= b_k h_k - 1, \\ c_k &= \left(\frac{a_k}{1 + a_k} \right)^{1/2}, \\ \rho_k^- &= \min(1, h_k(1 - c_k)), \\ \theta_k^- &= \frac{\rho_k^- - 1}{a_k}, \\ \theta_k^+ &= \frac{1}{\rho_k^-}, \\ \theta_k &= \max \left(\theta_k^-, \min \left(\theta_k^+, \frac{1 - b_k}{b_k} \right) \right), \\ \rho_k^+ &= \min \left(1, \frac{1}{b_k} \right), \\ \sigma_k &= 1 + \theta_k a_k, \\ \sigma_k^{(1-N)} &= |\sigma_k|^{\frac{1}{1-N}}, \end{aligned}$$

$$\tau_k = \begin{cases} \min\left(\rho_k^+ \sigma_k^{(1-N)}, \sigma_k\right), & \text{if } \theta_k \leq 0, \\ \rho_k^+ \min\left(\sigma_k^{(1-N)}, \frac{1}{\theta_k}\right), & \text{otherwise.} \end{cases}$$

If θ_k depends explicitly on \mathbf{B}_k , then both equations (10) and (11) must be used at each iteration to update the inverse Hessian estimates. However, when \mathbf{B}_k appears only through the product $\mathbf{B}_k \mathbf{s}_k$, the update can avoid direct use of (10), since:

$$\mathbf{B}_k \mathbf{s}_k = -\alpha_k \nabla f_k. \quad (12)$$

Among the methods in the Broyden class, the BFGS method is particularly effective for small- and medium-scale unconstrained optimization problems [62, 63]. However, its performance may deteriorate for ill-conditioned problems [64, 65]. To address this limitation, the *self-scaled* BFGS (SSBFGS) method was proposed by Oren and Luenberger [66]. In SSBFGS, the Hessian approximation \mathbf{B}_k is scaled by a positive scalar τ_k prior to the BFGS update, i.e.,

$$\mathbf{B}_{k+1} = \tau_k \left[\mathbf{B}_k - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{s}_k^\top \mathbf{B}_k}{\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k} \right] + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{y}_k^\top \mathbf{s}_k}. \quad (13)$$

The motivation behind scaling is to reduce the condition number of $\mathbf{H}_k^{1/2} \nabla^2 f_k \mathbf{H}_k^{1/2}$, which reflects the convergence rate. While initial results by Nocedal and Yuan [67] were discouraging for the scaling $\tau_k = 1/b_k$, later work by Al-Baali [68] showed promising results using the modified scaling $\tau_k = \min\{1, 1/b_k\}$. Al-Baali [69] extended the idea of scaling to other Broyden updates, proving favorable theoretical and numerical results for various $\theta_k \in [0, 1]$, while previously assuming $\tau_k \leq 1$. Additional scaling strategies have been developed for BFGS and other Broyden methods in works such as [70, 71, 72, 73, 74]. The self-scaled versions of the direct and inverse updates for general Broyden family methods—termed *SSBroyden* methods—are:

$$\mathbf{B}_{k+1} = \tau_k \left[\mathbf{B}_k - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{s}_k^\top \mathbf{B}_k}{\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k} + \theta_k (\mathbf{s}_k^\top \mathbf{B}_k \mathbf{s}_k) \mathbf{w}_k \mathbf{w}_k^\top \right] + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{y}_k^\top \mathbf{s}_k}, \quad (14)$$

$$\mathbf{H}_{k+1} = \frac{1}{\tau_k} \left[\mathbf{H}_k - \frac{\mathbf{H}_k \mathbf{y}_k \mathbf{y}_k^\top \mathbf{H}_k}{\mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k} + \phi_k (\mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k) \mathbf{v}_k \mathbf{v}_k^\top \right] + \frac{\mathbf{s}_k \mathbf{s}_k^\top}{\mathbf{y}_k^\top \mathbf{s}_k}. \quad (15)$$

In the context of Physics-Informed Neural Networks (PINNs) based on Multi-Layer Perceptrons (MLPs), recent research has begun to explore various families of SSBroyden updates and modifications to the loss formulation to address known challenges such as gradient conflicts and ill-conditioning. These second-order strategies significantly enhance both convergence speed and prediction accuracy compared to classical BFGS. In this work, we extend these findings to more complex problems beyond those presented in [53], employing not only SSBroyden but also other techniques proposed in the PINN literature. Additionally, we investigate whether these BFGS modifications enhance the performance of the recently introduced Kolmogorov-Arnold Networks (KANs) [3]. This study will further enable fair comparisons between KANs and MLPs using more advanced optimization algorithms. In the following section, we present a comprehensive study comparing the performance of different optimization methods with various line-search and trust-region strategies for solving PDEs. For each example, we analyze and discuss the performance of both PINNs and PIKANs.

2. Computational Experiments

In this section, we conduct various computational experiments to evaluate the performance of optimizers across a diverse range of steady-state and time-dependent PDEs. To ensure a comprehensive assessment, we select PDEs that encompass a wide spectrum of classes, including parabolic, hyperbolic, elliptic, and hyperbolic-parabolic equations. It is worth noting that the computations for BFGS and SSBroyden with Wolfe line-search are performed on an NVIDIA RTX A6000 GPU, while the computations for BFGS with backtracking and trust-region methods are conducted on an NVIDIA RTX 3090 GPU.

2.1. Burgers equation

To begin with, we consider the viscous Burgers' equation in a spatially periodic domain, given by,

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = \nu \frac{\partial^2 u}{\partial x^2}, \quad (16)$$

with an initial condition

$$u(x, 0) = -\sin(\pi x) \quad (17)$$

and periodic boundary conditions. The viscosity is given by $\nu = \frac{0.01}{\pi} \approx 0.003$, and the spatio-temporal domain for computing the solution is $(x, t) \in [-1, 1] \times [0, 1]$.

2.1.1. BFGS and SSBroyden optimization with Strong Wolfe line-search

A comparative study of the performance of L-BFGS, BFGS, SSBFGS and SSBroyden optimization with Strong Wolfe line-search for solving the Burgers' Equation has been conducted using seven different case studies. In all PINNs cases, the hidden layers utilize the hyperbolic tangent (\tanh) activation function. Model training is performed in two stages: first, the Adam optimizer is employed with a learning rate of 10^{-3} for 1000 iterations. Moreover, the adaptive sampling strategy (RAD) algorithm defined in [75] is employed to dynamically resample points in regions with high errors.

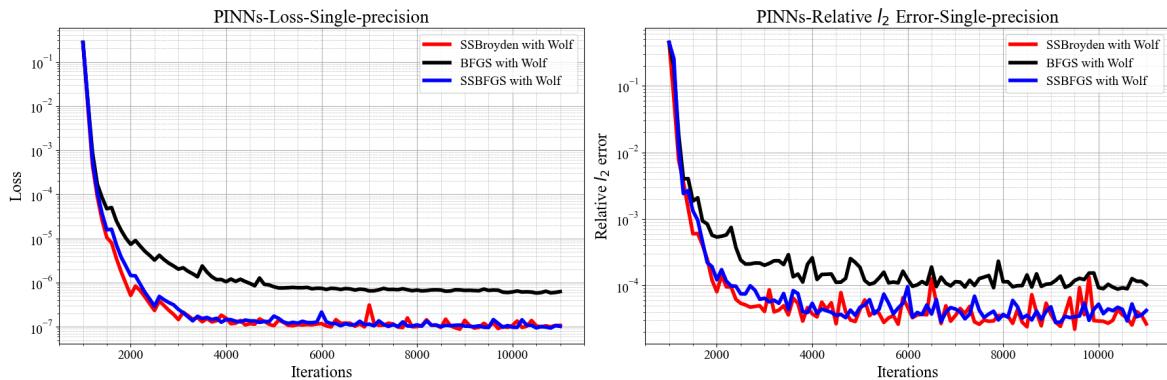


Figure 4: **PINNs with single precision for the Burgers equation:** The evolution of the loss function (left) and corresponding l_2 relative errors (right) over iterations for **Case 1**. After approximately 4,000 iterations, the loss function and relative error stabilize, showing no significant further improvement.

Case	Optimizer[# Iter., Line-search [#Iter.]	Relative l_2 error	Training time (s)	Total params
1	Adam [1000] + BFGS with Wolfe [10000]	1.04×10^{-4}	263	1,341
1	Adam [1000] + SSBroyden with Wolfe [10000]	3.60×10^{-5}	209	1,341
1	Adam [1000] + SSBFGS with Wolfe [10000]	4.04×10^{-5}	179	1,341

Table 1: **PINNs with single precision for the Burgers equation:** Relative L2 error and training time for solving the Burgers equation using **single precision** with different optimizers. A PINN with four layers, each containing 20 neurons, is trained for 1000 iterations using the Adam optimizer, followed by 10,000 iterations with the specified optimizer.

Figures 4 illustrate the loss function and relative error over iterations for **Case 1** in single precision. The results are obtained using a PINN with four layers, each comprising 20 neurons, trained with the Adam optimizer followed by 10,000 iterations of the specified optimizer. Table 1 provides a summary of the corresponding relative error and training time. Notably, SSBFGS and SSBroyden achieve a relative error of 10^{-5} , while BFGS converges to approximately 10^{-4} . It is observed that the loss function and relative error stabilize after roughly 4,000 iterations, indicating no significant further improvement.

The performance of PINNs with a similar network architecture and double precision is presented in Table 2 for **Case 2** and **Case 3**. The same network architecture is employed, consisting of four hidden layers, each with 20 neurons. After the initial Adam optimization stage, the training proceeds with BFGS, SSBFGS, and SSBroyden optimizers for two different iteration counts: **Case 2** utilizes 50,000 iterations, while **Case 3** uses 30,000 iterations. In **Case 4**, the same network structure as in **Case 2** and **Case 3** is used; however, training begins directly with second-order optimizers, BFGS, SSBFGS, and SSBroyden, instead of the Adam optimizer. The results, as shown in the table, demonstrate that increasing the iteration count to 50,000 significantly raises the training time. While this reduces

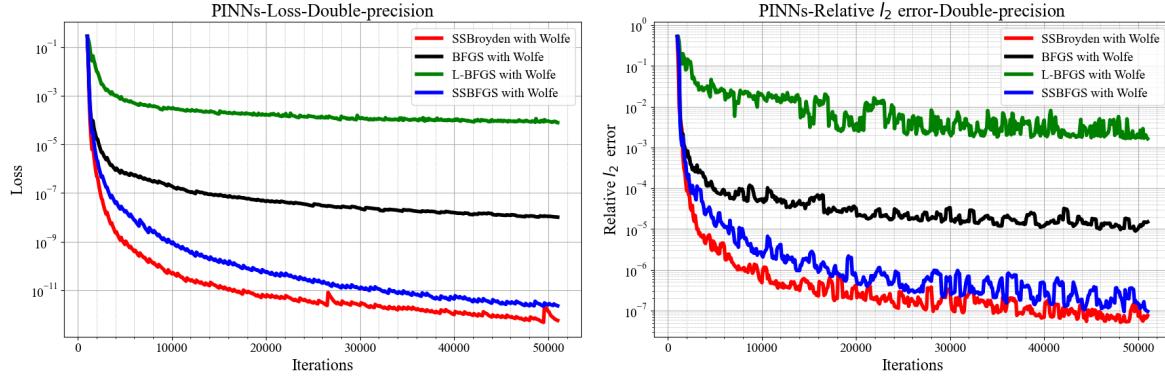


Figure 5: **Double-Precision PINNs for the Burgers equation:** The evolution of the loss function (left) and corresponding l_2 relative errors (right) over iterations for **Case 2**. After approximately 4,000 iterations, the loss function and relative error stabilize, showing no significant further improvement.

Case	Optimizer[# Iter., Line-search [#Iter.]]	Relative l_2 error	Training time (s)	Total params
2	Adam [1000] + BFGS with Wolfe [50000]	1.50×10^{-5}	1292	1,341
2	Adam [1000] + SSBroyden with Wolfe [50000]	7.57×10^{-8}	1354	1,341
2	Adam [1000] + SSBFGS with Wolfe [50000]	9.62×10^{-8}	1293	1,341
2	Adam [1000] + L-BFGS with Wolfe [50000]	2.05×10^{-3}	713	1,341
3	Adam [1000] + BFGS with Wolfe [30000]	2.21×10^{-5}	774	1,341
3	Adam [1000] + SSBroyden with Strong Wolfe [30000]	2.12×10^{-7}	819	1,341
4	BFGS with Wolfe [30000]	1.73×10^{-5}	863	1,341
4	SSBFGS with Wolfe [30000]	1.54×10^{-7}	711	1,341
4	SSBroyden with Wolfe [30000]	3.59×10^{-7}	862	1,341
5	Adam [1000] + BFGS with Wolfe [30000]	8.52×10^{-6}	3049	3,021
5	Adam [1000] + SSBroyden with Wolfe [30000]	1.62×10^{-8}	2812	3,021
5	Adam [1000] + SSBFGS with Wolfe [30000]	4.19×10^{-8}	2878	3,021

Table 2: **Double-Precision PINNs for the Burgers equation: Case 2:** A PINN with four layers, each containing 20 neurons, is trained for 1000 iterations using the Adam optimizer, followed by 50,000 iterations with BFGS, SSBroyden, SSBFGS, and L-BFGS optimizers with Strong Wolfe line-search. **Case 3:** The same network structure as **Case 2** is used, but the number of iterations with BFGS, and SSBroyden is reduced to 30,000. **Case 4:** The same network structure as in **Case 2** and **Case 3** is used; however, instead of starting with the Adam optimizer, the training begins directly with second-order optimizers, namely BFGS, SSBFGS, and SSBroyden. **Case 5:** A deeper PINN with eight layers (20 neurons per layer) is trained for 1000 iterations using Adam, followed by 30,000 iterations with BFGS and SSBroyden.

SSBFGS and SSBroyden's error to 10^{-8} , it does not result in any noticeable improvement in the error achieved by BFGS. In **Case 5**, a deeper network with eight hidden layers is used, trained for 30,000 iterations. This configuration further improves the error, though at the cost of a substantial increase in training time. Figure 5 illustrates the loss function and relative error over iterations for **Case 2** using four optimizers, BFGS, SSBFGS, SSBroyden, and L-BFGS. As summarized in Table 2, the relative errors achieved are 10^{-8} for SSBroyden, 10^{-5} for BFGS, and 10^{-3} for L-BFGS.

Additionally, Figure 6 compares the performance of BFGS and SSBroyden for PIKANs in **Case 6**. A PIKAN with four layers, each containing 20 neurons, is trained using Chebyshev polynomials of degree 3. As reported in Table 3, the error in this case is very close to that of **Case 2** in Table 2. Furthermore, a PIKAN with three layers, each containing 10 neurons, is trained using Chebyshev polynomials of degree 3 in **Case 7** and degree 5 in **Case 8**. The results demonstrate that SSBroyden consistently outperforms BFGS. Although the networks are trained for 50,000 iterations, minimal improvement is observed beyond approximately 30,000 iterations. As summarized in Table 3, increasing the polynomial degree from 3 to 5 does not improve the results, while the training time nearly doubles with the higher degree.

The conclusion highlights that the best error for the Burgers equation was achieved using SSBroyden with double precision. Across all seven case studies, SSBroyden consistently outperformed both BFGS and L-BFGS. A comparison between PINNs and PIKANs reveals that, for the same number of parameters, PIKANs did not perform as well. However, when trained with a higher number of parameters—at the cost of significantly increased training time—PIKANs can achieve results comparable to PINNs. In general, a comparison between **Case 2** and **Case 6** shows that, although both achieve

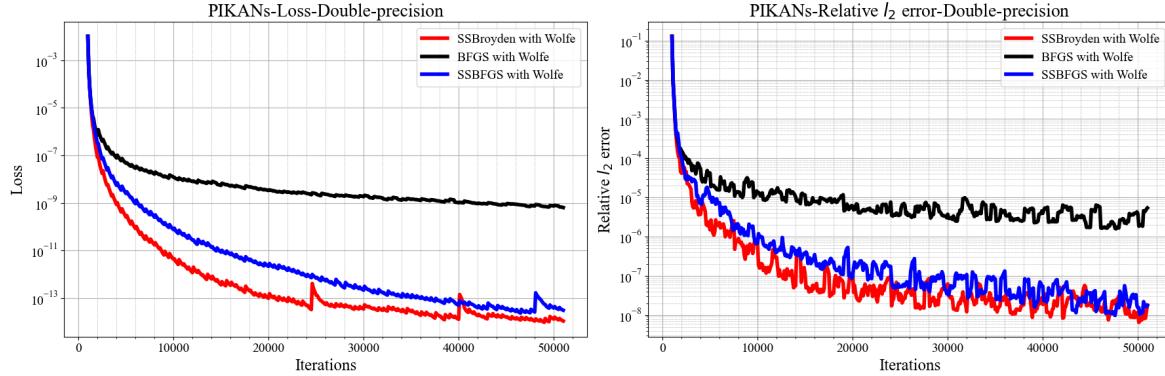


Figure 6: **Double-Precision PIKANs for the Burgers Equation:** Evolution of the loss function (left) and l_2 relative errors (right) for Case 6. The results compare the performance of PIKANs trained with SSBroyden and BFGS optimizers.

Case	Optimizer[# Iter., Line-search #Iter.]	Relative l_2 error	Training time (s)	Total params
6	Adam [1000] + BFGS with Wolfe [50000] (degree 3)	5.38×10^{-6}	9276	5,040
6	Adam [1000] + SSBFGS with Wolfe [50000] (degree 3)	1.77×10^{-8}	12796	5,040
6	Adam [1000] + SSBroyden with Wolfe [50000] (degree 3)	1.79×10^{-8}	14169	5,040
7	Adam [1000] + BFGS with Wolfe [50000] (degree 3)	3.18×10^{-5}	891	920
7	Adam [1000] + SSBroyden with Wolfe [50000] (degree 3)	2.41×10^{-6}	849	920
7	Adam [1000] + SSBFGS with Wolfe [50000] (degree 3)	4.21×10^{-6}	810	920
8	Adam [1000] + BFGS with Wolfe [50000] (degree 5)	4.29×10^{-5}	1364	1,380
8	Adam [1000] + SSBroyden with Wolfe [50000] (degree 5)	1.08×10^{-6}	1465	1,380

Table 3: **Double-Precision PIKANs for the Burgers Equation:** Relative l_2 error and training time for solving the Burgers equation using **double precision** with different optimizers. **Case 6** highlights a PIKAN with four layers of 20 neurons each and degree 3. A PIKAN with three layers, each containing 10 neurons, is trained using Chebyshev polynomials of degree 3 in **Case 7** and degree 5 in **Case 8**. The training consists of 1000 iterations with the Adam optimizer, followed by 50,000 iterations with BFGS and SSBroyden optimizers. **Case 6** demonstrates improved performance despite a significant increase in training time.

nearly identical errors, PINNs with SSBroyden demonstrate more efficient training times compared to PIKANs.

2.1.2. Performance of optimizers based on composition operator in the JAX ecosystem

Approximating solutions of PDEs with neural networks involves randomly initializing the network parameters using a normal probability distribution [76]. The random sampling of the parameters depends on the random seeds, which also defines the reproducibility of the results. The framework presented above is based on Tensorflow, which generates the random number using a stateful random number generators. A stateful generator has the following three characteristics:

1. It maintains an internal state that is updated after each random number generation.
2. The next random number depends on the current internal state.
3. The generator needs to be reseeded explicitly to reset or reproduce results.

Therefore, the main drawback with stateful random number generators can result in non-reproducibility as the internal state can cause issues when working in parallel systems. Secondly, it requires careful management in multithreaded contexts to avoid race conditions, which is very common in GPU based architectures, as stateful random number generator consider previous state for future number generation. To overcome these issues, the JAX framework [77] is based on a stateless random number generator which does not maintain any internal state between calls. Instead, it uses the provided input (such as a seed or key) to generate each random number independently, without relying on previous outputs. The state must be explicitly passed along with the call, and the RNG (Random Number Generator) does not store any internal data between calls. Especially, it becomes essential when true randomness or independent streams of numbers are required, e.g., parallel random number generation, GPU based programming models such as MPI + X ∈ OpenMP, CUDA, HIP, etc.

Therefore, to show the reproducibility of the proposed optimization method, in this paper we develop the same code in the JAX framework with the Optax library [79]. Optax is a library of

```

import optimistix as opx #← Rader et al. (2024)[74]
class BFGSTrustRegion(opx.AbstractBFGS):
    rtol: float
    atol: float
    norm: Callable = opx.max_norm # ← Maximum norm for convergence
    use_inverse: bool = False
    search: opx.AbstractSearch = opx.LinearTrustRegion() # ← Line search algorithm
    descent: opx.AbstractDescent = opx.NewtonDescent() # ← Newton descent algorithm
    verbose: frozenset[str] = frozenset()

```

Figure 7: A Python class for BFGS optimizers paired with trust-region line search and Newton descent. Equation (2.19) with L_∞ norm chosen as convergence criteria.

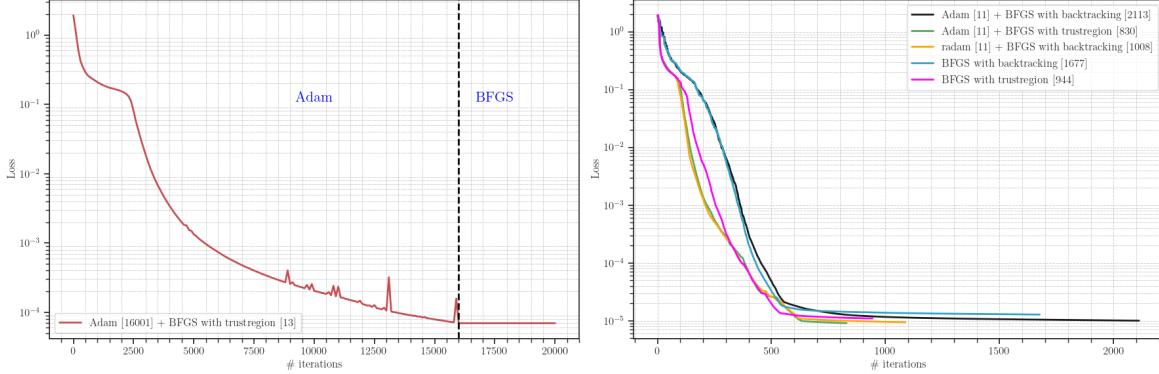


Figure 8: Convergence of PINNs for the viscous Burgers equation under various combinations of line search and BFGS algorithms for single-precision arithmetic (32-bit floats) with GPU fully saturated with model parameters to avoid the effect of latency on compute time: The left panel illustrates a scenario where the iteration count of the first-order Adam algorithm dominates, combined with BFGS using the trust-region approach. This pairing exhibits degeneracy due to limited numerical precision during BFGS iterations. The right panel presents the convergence behavior of BFGS paired with trust-region and backtracking line search methods. In this setup, Adam is applied for 11 iterations as a warmup. The BFGS-based optimizer requires significantly more iterations than first-order optimization methods like Adam and Rectified Adam (RAdam) [78].

Optimizer[# Iter., Line-search algorithm [#Iter.]]	Relative l_2 error	Training time (s)	Total params
Adam [16001] + BFGS with trust-region [13]	2.91×10^{-3}	127	3501
Adam [11] + BFGS with backtracking [2113]	7.25×10^{-4}	61	3501
Adam [11] + BFGS with trust-region [830]	6.42×10^{-4}	56	3501
RAdam [11] + BFGS with backtracking [1008]	7.67×10^{-4}	62	3501
BFGS with backtracking [1677]	2.48×10^{-3}	45	3501
BFGS with trust-region [944]	1.39×10^{-3}	49	3501

Table 4: [Corresponds to Figure 8] Performance metrics for the viscous Burgers equation using single-precision (32-bit) arithmetic for cases shown in Figure 8: The relative L_2 error, training time, and number of training parameters are evaluated for solving the viscous Burgers equation with various optimizers and combinations of line search algorithms. The convergence criteria are based on absolute and relative tolerances in norm of the gradient of the loss function, which is set as $[ATOL, RTOL] = [10^{-7}, 10^{-7}]$. This experiments shows that Adam combined with BFGS trust-region provides better accuracy (highlighted in bold fonts).

gradient transformations paired with composition operators (e.g., chain) that allow implementing many standard and new optimisers (e.g., RMSProp [80], Adam [81], Lion [82], RAdam [78], etc.) in just a single line of code. The compositional nature of Optax naturally supports recombining the same basic ingredients in custom optimisers. As we elaborated earlier, BFGS paired with a line search algorithm consists of the following components:

1. Relative tolerance: atol
2. Absolute tolerance: rtol
3. norm: L_1 , L_2 , L_∞ norm of gradient for convergence criteria.
4. Line search direction: Backtracking, Trust-region
5. Newton descent: Equation (2.19)

To customize an instance of the BFGS optimizer, the choices for items (1)–(5) should be made based on the specific problem being addressed. For this purpose, we employed the state-based Optimistix

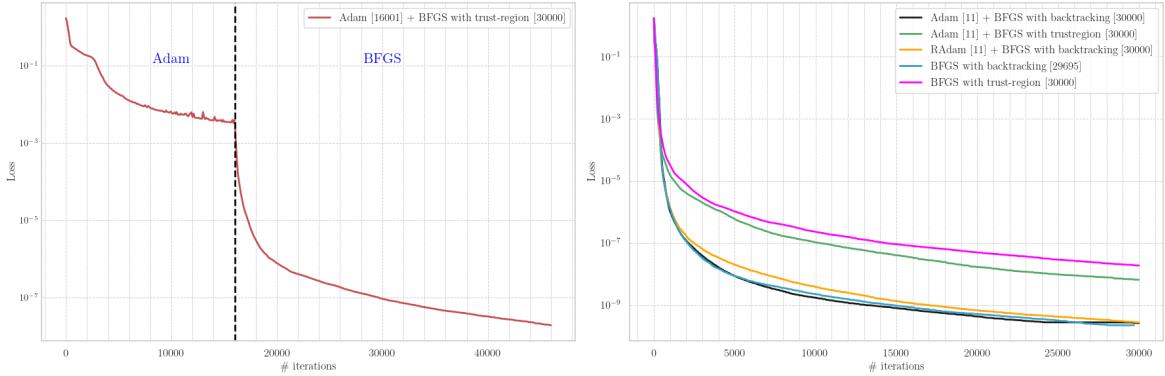


Figure 9: Convergence of PINNs for the viscous Burgers equation under various combinations of line search and BFGS algorithms for double-precision arithmetic (64-bit floats) with GPU fully saturated with model parameters to avoid the effect of latency on compute time: The left panel illustrates a scenario where the iteration count of the first-order Adam algorithm dominates, and combined with BFGS using the trust-region approach. This pairing unlike Figure 8 does not exhibit degeneracy during BFGS iterations due to double-precision and converges to much better loss value. The right panel presents the convergence behavior of BFGS paired with trust-region and backtracking line search methods. Inline to the Figure 8, in this setup as well, Adam is applied for 11 iterations as a warmup and then switched to BFGS algorithms. The BFGS-based optimizer requires significantly more iterations than first-order optimization methods like Adam and Rectified Adam (RAdam) [78].

Optimizer[# Iter., Line-search algorithm [#Iter.]]	Relative l_2 error	Training time (s)	Total params
Adam [16001] + BFGS with trust-region [30000]	2.8×10^{-5}	2030	3501
Adam [11] + BFGS with backtracking [30000]	3.0×10^{-6}	1607	3501
Adam [11] + BFGS with trust-region [30000]	1.3×10^{-5}	1726	3501
RAdam [11] + BFGS with backtracking [30000]	2.0×10^{-6}	1587	3501
BFGS with backtracking [29695]	4.0×10^{-6}	1555	3501
BFGS with trust-region [30000]	1.6×10^{-5}	1722	3501

Table 5: [Corresponds to Figure 9] Performance metrics for the viscous Burgers equation using double-precision (64-bit) arithmetic for cases shown in Figure 9: The relative L_2 error, training time, and number of training parameters are evaluated for solving the viscous Burgers equation with various optimizers and combinations of line search algorithms. The convergence criteria are based on absolute and relative tolerances in norm of the gradient of the loss function and number of iterations, and set as $[ATOL, RTOL] = \{[10^{-8}, 10^{-8}] \parallel (\# \text{ of steps} = 30000)\}$. This experiments shows that RAdam combined with BFGS with backtarcking linesearch provides better accuracy (highlighted in bold fonts). The computation is performed on Nvidia-GPU Card RTX-3090 with persistence memory usage of 76 % (Total 30 GB) and compute usage of 99% i.e. 550.44 GFLOPs of total theoretical peak of 556 GFLOPs.

library [83], which integrates seamlessly with the Optax library, to implement various variants of the BFGS optimizer by inheriting the `AbstractBFGS` class. These variants are tailored according to the selected line search algorithms, tolerance levels, convergence criteria, and types of Newton descent methods. For example, the BFGS optimizer with a trust-region line search algorithm and L_∞ -norm convergence is shown in code listing Figure 7. To evaluate the efficiency of the optimizers in terms of both accuracy and runtime on a GPU, we performed computational experiments in two distinct scenarios. In the first scenario, GPU saturation is ensured by setting the model parameters to a sufficiently large value, thus avoiding latency from impacting compute time. In the second scenario, saturation is achieved by increasing the number of collocation points.

In Scenario 1, the convergence history for single-precision (32-bit float) arithmetic is presented in Figure 8. This convergence is achieved using a neural network with 8 hidden layers, each containing 20 neurons, a tanh activation function, 200 randomly sampled data points for the initial and boundary conditions, and 10,000 collocation points to calculate the residual loss. The performance of five different optimizer combinations is summarized in Table 4. The left panel of Figure 8 shows a scenario where the iteration count of the first-order Adam algorithm dominates, combined with BFGS using the trust-region line search. The key observation from this setup is that when Adam is run for an extended period, the BFGS optimizer provides no additional advantages due to the degeneration of the loss function caused by single-precision floating-point representation [Row 1 of Table 4]. However, when Adam or RAdam is used with a small number of iterations (11 in this case) as a warm-up followed by BFGS with both line search algorithms, the error improves by an order of magnitude [Rows 2, 3, and 4 of Table 4]. In contrast, starting directly with the BFGS optimizer results in errors of a similar magnitude to those in the first scenario (left panel of 8). All runs in Table 4 are terminated once the

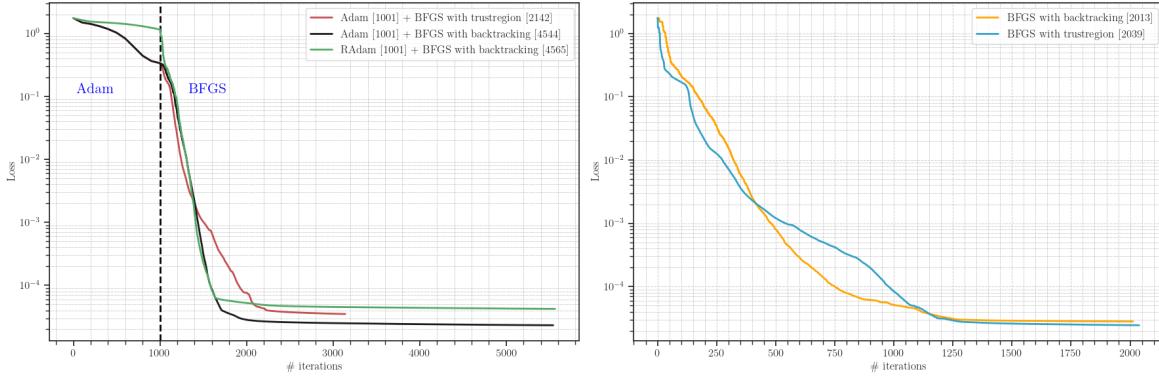


Figure 10: Convergence of PINNs for the viscous Burgers equation under various combinations of line search and BFGS algorithms for single-precision arithmetic (32-bit floats) with the GPU fully saturated by collocation points and data points (initial and boundary conditions) to mitigate the impact of latency on computation time: The left panel depicts a scenario where the iteration count for the first-order Adam and RAdam optimizers is fixed at 1001, followed by the second-order BFGS optimizer using the trust-region and backtracking line search algorithms. Similar to the cases shown in Figure 8, this setup also exhibits degeneracy during BFGS iterations due to limited numerical precision. The right panel shows the convergence behavior of the standalone BFGS optimizer with trust-region and backtracking line search methods, without any warmup by a first-order optimizer. In this configuration, the loss stagnates at 10^{-4} due to insufficient precision, particularly for the Hessian matrix inversion required for descent.

Optimizer[# Iter., Line-search algorithm [#Iter.]]	Relative l_2 error	Training time (s)	Total params
Adam [1001] + BFGS with trust-region [2142]	1.51×10^{-3}	30	1341
Adam [1001] + BFGS with backtracking [4544]	1.41×10^{-3}	43	1341
radam [1001] + BFGS with backtracking [4565]	1.43×10^{-3}	41	1341
BFGS with backtracking [2013]	1.65×10^{-3}	18	1341
BFGS with trust-region [2039]	1.1×10^{-3}	24	1341

Table 6: [Corresponding to Figure 10] Performance metrics for the viscous Burgers equation using single-precision (32-bit) arithmetic for the cases shown in Figure 10: The relative L_2 error, training time, and number of training parameters are reported for solving the viscous Burgers equation with various optimizers and combinations of line search algorithms. The convergence criteria are defined by absolute and relative tolerances on the gradient norm of the loss function, set as $[\text{ATOL}, \text{RTOL}] = [10^{-7}, 10^{-7}]$. The results demonstrate that BFGS with the trust-region line search achieves the highest accuracy (highlighted in bold) among all cases presented in Figure 10.

absolute and relative tolerances of $[\text{ATOL}, \text{RTOL}] = [10^{-7}, 10^{-7}]$ for the L_∞ norm of the gradient of the loss function are met. Consequently, the number of iterations reported for the BFGS optimizers corresponds to the step where the specified tolerance is reached. The runtime for each case is also recorded in Table 4. For all runs, the runtime is less than one minute, except for the run where Adam is used for a larger number of iterations Figure (8).

To address the issue of loss degeneration and demonstrate that it is caused by precision limitations, all cases presented in Figure 8 and Table 4 are rerun using double-precision arithmetic (64 bits float) but with same hyperparameters used for 8. The convergence criteria are based on absolute and relative tolerances in norm of the gradient of the loss function and number of iterations, and set as $[\text{ATOL}, \text{RTOL}] = \{[10^{-8}, 10^{-8}] \parallel (\# \text{ of steps} = 30000)\}$. The convergence plot with Adam having 16001 iterations and then using BFGS paired with trust-region linesearch algorithm is shown in left subfigure of Figure 9. Unlike single precision, the BFGS method achieves better accuracy in this case, converging to an L_2 error of 2.8×10^{-5} . Convergence results for other scenarios, similar to those in Figure 8, are displayed in the right subfigure of Figure 9. Performance metrics for these runs are provided in Table 5. Notably, the degeneration of the loss function is not observed and no longer impacts the convergence of the BFGS optimizer for any of the runs. Furthermore, relative L_2 errors for all the cases are reduced by two orders of magnitude compared to single-precision performance. For the remaining cases, the relative L_2 error is reduced to a minimum of 2×10^{-6} . However, runs in Table 5 with 30,000 BFGS iterations indicate that the optimizer failed to converge, likely due to the increased sensitivity of the tolerance limit in double precision.

In Scenario 2, where the GPU is fully utilized with collocation points and the PINN employs a smaller neural network, the convergence history for single and double precision arithmetic is shown in Figures 10 and 11, respectively. Convergence was achieved using a neural network with four hidden layers, each containing 20 neurons, a tanh activation function, 250 randomly sampled data points for

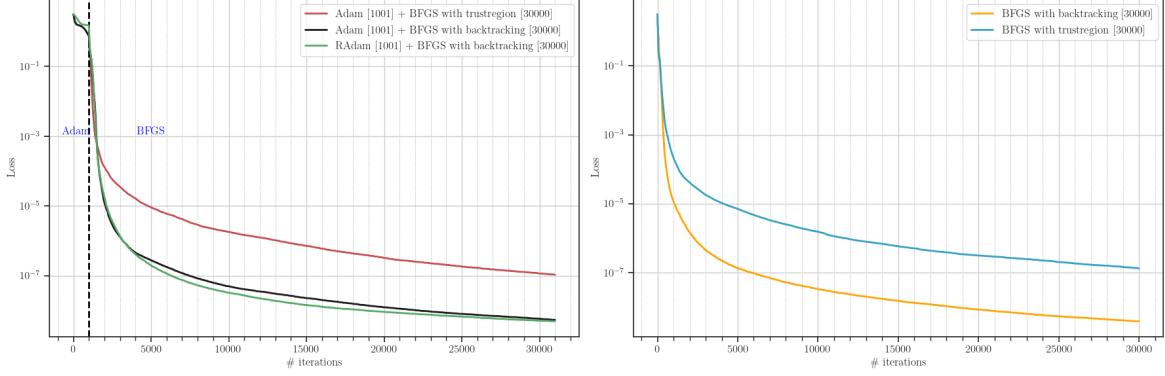


Figure 11: Convergence of PINNs for the viscous Burgers equation under various combinations of line search and BFGS algorithms for double-precision arithmetic (64-bit floats) with the GPU fully saturated by collocation points and data points (initial and boundary conditions) to minimize the impact of latency on computation time: The left panel illustrates a case where the iteration count for the first-order Adam and RAdam optimizers is fixed at 1001. Unlike Figure 10, this combination does not exhibit degeneracy during BFDS iterations due to the use of double-precision arithmetic and achieves a significantly lower loss value. The right panel presents the convergence behavior of the standalone BFGS optimizer with trust-region and backtracking line search methods, without warmup from a first-order optimizer. Notably, algorithms employing the backtracking line search method demonstrate better convergence rates and higher accuracy.

Optimizer[# Iter., Line-search [#Iters.]	Relative l_2 error	Training time (s)	Total params
Adam [1001] + BFDS with trust-region [30000]	4.4×10^{-5}	1141	1341
Adam [1001] + BFDS with backtracking [30000]	8×10^{-6}	1371	1341
RAdam [1001] + BFDS with backtracking [30000]	6×10^{-6}	1070	1341
BFGS with backtracking [30000]	9×10^{-6}	855	1341
BFGS with trust-region [30000]	4.9×10^{-5}	1423	1341

Table 7: [Corresponding to Figure 11] Performance metrics for the viscous Burgers equation using single-precision (64-bit) arithmetic for the cases shown in Figure 11: The relative L_2 error, training time, and number of training parameters are reported for solving the viscous Burgers equation with various optimizers and combinations of line search algorithms. The convergence criteria here is based on union of absolute and relative tolerance of norm of gradient and considered as: $[\text{ATOL}, \text{RTOL}] = \{[10^{-7}, 1 \times 10^{-7}] \parallel (\# \text{ of steps} = 30000)\}$. The results demonstrate that RAdam paired with BFDS with the backtracking line search achieves the highest accuracy (highlighted in bold) among all cases presented in Figure 11. The computation is performed on Nvidia-GPU Card RTX-3090 with persistence memory usage of 75 % (Total 24.56 GB) and compute usage of 99% i.e. 429 GFLOPs of total theoretical peak of 433.9 GFLOPs.

initial and boundary conditions, and 50,000 collocation points. The performance metrics, including the relative L_2 error and training time for single and double precision, are summarized in Tables 6 and 7, respectively. The interpretation of these results is similar to that of Scenario 1, though the overall accuracy is slightly reduced. However, in double precision, the BFDS optimizer failed to converge within 30,000 iterations, likely due to the limited representational capacity of the neural network.

2.2. Takeaways from 2.1

In Table 8, we provide our assessment of the most effective optimizer for PINN representation with double-precision arithmetic. The evaluation compares Wolfe and backtracking line searches alongside the trust-region approach, implemented in both TensorFlow and the JAX ecosystem. **Case 1*** corresponds to SS-Broyden with Wolfe line search, implemented in TensorFlow using double precision, while **Case 2*** refers to BFDS with backtracking line search, implemented in JAX.

Case	Optimizer	Training time (s)	Relative l_2 error	Total parameters, Verdict
1*	Adam [1000] + SS-Broyden with Wolfe [30000]	2812	1.62×10^{-8}	3,021, Winner
2*	RAdam [11] + BFDS with backtracking [30000]	1587	2×10^{-6}	3,501, Runner-up

Table 8: **Verdict on the choice of optimizer in double precision:** The winner and runner-up are determined based on the best relative L_2 error and training time for different optimizers applied to Burgers' equation. **Case 1*** represents SS-Broyden with Wolfe line search, implemented in TensorFlow using double precision, while **Case 2*** corresponds to BFDS with backtracking line search, implemented in JAX.

2.3. Allen-Cahn equation

Next, we consider the Allen-Cahn equation, given by

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

with the initial condition $u(x, 0) = x^2 \sin(2\pi x)$ and periodic boundary conditions:

$$u(t, -1) = u(t, 1), \quad u_x(t, -1) = u_x(t, 1),$$

where $\epsilon = 10^{-4}$, $\kappa = 5$, and the spatio-temporal domain for computing the solution is $[x, t] \in [-1, 1] \times [0, 1]$. The periodicity at the boundaries is enforced using interpolation polynomials $\{v^{(i)}(x)\}_{i=1}^n$ (see [84]) and defined as

$$v^{(i)}(x) = s_0^{(i)} + s_1^{(i)}(x - a)(b - x)(a + b - 2x) + (r_0^{(i)} + r_1^{(i)}x)(x - a)^2(x - b)^2, \quad (18)$$

where a and b represent the spatial domain boundaries, and the coefficients $\{s_0^{(i)}, s_1^{(i)}, r_0^{(i)}, r_1^{(i)}\}$ are defined for $i = 1, 2$ corresponding to $n = 2$. For all case studies involving the Allen-Cahn equation, the PINNs architecture utilizes the hyperbolic tangent activation function (\tanh) in its hidden layers. The training process begins with the Adam optimizer at $\kappa = 1$, incorporating a learning rate decay schedule. Following this initial phase, κ is increased to 5, and the model is further refined using the BFGS and SSBroyden algorithms. This training strategy starts with Adam on a simplified problem by reducing the parameter κ , and subsequently resumes with BFGS or SSBroyden at the original κ value until convergence. We found this strategy to be highly robust in achieving convergence to the global minimum. In contrast, maintaining $\kappa = 5$ throughout training occasionally led to stagnation in the loss function at relatively higher values for certain initializations, suggesting that the PINN converged to a stationary point distinct from the global minimum.

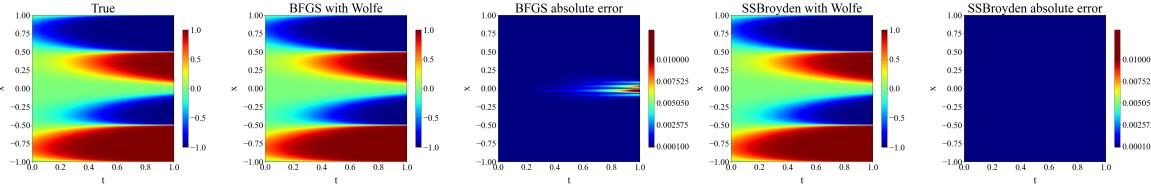


Figure 12: **Double-Precision PINNs for Allen-Cahn equation:** Comparison of Allen-Cahn equation predictions using PINNs optimized with BFGS and SSBroyden methods for **Case 3**. Absolute error plots are provided for each case, emphasizing the differences in error magnitudes between BFGS and SSBroyden.

Figure 12 presents contour plots comparing the PINNs solutions against highly accurate spectral solutions. For each case, the corresponding absolute errors for BFGS and SSBroyden are plotted as well. The numerical solution is computed using MATLAB’s Chebfun package [85] with 1000 Fourier modes per spatial dimension, combined with the ETDRK4 algorithm [86] for temporal integration and a time step of $dt = 10^{-5}$. Figure 13 illustrates the progression of the loss function over iterations for the BFGS and SSBroyden optimizers applied to PINNs with single precision. As summarized in Table 9, SSBroyden achieves a relative error of approximately 10^{-4} , while BFGS converges to around 10^{-3} . Notably, the loss function shows minimal change beyond 15,000 epochs.

Table 10 presents three cases, all of which share the same network architecture. The network consists of three hidden layers, each with 30 neurons. For **Case 2** and **Case 3**, the training process begins with the Adam optimizer, run for 5000 epochs at $\kappa = 1$, using a learning rate decay schedule. Following this initial phase, κ is increased to 5, and the model is further optimized using the BFGS and SSBroyden algorithms for 20,000 and 30,000 iterations, respectively. **Case 4** shows the results for the case that the network is trained using second optimizer BFGS and SSBroyden directly without ADAM optimizer. Figure 14 illustrates the evolution of the l_2 relative error over time and the loss function over iterations, comparing PINNs with double precision for **Case 3**.

We also demonstrate the performance of SSBroyden using PIKANs. The PIKANs architecture incorporates Chebyshev polynomials and B-spline basis function. The KAN architecture consists of four hidden layers, each containing 10 nodes, with cubic spline functions (order 3) defined over 10

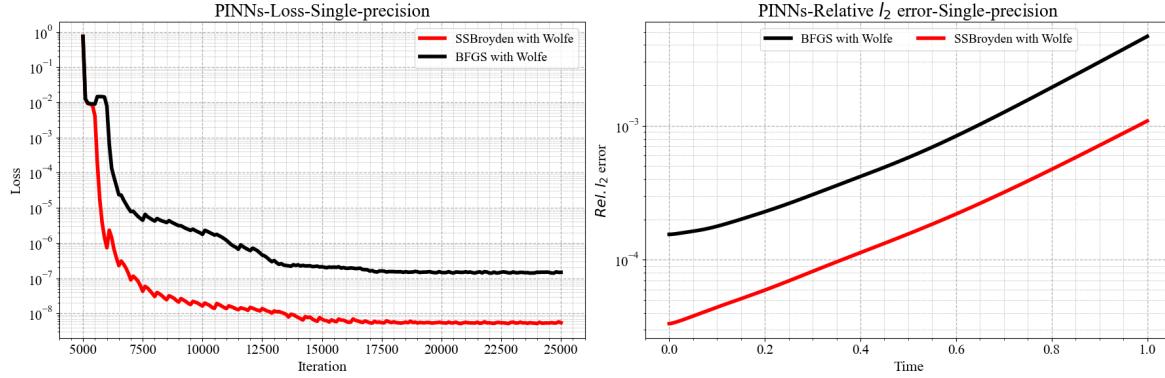


Figure 13: **Single-Precision PINNs for the Allen-Cahn equation:** Evolution of the loss function (left) and l_2 relative errors (right) for **Case 1**. The results compare the performance of PINNs trained with SSBroyden and BFGS optimizers with single precision.

Case	Optimizer[# Iter., Line-search [#Iter.]]	Relative l_2 error	Training Time (s)	Total Params
1	Adam [5000] + BFGS with Wolfe [20000]	1.97×10^{-3}	908	2,019
1	Adam [5000] + SSBroyden with Wolfe [20000]	$4.73e \times 10^{-4}$	762	2,019

Table 9: **PINNs with single precision for the Allen-Cahn equation:** Comparison of the loss function (left) and l_2 relative errors (right) for Allen-Cahn equation using SSBroyden and BFGS optimizers.

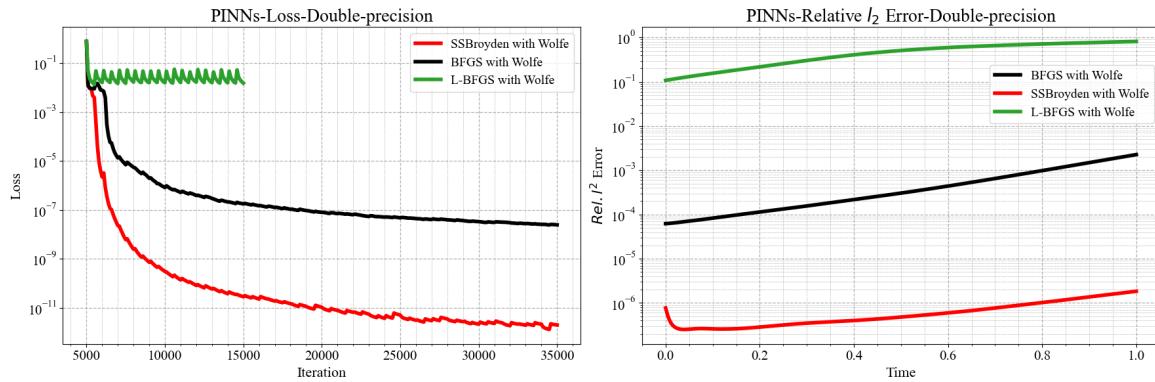


Figure 14: **Double-Precision PINNs for the Allen-Cahn equation:** Evolution of the loss function (left) and l_2 relative errors (right) for **Case 3**. The results compare the performance of PINNs trained with SSBroyden and BFGS optimizers with double precision.

Case	Optimizer[# Iter., Line-search [#Iter.]]	Relative l_2 error	Training time (s)	Total params
2	Adam [5000] + BFGS with Wolfe [20000]	7.59×10^{-4}	733	2,019
2	Adam [5000] + SSBroyden with Wolfe [20000]	1.15×10^{-6}	973	2,019
3	Adam [5000] + BFGS with Wolfe [30000]	9.84×10^{-4}	1493	2,019
3	Adam [5000] + SSBroyden with Wolfe [30000]	$9.43e \times 10^{-7}$	2000	2,019
4	BFGS with Wolfe [30000]	3.78×10^{-4}	1404	2,019
4	SSBroyden with Wolfe [30000]	1.28×10^{-6}	1488	2,019

Table 10: **Double-Precision PINNs for the Allen-Cahn equation:** **Case 2:** A PINN with three layers, each containing 30 neurons, is trained for 1000 iterations using the Adam optimizer, followed by 20,000 iterations with BFGS, SSBroyden, and L-BFGS optimizers with Wolfe line-search. **Case 3:** The same network structure as **Case 2** is used, but the number of iterations with BFGS and SSBroyden is reduced to 30,000. **Case 4:** The same network structure as in **Case 2** and **Case 3** is used; however, instead of starting with the Adam optimizer, the training begins directly with second-order optimizers, namely BFGS and SSBroyden.

grid points. Also, for Chebyshev polynomials, degree 5 is used. Figure 15 illustrates the evolution of the l_2 relative error over time and the loss function over iterations, comparing PIKANs with double precision for **Case 5** and **Case 6**. Table 11 summarizes the l_2 relative error for BFGS and SSBroyden optimizers applied to PINNs and PIKANs architectures with B-spline and Chebyshev polynomials. As

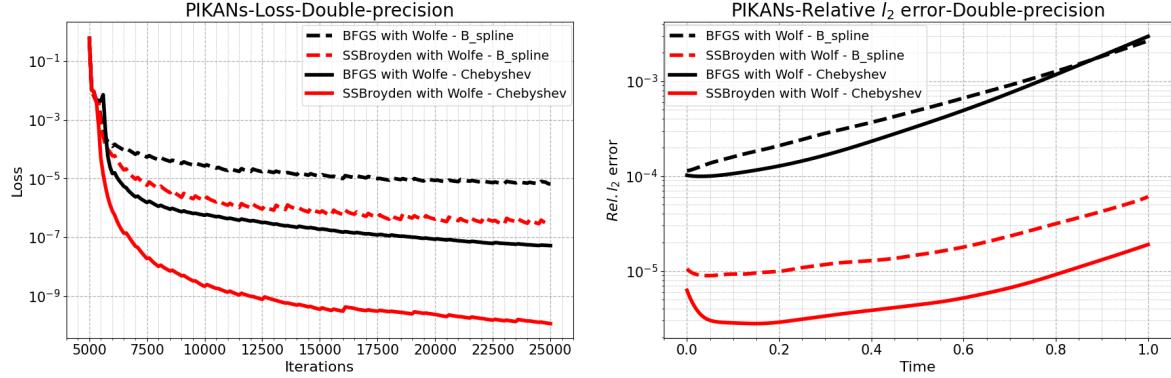


Figure 15: **Double-Precision PIKANs for the Allen-Cahn equation:** Evolution of the loss function (left) and L_2 relative errors (right) for **Case 5** (PIKAN using Chebyshev polynomials) and **Case 6** (PIKAN using B-spline polynomials). The results highlight the performance comparison of PIKANs trained using SSBroyden and BFGS optimizers.

Case	Optimizer[# Iter., Line-search #Iter.]	Relative l_2 error	Training time (s)	Total parameters
5	Adam [5000] + BFGS with Wolfe-Chebyshev	1.23×10^{-3}	1804.82	1,368
5	Adam [5000] + SSBroyden with Wolfe-Chebyshev	9.01×10^{-6}	2239.38	1,368
6	Adam [5000] + BFGS with Wolfe-Bspline	1.17×10^{-2}	11018.23	5,540
6	Adam [5000] + SSBroyden with Wolfe-Bspline	4.52×10^{-4}	12025.20	5,540

Table 11: **Double-Precision PIKANs for the Allen-Cahn equation: Case 5:** PIKANs with Chebyshev polynomials consisting of four hidden layers, each with 10 neurons, and polynomial degree 5. Training involves 5000 Adam steps followed by optimization using BFGS and SSBroyden. **Case 6:** PIKANs with B-splines consisting of four hidden layers with dimensions of 20, grid sizes of 10, and spline orders of 3. Training involves 5000 Adam steps followed by optimization using BFGS and SSBroyden. In all cases, SSBroyden consistently outperforms BFGS in terms of accuracy.

shown, the l_2 relative error obtained using SSBroyden for both KAN architectures is consistently lower than that obtained with BFGS. The results also confirm that PIKANs architectures with Chebyshev polynomials outperform those with B-splines in terms of both error and training time.

Overall, the results for the Allen-Cahn equation demonstrate that both PINNs and PIKANs with Chebyshev polynomials achieved a relative error of 10^{-6} . However, the training time to reach this accuracy is significantly shorter for PINNs (973 s) compared to PIKANs (1804 s). The l_2 relative error is computed across all time and domain points. In terms of training time, BFGS is generally faster than SSBroyden for all cases (PINNs and both PIKANs architectures), though the difference is minimal. Notably, there is a substantial difference in training times between PINNs and PIKANs with B-splines. Additionally, Chebyshev polynomials are shown to be considerably more efficient than B-splines in both accuracy and training time.

2.4. Kuramoto-Sivashinsky equation

In this example, we illustrate the effectiveness of the SSBroyden method for tackling spatio-temporal chaotic systems, with a particular focus on the one-dimensional Kuramoto–Sivashinsky equation. Known for its intricate spatial patterns and unpredictable temporal dynamics, the Kuramoto–Sivashinsky equation poses a significant challenge to conventional numerical approaches. The governing equation is given by:

$$\frac{\partial u}{\partial t} + \alpha u \frac{\partial u}{\partial x} + \beta \frac{\partial^2 u}{\partial x^2} + \gamma \frac{\partial^4 u}{\partial x^4} = 0, \quad (19)$$

where $\alpha = \frac{100}{16}$, $\beta = \frac{100}{16^2}$, and $\gamma = \frac{100}{16^4}$. The initial condition is defined as

$$u_0(x) = \cos(x)(1 + \sin(x)),$$

as presented in [87]. The solution domain is defined as $(t, x) \in [0, 1] \times [x_0 = 0, x_f = 2\pi]$. To evaluate the performance of BFGS and SSBroyden in solving the Kuramoto–Sivashinsky equation and addressing the complexity of this spatio-temporal problem, a time-marching strategy, which divides the temporal domain into smaller intervals, enables stable training and accurate predictions over time. However, it introduces computational overhead due to sequential training for each time window. The time domain

is divided into subdomains with a time increment of $dt = 0.05$, resulting in 20 time windows to train the equation over the interval $[0, 1]$. Moreover, we discuss the use of PINNs for 5 time windows over the interval $[0, 0.5]$ with $dt = 0.1$, as described in [26]. Periodic boundary conditions are implemented following the approach in [84], while the initial condition is softly enforced. The periodic nature of the problem is seamlessly encoded into the model by leveraging Fourier basis functions for the spatial domain. Therefore, the inputs to the PINN model are extended to incorporate these periodic boundary conditions. The total input to the neural network, X_{input} , is defined as:

$$X_{\text{input}} = \left[t, \cos\left(\frac{2\pi mx}{L_x}\right), \sin\left(\frac{2\pi mx}{L_x}\right) \right],$$

where $L_x = x_f - x_0$ is the length of the spatial domain, and M is the number of Fourier modes. Here, $M = 10$ is used.

In all case studies for the Kuramoto-Sivashinsky equation, the neural network architecture comprises five fully connected hidden layers, each containing 30 neurons, with the hyperbolic tangent (tanh) activation function applied across all hidden layers. To improve computational efficiency, the collocation points are dynamically resampled using the RAD method, with updates performed every 500 iterations. We explored two training scenarios. In the first scenario, training begins with the Adam optimizer. A learning rate decay schedule is implemented, starting at 5×10^{-3} and reducing by a factor of 0.98 every 1000 iterations. In the second scenario, training is initiated directly using a second optimizer, either BFGS or SSBroyden.

Figures 16 and 17 present the results for the Kuramoto-Sivashinsky equation for **Case 2**. Figure 16 compares the performance of BFGS and SSBroyden in predicting the Kuramoto-Sivashinsky equation using double-precision PINNs with 20 time windows, displayed through contour plots. The absolute error for each prediction is shown, highlighting a significant difference in error magnitude between the two optimizers, particularly around $t = 1.0$, where BFGS struggles to accurately predict the solution. Figure 17 illustrates the performance of BFGS and SSBroyden at three time steps: $t = 0.0$, $t = 0.5$, and $t = 1.0$ across the spatial domain for **Case 2**. The figure demonstrates the complexity of the Kuramoto-Sivashinsky equation solution as it evolves over time. Although predicting the solution around $t = 1.0$ is challenging, the network successfully predicts the solution, with the prediction and the true solution matching closely.

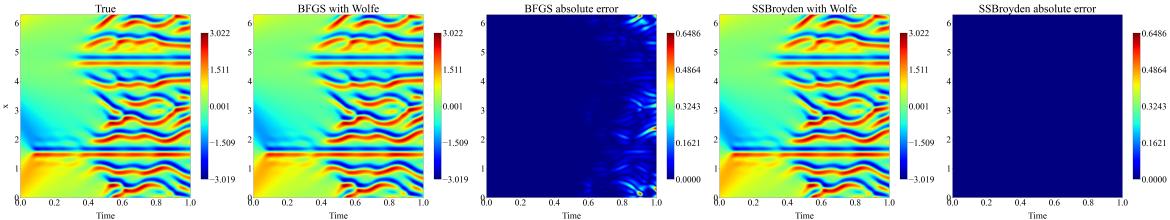


Figure 16: Double-Precision PINNs for the Kuramoto-Sivashinsky equation: Performance comparison of BFGS and SSBroyden in solution prediction for **Case 2**. Contour plots illustrate the absolute error for each method, revealing a notable difference in error magnitude between the two. The x -axis represents time, while the y -axis corresponds to the spatial domain.

Figure 18 illustrates the evolution of the loss function, which includes contributions from both the PDE residual and the initial conditions, across 20 time windows (**Case 2**). The time interval $[0, 1]$ is divided into 20 sequential windows, with each window trained independently using PINNs for 30,000 iterations. The comparison highlights the performance of the BFGS and SSBroyden optimization algorithms. Results show that SSBroyden consistently outperforms BFGS in all time slices, achieving faster convergence and delivering more accurate predictions across the entire domain. The figure also demonstrates that predicting the solution becomes increasingly challenging as we approach $t = 1$, as evidenced by higher loss values for the final time windows compared to the initial ones.

Furthermore, Figure 19 illustrates the relative errors obtained using BFGS and SSBroyden for the Kuramoto-Sivashinsky equation over time. The results are presented for two case studies, when the time interval $[0, 0.5]$ was divided into 5 windows (left), and when the time interval $[0, 1]$ was divided into 20 windows (right). The y -axis represents the relative error, while the x -axis represents time. As the complexity of the Kuramoto-Sivashinsky equation solution increases over time, predictions become more challenging, particularly after $t = 0.5$. This increased complexity leads to a noticeable rise in

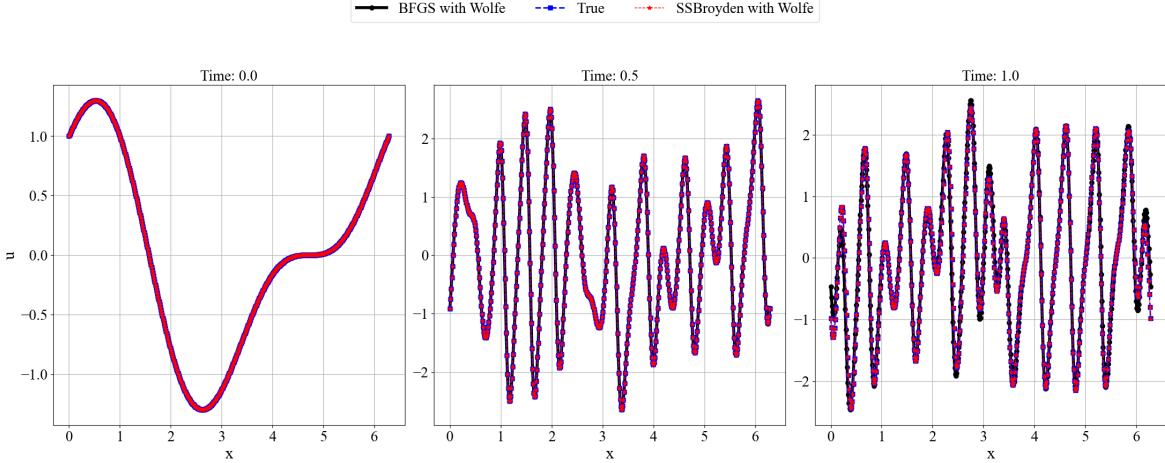


Figure 17: **Double-Precision PINNs for Kuramoto-Sivashinsky equation:** Comparison of numerical solution of the Kuramoto-Sivashinsky equation and predicted solution using BFGS and SSBroyden at three time steps $t = 0.0$, $t = 0.5$, and $t = 1.0$ for **Case 2**. In each plot, the x -axis represents the spatial domain, and the y -axis represents the solution of the Equation (19).

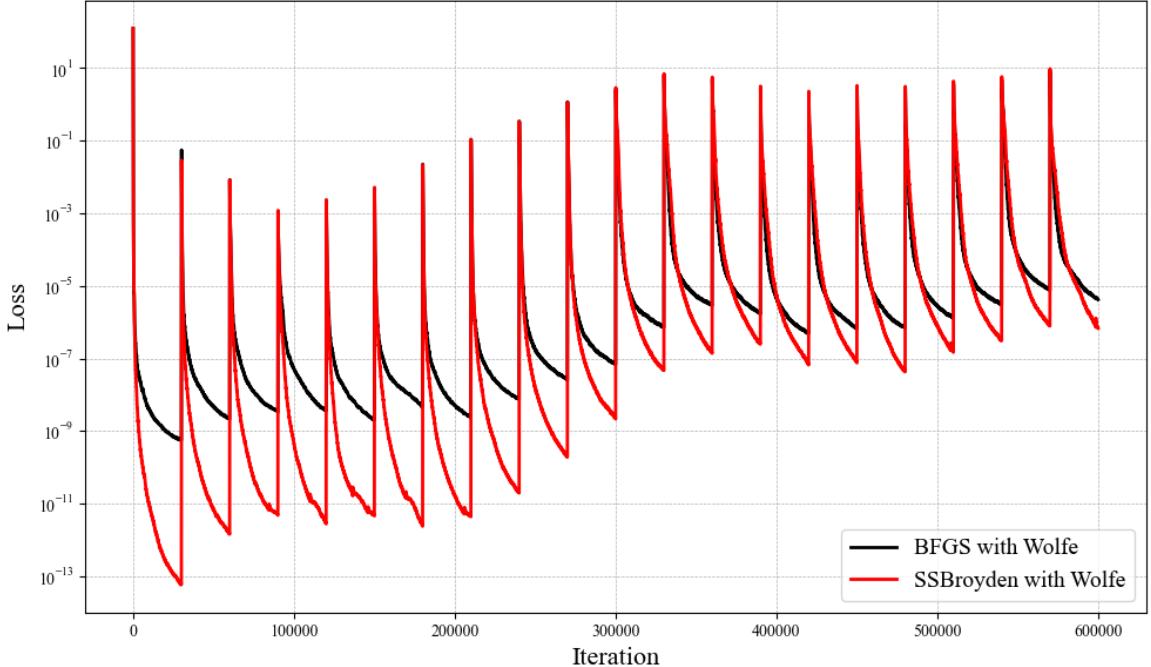


Figure 18: **Double-Precision PINNs for Kuramoto-Sivashinsky equation:** Comparison of loss functions for BFGS and SSBroyden over 20 time windows within the time interval $[0, 1]$ with $\Delta t = 0.05$, consisting of 20 separate PINNs (**Case 2**). Each PINNs was trained for 30,000 iterations, starting with the second optimizer BFGS or SSBroyden.

relative error. For the 5-window case, the relative error grows to approximately 10^{-4} at $t = 0.5$. For the 20-window case, the relative error increases to nearly 10^{-2} at $t = 1$.

Table 12 summarizes the relative l_2 errors and training times for the Kuramoto-Sivashinsky equation across three different cases, where PINNs are trained using double precision over 20 time windows on $[0, 1]$ and 5 time windows on the time interval $[0, 0.5]$. The architecture of the PINNs remains consistent across all cases. In **Case 1** and **Case 2**, the PINNs are trained directly using 20,000 and 30,000 iterations of the BFGS and SSBroyden optimizers, respectively, without prior Adam optimization. In **Case 3**, training involves 1,000 iterations of the Adam optimizer, followed by 30,000 iterations of BFGS and SSBroyden. The results highlight the superior accuracy and computational efficiency of the SSBroyden optimizer compared to BFGS across all cases. Moreover, the nearly identical results

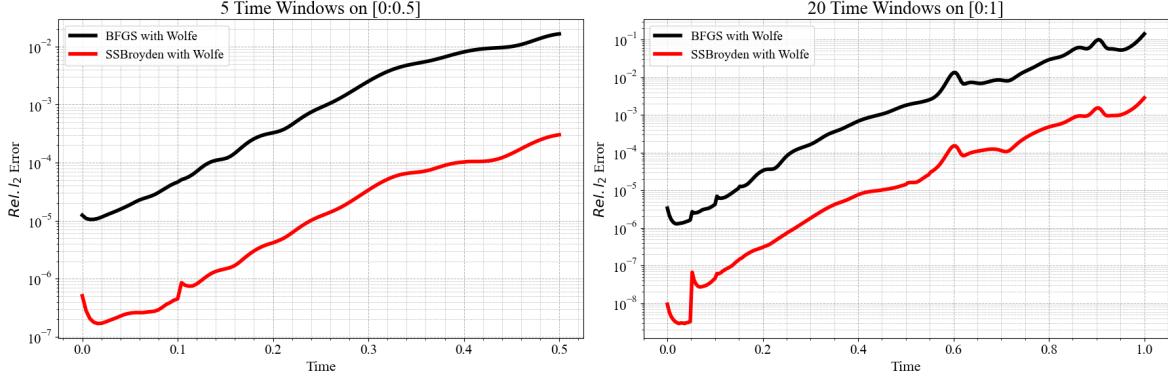


Figure 19: **Double-Precision PINNs for Kuramoto-Sivashinsky equation:** The evolution of l_2 relative errors is illustrated for two different time intervals and time window configurations in **Case 2**. Specifically, the plot on the left displays the l_2 relative errors over 5 time windows on the time interval $[0, 0.5]$, while the plot on the right shows the l_2 relative errors over 20 time windows on the time interval $[0, 1]$.

Case	Optimizer[# Iter., Line-search [#Iter.]]	Relative l_2 error	Training time (s)	Total parameters
1	BFGS with Wolfe [20000] (20 windows)	$8.75e \times 10^{-2}$	102177	4,411
1	SSBroyden with Wolfe [20000] (20 windows)	2.13×10^{-3}	113214	4,411
2	BFGS with Wolfe [30000] (20 windows)	3.65×10^{-2}	119286	4,411
2	SSBroyden with Wolfe [30000] (20 windows)	6.51×10^{-4}	130222	4,411
3	Adam [1000] + BFGS with Wolfe [30000] (20 windows)	6.15×10^{-2}	119353	4,411
3	Adam [1000] + SSBroyden with Wolfe [30000] (20 windows)	7.53×10^{-4}	127372	4,411
4	BFGS with Wolfe [20000] (5 windows)	7.54×10^{-3}	17793	4,411
4	SSBroyden with Wolfe [20000] (5 windows)	1.24×10^{-4}	19048	4,411
5	BFGS with Wolfe [30000] (5 windows)	2.39×10^{-3}	26697	4,411
5	SSBroyden with Wolfe [30000] (5 windows)	2.65×10^{-5}	30883	4,411

Table 12: **Double-Precision PINNs for the Kuramoto-Sivashinsky equation:** Summary of results comparing relative l_2 errors and training times for PINNs trained using BFGS and SSBroyden optimizers over 20 time windows on the time interval $[0, 1]$ and 5 time windows on the time interval $[0, 0.5]$. **Case 1:** The PINN architecture consists of five hidden layers with 30 neurons each. Optimization is performed directly using BFGS and SSBroyden for 30,000 iterations. **Case 2:** The same architecture as **Case 1** is used. Training includes 30,000 iterations of BFGS and SSBroyden optimization. **Case 3:** The architecture is consistent with **Case 1** and **Case 2**. Training consists of 1,000 steps of Adam optimization, followed by 30,000 iterations of BFGS and SSBroyden. Similarly, for **Case 4** and **Case 5**, the architecture remains unchanged, i.e., it employs five hidden layers with 30 neurons each and utilizes five time windows over the time interval $[0, 0.5]$, with training conducted for 20,000 and 30,000 iterations of BFGS and SSBroyden with Strong Wolfe line-search, respectively.

for **Case 2** and **Case 3** suggest that starting optimization directly with BFGS and SSBroyden or including an initial phase of Adam optimization yields similar error levels. Similarly, for **Case 4** and **Case 5**, the architecture remains unchanged, employing five hidden layers with 30 neurons each and utilizing five time windows over the time interval $[0, 0.5]$. Training is conducted for 20,000 and 30,000 iterations of BFGS and SSBroyden, respectively.

Table 13 summarizes the performance of PIKANs using Chebyshev polynomials of degree 3 with two architectures: one consisting of five hidden layers with 10 neurons per layer (**Case 6**) and another with five hidden layers and 20 neurons per layer (**Case 7**). The table highlights the relative accuracy and efficiency of each configuration. It confirms that while increasing the number of neurons to 20 significantly raises the number of parameters and training time, it does not substantially reduce the error. Figure 20 presents the results for **Case 6**. The figure illustrates that the SSBroyden optimizer successfully captures the solution, whereas the BFGS optimizer fails to achieve accurate predictions, demonstrating the superior performance of SSBroyden in this scenario. Comparing the prediction figures obtained using BFGS and SSBroyden with the true solution, it is evident that although the relative error reported in the table for both optimizers is approximately 10^{-1} , SSBroyden's predictions are significantly more accurate than those of BFGS.

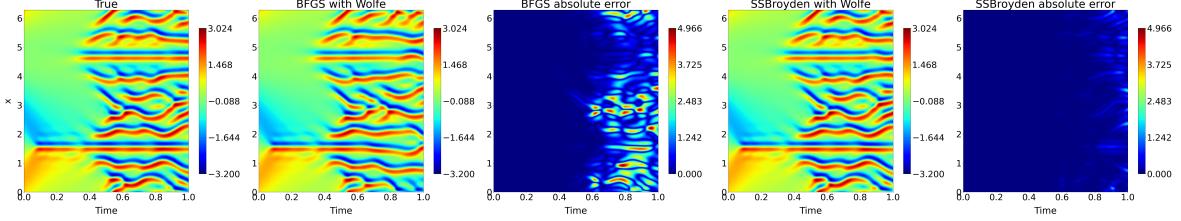


Figure 20: **Double-Precision PIKANs for Kuramoto-Sivashinsky equation:** Evolution of the loss function (left) and l_2 relative errors (right) for **Case 6**. The results compare the performance of PIKANs trained using SSBroyden and BFGS optimizers.

Case	Optimizer[# Iter., Line-search [#Iter.]]	Relative l_2 error	Training time (s)	Total parameters
6	BFGS with Wolfe [20000]-Chebyshev (degree 3)	7.73×10^{-1}	50387	2,480
6	SSBroyden with Wolfe [20000]-Chebyshev (degree 3)	1.27×10^{-1}	54130	2,480
7	BFGS with Wolfe [20000]-Chebyshev (degree 3)	1.68×10^{-1}	219627	8,160
7	SSBroyden with Wolfe [20000]-Chebyshev (degree 3)	1.55×10^{-1}	238393	8,160

Table 13: **Double-Precision PIKANs for Kuramoto-Sivashinsky equation:** Summary of results highlighting the relative L_2 errors and training durations for PIKANs optimized with BFGS and SSBroyden methods. **Case 6:** PIKANs with Chebyshev polynomials of degree 3 consist of five hidden layers, each containing 10 neurons. Training is initiated directly using BFGS and SSBroyden optimizers. **Case 7:** This PIKAN is trained similarly to **Case 6**, utilizing Chebyshev polynomials of degree 3 and consisting of five hidden layers, each with 20 neurons.

2.5. Ginzburg-Landau equation

In this section, we illustrate the effectiveness of the two-dimensional Ginzburg-Landau equation, expressed as:

$$\frac{\partial A}{\partial t} = \epsilon \nabla^2 A + (\nu - \gamma |A|^2) A, \quad (20)$$

where A is a complex-valued function, and ϵ , ν , and γ are constant coefficients. Here, ϵ and ν are real numbers, while γ is a complex constant ($\gamma \in \mathbb{C}$). Representing A in terms of its real and imaginary components, $A = u + iv$, where u and v are real-valued functions and i is the imaginary unit, we derive the following system of PDEs:

$$\frac{\partial u}{\partial t} = \epsilon \nabla^2 u + \nu u - (u^2 + v^2) (\text{Re}(\gamma)u - \text{Im}(\gamma)v), \quad (21)$$

$$\frac{\partial v}{\partial t} = \epsilon \nabla^2 v + \nu v - (u^2 + v^2) (\text{Re}(\gamma)v + \text{Im}(\gamma)u), \quad (22)$$

where $\text{Re}(\cdot)$ and $\text{Im}(\cdot)$ denotes the real and the imaginary parts, respectively. Specific values for these coefficients are $\epsilon = 0.004$, $\nu = 10$, $\gamma = 10 + 15i$, which are the same values as the ones chosen in [88]. The initial condition is

$$A_0(x, y) = 10(y + ix)e^{-25(x^2 + y^2)}.$$

The solution domain is defined as $(t, x) \in [0, 1] \times [-1, 1]^2$. Periodic boundary conditions are enforced using a hard-enforcement method with polynomials defined in (18) for $n = 4$, applied to each spatial variable x and y . Initial conditions, on the other hand, are applied using a soft-enforcement method. Consequently, the loss function integrates contributions from the PDE residuals and the initial conditions. The temporal domain is divided into 5 time windows, each of length $\Delta t = 0.2$, with a separate PINN trained for each window. A fully connected neural network is designed to predict the solution fields u and v . The network architecture consists of five dense layers, each with 30 neurons, where all hidden layers utilize the hyperbolic tangent (\tanh) activation function. The output layer contains two neurons, providing simultaneous predictions for u and v . To enhance training efficiency, the model incorporates the adaptive sampling strategy (RAD) algorithm.

Figures 21 compare the numerical solution of the Ginzburg-Landau equation with the predicted solutions obtained using double-precision PINNs trained with BFGS and SSBroyden, presented as contour plots for **Case 3**. The absolute errors for both optimizers are shown for (a) the real part u and (b) the imaginary part v , demonstrating the superior accuracy of SSBroyden over BFGS. The numerical

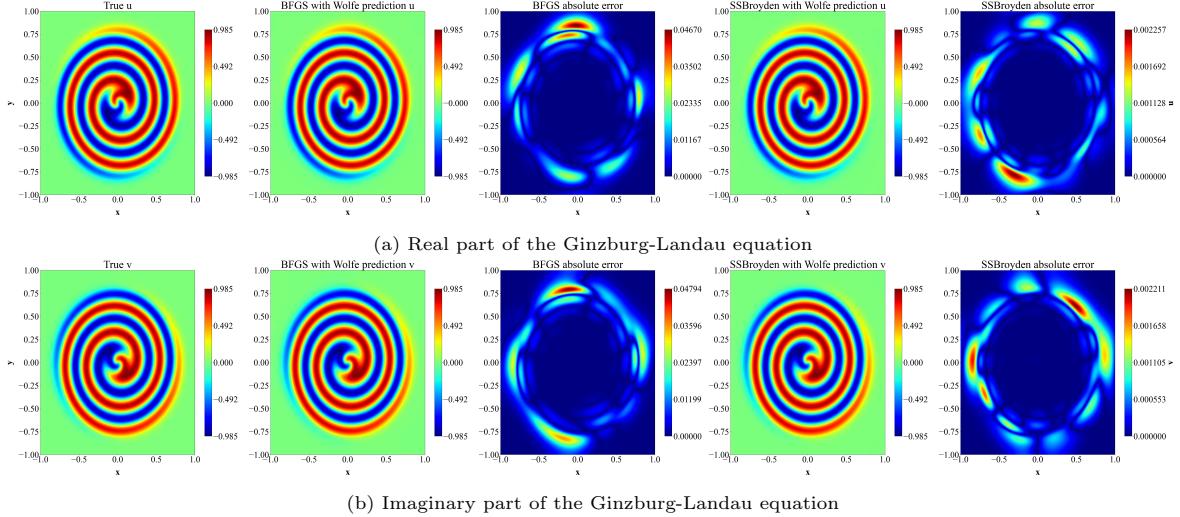


Figure 21: Double-Precision PINNs for Ginzburg-Landau equation: Comparison of the numerical solution for the Ginzburg-Landau equation with predictions obtained using double-precision PIKANs with 5 time windows, each of length $\Delta t = 0.2$. (a) Real part. (b) Imaginary part. In each case, the absolute error is also plotted, demonstrating the superior performance of SSBroyden compared to BFGS for both the real and imaginary components.

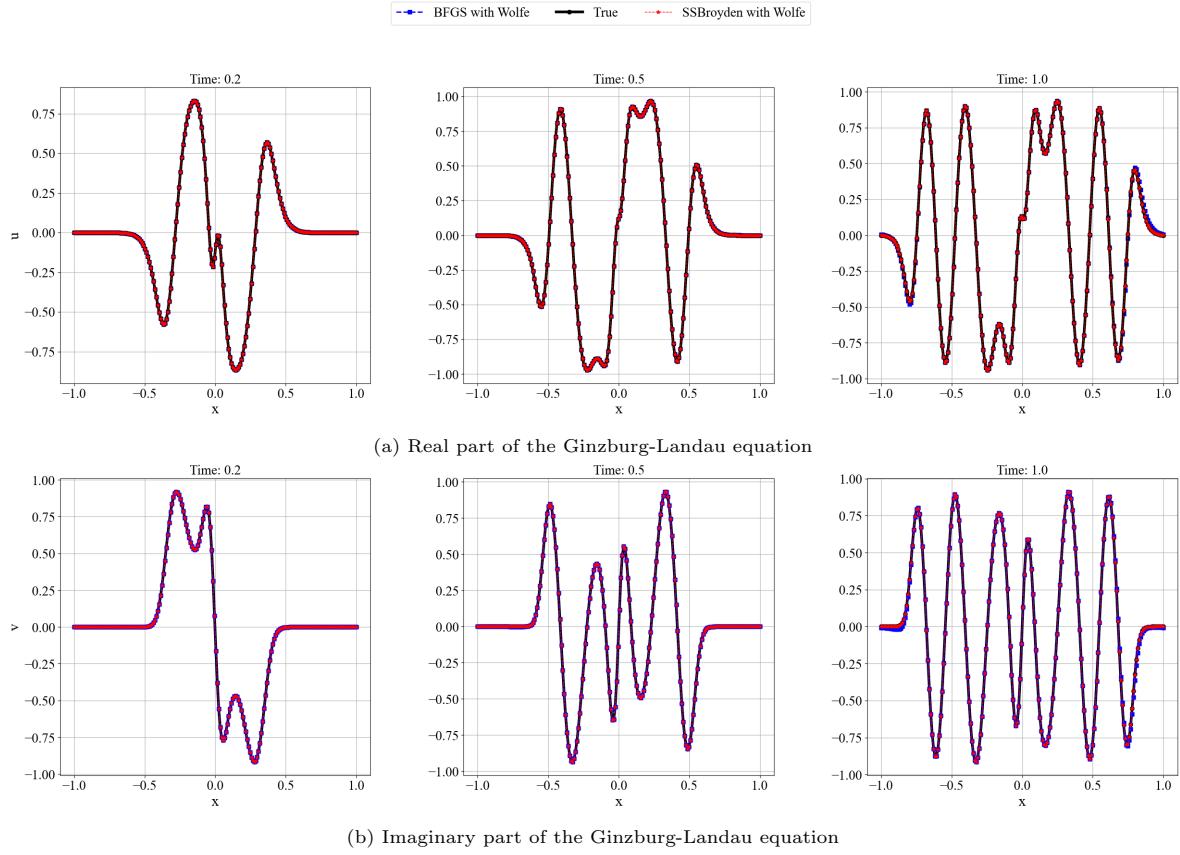


Figure 22: Double-Precision PINNs for Ginzburg-Landau equation: Comparison of the numerical solution and the predicted solution using BFGS and SSBroyden. (a) real part. (b) imaginary part. Results are shown at time steps $t = 0.2$, $t = 0.5$, and $t = 1.0$ obtained using double-precision PIKANs with 5 time windows, each of length $\Delta t = 0.2$.

solution is computed using Chebfun with 200 Fourier modes in each spatial direction, combined with the EDTRK4 algorithm for time integration, employing a step size of $dt = 10^{-5}$. Additionally, Figures 22 illustrate the performance of BFGS and SSBroyden at three distinct time steps: $t = 0.0$, $t = 0.5$, and $t = 1.0$, across the spatial domain for **Case 3**. The plots display (a) the real fields u and (b)

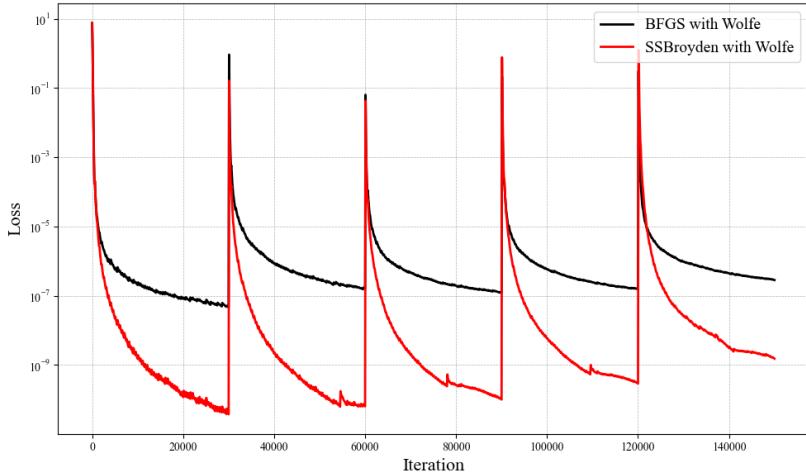


Figure 23: **Double-Precision PINNs for Ginzburg-Landau equation:** Comparison of loss functions for BFGS and SSBroyden over 5 time windows within the time interval $[0, 1]$ with $\Delta t = 0.2$, consisting of 5 separate PINNs (**Case 2**). Each PINNs was trained for 30,000 iterations, starting with the second optimizer BFGS or SSBroyden.

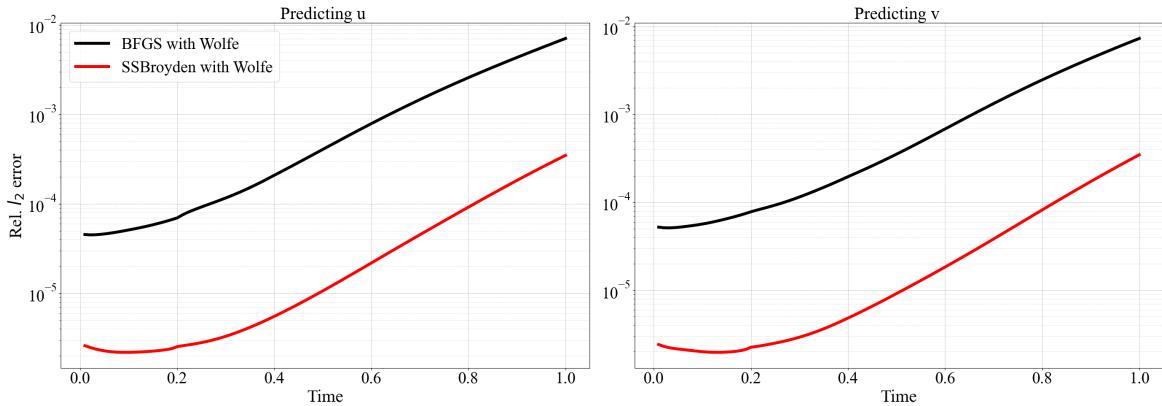


Figure 24: **Double-Precision PINNs for Ginzburg-Landau equation:** The l_2 relative error for the Ginzburg-Landau equation, computed using BFGS and SSBroyden over 5 time windows within the time interval $[0, 1]$ with $\Delta t = 0.2$ for **Case 3**.

Case	Optimizer[# Iter., Line-search #Iter.]	Relative l_2 error	Training time (s)	Total parameters
1	BFGS with Wolfe [20000]	2.73×10^{-2}	21265	3,184
1	SSBroyden with Wolfe [20000]	7.68×10^{-3}	21833	3,184
2	BFGS with Wolfe [20000]	1.05×10^{-2}	41159	6,904
2	SSBroyden with Wolfe [20000]	7.80×10^{-4}	47131	6,904
3	BFGS with Wolfe [30000]	7.19×10^{-3}	61178	6,904
3	SSBroyden with Wolfe [30000]	3.48×10^{-4}	64883	6,904

Table 14: **Double-Precision PINNNs for the Ginzburg-Landau equation:** Comparison of relative l_2 errors, training times, and total parameters for PINNs trained using BFGS and SSBroyden optimizers. **Case 1:** The network architecture consists of five dense layers, each with 30 neurons, designed to predict the solution fields u and v . The temporal domain is divided into 5 time windows, each of length $\Delta t = 0.2$, with a separate PINN trained for each window. Training is performed using 20,000 iterations of BFGS and SSBroyden optimizers. **Case 2:** An extended network architecture with eight dense layers, each containing 30 neurons, is used to predict the solution fields u and v . The temporal domain is divided into 5 time windows, and the PINNs are trained using 20,000 iterations of BFGS and SSBroyden optimizers. **Case 3:** The same architecture as **Case 1** is trained using 30,000 iterations of BFGS and SSBroyden optimizers. The results highlight that SSBroyden consistently outperforms BFGS across all cases. It is worth noting that the error reported in this table represents the average error for the prediction of the real part u and the imaginary part v .

the imaginary fields v , highlighting the evolution and increasing complexity of the solution over time. These results further emphasize the effectiveness of SSBroyden in capturing the solution dynamics accurately.

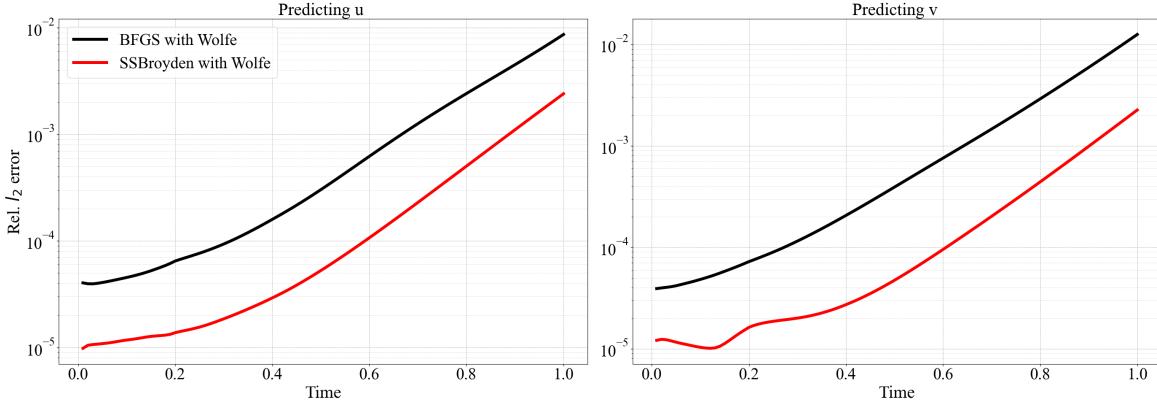


Figure 25: **Double-Precision PIKANs for the Ginzburg-Landau equation:** The l_2 relative error over time for the Ginzburg-Landau equation (**Case 4**), computed using BFGS and SSBroyden optimizers with PIKAN architectures using Chebyshev polynomials.

Case	Optimizer[# Iter., Line-search [#Iter.]]	Relative l_2 error	Training time (s)	Total parameters
4	BFGS with Wolfe-Chebyshev[30000](degree 3)	7.88×10^{-2}	14614	12,152
4	SSBroyden with Wolfe-Chebyshev[30000](degree 3)	1.10×10^{-3}	26152	12,152
5	BFGS with Wolfe-Chebyshev[30000](degree 5)	1.06×10^{-2}	44151	18,212
5	SSBroyden with Wolfe-Chebyshev[30000](degree 5)	2.33×10^{-3}	46625	18,212

Table 15: **Double-Precision PIKANs for the Ginzburg-Landau equation:** Comparison of relative l_2 errors, training times, and total parameters for PIKANs using Chebyshev polynomials of degree 3 (**Case 4**) and degree 5 (**Case 5**). Both models employ four dense layers with 30 neurons each and are trained for 30,000 iterations using BFGS and SSBroyden optimizers. The results show that SSBroyden consistently outperforms BFGS. All values are computed across all data points and time steps, averaged over the real (u) and imaginary (v) components.

Figure 23 presents the evolution of the loss function, which incorporates contributions from both the PDE residual and the initial conditions, across 5 time windows for **Case 3**. The comparison highlights the performance of the BFGS and SSBroyden optimization algorithms. The results indicate that SSBroyden consistently outperforms BFGS across all time slices, achieving faster convergence and providing more accurate predictions throughout the domain.

Table 14 summarizes the relative l_2 error, training time, and total number of training parameters for solving the Ginzburg-Landau equation using different optimizers. The results are computed across all data points and time steps, averaged over the real part (u) and the imaginary part (v) for three cases. In **Case 1**, the network architecture consists of five dense layers, each with 30 neurons, designed to predict the solution fields u and v . In contrast, for **Case 2** and **Case 3**, an extended network architecture with eight dense layers, each containing 30 neurons, is utilized. **Case 2** and **Case 3** are trained using 20,000 and 30,000 iterations of the BFGS and SSBroyden optimizers, respectively. Figure 24 illustrates the evolution of the l_2 relative error over time for both the real part (u) and the imaginary part (v) for **Case 3**. While the errors increase as the prediction complexity grows over time, SSBroyden consistently achieves lower errors compared to BFGS for both components.

Table 15 provides a summary of the relative l_2 error, training time, and total number of training parameters for solving the Ginzburg-Landau equation using PIKANs with Chebyshev polynomials of degree 3 (**Case 4**) and degree 5 (**Case 5**). The results are computed across all data points and time steps, averaged over the real part (u) and the imaginary part (v). In both cases, the network architecture consists of four dense layers, each with 30 neurons, and the models are trained using 30,000 iterations of the BFGS and SSBroyden optimizers. Figure 25 shows the evolution of the l_2 relative error over time for both the real part (u) and the imaginary part (v) for **Case 5**. Although the errors increase as the prediction complexity grows over time, SSBroyden consistently achieves lower errors compared to BFGS for both components.

3. Summary

In this study, we conducted a comprehensive comparison of multiple second-order optimizers for PINNs and PIKANs, focusing on BFGS, SSBFGS, SSBroyden, and L-BFGS with Wolfe line-search

conditions, as well as BFGS with Backtracking line search and trust-region methods. Utilizing the `optax` and `optimistix` library in JAX, we systematically evaluated these optimization techniques across a range of PDEs. We note that we obtained state-of-the-art results by only employing the best optimizers without fine-tuning our loss functions with self-adaptive or attention-based weights or any other enhancements. We simply pursued a straightforward application of good optimization solvers in order to demonstrate that the big bottleneck for PINNs or PIKANs is the optimization error.

Our initial investigation focused on the Burgers equation, where we evaluated the impact of combining first-order optimizers such as Adam with second-order optimizers such as BFGS and SSBroyden. We then extended our analysis to more challenging PDEs, including the Allen–Cahn equation, the Kuramoto–Sivashinsky equation, and the Ginzburg–Landau equation, providing a detailed evaluation of efficiency and accuracy. The results demonstrate that advanced quasi-Newton methods, particularly SSBroyden, SSBFGS and BFGS with Wolfe line-search, significantly improve the convergence rate and accuracy of PINNs, especially for complex and stiff PDEs. Across all cases, SSBroyden consistently outperformed SSBFGS and BFGS, achieving faster convergence and exhibiting greater robustness in handling complex optimization landscapes. These findings highlight the effectiveness of quasi-Newton methods in accelerating the training of PINNs and enhancing numerical stability. Furthermore, we extended our analysis to PIKANs, replacing traditional multilayer perceptron architectures with KANs utilizing Chebyshev polynomials. In this setting, SSBroyden continued to demonstrate superior optimization performance, reinforcing its robustness across diverse network architectures. Additionally, our study emphasized the importance of implementing PINNs and PIKANs with double-precision arithmetic, which improves numerical stability and enhances optimization efficiency across all tested scenarios.

In Appendix [Appendix B](#), we provide a detailed comparison of the number of iterations required to minimize the Rosenbrock function using various optimizers. This analysis highlights how increasing the dimensionality affects the convergence behavior of second-order methods such as BFGS and SSBroyden. We observed that these methods remain efficient across dimensions, while first-order optimizers like ADAM require significantly more iterations and fail to achieve comparable accuracy. The errors reported in this study show that using SSBroyden with the Wolfe line search leads to significantly lower relative \mathcal{L}_2 errors in all benchmark problems tested compared to the state-of-the-art SOAP method [56]. For instance, in the Burgers equation, our approach reduces the error from 8.06×10^{-6} to 4.19×10^{-8} . In the Kuramoto–Sivashinsky equation, we achieve an error of 2.65×10^{-5} over the full time interval $t \in [0, 1]$, whereas SOAP reports 3.86×10^{-2} over a shorter interval $t \in [0, 0.8]$. Similar improvements are observed for the Allen–Cahn and Ginzburg–Landau equations (see Table [16](#)).

Our study demonstrates that SSBFGS and SSBroyden with Wolfe line-search are effective and reliable optimizers for training PINNs and PIKANs. Their capacity to handle complex PDEs and their strong convergence properties make them well-suited for advanced applications in scientific machine learning applications. Future research should focus on incorporating domain-specific adaptations to further improve the accuracy, efficiency, and generalization of PINNs and PIKANs, particularly in high-dimensional and multi-scale problem settings.

Benchmark	\mathcal{L}_2 Error (SOAP [56])	\mathcal{L}_2 Error (SSbroyden with Wolfe)
Burgers Equation	8.06×10^{-6}	4.19×10^{-8}
Allen–Cahn Equation	3.48×10^{-6}	9.43×10^{-7}
Kuramoto–Sivashinsky Equation	3.86×10^{-2} on $t \in [0, 0.8]$	2.65×10^{-5} on $t \in [0, 1]$
Ginzburg–Landau Equation	4.78×10^{-3}	2.33×10^{-3}

Table 16: **Benchmark comparison of relative \mathcal{L}_2 errors.** Results from SSBFGS with Wolfe line-search are compared against SOAP [56]. SSBFGS demonstrates improved accuracy across all PDE benchmarks.

Acknowledgments

This research was primarily supported as part of the AIM for Composites, an Energy Frontier Research Center funded by the U.S. Department of Energy (DOE), Office of Science, Basic Energy Sciences (BES), under Award #DE-SC0023389 (computational studies, data analysis). Additional funding was provided by the DOE-MMICS SEA-CROGS DE-SC0023191 award, the MURI/AFOSR FA9550-20-1-0358 project, and a grant for GPU Cluster for Neural PDEs and Neural Operators to

Support MURI Research and Beyond, under Award #FA9550-23-1-0671. JFU is supported by the pre-doctoral fellowship ACIF 2023, cofunded by Generalitat Valenciana and the European Union through the European Social Fund. JFU acknowledges the support through the grant PID2021-127495NB-I00 funded by MCIN/AEI/10.13039/501100011033 and by the European Union, and the Astrophysics and High Energy Physics programme of the Generalitat Valenciana ASFAE/2022/026 funded by MCIN and the European Union NextGenerationEU (PRTR-C17.I1).

Appendix A. Lorenz system

Considering the the Lorenz system, as described in [89]. It is governed by the following set of coupled ordinary differential equations:

$$\frac{dx}{dt} = \sigma(y - x), \quad (\text{A.1})$$

$$\frac{dy}{dt} = x(\rho - z) - y, \quad (\text{A.2})$$

$$\frac{dz}{dt} = xy - \beta z, \quad (\text{A.3})$$

where the parameters of the Lorenz system are given as $\sigma = 10$, $\rho = 28$, and $\beta = \frac{8}{3}$. The system continues to display chaotic behavior [90]. The solution is computed up to $t = 20$, starting from the initial condition $(x(0), y(0), z(0)) = (1, 1, 1)$.

The neural network architecture consists of three hidden layers, each containing 30 neurons. The hyperbolic tangent (\tanh) activation function is applied across all hidden layers to capture the system's inherent nonlinearity. Training begins with the Adam optimizer, following an exponential decay schedule with an initial learning rate of 5×10^{-3} , which decays by a factor of 0.98 every 1000 iterations. Subsequently, the BFGS and SSBroyden optimizers, incorporating Wolfe line-search, are employed for further optimization.

To address the extended time horizon $[0, 20]$, the time domain is partitioned into 40 time windows, each of length $\Delta t = 0.5$. A separate PINNs is trained for each window, using the final state of the previous window as the initial condition for the next. Figure A.26 compares the numerical solution of the Lorenz system with PINNs predictions obtained via BFGS and SSBroyden over 40 time windows in $[0, 20]$. The plots show $x(t)$ (left), $y(t)$ (middle), and $z(t)$ (right). Table A.17 summarizes the L_2 relative errors for predicting $x(t)$, $y(t)$, and $z(t)$ across all 40 time windows and reports the corresponding training times for BFGS and SSBroyden.

Appendix B. Multi-dimensional Rosenbrock function

Here, we consider the non-quadratic objective functions Rosenbrock, which is commonly used to evaluate optimization algorithms, defined in two-dimensions as:

$$f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (x_1 - 1)^2.$$

The SSBroyden method with Strong Wolfe line-search demonstrates performance comparable to SSFBS and BFGS with Strong Wolfe line-search, with the optimization starting from $x_i = 0.5$. When applied to challenging optimization problems, including the Rosenbrock function—characterized by a narrow, curved valleys and multi-modal landscapes- SSBroyden efficiently converges to the global minimum with accuracy similar to BFGS.

Figure B.27 presents a comparison of optimization performance on the Rosenbrock function across four different dimensionalities: 2D, 5D, 10D, and 20D. To ensure a fair comparison, all optimizers were initialized with the same starting point ($\mathbf{x}_0 = [0.5, \dots, 0.5]$), and the stopping criterion was uniformly defined as the gradient norm falling below 10^{-6} or a maximum of 5000 iterations. Both methods used the same initial inverse Hessian approximation ($H_0 = I$) and consistent optimization parameters. The plots highlight the convergence behavior of each method across increasing dimensions, illustrating how the number of iterations and the rate of loss decay vary with problem size. The number of iterations required to reach the global minimum increases as the dimensionality of the problem grows.

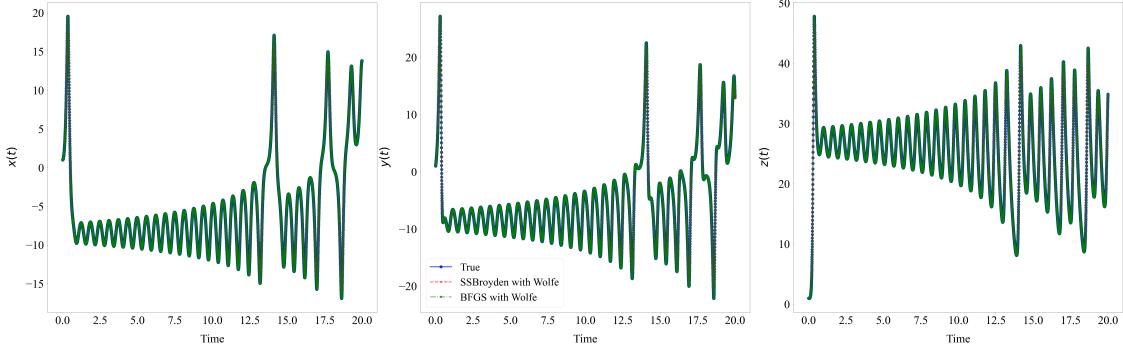


Figure A.26: **Double-Precision PINNs for the Lorenz system:** Comparison of PINNs predictions using BFGS and SSBroyden with Strong Wolfe line-search against the numerical solution for the Lorenz system. The plots display $x(t)$, $y(t)$, and $z(t)$ over time.

Metric	BFGS with Wolfe	SSBroyden with Wolfe
Relative l_2 error for $x(t)$	6.74×10^{-4}	9.65×10^{-5}
Relative l_2 error for $y(t)$	9.81×10^{-4}	1.40×10^{-4}
Relative l_2 error for $z(t)$	4.15×10^{-4}	5.94×10^{-5}
Training time (s)	4541.61	5022.95

Table A.17: **Double-Precision PINNs for the Lorenz system:** Double-Precision PINN results for the Lorenz system: L_2 Relative Errors and training times for PINN predictions obtained using BFGS and SSBroyden over 40 time windows in $[0, 20]$.

However, for the 20-dimensional case, SSBFGS and SSBroyden requires more iterations to reach global optimization compared to BFGS. While, for the 2D case, BFGS, SSBFGS and SSBroyden require approximately the same number of iterations to achieve the global minimum.

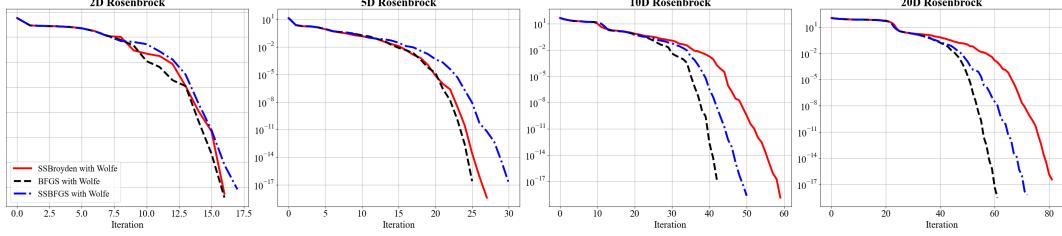


Figure B.27: **Convergence comparison** of BFGS, SS Broyden, and SSBFGS with Strong Wolfe line-search on the Rosenbrock function across multiple dimensions (2D, 5D, 10D, and 20D).

Dimension	Optimizer, Line-search	Iterations	Optimized parameters	Final Loss
2	BFGS with Wolfe	17	[1.0, 1.0]	1.50×10^{-15}
2	SSBroyden with Wolfe	17	[1.0, 1.0]	2.97×10^{-15}
2	SSBFGS with Wolfe	19	[1.0, 1.0]	3.42×10^{-20}
2	Gradient Descent	5000	[0.9689, 0.9386]	9.71×10^{-4}
2	Adam	3899	[1.0, 1.0]	1.24×10^{-12}
5	BFGS with Wolfe	26	[1.0, 1.0, 1.0, 1.0, 1.0]	2.56×10^{-17}
5	SSBroyden with Wolfe	27	[1.0, 1.0, 1.0, 1.0, 1.0]	7.79×10^{-17}
5	SSBFGS with Wolfe	31	[1.0, 1.0, 1.0, 1.0, 1.0]	1.44×10^{-17}
5	Gradient Descent	5000	[0.9970, 0.9939, 0.9879, 0.9759, 0.9522]	7.78×10^{-4}
5	Adam	5000	[1.0, 1.0, 1.0, 1.0, 1.0]	2.77×10^{-11}
10	BFGS with Wolfe	43	[1.0, ..., 1.0]	1.53×10^{-17}
10	SSBroyden with Wolfe	57	[1.0, ..., 1.0]	1.31×10^{-15}
10	SSBFGS with Wolfe	49	[1.0, ..., 1.0]	5.19×10^{-17}
10	Gradient Descent	5000	[0.9999, ..., 0.9421]	1.15×10^{-3}
10	Adam	5000	[0.9998, 1.0011, ..., 1.0002]	3.77×10^{-3}
20	BFGS with Wolfe	60	[1.0, ..., 1.0]	6.35×10^{-16}
20	SSBroyden with Wolfe	81	[1.0, ..., 1.0]	1.40×10^{-16}
20	SSBFGS with Wolfe	70	[1.0, ..., 1.0]	4.90×10^{-16}
20	Gradient Descent	5000	[1.0, ..., 0.9312]	1.64×10^{-3}
20	Adam	5000	[1.0, ..., 0.9999]	6.57×10^{-7}

Table B.18: Final optimization results for the Rosenbrock function using various methods across dimensions 2D, 5D, 10D, and 20D. All quasi-Newton methods recover the global minimum to machine precision. First-order methods (GD, Adam) either converge slowly or stall at suboptimal values in higher dimensions.

References

- [1] S. Cuomo, V. S. Di Cola, F. Giampaolo, G. Rozza, M. Raissi, F. Piccialli, Scientific machine learning through physics-informed neural networks: Where we are and what's next, *Journal of Scientific Computing* (2022) 88.
- [2] M. Raissi, P. Perdikaris, G. E. Karniadakis, Machine learning of linear differential equations using gaussian processes, *Journal of Computational Physics* (2017) 683–693.
- [3] Z. Liu, Y. Wang, S. Vaidya, F. Ruehle, J. Halverson, M. Soljacic, T. Y. Hou, M. Tegmark, [Kan: Kolmogorov-arnold networks](#), ArXiv (2024).
URL <https://api.semanticscholar.org/CorpusID:269457619>
- [4] K. Shukla, J. D. Toscano, Z. Wang, Z. Zou, G. E. Karniadakis, A comprehensive and fair comparison between mlp and kan representations for differential equations and operator networks, *Computer Methods in Applied Mechanics and Engineering* (2024) 117290.
- [5] T. Ji, Y. Hou, D. Zhang, A comprehensive survey on kolmogorov arnold networks (kan), arXiv preprint arXiv:2407.11075 (2024).
- [6] D. Zhang, L. Lu, L. Guo, G. E. Karniadakis, Quantifying total uncertainty in physics-informed neural networks for solving forward and inverse stochastic problems, *Journal of Computational Physics* 397 (2019) 108850.
- [7] A. D. Jagtap, E. Kharazmi, G. E. Karniadakis, Conservative physics-informed neural networks on discrete domains for conservation laws: Applications to forward and inverse problems, *Computer Methods in Applied Mechanics and Engineering* 365 (2020) 113028.
- [8] A. D. Jagtap, G. E. Karniadakis, Extended physics-informed neural networks (xpinns): A generalized space-time domain decomposition based deep learning framework for nonlinear partial differential equations, *Communications in Computational Physics* 28 (2020).
- [9] K. Shukla, A. D. Jagtap, G. E. Karniadakis, Parallel physics-informed neural networks via domain decomposition, *Journal of Computational Physics* 447 (2021) 110683.
- [10] E. Kiyani, K. Shukla, G. E. Karniadakis, M. Karttunen, A framework based on symbolic regression coupled with extended physics-informed neural networks for gray-box learning of equations of motion from data, *Computer Methods in Applied Mechanics and Engineering* 415 (2023) 116258.
- [11] H. Gao, L. Sun, J.-X. Wang, Phygeonet: Physics-informed geometry-adaptive convolutional neural networks for solving parameterized steady-state pdes on irregular domain, *Journal of Computational Physics* (2021) 110079.
- [12] P. Ren, C. Rao, Y. Liu, J.-X. Wang, H. Sun, Phycrnet: Physics-informed convolutional-recurrent network for solving spatiotemporal pdes, *Computer Methods in Applied Mechanics and Engineering* (2022) 114399.
- [13] W. Li, C. Zhang, C. Wang, H. Guan, D. Tao, Revisiting pinns: Generative adversarial physics-informed neural networks and point-weighting method, arXiv preprint arXiv:2205.08754 (2022).
- [14] J.-Z. Peng, Y. Hua, N. Aubry, Z.-H. Chen, M. Mei, W.-T. Wu, Data and physics-driven modeling for fluid flow with a physics-informed graph convolutional neural network, *Ocean Engineering* 301 (2024) 117551.
- [15] E. Rodrigo-Bonet, N. Deligiannis, Physics-guided graph convolutional deep equilibrium network for environmental data, in: 2024 32nd European Signal Processing Conference (EUSIPCO), IEEE, 2024, pp. 987–991.
- [16] L. Yang, X. Meng, G. E. Karniadakis, B-pinns: Bayesian physics-informed neural networks for forward and inverse pde problems with noisy data, *Journal of Computational Physics* (2021) 109913.

- [17] O. Graf, P. Flores, P. Protopapas, K. Pichara, Error-aware b-pinns: Improving uncertainty quantification in bayesian physics-informed neural networks, arXiv preprint arXiv:2212.06965 (2022).
- [18] E. Kianiharchegani, Data-driven exploration of coarse-grained equations: Harnessing machine learning, Ph.D. thesis, The University of Western Ontario (Canada) (2023).
- [19] E. Kiyani, M. Kooshkbaghi, K. Shukla, R. B. Koneru, Z. Li, L. Bravo, A. Ghoshal, G. E. Karniadakis, M. Karttunen, Characterization of partial wetting by cmas droplets using multiphase many-body dissipative particle dynamics and data-driven discovery based on pinns, Journal of Fluid Mechanics 985 (2024) A7.
- [20] V. Dwivedi, B. Srinivasan, Physics informed extreme learning machine (pielm)—a rapid method for the numerical solution of partial differential equations, Neurocomputing (2020) 96–118.
- [21] E. Kharazmi, Z. Zhang, G. E. Karniadakis, hp-vpinns: Variational physics-informed neural networks with domain decomposition, Computer Methods in Applied Mechanics and Engineering (2021) 113547.
- [22] T. De Ryck, S. Mishra, Error analysis for physics-informed neural networks (pinns) approximating kolmogorov pdes, Advances in Computational Mathematics (2022) 79.
- [23] S. Mishra, R. Molinaro, Estimates on the generalization error of physics-informed neural networks for approximating a class of inverse problems for pdes, IMA Journal of Numerical Analysis (2022) 981–1022.
- [24] Y. Shin, J. Darbon, G. E. Karniadakis, On the convergence of physics informed neural networks for linear second-order elliptic and parabolic type pdes, Communications in Computational Physics 28 (5) (11 2020). [doi:10.4208/cicp.oa-2020-0193](https://doi.org/10.4208/cicp.oa-2020-0193).
- [25] Y. Shin, J. Darbon, G. E. Karniadakis, On the convergence and generalization of physics informed neural networks, arXiv e-prints (2020) arXiv–2004.
- [26] S. Wang, S. Sankaran, P. Perdikaris, Respecting causality is all you need for training physics-informed neural networks, arXiv preprint arXiv:2203.07404 (2022).
- [27] S. Wang, X. Yu, P. Perdikaris, When and why pinns fail to train: A neural tangent kernel perspective, Journal of Computational Physics 449 (2022) 110768.
- [28] S. Wang, H. Wang, P. Perdikaris, [On the eigenvector bias of fourier feature networks: From regression to solving multi-scale pdes with physics-informed neural networks](#), Computer Methods in Applied Mechanics and Engineering 384 (2021) 113938. [doi:<https://doi.org/10.1016/j.cma.2021.113938>](https://doi.org/10.1016/j.cma.2021.113938).
URL <https://www.sciencedirect.com/science/article/pii/S0045782521002759>
- [29] L. D. McClenney, U. M. Braga-Neto, [Self-adaptive physics-informed neural networks](#), Journal of Computational Physics 474 (2023) 111722. [doi:<https://doi.org/10.1016/j.jcp.2022.111722>](https://doi.org/10.1016/j.jcp.2022.111722).
URL <https://www.sciencedirect.com/science/article/pii/S0021999122007859>
- [30] J. Bai, G.-R. Liu, A. Gupta, L. Alzubaidi, X.-Q. Feng, Y. Gu, [Physics-informed radial basis network \(pirbn\): A local approximating neural network for solving nonlinear partial differential equations](#), Computer Methods in Applied Mechanics and Engineering 415 (2023) 116290. [doi:<https://doi.org/10.1016/j.cma.2023.116290>](https://doi.org/10.1016/j.cma.2023.116290).
URL <https://www.sciencedirect.com/science/article/pii/S0045782523004140>
- [31] N. Jha, E. Mallik, [Gpinn with neural tangent kernel technique for nonlinear two point boundary value problems](#), Neural Processing Letters 56 (2024) 192. [doi:\[10.1007/s11063-024-11644-7\]\(https://doi.org/10.1007/s11063-024-11644-7\)](https://doi.org/10.1007/s11063-024-11644-7).
URL <https://doi.org/10.1007/s11063-024-11644-7>
- [32] S. J. Anagnostopoulos, J. D. Toscano, N. Stergiopoulos, G. E. Karniadakis, Learning in PINNs: Phase transition, total diffusion, and generalization, arXiv e-prints (2024). [arXiv:2403.18494](https://arxiv.org/abs/2403.18494).

- [33] N. Tishby, F. C. Pereira, W. Bialek, The information bottleneck method, arXiv e-prints (2000) physics/0004057doi:[10.48550/arXiv.physics/0004057](https://doi.org/10.48550/arXiv.physics/0004057).
- [34] J. Lee, L. Xiao, S. S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, J. Pennington, [Wide neural networks of any depth evolve as linear models under gradient descent*](#), Journal of Statistical Mechanics: Theory and Experiment (2020) 124002doi:[10.1088/1742-5468/abc62b](https://doi.org/10.1088/1742-5468/abc62b).
URL <https://dx.doi.org/10.1088/1742-5468/abc62b>
- [35] J. M. Cohen, S. Kaur, Y. Li, J. Z. Kolter, A. Talwalkar, [Gradient descent on neural networks typically occurs at the edge of stability](#), ArXiv abs/2103.00065 (2021).
URL <https://api.semanticscholar.org/CorpusID:232076011>
- [36] J. M. Cohen, B. Ghorbani, S. Krishnan, N. Agarwal, S. Medapati, M. Badura, D. Suo, D. Cardoze, Z. Nado, G. E. Dahl, J. Gilmer, Adaptive Gradient Methods at the Edge of Stability, arXiv e-prints (2022). doi:[10.48550/arXiv.2207.14484](https://doi.org/10.48550/arXiv.2207.14484).
- [37] H. Robbins, S. Monro, A stochastic approximation method, *The annals of mathematical statistics* (1951) 400–407.
- [38] P. Jain, S. M. Kakade, R. Kidambi, P. Netrapalli, A. Sidford, Parallelizing stochastic gradient descent for least squares regression: mini-batching, averaging, and model misspecification, *Journal of machine learning research* 18 (2018) 1–42.
- [39] D. F. Shanno, Conditioning of quasi-newton methods for function minimization, *Mathematics of computation* 24 (1970) 647–656.
- [40] J. Hu, B. Jiang, L. Lin, Z. Wen, Y.-x. Yuan, Structured quasi-newton methods for optimization with orthogonality constraints, *SIAM Journal on Scientific Computing* 41 (2019) A2239–A2269.
- [41] J. E. Dennis Jr, J. Jorge, Moré, [quasi-newton methods, motivation and theory](#), *SIAM review* 19 (1977) 46–89.
- [42] J. Larson, M. Menickelly, S. M. Wild, Derivative-free optimization methods, *Acta Numerica* 28 (2019) 287–404.
- [43] O. Kramer, D. E. Ciaurri, S. Koziel, Derivative-free optimization, in: *Computational optimization, methods and algorithms*, Springer, 2011, pp. 61–83.
- [44] N. M. Nawi, M. R. Ransing, R. S. Ransing, An improved learning algorithm based on the broyden-fletcher-goldfarb-shanno (bfsgs) method for back propagation neural networks, in: *Sixth International Conference on Intelligent Systems Design and Applications*, Vol. 1, IEEE, 2006, pp. 152–157.
- [45] D. R. S. Saputro, P. Widyaningsih, Limited memory broyden-fletcher-goldfarb-shanno (l-bfgs) method for the parameter estimation on geographically weighted ordinal logistic regression model (gwolr), in: *AIP conference proceedings*, Vol. 1868, AIP Publishing, 2017, pp. 040009–1–040009–9.
- [46] C. G. Broyden, The Convergence of a Class of Double-rank Minimization Algorithms 1. General Considerations, *IMA Journal of Applied Mathematics* (1970) 76–90doi:[10.1093/imamat/6.1.76](https://doi.org/10.1093/imamat/6.1.76).
- [47] R. Fletcher, A new approach to variable metric algorithms, *The Computer Journal* 13 (1970) 317–322. doi:[10.1093/comjnl/13.3.317](https://doi.org/10.1093/comjnl/13.3.317).
- [48] D. Goldfarb, [A family of variable-metric methods derived by variational means](#), *Mathematics of Computation* 24 (1970) 23–26.
URL <https://api.semanticscholar.org/CorpusID:790344>
- [49] D. F. Shanno, [Conditioning of quasi-newton methods for function minimization](#), *Mathematics of Computation* 24 (1970) 647–656.
URL <https://api.semanticscholar.org/CorpusID:7977144>
- [50] D. C. Liu, J. Nocedal, [On the limited memory bfsgs method for large scale optimization](#), *Mathematical Programming* 45 (1989) 503–528.
URL <https://api.semanticscholar.org/CorpusID:5681609>

- [51] P. Rathore, W. Lei, Z. Frangella, L. Lu, M. Udell, Challenges in Training PINNs: A Loss Landscape Perspective, arXiv e-prints (2024) arXiv:2402.01868 [doi:10.48550/arXiv.2402.01868](https://doi.org/10.48550/arXiv.2402.01868).
- [52] A. Jnini, F. Vella, M. Zeinhofer, Gauss-Newton Natural Gradient Descent for Physics-Informed Computational Fluid Dynamics, arXiv e-prints (2024) arXiv:2402.10680 [doi:10.48550/arXiv.2402.10680](https://doi.org/10.48550/arXiv.2402.10680).
- [53] J. F. Urbán, P. Stefanou, J. A. Pons, Unveiling the optimization process of physics informed neural networks: How accurate and competitive can pinns be?, arXiv preprint arXiv:2405.04230 (2024).
- [54] M. Al-Baali, E. Spedicato, F. Maggioni, Broyden's quasi-newton methods for a nonlinear system of equations and unconstrained optimization: a review and open problems, Optimization Methods and Software (2014) 937–954 [doi:10.1080/10556788.2013.856909](https://doi.org/10.1080/10556788.2013.856909).
- [55] M. Al-Baali, H. Khalfan, Wide interval for efficient self-scaling quasi-newton algorithms, Optimization Methods and Software (2005) 679–691 [doi:10.1080/10556780410001709448](https://doi.org/10.1080/10556780410001709448).
- [56] S. Wang, A. K. Bhartari, B. Li, P. Perdikaris, Gradient alignment in physics-informed neural networks: A second-order optimization perspective, arXiv preprint arXiv:2502.00604 (2025).
- [57] P. Niyogi, F. Girosi, Generalization bounds for function approximation from scattered noisy data, Advances in Computational Mathematics 10 (1999) 51–80.
- [58] L. Bottou, O. Bousquet, The tradeoffs of large scale learning, Advances in neural information processing systems 20 (2007).
- [59] H. Li, Z. Xu, G. Taylor, C. Studer, T. Goldstein, Visualizing the loss landscape of neural nets, Advances in neural information processing systems 31 (2018).
- [60] K. Kawaguchi, Deep learning without poor local minima, Advances in neural information processing systems 29 (2016).
- [61] L. Armijo, Minimization of functions having lipschitz continuous first partial derivatives, Pacific Journal of mathematics 16 (1966) 1–3.
- [62] J. Nocedal, Theory of algorithms for unconstrained optimization, Acta Numerica 1 (1992) 199–242. [doi:10.1017/S0962492900002270](https://doi.org/10.1017/S0962492900002270).
- [63] R. Fletcher, An overview of unconstrained optimization, in: C. A. Floudas, P. M. Pardalos (Eds.), Algorithms for Continuous Optimization: The State of the Art, Springer Netherlands, Dordrecht, 1994, pp. 109–143. [doi:doi.org/10.1007/978-94-009-0369-2_5](https://doi.org/10.1007/978-94-009-0369-2_5).
- [64] M. J. D. Powell, How bad are the bfgs and dfp methods when the objective function is quadratic?, Mathematical Programming (1986) 34–47.
- [65] R. H. Byrd, D. C. Liu, J. Nocedal, On the behavior of broyden's class of quasi-newton methods, SIAM Journal on Optimization (1992) 533–557 [doi:10.1137/0802026](https://doi.org/10.1137/0802026).
- [66] S. S. Oren, D. G. Luenberger, Self-scaling variable metric (ssvm) algorithms, Management Science (1974) 845–862 [doi:10.1287/mnsc.20.5.845](https://doi.org/10.1287/mnsc.20.5.845).
- [67] J. Nocedal, Y.-x. Yuan, Analysis of a self-scaling quasi-newton method, Mathematical Programming (1993) 19–37 [doi:10.1007/BF01582136](https://doi.org/10.1007/BF01582136).
URL <https://doi.org/10.1007/BF01582136>
- [68] M. Al-Baali, Analysis of a family of self-scaling quasi-newton methods, Dept. of Mathematics and Computer Science, United Arab Emirates University, Tech. Report (1993).
- [69] M. Al-Baali, Numerical experience with a class of self-scaling quasi-newton algorithms, Journal of Optimization Theory and Applications (1998) 533–553 [doi:10.1023/A:1022608410710](https://doi.org/10.1023/A:1022608410710).
URL <https://doi.org/10.1023/A:1022608410710>

- [70] M. Contreras, R. A. Tapia, [Sizing the bfgs and dfp updates: Numerical study](#), Journal of Optimization Theory and Applications (1993) 93–108doi:[10.1007/BF00940702](https://doi.org/10.1007/BF00940702).
URL <https://doi.org/10.1007/BF00940702>
- [71] L. Lukšan, [Computational experience with known variable metric updates](#), Journal of Optimization Theory and Applications (1994) 27–47doi:[10.1007/BF02191760](https://doi.org/10.1007/BF02191760).
URL <https://doi.org/10.1007/BF02191760>
- [72] H. Wolkowicz, Q. Zhao, [An all-inclusive efficient region of updates for least change secant methods](#), SIAM Journal on Optimization (1995) 172–191doi:[10.1137/0805009](https://doi.org/10.1137/0805009).
URL <https://doi.org/10.1137/0805009>
- [73] M. Al-Baali, H. Khalfan, [An overview of some practical quasi-newton methods for unconstrained optimization](#), Sultan Qaboos University Journal for Science 12 (2007) 199–209.
URL <https://api.semanticscholar.org/CorpusID:125566155>
- [74] H. Yabe, H. J. Martinez, R. A. Tapia, On sizing and shifting the bfgs update within the sized-broyden family of secant updates, SIAM Journal on Optimization 15 (2004) 139–160.
- [75] C. Wu, M. Zhu, Q. Tan, Y. Kartha, L. Lu, A comprehensive study of non-adaptive and residual-based adaptive sampling for physics-informed neural networks, Computer Methods in Applied Mechanics and Engineering (2023) 115671doi:<https://doi.org/10.1016/j.cma.2022.115671>.
- [76] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, in: Proceedings of the thirteenth international conference on artificial intelligence and statistics, JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.
- [77] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, Q. Zhang, JAX: Composable transformations of python+numpy programs, <https://github.com/google/jax>, deepMind, Google Research (2018).
- [78] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, J. Han, On the variance of the adaptive learning rate and beyond, arXiv preprint arXiv:1908.03265 (2019).
- [79] I. Babuschkin, K. Baumli, A. Bell, S. Bhupatiraju, J. Bruce, P. Buchlovsky, D. Budden, T. Cai, A. Clark, I. Danihelka, et al., [The deepmind jax ecosystem](#), DeepMind (2020).
URL <http://github.com/google-deepmind>
- [80] G. Hinton, [Neural networks for machine learning \(lecture 6e\)](#), Online lecture, coursera (2012).
URL <https://www.coursera.org/learn/neural-networks>
- [81] D. P. Kingma, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014).
- [82] X. Chen, et al., [Symbolic discovery of optimization algorithms](#), arXiv preprint arXiv:2302.06675 (2023).
URL <https://arxiv.org/abs/2302.06675>
- [83] J. Rader, T. Lyons, P. Kidger, Optimistix: modular optimisation in jax and equinox, arXiv:2402.09983 (2024).
- [84] S. Dong, N. Ni, [A method for representing periodic functions and enforcing exactly periodic boundary conditions with deep neural networks](#), Journal of Computational Physics (2021) 110242doi:<https://doi.org/10.1016/j.jcp.2021.110242>.
URL <https://www.sciencedirect.com/science/article/pii/S0021999121001376>
- [85] Z. Battles, L. N. Trefethen, [An extension of matlab to continuous functions and operators](#), SIAM J. Sci. Comput. (2004) 1743–1770.
URL <https://api.semanticscholar.org/CorpusID:14283334>
- [86] A.-K. Kassam, L. N. Trefethen, Fourth-order time-stepping for stiff pdes, SIAM Journal on Scientific Computing (2005) 1214–1233doi:[10.1137/S1064827502410633](https://doi.org/10.1137/S1064827502410633).

- [87] S. Wang, S. Sankaran, P. Perdikaris, [Respecting causality for training physics-informed neural networks](#), Computer Methods in Applied Mechanics and Engineering (2024) 116813doi:<https://doi.org/10.1016/j.cma.2024.116813>.
URL <https://www.sciencedirect.com/science/article/pii/S0045782524000690>
- [88] S. Wang, B. Li, Y. Chen, P. Perdikaris, PirateNets: Physics-informed Deep Learning with Residual Adaptive Networks, arXiv e-prints (2024). doi:[10.48550/arXiv.2402.00326](https://arxiv.org/abs/2402.00326).
- [89] E. N. Lorenz, Deterministic Nonperiodic Flow., Journal of the Atmospheric Sciences (1963) 130–148doi:[10.1175/1520-0469\(1963\)020<0130:DNF>2.0.CO;2](https://doi.org/10.1175/1520-0469(1963)020<0130:DNF>2.0.CO;2).
- [90] M. W. Hirsch, S. Smale, R. L. Devaney, Differential equations, dynamical systems, and an introduction to chaos, Academic press, 2013.