

# FORWARD PINN: UNIFIED COMPUTATION OF SOLUTIONS AND DERIVATIVES IN A SINGLE FORWARD PASS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Physics-Informed Neural Networks (PINNs) solve partial differential equations by embedding physical laws into their training process. A computational bottleneck, however, limits conventional PINNs. They rely on multiple backward passes to compute derivatives sequentially, a process that is memory-intensive and fails to leverage the parallelism of modern GPUs. To address this, we introduce Forward PINN, a framework that breaks from this computational architecture. Instead of relying on the backward pass, our innovation is to redesign the forward pass itself to perform differentiation. By exploiting the mathematical properties of specific activation functions, we unify the computation. The network’s output and all its necessary partial derivatives—first, second, and higher-order—are calculated concurrently within a single forward propagation. This approach effectively eliminates the need for multiple backward passes for derivative computation. We validated this new architecture on the two benchmark PDE problems: the two-dimensional heat equation and the anisotropic wave equation. The experimental results show that Forward PINN achieves accuracy comparable to its conventional counterparts while delivering performance speedups of 1.55× and 1.8× on the respective benchmark problems.

## 1 INTRODUCTION

Physics-Informed Neural Networks (PINNs) (Raissi et al., 2019) have emerged as a powerful paradigm for solving partial differential equations (PDEs) by embedding physical laws directly into neural network training objectives. The fundamental principle of PINNs relies on automatic differentiation to compute partial derivatives of the network output with respect to its inputs, enabling the enforcement of governing equations as soft constraints during training (Karniadakis et al., 2021).

However, the conventional implementation of PINNs faces significant computational challenges that limit their scalability and efficiency. The standard approach requires multiple sequential operations: first, a forward pass computes the network output  $u_\theta(\mathbf{x})$ ; subsequently, automatic differentiation is invoked to compute first-order derivatives  $\frac{\partial u}{\partial x_i}$ ; finally, additional backward passes compute second-order derivatives  $\frac{\partial^2 u}{\partial x_i \partial x_j}$  required by most PDEs (Lu et al., 2021). This sequential computation pattern presents critical limitations: **memory inefficiency** from maintaining separate computational graphs (Baydin et al., 2017), **sequential execution bottlenecks** that underutilize GPU cores and can exacerbate gradient flow issues during training (Wang et al., 2022), and **redundant graph traversals** leading to inefficient memory access patterns.

Recent efforts have explored various optimization strategies. Wang et al. (2021) proposed adaptive sampling techniques, while Krishnapriyan et al. (2021) investigated curriculum learning approaches. Self-adaptive approaches have addressed loss function balancing challenges (McClenny et al., 2023), demonstrating improved training stability through dynamic weight adjustment mechanisms. However, existing optimization methods do not address the fundamental computational inefficiency of derivative computation. More directly relevant, Bischof et al. (2008) explored forward-mode automatic differentiation for scientific computing, demonstrating theoretical advantages when the number of inputs exceeds outputs. However, these works exhibit significant limitations for PINN applications: they primarily focus on first-order derivatives while insufficiently addressing second-order

054 derivatives and higher-order derivatives essential for most PDE formulations, and lack specific ap-  
 055 plications to PINNs.

056 A fundamental gap remains: no prior work has successfully eliminated the sequential derivative  
 057 computation bottleneck while maintaining the mathematical rigor and generality of the original  
 058 PINN formulation.

059 In this paper, we introduce **Forward PINN**, a novel approach that reformulates the structure of  
 060 forward propagation in PINNs. Our key insight is that by leveraging the mathematical properties  
 061 of specific activation functions, we can compute network outputs and all required derivatives (first,  
 062 second or higher derivatives) simultaneously in a single forward pass. This eliminates the need  
 063 for multiple backward passes and their associated computational graphs, resulting in substantial  
 064 improvements in both memory efficiency and computational speed.

065 Our main contributions are: (1) A mathematical framework that enables simultaneous calculation  
 066 of solutions and derivatives in a single forward propagation process; (2) Experiments on two bench-  
 067 mark PDE problems demonstrating that Forward PINN achieves comparable accuracy to conven-  
 068 tional PINNs while significantly reducing training time.

## 071 2 METHOD

072 Our Forward PINN architecture employs the Exponential Linear Unit (ELU) (Clevert et al., 2015)  
 073 function as the activation function and achieves computational acceleration by rewriting the forward  
 074 function to compute both solutions and partial derivatives in a single forward pass. The specific  
 075 mathematical principles are as follows:

### 078 2.1 PROPERTY OF ELU FUNCTION

079 The Exponential Linear Unit (ELU) (Clevert et al., 2015) function has a property. The first-order  
 080 and second-order derivatives of this function can be directly represented by its own results, without  
 081 the need to input  $x$ . Denote the activation state as  $act$ . The specific representation is as follows:

$$\begin{aligned}
 \text{ELU}'(x) &= \begin{cases} 1 & \text{if } x \geq 0 \\ \alpha e^x & \text{if } x < 0 \end{cases} = act + (1 - act) \times (\text{ELU}(x) + \alpha) \\
 \text{ELU}''(x) &= \begin{cases} 0 & \text{if } x \geq 0 \\ \alpha e^x & \text{if } x < 0 \end{cases} = (1 - act) \times (\text{ELU}(x) + \alpha) \\
 act &= \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}
 \end{aligned}$$

091 Thus, the results of the first-order and second-order derivatives at each neuron can be directly repre-  
 092 sented through the results of the forward propagation. This property is used in our algorithm.

### 094 2.2 CALCULATE DERIVATIVES USING FORWARD PROPAGATION

095 The simple feedforward neural network shown in Figure 1 is used to illustrate our algorithm. The  
 096 feedforward neural network is a multilayer perceptron with two hidden layers.  $h_{11}, h_{12}, \dots, h_{23}$  are  
 097 the results of the activation function at each neuron. The activation function is set as ELU, which is  
 098 mentioned in section 2.1. Here,  $[h_{11}, h_{12}, h_{13}]^T$  is denoted as  $\mathbf{h}_1$ ,  $[h_{21}, h_{22}, h_{23}]^T$  is denoted as  $\mathbf{h}_2$ ,  
 099  $[x_1, x_2]^T$  is denoted as  $\mathbf{x}$  and  $[y_1, y_2]^T$  is denoted as  $\mathbf{y}$ . The parameters for the connections between  
 100 layers are defined as  $\mathbf{W}_1, \mathbf{b}_1$  for the first hidden layer,  $\mathbf{W}_2, \mathbf{b}_2$  for the second hidden layer,  $\mathbf{W}_3, \mathbf{b}_3$  for  
 101 the output layer. The process of forward propagation can be represented as follows:

$$102 \mathbf{h}_1 = \text{ELU}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \quad \mathbf{h}_2 = \text{ELU}(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2) \quad \mathbf{y} = \mathbf{W}_3 \mathbf{h}_2 + \mathbf{b}_3$$

#### 105 2.2.1 FIRST-ORDER DERIVATIVE

106 Take  $\frac{\partial \mathbf{y}}{\partial x_1}$  as an example. The first-order derivative of the ELU function at each neuron needs to  
 107 be obtained. Denote the first-order derivative at each neuron as  $h'_{11} \sim h'_{23}$ . Additionally, denote

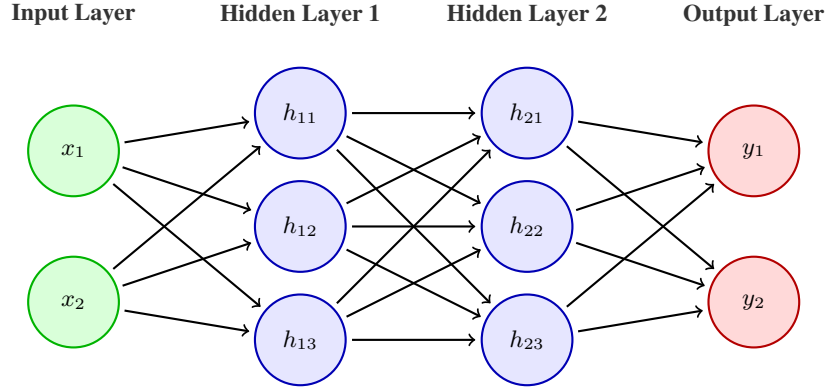


Figure 1: A multilayer perceptron architecture used to illustrate how to compute all the derivatives in a single forward pass. The network consists of an input layer with 2 neurons, two hidden layers with 3 neurons each, and an output layer with 2 neurons.

$[h'_{11}, h'_{12}, h'_{13}]^T$  as  $\mathbf{h}'_1$  and  $[h'_{21}, h'_{22}, h'_{23}]^T$  as  $\mathbf{h}'_2$ , the following can be derived:

$$\mathbf{h}'_1 = \text{ELU}'(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \quad \mathbf{h}'_2 = \text{ELU}'(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2)$$

$$\frac{\partial \mathbf{y}}{\partial x_1} = \mathbf{W}_3 \times \frac{\partial \mathbf{h}_2}{\partial x_1} = \mathbf{W}_3 \times (\mathbf{h}'_2 \odot (\mathbf{W}_2 \times \frac{\partial \mathbf{h}_1}{\partial x_1})) = \mathbf{W}_3 \times (\mathbf{h}'_2 \odot (\mathbf{W}_2 \times (\mathbf{h}'_1 \odot (\mathbf{W}_1 \times \begin{pmatrix} 1 \\ 0 \end{pmatrix}))))$$

**Note:**  $\odot$  is the Hadamard Product (Hadamard, 1893), specifically:

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \odot \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} x_1 y_1 \\ x_2 y_2 \\ x_3 y_3 \end{pmatrix}$$

As is shown in section 2.1,  $\mathbf{h}_1$  and  $\mathbf{h}_2$  can be used to represent  $\mathbf{h}'_1$  and  $\mathbf{h}'_2$ . Specifically, denote the activation state at each neuron as  $act_{11} \sim act_{23}$ , denote  $[act_{11}, act_{12}, act_{13}]^T$  as  $\mathbf{act}_1$  and  $[act_{21}, act_{22}, act_{23}]^T$  as  $\mathbf{act}_2$ , then  $\mathbf{h}'_1$  and  $\mathbf{h}'_2$  can be represented as follows:

$$\mathbf{h}'_1 = \mathbf{act}_1 + \left( \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \mathbf{act}_1 \right) \odot \left( \mathbf{h}_1 + \begin{pmatrix} \alpha \\ \alpha \end{pmatrix} \right) \quad \mathbf{h}'_2 = \mathbf{act}_2 + \left( \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \mathbf{act}_2 \right) \odot \left( \mathbf{h}_2 + \begin{pmatrix} \alpha \\ \alpha \end{pmatrix} \right)$$

In this way, the process of computing the first-order derivative of the output with respect to the input can be regarded as a kind of forward propagation process with the input being the standard basis vectors. Specifically, for a multilayer perceptron with  $n$  hidden layers, if the result of the  $i$ -th hidden layer is denoted as  $\mathbf{h}_i$ , the transfer function from the  $i$ -th layer to the  $(i+1)$ -th layer to compute the first-order derivative of the output to the input  $x_1$  can be defined as:

$$\frac{\partial \mathbf{h}_{i+1}}{\partial x_1} = \mathbf{h}'_{i+1} \odot (\mathbf{W}_{i+1} \times \frac{\partial \mathbf{h}_i}{\partial x_1})$$

Denote  $\frac{\partial \mathbf{h}_i}{\partial x_1}$  as  $\mathbf{s}_i$ , then the transfer function can be represented as:

$$\mathbf{s}_{i+1} = \mathbf{h}'_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{s}_i)$$

### 2.2.2 SECOND-ORDER DERIVATIVE

Chain rule is used to construct a recursive formula, thereby we can build a model for calculating the second-order derivatives in a forward propagation way. Here, take  $\frac{\partial^2 \mathbf{y}}{\partial x_1^2}$  as an example. With chain rule, the following can be derived:

$$\frac{\partial^2 \mathbf{y}}{\partial x_1^2} = \mathbf{W}_3 \times \frac{\partial^2 \mathbf{h}_2}{\partial x_1^2}$$

To calculate  $\frac{\partial^2 \mathbf{h}_2}{\partial x_1^2}$ , the second-order derivative of the ELU function at each neuron needs to be obtained. Similar to the representation used when computing the first-order derivative, the second-order derivative at each neuron is denoted as  $h''_{11} \sim h''_{23}$ . Additionally  $[h''_{11}, h''_{12}, h''_{13}]^T$  is denoted as  $\mathbf{h}''_1$ ,  $[h''_{21}, h''_{22}, h''_{23}]^T$  is denoted as  $\mathbf{h}''_2$ , then the following can be derived:

$$\mathbf{h}''_1 = \text{ELU}''(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \quad \mathbf{h}''_2 = \text{ELU}''(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2)$$

Shown in section 2.1,  $\mathbf{h}'_1$  and  $\mathbf{h}'_2$  can be represented with  $\mathbf{h}_1$  and  $\mathbf{h}_2$ , To be specific:

$$\mathbf{h}'_1 = \left( \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} - \text{act}_1 \right) \odot (\mathbf{h}_1 + \begin{pmatrix} \alpha \\ \alpha \\ \alpha \end{pmatrix}) \quad \mathbf{h}'_2 = \left( \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} - \text{act}_2 \right) \odot (\mathbf{h}_2 + \begin{pmatrix} \alpha \\ \alpha \\ \alpha \end{pmatrix})$$

As is shown in section 2.2.1:

$$\frac{\partial \mathbf{h}_2}{\partial x_1} = \mathbf{h}'_2 \odot (\mathbf{W}_2 \times \frac{\partial \mathbf{h}_1}{\partial x_1}) = \text{ELU}'(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2) \odot (\mathbf{W}_2 \times \frac{\partial \mathbf{h}_1}{\partial x_1})$$

According to the Leibniz formula, the following can be derived:

$$\frac{\partial^2 \mathbf{h}_2}{\partial x_1^2} = \text{ELU}''(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2) \odot (\mathbf{W}_2 \times \frac{\partial \mathbf{h}_1}{\partial x_1}) \odot (\mathbf{W}_2 \times \frac{\partial \mathbf{h}_1}{\partial x_1}) + \text{ELU}'(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2) \odot (\mathbf{W}_2 \times \frac{\partial^2 \mathbf{h}_1}{\partial x_1^2})$$

Denote  $\frac{\partial^2 \mathbf{h}_i}{\partial x_1^2}$  as  $\mathbf{t}_i$ , and  $\frac{\partial \mathbf{h}_i}{\partial x_1}$  as  $\mathbf{s}_i$ , then the formula above can be denoted as follows:

$$\mathbf{t}_2 = \mathbf{h}''_2 \odot (\mathbf{W}_2 \times \mathbf{s}_1) \odot (\mathbf{W}_2 \times \mathbf{s}_1) + \mathbf{h}'_2 \odot (\mathbf{W}_2 \times \mathbf{t}_1)$$

Generally, the transfer function from the  $i$ -th layer to the  $(i+1)$ -th layer to compute the second-order derivative of the output to the input  $x_1$  can be defined as:

$$\mathbf{t}_{i+1} = \mathbf{h}''_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{s}_i) \odot (\mathbf{W}_{i+1} \times \mathbf{s}_i) + \mathbf{h}'_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{t}_i)$$

Similarly, To calculate  $\frac{\partial^2 \mathbf{y}}{\partial x_1 \partial x_2}$ , the following can be derived:

$$\frac{\partial^2 \mathbf{h}_2}{\partial x_1 \partial x_2} = \text{ELU}''(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2) \odot (\mathbf{W}_2 \times \frac{\partial \mathbf{h}_1}{\partial x_2}) \odot (\mathbf{W}_2 \times \frac{\partial \mathbf{h}_1}{\partial x_1}) + \text{ELU}'(\mathbf{W}_2 \mathbf{h}_1 + \mathbf{b}_2) \odot (\mathbf{W}_2 \times \frac{\partial^2 \mathbf{h}_1}{\partial x_1 \partial x_2})$$

Denote  $\frac{\partial^2 \mathbf{h}_i}{\partial x_1 \partial x_2}$  as  $\mathbf{t}_i$ ,  $\frac{\partial \mathbf{h}_i}{\partial x_1}$  as  $\mathbf{s}_{i1}$ ,  $\frac{\partial \mathbf{h}_i}{\partial x_2}$  as  $\mathbf{s}_{i2}$ , then the transfer function from the  $i$ -th layer to the  $(i+1)$ -th layer to compute the second-order derivative of the output to the input  $x_1$  and  $x_2$  can be defined as:

$$\mathbf{t}_{i+1} = \mathbf{h}''_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{s}_{i2}) \odot (\mathbf{W}_{i+1} \times \mathbf{s}_{i1}) + \mathbf{h}'_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{t}_i)$$

### 2.2.3 SUMMARY OF THE FORWARD ALGORITHM

Above all, the process of computing first-order derivative and second-order derivative can be regarded as a kind of forward propagation. Specifically, define the output of the  $i$ -th hidden layer as  $\mathbf{h}_i$ . And denote  $\frac{\partial \mathbf{h}_i}{\partial x_j}$  as  $\mathbf{s}_{ij}$ ,  $\frac{\partial^2 \mathbf{h}_i}{\partial x_j \partial x_k}$  as  $\mathbf{t}_{ijk}$ , the forward propagation between the  $i$ -th layer and the  $(i+1)$ -th layer ( $i \geq 0$ ) can be defined as:

$$\mathbf{h}_{i+1} = \text{ELU}(\mathbf{W}_{i+1} \mathbf{h}_i + \mathbf{b}_{i+1})$$

$$\mathbf{h}'_{i+1} = \text{act}_{i+1} + (\mathbf{1}_{i+1} - \text{act}_{i+1}) \odot (\mathbf{h}_{i+1} + \alpha \times \mathbf{1}_{i+1})$$

$$\mathbf{h}''_{i+1} = (\mathbf{1}_{i+1} - \text{act}_{i+1}) \odot (\mathbf{h}_{i+1} + \alpha \times \mathbf{1}_{i+1})$$

$$\mathbf{s}_{(i+1)j} = \mathbf{h}'_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{s}_{ij})$$

$$\mathbf{t}_{(i+1)jk} = \mathbf{h}''_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{s}_{ij}) \odot (\mathbf{W}_{i+1} \times \mathbf{s}_{ik}) + \mathbf{h}'_{i+1} \odot (\mathbf{W}_{i+1} \times \mathbf{t}_{ijk})$$

$\mathbf{1}_i$  is the vector of ones, of which the dimension is the same as  $\mathbf{h}_i$ . Specifically, for  $i = 0$ ,  $\mathbf{h}_0 = \mathbf{x}$ ,  $\mathbf{s}_{0j} = \mathbf{e}_j$ ,  $\mathbf{t}_{0jk} = \mathbf{0}$ .  $\mathbf{x}$  is the input of the MLP,  $\mathbf{e}_j$  is the standard basis vector with 1 in the  $j$ -th position and 0 elsewhere.  $\mathbf{0}$  is the zero vector (all elements are 0).

With these formulas, we can calculate the output, the first-order derivative and the second-order derivative at the same time, just using forward propagation.

**Remark 1:** The forward propagation framework presented above is not limited to the ELU activation function. Any activation function with known analytical expressions for its first-order derivative  $h'_{i+1}$  and second-order derivative  $h''_{i+1}$  can be incorporated into this framework, enabling efficient derivative computation.

**Remark 2:** While we have demonstrated the computation of first- and second-order derivatives, the proposed method can be extended to higher-order derivatives. By applying the chain rule recursively to the second-order derivative terms  $t_{(i+1)jk}$ , one can derive the update formulas for third-order and higher-order derivatives following the same forward propagation principle.

### 2.3 ADVANTAGE ANALYSIS

The conventional Physics-Informed Neural Network(PINN)(Raissi et al., 2019) employs a fully-connected neural network to approximate the solution  $u$ , and record the computation graph of forward propagation at the same time. Then the computation graph is used to calculate the first-order derivative. Similarly, the computation graph is recorded to enable second-order differentiation while computing the first-order derivatives. The overall computation graph is shown in Figure 2 (a). The computation graph illustrates a propagation process for a multilayer perceptron with three hidden layers when the conventional PINN approach is employed. It shows how the output, first-order derivatives and second-order derivatives are computed through forward and backward propagation.

The Forward PINN also employs a fully-connected neural network to approximate the solution  $u$ . However, it doesn't need computation graph to calculate first-order or second-order derivatives. Derivatives can be calculated just using forward propagation. The overall computation graph is shown in Figure 2 (b).

Comparing the two computation graphs, the advantages of Forward PINN are as follows:

1. The propagation path to calculate second-order derivative is much shorter, approximately half the length of that required by conventional algorithms.
2. Forward PINN stores only a single computational graph for network-parameter optimization, whereas the conventional algorithm must maintain at least three separate graphs: one for calculating first-order derivatives, one for calculating second-order derivatives and one for parameters updation. This approach substantially reduces memory consumption.
3. Forward PINN is ideally suited to the GPU's parallel architecture(NVIDIA Corporation, 2020): its computational density(the ratio of arithmetic operations to global memory accesses) is significantly higher. Because the output, first-order derivatives and second-order derivatives are computed simultaneously within the same forward kernel.

In contrast, the conventional PINN compute these quantities sequentially, which forces multiple, memory-intensive passes over the network. This sequential pattern under-utilizes GPU cores and incurs repeated global-memory traffic, making it inherently less efficient on modern accelerators. As noted by Hooker (2021), the evolution of machine learning algorithms has been significantly shaped by hardware constraints, particularly the parallel architecture of GPUs, emphasizing the importance of designing algorithms that align with hardware capabilities rather than fighting against them.

## 3 EXPERIMENT DESIGN

The experiments were conducted in two scenarios: one based on the two-dimensional heat equation, and the other based on the anisotropic wave equation.

### 3.1 SCENARIO ONE

Scenario One is based on two-dimensional heat equation(Raissi et al., 2019; Lu et al., 2021), the physical constraints can be described as follows:

$$\frac{\partial u}{\partial t} = \alpha \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right)$$

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

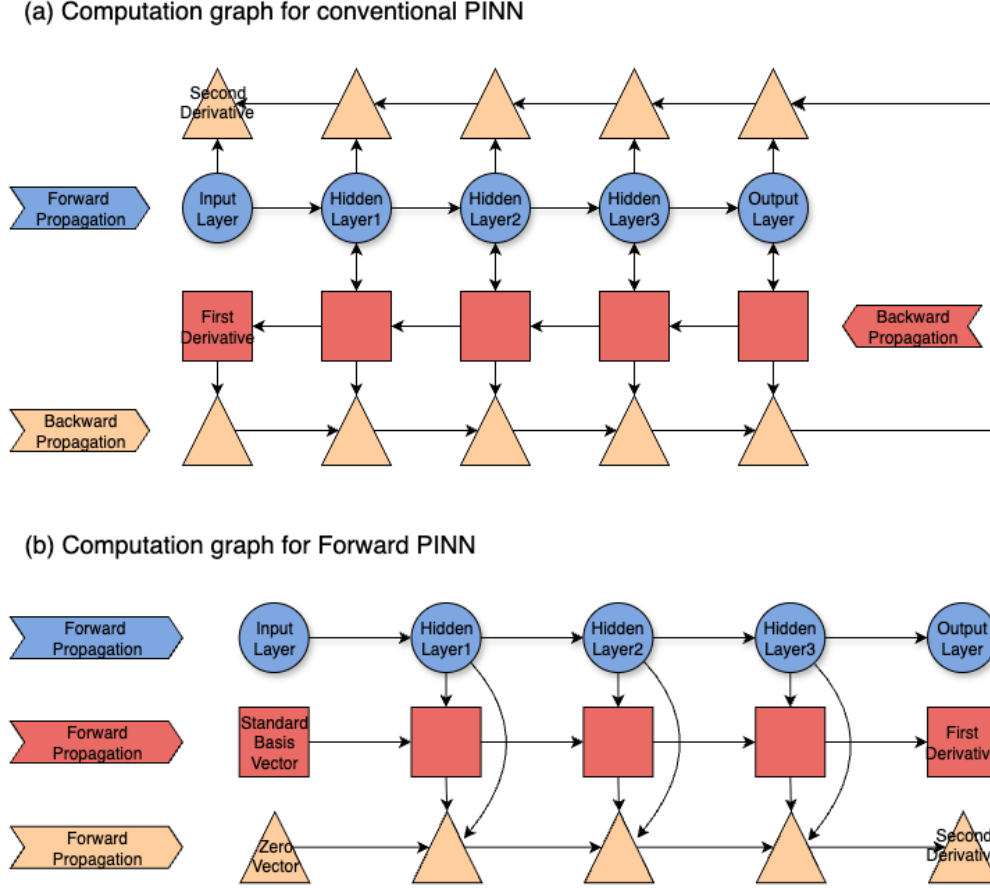


Figure 2: Computation graph to help analyze the advantages: (a) for conventional PINN and (b) for Forward PINN. Circles denote the nodes during the forward process to compute solutions. Squares represent the nodes where the first-order derivatives are computed. Triangles mark the nodes involved in computing the second-order derivatives. The computation graph ignores the details of operations, instead emphasizes the routes taken by forward or backward propagation. Arrows between nodes indicate the direction of data flow.

$u(t,x,y)$  is the temperature distribution function,  $\alpha$  is the thermal diffusivity constant, which is set as 0.01 in this experiment. The domain is defined as  $(t, x, y) \in [0, 1] \times [0, 1] \times [0, 1]$ .

Experimental data were generated using the analytical solution for diffusion from a Gaussian heat source (Polyanin & Zaitsev, 2002; Özışık, 1993). The expression for the solution is as follows:

$$u(t, x, y) = \frac{1}{2\pi\sigma_t^2} \exp\left(-\frac{(x-x_0)^2 + (y-y_0)^2}{2\sigma_t^2}\right)$$

$(x_0, y_0) = (0.5, 0.5)$  is the center position of the heat source,  $\sigma_t = \sqrt{\sigma_0^2 + 2\alpha t}$  is the time-dependent diffusion width,  $\sigma_0 = 0.1$  is the initial width of the heat source.

Specifically, 20 time points are uniformly selected in the interval (0, 1). For each time point, a  $20 \times 20$  grid is constructed to cover the region  $(0, 1) \times (0, 1)$ , resulting in 8,000 spatiotemporal points. For each spatiotemporal point, the analytical solution of the Gaussian heat source diffusion is used as the ground truth, and maximum normalization is applied to construct the training set.

### 3.2 SCENARIO TWO

Scenario Two is based on anisotropic wave equation (Raissi et al., 2019; McClenny et al., 2023), the physical constraints can be described as follows:

$$\frac{\partial^2 u}{\partial t^2} = c_1^2 \times \frac{\partial^2 u}{\partial x^2} + c_2^2 \times \frac{\partial^2 u}{\partial y^2} + c_3^2 \times \frac{\partial^2 u}{\partial x \partial y}$$

$u(t, x, y)$  is the wave field function,  $c_1 = 1.0$  is the wave speed parameter in the  $x$  direction,  $c_2 = 0.8$  is the wave speed parameter in the  $y$  direction, and  $c_3 = 0.3$  is the coupling coefficient for the cross term. The domain is defined as  $(t, x, y) \in [0, 1] \times [-1, 1] \times [-1, 1]$ .

Experimental data were generated using the analytical solution for anisotropic wave propagation from a Gaussian wave packet (Brillouin, 1946; Gabor, 1946). The mathematical expression for the solution is as follows:

$$u(t, x, y) = \exp\left(-\frac{1}{2}r^2\right) \times \sin(\omega t - k_1 x - k_2 y - k_3 xy)$$

where  $r = \sqrt{(c_1 x)^2 + (c_2 y)^2 + c_3 xy}$  is the anisotropic distance function,  $\omega = \sqrt{c_1^2 + c_2^2 + c_3^2}$  is the angular frequency,  $k_1 = c_1$ ,  $k_2 = c_2$ ,  $k_3 = c_3$  are the wave numbers in different directions.

Similar to Scenario One, 20 time points are uniformly selected in the interval  $(0, 1)$ . For each time point, a  $20 \times 20$  grid is constructed to cover the region  $(-1, 1) \times (-1, 1)$ , resulting in 8,000 spatiotemporal points. For each spatiotemporal point, the analytical solution of the anisotropic wave propagation is used as the ground truth, and maximum normalization is applied to construct the training set.

### 3.3 COMPARISON EXPERIMENT DESIGN

To compare the performance of the conventional PINN and the proposed Forward PINN, both methods are evaluated under the two scenarios described. For each scenario, the same feedforward neural network, training data, and optimizer settings are used to ensure a fair comparison between methods.

Specifically, a feedforward neural network with an input dimension of 3, an output dimension of 1, four hidden layers each containing 64 neurons, and ELU activation functions is used to fit the input and output of the training set. Both methods use the Adam optimizer with a learning rate of 0.001.

In the PyTorch framework, both the conventional PINN and our proposed Forward PINN methods were implemented and trained on a single NVIDIA A100 GPU for 50 iterations. For the Standard PINN, automatic differentiation provided by PyTorch was used to compute the required first and second-order derivatives during training. This involves recording the computation graph during the forward pass and performing multiple backward passes to obtain the derivatives. For the Forward PINN, the forward propagation formulas described in Section 2.2 were implemented to compute the output, first-order, and second-order derivatives simultaneously in a single forward pass, without relying on PyTorch’s autograd for derivatives. The pseudocode for both algorithms is in Appendix B.

## 4 RESULTS

This section presents a comprehensive evaluation of the Forward PINN framework. The empirical analysis is structured to validate three critical aspects of the proposed method: the numerical accuracy of its derivative computations, its training efficiency relative to conventional PINNs, and the quality of the final solutions it produces.

### 4.1 DERIVATIVE COMPUTATION ACCURACY VALIDATION

A preliminary validation was conducted to rigorously assess the numerical accuracy of the proposed forward-mode derivative computation. At each step of the training process, the first and second-order partial derivatives computed via the Forward PINN algorithm were quantitatively compared against those obtained from PyTorch’s established automatic differentiation (autograd) engine.

The comparative analysis revealed excellent agreement between the two methods. For derivative values typically in the range of  $10^{-2}$  to  $10^{-3}$ , the discrepancies between our Forward PINN method and PyTorch’s `autograd` were consistently on the order of  $10^{-9}$ , demonstrating machine-precision accuracy. This outcome not only confirms the mathematical integrity and correct implementation of our approach but also validates that our method produces results virtually identical to the established automatic differentiation framework.

#### 4.2 TRAINING EFFICIENCY COMPARISON

The training efficiency of the Forward PINN and the conventional PINN was compared in terms of convergence speed and training time. For the convergence speed, the loss curves of both methods for two scenarios in the first 50 training iterations are shown in Figure 3. The loss curves indicate that both methods converge at a similar speed, confirming the correctness of our algorithm.

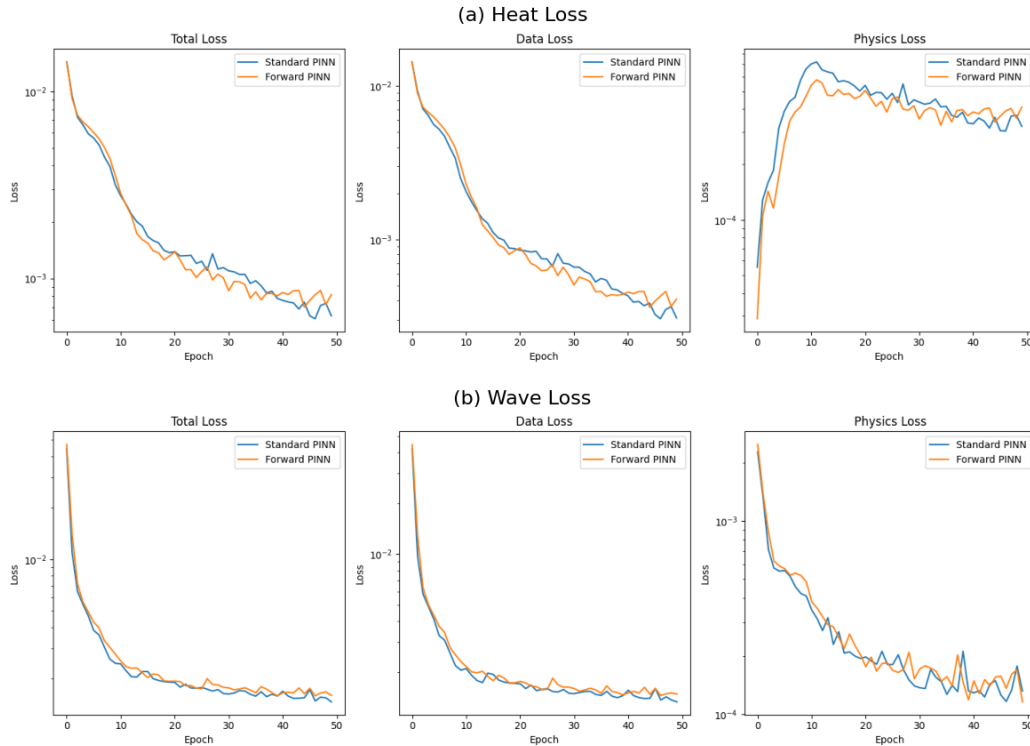


Figure 3: Loss curves comparing the two methods: (a) for Scenario One and (b) for Scenario Two. Both methods demonstrate similar convergence behavior.

For the training time, the training time for each epoch and the total training time for both methods across two scenarios are shown in Figure 4

For Scenario One, the total training time for the Forward PINN is 32.32seconds, while the conventional PINN takes 50.24 seconds, indicating a significant speedup of approximately 1.55 times. For Scenario Two, the total training time for the Forward PINN is 42.85 seconds, while the conventional PINN takes 77.21 seconds, indicating a significant speedup of approximately 1.8 times. The results demonstrate that the Forward PINN has a significant advantage in training efficiency.

#### 4.3 TRAINING EFFECT COMPARISON

To ensure that the observed efficiency gains did not compromise the quality of the final solution, the predictive accuracy of both methods was evaluated against the known analytical solutions. The

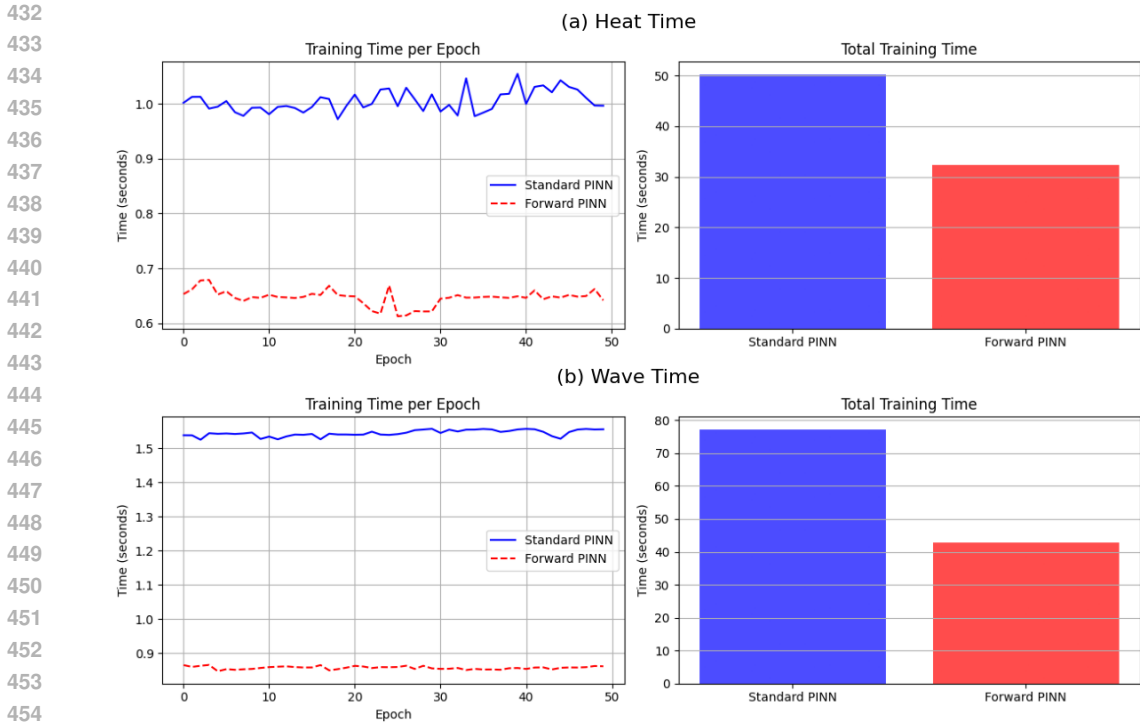


Figure 4: Training time for each epoch and total training time for both methods: (a) for Scenario One and (b) for Scenario Two. The Forward PINN is approximately 1.5 times faster than the conventional PINN in Scenario One and 1.8 times faster in Scenario Two.

temporal evolution of the prediction error was visualized through animated GIFs (available in the supplementary materials), with key time-slice comparisons provided in Appendix C.

This qualitative analysis confirms that both the Forward PINN and the conventional PINN produce solutions of a similar high fidelity when compared to the ground truth. The error distributions over the spatiotemporal domain are visually indistinguishable between the two methods, affirming that the Forward PINN architecture successfully learns the physical dynamics with an accuracy comparable to its traditional counterpart.

## 5 CONCLUSION

In this paper, we addressed a critical computational bottleneck in conventional Physics-Informed Neural Networks (PINNs) arising from their reliance on sequential backward passes for derivative computation. We introduced **Forward PINN**, a novel framework that reformulates the network’s forward pass to compute the solution and its required partial derivatives simultaneously. By leveraging the analytical properties of activation functions like ELU, our method eliminates the need for multiple, memory-intensive autograd calls, creating a more efficient, unified computational pipeline.

Our experimental results on the 2D heat and anisotropic wave equations validate the effectiveness of our approach. We demonstrated that Forward PINN achieves accuracy and convergence rates comparable to conventional PINNs while delivering significant performance improvements, with training speedups of up to  $1.8\times$ . These findings confirm that our method successfully reduces computational overhead without compromising the accuracy or training stability of the PINN framework.

The primary implication of this work is a more scalable and hardware-efficient formulation of PINNs. By reducing memory footprint and aligning the computational pattern with the parallel architecture of modern GPUs, Forward PINN paves the way for tackling larger and more complex scientific machine learning problems.

## REFERENCES

- 486  
487  
488 Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind.  
489 Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research*,  
490 18(1):5595–5637, 2017.
- 491 Christian H Bischof, Alan Carle, Peyvand Khademi, and Andrew Mauer. Automatic differentiation  
492 for inverse problems in x-ray astronomy. *Computational Science–ICCS 2008*, pp. 584–593, 2008.
- 493  
494 Leon Brillouin. *Wave propagation and group velocity*. Academic Press, 1946. Classic work on  
495 wave packet propagation.
- 496 Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network  
497 learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 4(5):11, 2015.
- 498  
499 Dennis Gabor. Theory of communication. part 1: The analysis of information. *Journal of the*  
500 *Institution of Electrical Engineers-Part III: Radio and Communication Engineering*, 93(26):429–  
501 441, 1946. Foundation of Gaussian wave packet theory.
- 502 Jacques Hadamard. Résolution d’une question relative aux déterminants. *Bulletin des Sciences*  
503 *Mathématiques*, 17:240–246, 1893. Original work introducing the Hadamard product (element-  
504 wise multiplication).
- 505  
506 Sara Hooker. The hardware lottery. *Communications of the ACM*, 64(12):58–65, 2021.
- 507 George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang.  
508 Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021.
- 509  
510 Aditi Krishnapriyan, Amir Gholami, Shandian Zhe, Robert Kirby, and Michael W Mahoney. Char-  
511 acterizing possible failure modes in physics-informed neural networks. *Advances in Neural In-*  
512 *formation Processing Systems*, 34:26548–26560, 2021.
- 513 Lu Lu, Xuhui Meng, Zhiping Mao, and George Em Karniadakis. Deepxde: A deep learning library  
514 for solving differential equations. *SIAM Review*, 63(1):208–228, 2021.
- 515  
516 Levi McClenny, Mulugeta A Haile, and Ulisses M Braga-Neto. Self-adaptive physics-informed  
517 neural networks. *Journal of Computational Physics*, 474:111722, 2023.
- 518 NVIDIA Corporation. Cuda c++ programming guide. *NVIDIA Developer Documentation*, 2020.  
519 Official CUDA programming guide for GPU parallel computing.
- 520  
521 M Necati Özışık. *Heat conduction*. John Wiley & Sons, 2nd edition, 1993.
- 522  
523 Andrei D Polyaniin and Valentin F Zaitsev. *Handbook of linear partial differential equations for*  
524 *engineers and scientists*. Chapman and Hall/CRC, 2002. Contains analytical solutions for heat  
525 equations with various source terms.
- 526  
527 Maziar Raissi, Paris Perdikaris, and George Em Karniadakis. Physics-informed neural networks:  
528 A deep learning framework for solving forward and inverse problems involving nonlinear partial  
529 differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- 530  
531 Sifan Wang, Yujun Teng, and Paris Perdikaris. Understanding and mitigating gradient flow patholo-  
532 gies in physics-informed neural networks. *SIAM Journal on Scientific Computing*, 43(5):A3055–  
533 A3081, 2021.
- 534  
535 Sifan Wang, Xinling Yu, and Paris Perdikaris. When and why pinns fail to train: A neural tangent  
536 kernel perspective. *Journal of Computational Physics*, 449:110768, 2022.
- 537  
538  
539

## 540 A LLM USAGE DISCLOSURE

541 In accordance with ICLR 2026 guidelines, we disclose the use of (LLMs) in this research:

542 **Writing Assistance:** LLMs were used to assist in drafting portions of the abstract and introduction  
 543 sections. The authors provided the core research ideas, technical content, and scientific contribu-  
 544 tions, while LLMs helped with language refinement, structure organization, and clarity improve-  
 545 ment. 546  
 547

548 **Code Development:** LLMs were utilized to generate code components for:

- 549 • Implementation of the conventional PINN baseline model
- 550 • Data visualization scripts for comparative analysis between Forward PINN and conven-  
 551 tional PINN
- 552 • Utility functions for performance benchmarking and result plotting

553 All generated code was thoroughly reviewed, tested and modified by the authors to ensure correct-  
 554 ness and alignment with the research objectives. The core algorithmic contributions, mathematical  
 555 derivations, and experimental design were conceived and implemented by the human authors. 556  
 557

## 558 B PSEUDOCODE FOR BOTH ALGORITHMS

### 559 Conventional PINN:

- 560 1. Initialize neural network.
- 561 2. For each training iteration:
  - 562 (a) Perform forward propagation to compute the network output  $u_\theta(\mathbf{x})$ .
  - 563 (b) Use automatic differentiation (torch.autograd.grad() function, with both create\_graph  
 564 and retain\_graph parameters set to True) to compute the first-order derivatives  $\frac{\partial u}{\partial t}$ ,  $\frac{\partial u}{\partial x}$ ,  
 565 and  $\frac{\partial u}{\partial y}$ . Here manual optimization is used to calculate the three derivatives simulta-  
 566 neously.
  - 567 (c) Use automatic differentiation again to compute the required second-order derivatives
  - 568 (d) Compute the loss function using the solutions and derivatives.
  - 569 (e) Update parameters  $\theta$  using the optimizer.

### 570 Forward PINN:

571 Note: To reduce memory usage, the pseudocode occasionally reuses the same variable name on both  
 572 sides of the assignment, here the equals sign denotes programming assignment, not mathematical  
 573 equality. 574

- 575 1. Initialize neural network.
- 576 2. Redefine forward function of the network to compute output, first-order derivatives, and  
 577 second-order derivatives using the formulas described in Section 2.2 at the same time:
  - 578 (a) Get the input and reshape them into the vector form  $[t, x, y]^T$ .
  - 579 (b) Create input for first-order derivative calculation:

$$580 \mathbf{e}_1 = [1, 0, 0]^T, \quad \mathbf{e}_2 = [0, 1, 0]^T, \quad \mathbf{e}_3 = [0, 0, 1]^T$$

581 where  $\mathbf{e}_1$  is for  $\frac{\partial u}{\partial t}$ ,  $\mathbf{e}_2$  is for  $\frac{\partial u}{\partial x}$ , and  $\mathbf{e}_3$  is for  $\frac{\partial u}{\partial y}$ .

- 582 (c) Create input for second-order derivative calculation:  $\mathbf{z} = [0, 0, 0]^T$
- 583 (d) Let  $\mathbf{u} = [t, x, y]^T$ ,  $\mathbf{u}_t = \mathbf{e}_1$ ,  $\mathbf{u}_x = \mathbf{e}_2$ ,  $\mathbf{u}_y = \mathbf{e}_3$ ,  $\mathbf{u}_{xx} = \mathbf{u}_{xy} = \mathbf{u}_{yy} = \mathbf{u}_{tt} = \mathbf{z}$ .
- 584 (e) For each layer  $i$  in the network:
  - 585 i. Get the weight matrix  $\mathbf{W}_i$  and bias vector  $\mathbf{b}_i$  for the layer.
  - 586 ii. Compute the pre-activation values:  $\mathbf{z}_i = \mathbf{W}_i \mathbf{u} + \mathbf{b}_i$
  - 587 iii. Apply ELU activation:  $\mathbf{u} = \text{ELU}(\mathbf{z}_i)$

594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

iv. Construct activation state vector by element-wise comparison with zero:

$$\mathbf{act}_i[j] = \begin{cases} 1 & \text{if } z_i[j] \geq 0 \\ 0 & \text{if } z_i[j] < 0 \end{cases}$$

where  $j$  indexes each dimension of the layer output.

v. Compute first-order derivative coefficients:

$$\mathbf{u}'_i = \mathbf{act}_i + (\mathbf{1} - \mathbf{act}_i) \odot (\mathbf{u} + \alpha \cdot \mathbf{1})$$

vi. Compute second-order derivative coefficients:

$$\mathbf{u}''_i = (\mathbf{1} - \mathbf{act}_i) \odot (\mathbf{u} + \alpha \cdot \mathbf{1})$$

vii. Multiply first-order derivatives with weight matrix:

$$\mathbf{u}_t = (\mathbf{W}_i \times \mathbf{u}_t), \quad \mathbf{u}_x = (\mathbf{W}_i \times \mathbf{u}_x), \quad \mathbf{u}_y = (\mathbf{W}_i \times \mathbf{u}_y)$$

viii. Update second-order derivatives:

$$\mathbf{u}_{tt} = \mathbf{u}''_i \odot \mathbf{u}_t \odot \mathbf{u}_t + \mathbf{u}'_i \odot (\mathbf{W}_i \times \mathbf{u}_{tt})$$

$$\mathbf{u}_{xx} = \mathbf{u}''_i \odot \mathbf{u}_x \odot \mathbf{u}_x + \mathbf{u}'_i \odot (\mathbf{W}_i \times \mathbf{u}_{xx})$$

$$\mathbf{u}_{yy} = \mathbf{u}''_i \odot \mathbf{u}_y \odot \mathbf{u}_y + \mathbf{u}'_i \odot (\mathbf{W}_i \times \mathbf{u}_{yy})$$

$$\mathbf{u}_{xy} = \mathbf{u}''_i \odot \mathbf{u}_x \odot \mathbf{u}_y + \mathbf{u}'_i \odot (\mathbf{W}_i \times \mathbf{u}_{xy})$$

ix. Update the first-order derivatives:

$$\mathbf{u}_t = \mathbf{u}'_i \odot \mathbf{u}_t$$

$$\mathbf{u}_x = \mathbf{u}'_i \odot \mathbf{u}_x$$

$$\mathbf{u}_y = \mathbf{u}'_i \odot \mathbf{u}_y$$

(f) Extract final outputs: solution  $\mathbf{u}$ , and derivatives  $\mathbf{u}_t, \mathbf{u}_x, \mathbf{u}_y, \mathbf{u}_{xx}, \mathbf{u}_{yy}, \mathbf{u}_{xy}, \mathbf{u}_{tt}$ .

3. For each training iteration:

- (a) Call the redefined forward function to get solutions and all required derivatives.
- (b) Compute the physics-informed loss function using the solutions and derivatives.
- (c) Update network parameters using backward propagation.

## C THE SLICES OF THE GIF ANIMATION WHICH IS CREATED TO SHOW THE EVOLUTION OF THE PREDICTED SOLUTION ERROR OVER TIME

These two figures demonstrate that Forward PINN and conventional PINN achieve very similar performance on the test set, providing additional evidence for the correctness of our proposed algorithm.

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

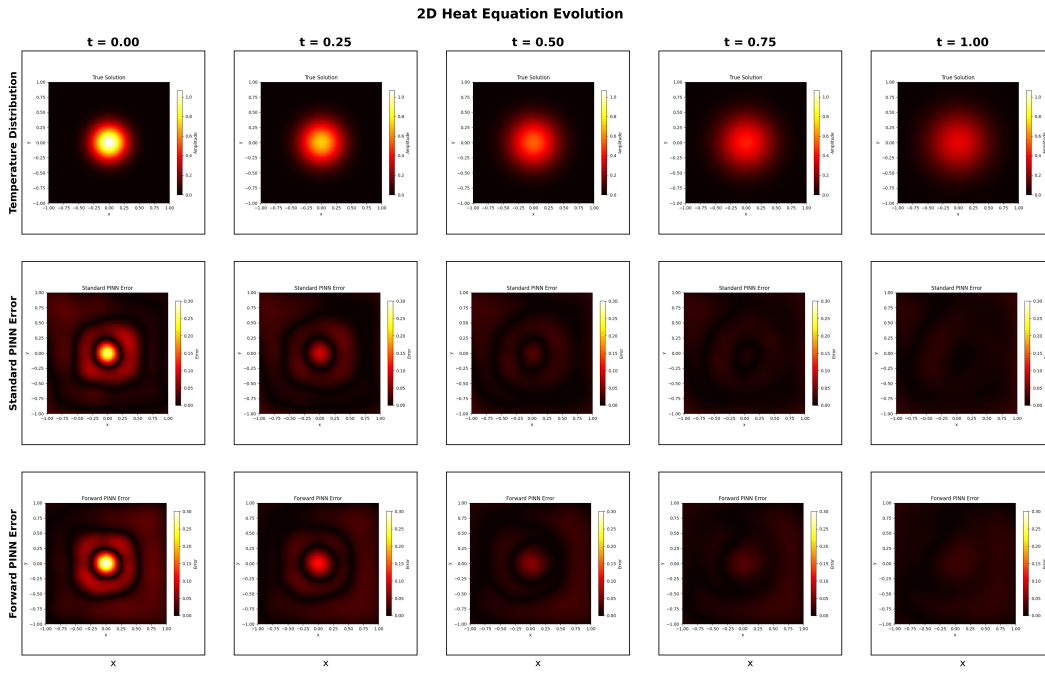


Figure 5: Slices of the GIF animation showing the predicted solution error at different time points for Scenario One. The ground truth is shown in the first row, the predicted solution error by the Standard PINN is shown in the second row, and the predicted solution error by the Forward PINN is shown in the third row. The error is shown in the fourth row. The error is calculated as the absolute difference between the predicted solution and the ground truth.

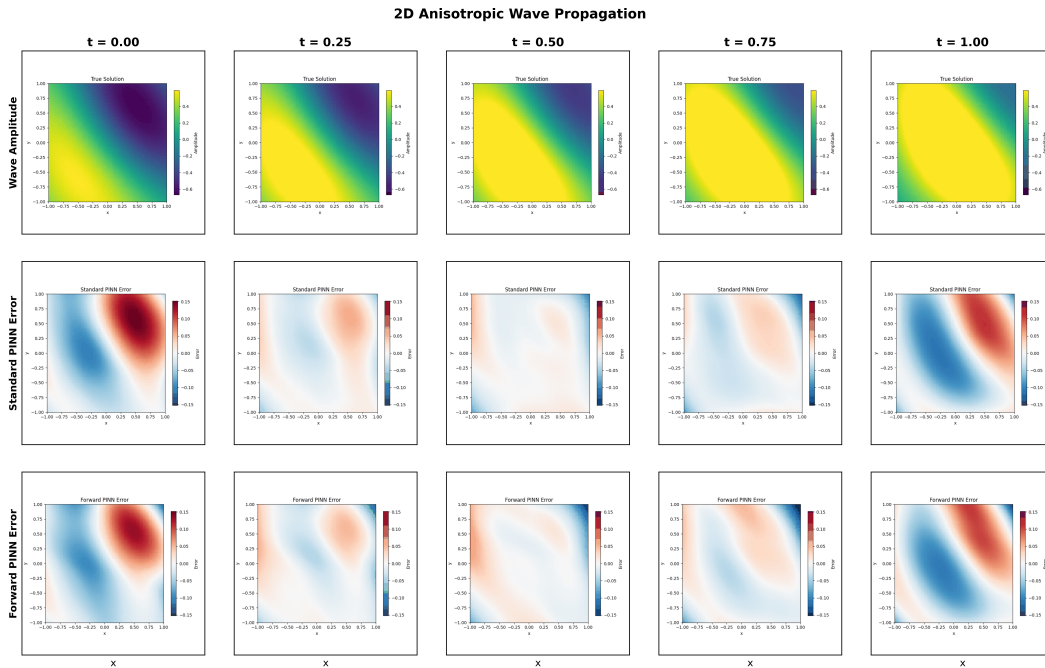


Figure 6: Slices of the GIF animation showing the predicted solution error at different time points for Scenario Two. The visualization is presented in the same way as in Figure 5.