



OPEN ACCESS

EDITED BY

Lev Shchur,
National Research University Higher School of
Economics, Russia

REVIEWED BY

Ali Mehri,
Babol Noshirvani University of
Technology, Iran
Zhao Li,
Chengdu University, China

*CORRESPONDENCE

Min He,
✉ hemin20052005@163.com

RECEIVED 07 January 2025

ACCEPTED 13 February 2025

PUBLISHED 19 March 2025

CITATION

Xie Q, He M and Lu Z (2025) Application of
computer vision based nonlinear physical
system dynamic behavior analysis in
education.

Front. Phys. 13:1556601.

doi: 10.3389/fphy.2025.1556601

COPYRIGHT

© 2025 Xie, He and Lu. This is an open-access
article distributed under the terms of the
[Creative Commons Attribution License \(CC
BY\)](#). The use, distribution or reproduction in
other forums is permitted, provided the
original author(s) and the copyright owner(s)
are credited and that the original publication
in this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

Application of computer vision based nonlinear physical system dynamic behavior analysis in education

Qiuen Xie¹, Min He^{2*} and Zhang Lu³

¹School of Foreign Studies, Jiangxi University of Science and Technology, Ganzhou, Jiangxi, China,

²School of Social Management, Jiangxi College of Applied Technology, Ganzhou, Jiangxi, China,

³Jiangmen Polytechnic, Guangdong Jiangmen, Jiangmen, China

Introduction: The dynamic behavior analysis of nonlinear physical systems plays a critical role in understanding complex processes across various domains, including education, where interactive simulations of such systems can enhance conceptual learning. Traditional modeling techniques for nonlinear systems often fail to capture their high-dimensional, multi-scale, and chaotic nature due to oversimplified assumptions or reliance on linear approximations.

Methods: In this study, we present a novel framework leveraging computer vision and advanced neural architectures to analyze the dynamic behaviors of nonlinear physical systems. The proposed Physics-Informed Nonlinear Dynamics Network (PNDN) integrates data-driven embeddings with physics-based constraints, offering a robust solution for capturing intricate dynamics and ensuring adherence to physical principles.

Results: Experimental results highlight the model's superior performance in reconstructing and predicting nonlinear system behaviors under diverse conditions, establishing its utility for real-time educational simulations.

Discussion: This approach bridges the gap between computational modeling and educational innovation, providing learners with interactive tools to explore complex physical phenomena.

KEYWORDS

nonlinear physical systems, dynamic behavior analysis, computer vision, education, physics-informed neural networks

1 Introduction

The integration of computer vision in analyzing the dynamic behavior of nonlinear physical systems represents a significant advancement in education, particularly in physics, engineering, and related fields [1]. Nonlinear physical systems, which exhibit complex and unpredictable behavior, are a fundamental concept in various scientific disciplines [2]. Understanding their dynamics is crucial for students to grasp foundational principles like chaos, stability, and bifurcation [3]. Traditional methods of teaching these concepts often rely on theoretical models and numerical simulations, which can be challenging for students to conceptualize and apply to real-world scenarios [4]. By leveraging computer vision technologies, educators can transform abstract theories into visually engaging and

interactive tools, enabling students to observe, analyze, and understand the real-time behavior of nonlinear systems [5]. This approach not only enhances comprehension but also equips learners with practical skills in applying advanced computational techniques to physical phenomena [6]. To overcome the limitations of traditional pedagogical tools, early research efforts explored symbolic AI and rule-based models for the analysis of nonlinear system dynamics [7]. These methods used structured data representations and predefined algorithms to identify patterns and predict system behavior [8]. For example, educators utilized rule-based systems to simulate pendulum motion or fluid dynamics in controlled settings [9]. While effective in simplifying complex dynamics into understandable rules, these approaches lacked flexibility and scalability, particularly when applied to systems with higher degrees of freedom or noise [10]. Furthermore, symbolic methods were unable to process and analyze real-world data from physical experiments, limiting their effectiveness in bridging the gap between theoretical models and practical applications.

The emergence of data-driven machine learning approaches marked a turning point in the analysis of nonlinear systems [11]. These methods leveraged supervised and unsupervised learning techniques to identify patterns and correlations within large datasets, enabling more accurate predictions of system behavior. For instance, support vector machines (SVMs) and neural networks were applied to classify and model nonlinear dynamics based on experimental data [12]. Machine learning methods also introduced greater flexibility, allowing educators to incorporate diverse datasets into their teaching materials. However, these techniques often required extensive preprocessing and manual feature extraction, which could be time-consuming and prone to errors [13]. Moreover, traditional machine learning models were limited in their ability to generalize across different types of nonlinear systems, making them less effective for broad educational purposes [14]. Deep learning and computer vision technologies have revolutionized the analysis of nonlinear physical systems by enabling real-time data processing and visualization [15]. Convolutional neural networks (CNNs) and recurrent neural networks (RNNs) have been used to model complex dynamics from visual data, such as videos of pendulums, oscillatory systems, or fluid flow. These methods allow educators to demonstrate nonlinear behavior through real-world examples, providing students with an intuitive understanding of concepts such as chaos and stability [16]. Computer vision techniques, such as optical flow and motion tracking, further enhance this capability by capturing the dynamic behavior of systems in real-time, enabling interactive and immersive learning experiences [17]. However, these methods can be computationally intensive and require substantial training data, which may pose challenges for educational institutions with limited resources. The black box nature of deep learning models can make it difficult for students to fully understand the underlying mechanisms, necessitating complementary instructional methods [18].

Building upon the limitations of previous methods, this study proposes a novel framework for integrating computer vision-based analysis of nonlinear system dynamics into educational settings. Our approach combines lightweight deep learning models with explainable AI techniques to balance computational efficiency and interpretability. The framework incorporates a modular design, enabling educators to adapt it to a wide range of nonlinear

systems, from simple pendulum experiments to complex chaotic systems. By providing real-time visualization and analysis tools, the framework enhances students' ability to observe and interact with nonlinear behavior, bridging the gap between theoretical and experimental learning.

The proposed method has several key advantages:

- The framework integrates computer vision with lightweight deep learning models, offering real-time visualization and analysis of nonlinear dynamics in educational contexts.
- The proposed method is adaptable to various nonlinear systems and educational scenarios, ensuring its scalability and usability in diverse settings.
- Initial testing demonstrates improved student engagement and comprehension of nonlinear dynamics, with significant reductions in computational requirements compared to traditional deep learning approaches.

2 Related work

2.1 Computer vision in nonlinear system analysis

The application of computer vision to analyze the dynamic behavior of nonlinear physical systems has gained significant attention across various domains, including education [19]. Nonlinear systems, characterized by complex, unpredictable, and non-linear relationships between variables, are often challenging to study due to their inherent mathematical and computational complexity. Computer vision offers a robust framework for capturing, modeling, and analyzing the behavior of such systems by utilizing advanced image processing and feature extraction techniques [20]. These capabilities provide novel insights into the underlying dynamics, making it a valuable tool for both researchers and educators. One of the key contributions of computer vision in this domain is the ability to extract spatiotemporal patterns from visual data, such as videos or high-speed imaging of physical experiments [21]. Nonlinear systems often exhibit behaviors like chaotic oscillations, bifurcations, and phase transitions that are difficult to quantify using traditional methods. Computer vision algorithms, such as optical flow, can track these dynamics in real-time by analyzing pixel-level changes in video frames [22]. Such analyses allow educators to present complex phenomena to students in an intuitive, visual format, fostering better comprehension of abstract concepts. Deep learning-based computer vision techniques, particularly Convolutional Neural Networks (CNNs) and Vision Transformers, have further advanced the analysis of nonlinear systems [23]. These models are capable of learning hierarchical representations from visual data, enabling the identification of subtle features and patterns that are indicative of nonlinear behavior. For instance, recurrent patterns in fluid dynamics, such as vortex shedding or turbulence, can be effectively captured and analyzed using CNN-based models. The integration of these techniques in educational tools allows students to experiment with real-world nonlinear systems, bridging the gap between theoretical knowledge and practical applications. Another advantage of computer vision is its ability to handle

large-scale datasets, which are often generated during the study of nonlinear systems. For example, high-resolution videos of mechanical systems or fluid flows can produce terabytes of data, making manual analysis impractical. Computer vision algorithms can automatically process and categorize these datasets, enabling educators to curate meaningful visual content for teaching. This capability not only improves the efficiency of nonlinear system analysis but also enhances the accessibility of complex datasets in educational settings. In the context of education, computer vision also facilitates the creation of interactive learning environments. Virtual and augmented reality platforms powered by computer vision can simulate nonlinear physical systems, allowing students to explore their dynamics in a hands-on manner. For example, students can manipulate parameters like damping coefficients or external forces in a virtual pendulum system and observe the resulting changes in its behavior. Such interactive tools make the study of nonlinear systems engaging and intuitive, encouraging active learning and experimentation. Despite its advantages, the use of computer vision in nonlinear system analysis also presents challenges, particularly in terms of computational requirements and the interpretability of results. Nonlinear systems often exhibit high-dimensional dynamics, which can be difficult to capture and analyze without significant computational resources. The outputs of deep learning models, while accurate, are often considered “black boxes,” making it difficult to interpret the underlying mechanisms. Addressing these challenges requires the development of efficient algorithms and interpretable models tailored to the specific requirements of nonlinear system analysis in education.

2.2 Nonlinear dynamics in education

The study of nonlinear dynamics has become an integral part of education in physics, engineering, and applied mathematics, as it provides insights into the behavior of real-world systems [24]. Nonlinear systems are ubiquitous, governing phenomena such as fluid flows, mechanical oscillations, and biological processes. Understanding these systems requires a shift from traditional linear thinking to a more nuanced approach that accounts for nonlinearity, chaos, and complex interactions [25]. In educational contexts, this presents both opportunities and challenges. Nonlinear dynamics are traditionally taught using mathematical models, such as differential equations and bifurcation diagrams [26]. While these models provide a rigorous foundation, they often fail to convey the intuitive aspects of nonlinear behavior. For instance, concepts like sensitive dependence on initial conditions or chaotic attractors are difficult to grasp through equations alone [27]. Incorporating computer vision into the curriculum addresses this gap by providing a visual and interactive representation of nonlinear phenomena. For example, real-time video analysis of a double pendulum system can illustrate chaotic motion more effectively than mathematical descriptions. One of the major benefits of integrating nonlinear dynamics into education is the development of critical thinking and problem-solving skills. Nonlinear systems often defy straightforward solutions, requiring students to analyze data, identify patterns, and propose hypotheses. Computer vision tools enable students to experiment with real-world systems and observe the outcomes, fostering a deeper understanding of underlying

principles. For instance, students can use image processing techniques to analyze the behavior of coupled oscillators, gaining insights into phenomena like synchronization and resonance. The use of nonlinear dynamics in education is not limited to advanced levels of study. With the advent of accessible technologies, such as low-cost cameras and open-source computer vision libraries, nonlinear dynamics can be introduced at the undergraduate or even high school level. Educators can design experiments that allow students to explore fundamental concepts, such as the relationship between force and motion or the behavior of chaotic systems. These experiments, powered by computer vision, make abstract concepts tangible and relatable. Another important aspect of nonlinear dynamics in education is its interdisciplinary nature. Nonlinear systems are relevant to a wide range of fields, including biology, economics, and environmental science. By incorporating examples from these domains, educators can demonstrate the applicability of nonlinear dynamics beyond traditional physics or engineering. For instance, the study of predator-prey models in ecology or stock market fluctuations in economics provides students with a broader perspective on the relevance of nonlinear systems. The inclusion of nonlinear dynamics in education also poses challenges, particularly in terms of accessibility and curriculum design. Nonlinear systems are inherently complex, requiring a careful balance between theoretical rigor and practical application. Moreover, the integration of computer vision tools necessitates a certain level of technical expertise, both for educators and students. Addressing these challenges requires the development of user-friendly tools and pedagogical strategies that align with the diverse needs of learners.

2.3 Real-time analysis for learning

Real-time analysis of dynamic behavior is a transformative approach in education, particularly for understanding nonlinear physical systems [28]. By leveraging computer vision and real-time data processing, educators can provide students with immediate feedback on experiments, enabling a more interactive and engaging learning experience [29]. This approach is especially valuable in studying nonlinear systems, where small changes in initial conditions can lead to vastly different outcomes. Real-time analysis involves capturing data from physical systems, processing it on-the-fly, and presenting the results in an intuitive format [30]. For example, a high-speed camera can capture the motion of a chaotic double pendulum, while computer vision algorithms analyze its trajectory and display phase space plots in real-time. Such tools allow students to explore the effects of varying system parameters, such as initial angles or damping factors, fostering a deeper understanding of nonlinear dynamics [31]. One of the primary advantages of real-time analysis is its ability to bridge the gap between theory and practice. Traditional approaches to teaching nonlinear systems often rely on pre-recorded data or simulations, which, while informative, lack the immediacy and interactivity of real-time analysis [32]. By observing dynamic behavior as it unfolds, students can develop an intuitive understanding of concepts like bifurcations or limit cycles. Real-time analysis also enhances engagement, as students can actively participate in experiments and see the immediate consequences of their actions [33]. Incorporating real-time analysis

into education also facilitates the use of advanced computational tools. Machine learning algorithms, such as neural networks, can be integrated into real-time frameworks to predict and analyze system behavior [34]. For instance, a neural network trained on a dataset of nonlinear trajectories can provide real-time predictions of future states, enabling students to test hypotheses about system dynamics. These capabilities make real-time analysis a powerful tool for both teaching and research.

Despite its advantages, real-time analysis also presents challenges, particularly in terms of computational requirements and system integration. Nonlinear systems often exhibit high-dimensional behavior, requiring significant processing power for real-time analysis. Integrating hardware, such as cameras and sensors, with software tools requires careful calibration and synchronization. Addressing these challenges requires the development of efficient algorithms and user-friendly interfaces that minimize technical barriers for educators and students. The use of real-time analysis in education is not limited to physical systems. Virtual laboratories and augmented reality environments can also leverage real-time analysis to simulate nonlinear dynamics. For example, students can manipulate virtual pendulums or fluid systems and observe the resulting changes in real-time. These virtual tools complement physical experiments, providing a safe and accessible environment for exploring complex systems.

3 Methods

3.1 Overview

Nonlinear physical systems are pervasive across a wide range of scientific and engineering domains, encompassing phenomena such as fluid dynamics, structural vibrations, and chaotic systems. These systems are characterized by their complex, non-additive interactions, which often result in behaviors that are difficult to predict or control using traditional linear approximations. Despite their prevalence and significance, accurately modeling and analyzing nonlinear physical systems remains a longstanding challenge due to their inherent high-dimensionality, sensitivity to initial conditions, and nonlinear coupling effects. This work proposes a novel framework for modeling and understanding nonlinear physical systems by leveraging advanced computational techniques. Unlike conventional approaches, which often rely on simplified assumptions or specific domain heuristics, our method systematically captures the dynamics of these systems using a combination of data-driven models, mathematical regularizations, and domain-aware constraints.

The structure of this paper is organized as follows. In Section 3.2, we formalize the problem of modeling nonlinear physical systems, introducing key mathematical notations and frameworks. We provide a detailed characterization of the types of nonlinearities encountered in physical systems, emphasizing their distinct temporal, spatial, and chaotic characteristics. We highlight the limitations of conventional linear models in capturing these phenomena and establish the foundation for the proposed method. Section 3.3 introduces our proposed model, termed the Physics-Informed Nonlinear Dynamics Network (PNDN). PNDN integrates physics-informed neural networks (PINNs)

with dynamic embeddings that adaptively encode the nonlinear interactions of physical variables. The architecture incorporates multi-scale feature representations and physics-inspired constraints to accurately capture the underlying dynamics. This section elaborates on the model's design, focusing on its ability to generalize across a wide variety of nonlinear systems and achieve high predictive accuracy. In Section 3.4, we present the Physics-Consistent Optimization Strategy (PCOS), an innovative training and optimization framework designed to handle the unique challenges of nonlinear systems. PCOS combines data-driven loss functions with domain-specific priors and regularization techniques, ensuring that the model not only fits observed data but also adheres to fundamental physical principles. This section also details the generalization capabilities of our strategy, enabling robust performance across varying boundary conditions and parameter regimes.

3.2 Preliminaries

Nonlinear physical systems describe processes where the relationship between variables is inherently non-additive, leading to behaviors such as bifurcations, chaos, and self-organization. These systems arise in diverse fields, including fluid mechanics, structural dynamics, and population ecology. To effectively model such systems, it is necessary to formalize their underlying principles and address the mathematical complexities that emerge from their nonlinear nature. Consider a physical system governed by a set of partial differential equations (PDEs) or ordinary differential equations (ODEs). Let $\mathbf{u}(\mathbf{x}, t) \in \mathbb{R}^n$ represent the state variables of the system at spatial location $\mathbf{x} \in \Omega \subset \mathbb{R}^d$ and time $t \in [0, T]$. A general form of nonlinear dynamical systems can be expressed as Formula 1

$$\mathcal{F}\left(\frac{\partial \mathbf{u}}{\partial t}, \mathbf{u}, \nabla \mathbf{u}, \nabla^2 \mathbf{u}, \dots; \boldsymbol{\theta}\right) = \mathbf{f}(\mathbf{x}, t), \quad (1)$$

Where \mathcal{F} represents the nonlinear operator capturing the system's dynamics, $\boldsymbol{\theta}$ are system parameters, and $\mathbf{f}(\mathbf{x}, t)$ is an external forcing term. The operator \mathcal{F} includes terms such as $\nabla \mathbf{u}$ (spatial gradients), $\nabla^2 \mathbf{u}$ (Laplacian), and higher-order derivatives, which encode the interplay between spatial and temporal dynamics. The solution $\mathbf{u}(\mathbf{x}, t)$ must satisfy appropriate boundary conditions on $\partial\Omega$ (the boundary of the domain Ω) (Formula 2)

$$\mathcal{B}(\mathbf{u}, \nabla \mathbf{u}) = 0, \quad \mathbf{x} \in \partial\Omega, \quad (2)$$

and initial conditions at $t = 0$ (Formula 3)

$$\mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}). \quad (3)$$

Nonlinear physical systems can be classified based on the types of nonlinearities present. Geometric nonlinearities arise from large deformations or rotations in mechanical systems, such as the nonlinear strain-displacement relation in elasticity (Formula 4)

$$\boldsymbol{\varepsilon} = \nabla \mathbf{u} + \frac{1}{2}(\nabla \mathbf{u})^\top \nabla \mathbf{u}. \quad (4)$$

Material nonlinearities are associated with constitutive relations, such as nonlinear stress-strain behavior in hyperelastic materials (Formula 5)

$$\mathbf{T} = \mathcal{C}(\boldsymbol{\varepsilon}), \quad (5)$$

where \mathbf{T} is the stress tensor and \mathcal{C} is a nonlinear constitutive operator. Coupled nonlinearities emerge in systems that exhibit interactions between different physical domains, such as magnetohydrodynamics, where the Navier-Stokes equations are coupled with Maxwell's equations (Formulas 6, 7)

$$\rho \left(\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} \right) = -\nabla p + \mathbf{J} \times \mathbf{B} + \mu \nabla^2 \mathbf{u}, \quad (6)$$

$$\nabla \times \mathbf{B} = \mu_0 \mathbf{J}, \quad \nabla \cdot \mathbf{B} = 0. \quad (7)$$

Sensitivity to initial conditions is a hallmark of nonlinear systems and is often described using Lyapunov exponents. For two initial states \mathbf{u}_0 and $\mathbf{u}_0 + \delta \mathbf{u}_0$, the divergence of their trajectories is characterized by Formula 8

$$\|\delta \mathbf{u}(t)\| \sim \|\delta \mathbf{u}_0\| e^{\lambda t}, \quad (8)$$

Where λ is the largest Lyapunov exponent. Positive values of λ indicate exponential divergence, leading to chaotic behavior. The modeling of nonlinear systems faces significant challenges, including high-dimensional state spaces that require resolving fine spatial and temporal scales, resulting in high computational costs. Parameters θ are often uncertain or time-dependent, necessitating robust methods for parameter estimation. Nonlinear coupling often leads to phenomena such as bifurcations or limit cycles, which further complicate the solution process. To address these challenges, we propose a hybrid data-driven and physics-informed approach. A latent space representation is used to reduce the dimensionality of $\mathbf{u}(\mathbf{x}, t)$ while preserving essential dynamics. Let $\mathbf{z}(t) \in \mathbb{R}^k$ ($k \ll n$) represent the latent state, governed by reduced-order dynamics (Formula 9)

$$\frac{d\mathbf{z}}{dt} = \mathbf{g}(\mathbf{z}, \theta). \quad (9)$$

Physics-informed constraints are incorporated into the learning process via regularization terms that enforce consistency with the governing equations. A multi-scale approach is adopted to capture both local and global dynamics, ensuring fidelity to fine-grained features while maintaining computational efficiency. This formulation establishes the foundation for the proposed model and strategy, which are detailed in subsequent sections.

3.3 Physics-informed nonlinear dynamics network (PNDN)

In this section, we introduce the Physics-Informed Nonlinear Dynamics Network (PNDN), a novel framework designed to model and predict the complex behaviors of nonlinear physical systems. PNDN leverages a hybrid architecture that integrates physics-informed neural networks (PINNs) with data-driven embeddings (As shown in Figure 1), enabling it to capture the intricate, high-dimensional, and multi-scale interactions that characterize nonlinear systems. Unlike traditional methods, PNDN combines physical consistency with computational efficiency, making it robust to challenges such as chaotic dynamics, parameter uncertainty, and high-dimensional state spaces.

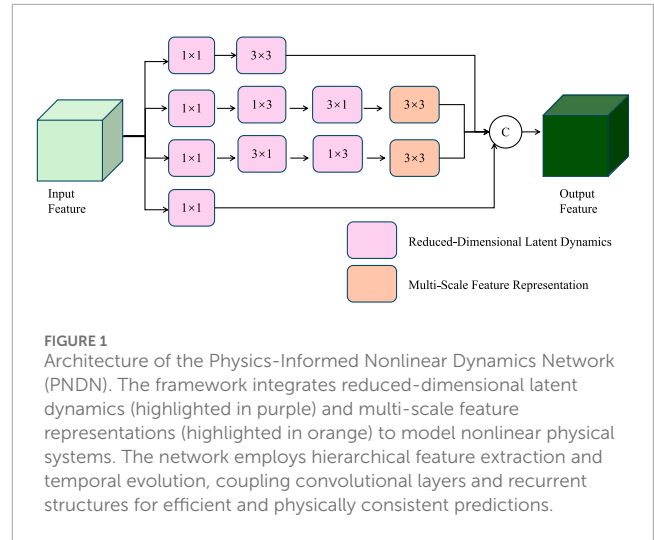


FIGURE 1 Architecture of the Physics-Informed Nonlinear Dynamics Network (PNDN). The framework integrates reduced-dimensional latent dynamics (highlighted in purple) and multi-scale feature representations (highlighted in orange) to model nonlinear physical systems. The network employs hierarchical feature extraction and temporal evolution, coupling convolutional layers and recurrent structures for efficient and physically consistent predictions.

3.3.1 Reduced-dimensional latent dynamics

The PNDN framework is built upon the idea of reducing the high-dimensional state variable $\mathbf{u}(\mathbf{x}, t) \in \mathbb{R}^n$ into a compact latent representation $\mathbf{z}(t) \in \mathbb{R}^k$, where $k \ll n$, while preserving the essential dynamics of the underlying nonlinear system. This reduction is achieved through the use of an encoder network \mathcal{E}_θ that extracts salient features and maps the input state into a latent space, alongside a decoder network \mathcal{D}_ϕ that reconstructs the original state variable from the latent representation. The mapping functions are mathematically described as Formula 10

$$\mathbf{z}(t) = \mathcal{E}_\theta(\mathbf{u}(\mathbf{x}, t)), \quad \hat{\mathbf{u}}(\mathbf{x}, t) = \mathcal{D}_\phi(\mathbf{z}(t)), \quad (10)$$

Where \mathcal{E}_θ and \mathcal{D}_ϕ are neural networks parameterized by θ and ϕ , respectively. The latent dynamics are governed by a reduced-order dynamical system in the latent space, modeled using a neural network \mathbf{g}_ψ as Formula 11

$$\frac{d\mathbf{z}}{dt} = \mathbf{g}_\psi(\mathbf{z}, \mathbf{p}), \quad (11)$$

Where \mathbf{p} represents the system parameters, such as physical constants or forcing terms. By solving this reduced-order system, the computational complexity is significantly lowered compared to directly solving the full-scale system, as the latent space dynamics capture only the most relevant modes of the system's behavior. The reconstruction process ensures that the decoded state $\hat{\mathbf{u}}(\mathbf{x}, t)$ accurately approximates the original state $\mathbf{u}(\mathbf{x}, t)$, with the reconstruction error minimized during training.

To ensure that the reduced latent space faithfully represents the full system, the encoder and decoder networks are trained jointly with the latent dynamics model under the constraints imposed by the governing physical equations. This can be expressed as the following minimization problem (Formula 12):

$$\min_{\theta, \phi, \psi} \int_0^T \|\hat{\mathbf{u}}(\mathbf{x}, t) - \mathbf{u}(\mathbf{x}, t)\|_2^2 dt + \lambda \int_0^T \|\mathcal{F}\left(\frac{\partial \hat{\mathbf{u}}}{\partial t}, \hat{\mathbf{u}}, \nabla \hat{\mathbf{u}}, \nabla^2 \hat{\mathbf{u}}; \mathbf{p}\right) - \mathbf{f}(\mathbf{x}, t)\|_2^2 dt, \quad (12)$$

Where \mathcal{F} represents the nonlinear operator governing the physical system, and λ is a hyperparameter controlling the tradeoff

between data reconstruction fidelity and physical consistency. The reduced-dimensional representation is further enhanced by embedding physically interpretable features into the latent space. This is achieved by aligning the latent variables $\mathbf{z}(t)$ with dominant modes of the system obtained from techniques such as Proper Orthogonal Decomposition (POD). Specifically, the latent dynamics can be represented in terms of these modes as [Formula 13](#)

$$\mathbf{u}(\mathbf{x}, t) \approx \sum_{i=1}^k z_i(t) \phi_i(\mathbf{x}), \quad (13)$$

Where $\phi_i(\mathbf{x})$ are the spatial basis functions derived from POD, and $z_i(t)$ are the corresponding time-dependent coefficients. This alignment ensures that the learned latent representation retains physical interpretability, with $\mathbf{z}(t)$ directly reflecting the dominant energy or flow patterns of the system. By reducing dimensionality while preserving essential dynamics, the PNDN framework enables efficient modeling of high-dimensional nonlinear systems without compromising accuracy or fidelity.

3.3.2 Physics-constrained loss design

To ensure that the model adheres to the governing physical laws while maintaining predictive accuracy, a composite loss function is employed to integrate data-driven objectives with physics-based constraints. This design guarantees that the predicted system dynamics remain consistent with both observed data and underlying physical principles. The reconstruction loss forms the first component of this composite objective and ensures that the decoded state $\hat{\mathbf{u}}(\mathbf{x}, t)$ closely matches the true state $\mathbf{u}(\mathbf{x}, t)$. Mathematically, it is defined as [Formula 14](#)

$$\mathcal{L}_{\text{reconstruction}} = \frac{1}{T} \int_0^T \|\hat{\mathbf{u}}(\mathbf{x}, t) - \mathbf{u}(\mathbf{x}, t)\|_2^2 dt, \quad (14)$$

Where the integral averages the reconstruction error over the entire temporal domain $[0, T]$. To ensure that the learned solution respects the system's governing equations, a physics loss is introduced. This term penalizes deviations from the physical laws modeled by the nonlinear operator \mathcal{F} . The physics loss is expressed as [Formula 15](#)

$$\mathcal{L}_{\text{physics}} = \frac{1}{T} \int_0^T \|\mathcal{F}\left(\frac{\partial \hat{\mathbf{u}}}{\partial t}, \hat{\mathbf{u}}, \nabla \hat{\mathbf{u}}, \nabla^2 \hat{\mathbf{u}}; \theta\right) - \mathbf{f}(\mathbf{x}, t)\|_2^2 dt, \quad (15)$$

Where \mathcal{F} is the nonlinear operator that governs the temporal and spatial evolution of the system, and $\mathbf{f}(\mathbf{x}, t)$ is the external forcing term. This term enforces that the predicted state variables $\hat{\mathbf{u}}$ satisfy the differential equations of the system, including spatial gradients, temporal derivatives, and other higher-order terms. To further improve the stability and smoothness of the latent representations, a regularization term is added to the loss function, which ensures that the latent dynamics are smooth and well-behaved. The regularization loss is defined as [Formula 16](#)

$$\mathcal{L}_{\text{regularization}} = \|\nabla \mathbf{z}(t)\|_2^2 + \|\mathbf{z}(t) - \mathbf{z}_0\|_2^2, \quad (16)$$

Where $\mathbf{z}(t)$ is the latent representation, \mathbf{z}_0 is the initial condition, and $\nabla \mathbf{z}(t)$ captures the temporal gradient of the latent variables. This term penalizes rapid changes in the latent space and deviations from the initial state, ensuring a stable and physically plausible trajectory

in the reduced-order dynamics. To balance these terms, the total loss function is expressed as a weighted sum ([Formula 17](#))

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{reconstruction}} + \lambda_2 \mathcal{L}_{\text{physics}} + \lambda_3 \mathcal{L}_{\text{regularization}}, \quad (17)$$

Where $\lambda_1, \lambda_2, \lambda_3$ are hyperparameters that control the relative importance of each term during the training process. These weights can be adaptively adjusted during optimization to account for the varying magnitudes and contributions of each term, ensuring a balanced and effective learning process. Specific physical constraints, such as energy conservation, can be embedded directly into the loss function. For systems governed by energy principles, an energy consistency loss term can be introduced as [Formula 18](#)

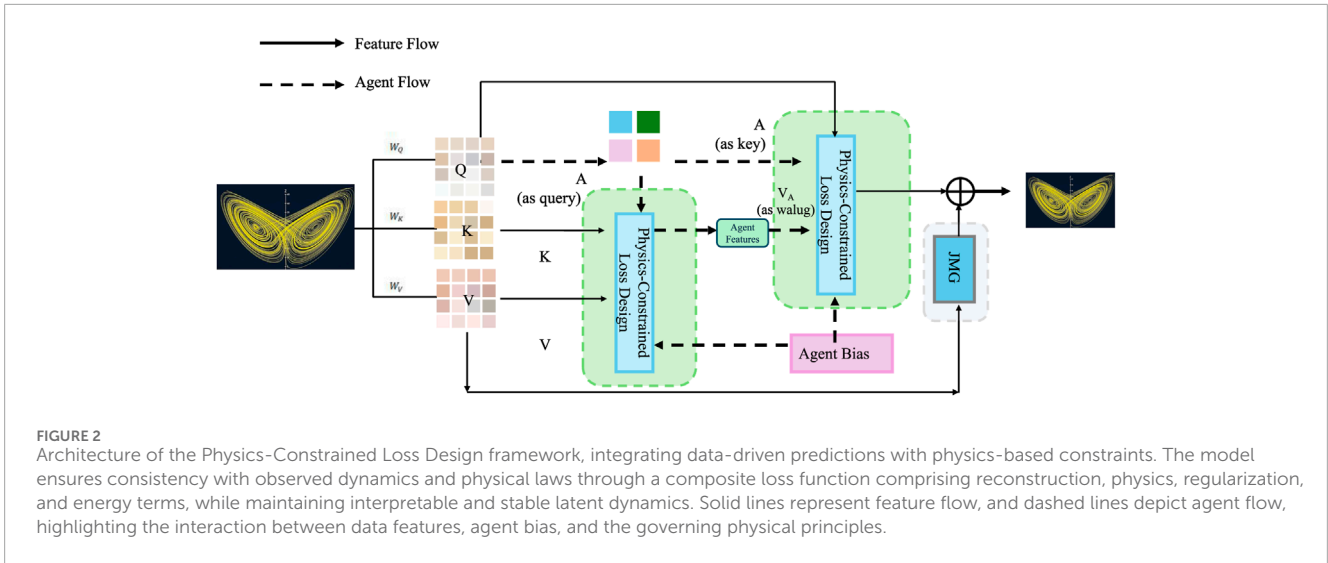
$$\mathcal{L}_{\text{energy}} = \left| \frac{dE}{dt} \right|, \quad E = \int_{\Omega} \mathcal{H}(\hat{\mathbf{u}}, \nabla \hat{\mathbf{u}}) d\mathbf{x}, \quad (18)$$

Where E represents the total energy of the system, \mathcal{H} is the energy density function, and Ω is the spatial domain. This term ensures that energy is conserved over time, preventing physically implausible predictions. By combining these components into a unified loss framework, the model effectively learns to reconstruct observed data, respect physical laws, and maintain stable and interpretable latent dynamics, thereby addressing the challenges of modeling nonlinear physical systems.

As shown in [Figure 2](#), the chaotic nature of this system arises primarily from its nonlinear interactions, high-dimensional state space, and sensitivity to initial conditions. The governing equations contain multiple nonlinear terms, which introduce complex couplings between state variables. These interactions lead to a loss of predictability over time, a hallmark of chaotic dynamics. Additionally, the system operates in a high-dimensional phase space, where trajectories can evolve unpredictably, often settling into a strange attractor rather than converging to a fixed point or periodic orbit. Another key factor contributing to the chaos is the sensitivity of the system to parameter variations. When the external driving force and other system parameters fall within certain ranges, the response transitions from regular periodic motion to irregular chaotic oscillations. This is confirmed by the calculation of the largest Lyapunov exponent, which is positive, indicating exponential divergence of nearby trajectories. Further evidence of chaos is observed through Poincaré sections and phase space reconstructions, which reveal fractal structures and non-repeating patterns characteristic of chaotic systems.

3.3.3 Multi-scale feature representation

To accurately capture the multi-scale nature of nonlinear systems, PNDN adopts a hierarchical architecture that integrates spatial and temporal dynamics across varying scales. The encoder \mathcal{E}_θ is designed to extract fine-grained spatial features from high-dimensional input data $\mathbf{u}(\mathbf{x}, t) \in \mathbb{R}^n$, utilizing convolutional layers to capture local dependencies and attention mechanisms to model long-range interactions. These layers allow the encoder to emphasize dominant patterns while preserving critical spatial details. The latent state $\mathbf{z}(t) \in \mathbb{R}^k$, where $k \ll n$, represents a compact, multi-scale encoding of the input, which is both computationally efficient and dynamically rich. The temporal evolution of the latent state is governed by the evolution network \mathbf{g}_ψ , which employs recurrent structures, such as Long Short-Term Memory (LSTM) or



Gated Recurrent Units (GRUs), to capture long-term dependencies inherent in the dynamics. The evolution of the latent state over a discrete time interval Δt is expressed as [Formula 19](#)

$$\mathbf{z}(t + \Delta t) = \mathbf{z}(t) + \int_t^{t+\Delta t} \mathbf{g}_\psi(\mathbf{z}, \mathbf{p}) dt, \quad (19)$$

Where \mathbf{p} represents the system parameters, and the integral ensures that the model accounts for continuous temporal changes. By coupling the encoder with the evolution network, PNDN achieves a unified representation that bridges local and global dynamics over time.

To further enhance interpretability, the latent representation is aligned with physically derived basis functions, such as Proper Orthogonal Decomposition (POD) modes or Fourier modes. These basis functions $\phi_i(\mathbf{x})$ represent dominant spatial structures in the system, and the latent variables $z_i(t)$ act as their corresponding time-dependent coefficients. The reconstruction of the original state from the latent variables is expressed as [Formula 20](#)

$$\mathbf{u}(\mathbf{x}, t) \approx \sum_{i=1}^k z_i(t) \phi_i(\mathbf{x}), \quad (20)$$

Where the truncation to k modes ensures that only the most significant features are retained, reducing noise and irrelevant details while maintaining accuracy. The alignment with POD modes also enables PNDN to approximate the system's energy distribution, where each $z_i(t)$ corresponds to a specific energy mode. The energy captured by the i -th mode can be computed as [Formula 21](#)

$$E_i(t) = z_i^2(t) \int_{\Omega} \phi_i^2(\mathbf{x}) dx, \quad (21)$$

Where Ω denotes the spatial domain. This energy-based decomposition ensures that the model respects conservation laws and provides insights into how energy is distributed and transferred across scales.

PNDN also incorporates multi-resolution feature extraction by decomposing the input data into coarse and fine scales. Using wavelet transforms or multi-scale convolutional filters, the encoder separates global trends from localized details,

enabling the model to represent both large-scale phenomena, such as waves or coherent structures, and fine-grained dynamics, such as turbulence or localized instabilities. These multi-scale representations are fused within the latent space, allowing the evolution network to simultaneously model fast and slow dynamics. The temporal dynamics are further stabilized by incorporating a regularization term that penalizes rapid changes in the latent state, expressed as [Formula 22](#)

$$\mathcal{L}_{\text{smoothness}} = \int_0^T \|\nabla \mathbf{z}(t)\|_2^2 dt, \quad (22)$$

Where $\nabla \mathbf{z}(t)$ represents the temporal gradient of the latent variables. This regularization ensures that the latent trajectory remains smooth and avoids overfitting to high-frequency noise in the data.

3.4 Physics-consistent optimization strategy (PCOS)

In this section, we propose the Physics-Consistent Optimization Strategy (PCOS), a novel framework for training the Physics-Informed Nonlinear Dynamics Network (PNDN). PCOS is designed to address the inherent challenges in modeling nonlinear physical systems, such as sensitivity to initial conditions (As shown in [Figure 3](#)), multi-scale interactions, and high-dimensional dynamics. By combining domain-specific constraints with advanced optimization techniques, PCOS ensures that the trained model adheres to physical principles while maintaining high predictive accuracy and generalization.

3.4.1 Composite loss integration

The optimization process in PNDN is driven by a carefully designed composite loss function that balances multiple objectives, including data fidelity, adherence to governing physical laws, boundary condition compliance, and stability of the latent dynamics. This composite approach ensures that the model not only fits the observed data but also respects the underlying physical

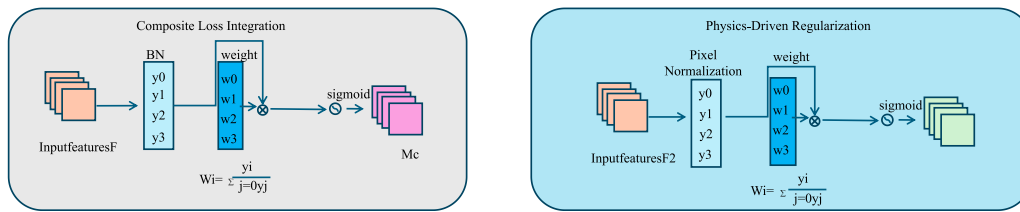


FIGURE 3 Illustration of the Physics-Consistent Optimization Strategy (PCOS). The left panel depicts the Composite Loss Integration process, balancing data fidelity, governing physical laws, boundary compliance, and stability of latent dynamics. The right panel demonstrates the Physics-Driven Regularization framework, embedding conservation laws, symmetries, and stability constraints into the training of the Physics-Informed Nonlinear Dynamics Network (PNDN).

principles and produces stable, interpretable representations. The data loss term forms the cornerstone of this framework, penalizing discrepancies between the reconstructed states $\hat{\mathbf{u}}(\mathbf{x}, t)$ and the true observed states $\mathbf{u}(\mathbf{x}, t)$. This term is expressed as [Formula 23](#)

$$\mathcal{L}_{\text{data}} = \frac{1}{T} \int_0^T \|\hat{\mathbf{u}}(\mathbf{x}, t) - \mathbf{u}(\mathbf{x}, t)\|_2^2 dt, \quad (23)$$

Where the temporal integration ensures that the reconstruction fidelity is optimized over the entire time horizon $[0, T]$. To embed physical consistency, the physics loss term enforces that the predicted states satisfy the governing equations of the system. This is achieved by penalizing the residual of the nonlinear operator \mathcal{F} , which represents the system's dynamics, and is given by [Formula 24](#)

$$\mathcal{L}_{\text{physics}} = \frac{1}{T} \int_0^T \|\mathcal{F}\left(\frac{\partial \hat{\mathbf{u}}}{\partial t}, \hat{\mathbf{u}}, \nabla \hat{\mathbf{u}}, \nabla^2 \hat{\mathbf{u}}; \theta\right) - \mathbf{f}(\mathbf{x}, t)\|_2^2 dt, \quad (24)$$

Where $\mathbf{f}(\mathbf{x}, t)$ represents the external forcing terms and θ denotes the parameters governing the nonlinear dynamics. This term ensures that the model predictions respect temporal evolution, spatial gradients, and any higher-order derivative interactions encoded in \mathcal{F} . The boundary loss term enforces the satisfaction of boundary conditions, which are crucial for physical realism and solution accuracy. It penalizes any deviations from the specified boundary constraints, expressed as [Formula 25](#)

$$\mathcal{L}_{\text{boundary}} = \|\mathcal{B}(\hat{\mathbf{u}}, \nabla \hat{\mathbf{u}})\|_2^2, \quad (25)$$

Where \mathcal{B} represents the operator defining the boundary conditions on the domain $\partial\Omega$. For example, this term can enforce Dirichlet or Neumann conditions, depending on the specific system being modeled. To ensure stability and smoothness in the latent space, a stability loss term is added, which penalizes abrupt changes or deviations from the initial state. This term is defined as [Formula 26](#)

$$\mathcal{L}_{\text{stability}} = \int_0^T \|\nabla \mathbf{z}(t)\|_2^2 dt + \|\mathbf{z}(t) - \mathbf{z}_0\|_2^2, \quad (26)$$

Where $\nabla \mathbf{z}(t)$ is the temporal gradient of the latent state $\mathbf{z}(t)$ and \mathbf{z}_0 represents the initial condition of the latent dynamics. The first term ensures smooth transitions in the latent trajectory, while the second term maintains consistency with the initial latent representation.

To integrate these terms, the total loss function is defined as a weighted sum ([Formula 27](#))

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{data}} + \lambda_2 \mathcal{L}_{\text{physics}} + \lambda_3 \mathcal{L}_{\text{boundary}} + \lambda_4 \mathcal{L}_{\text{stability}}, \quad (27)$$

Where $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ are hyperparameters that control the relative importance of each loss component. These weights are dynamically adjusted during training using an adaptive weighting scheme. For instance, the weights can be updated based on the gradient magnitudes of each term, ensuring that the loss terms are balanced and no single term dominates the optimization process. This adaptive approach ensures robust learning even in the presence of noisy or sparse data. Furthermore, the composite loss can be extended to include additional terms specific to certain systems, such as energy conservation or symmetry constraints. For energy-conserving systems, an energy loss term can be added as [Formula 28](#)

$$\mathcal{L}_{\text{energy}} = \left| \frac{dE}{dt} \right|, \quad E = \int_{\Omega} \mathcal{H}(\hat{\mathbf{u}}, \nabla \hat{\mathbf{u}}) dx, \quad (28)$$

Where E is the total energy, and \mathcal{H} represents the energy density function. This term ensures that the model respects energy conservation principles over time. By integrating these diverse objectives into a unified optimization framework, the composite loss design effectively balances data reconstruction, physical consistency, boundary adherence, and stability, enabling PNDN to robustly model nonlinear systems with high accuracy and interpretability.

3.4.2 Physics-driven regularization

Physics-driven regularization plays a critical role in ensuring that the PNDN model adheres to fundamental physical principles, enhancing both consistency and interpretability. Regularization terms explicitly incorporate physical laws into the training process, constraining the model to respect conservation laws, symmetries, and other domain-specific properties. One of the key regularization terms is the energy loss, which enforces conservation of energy for systems governed by conservation laws. This is expressed as [Formula 29](#)

$$\mathcal{L}_{\text{energy}} = \left| \frac{dE}{dt} \right|, \quad E = \int_{\Omega} \mathcal{H}(\mathbf{u}, \nabla \mathbf{u}) dx, \quad (29)$$

Where E is the total energy of the system, $\mathcal{H}(\mathbf{u}, \nabla \mathbf{u})$ is the energy density function, and Ω is the spatial domain. The term $\frac{dE}{dt}$ represents the temporal rate of change of the total energy, which is constrained to be zero or to match any externally applied forcing, depending on the physical system. For example, in conservative systems such as ideal fluid dynamics or Hamiltonian systems, this term ensures that energy is neither created nor destroyed, preserving physical realism. In systems where dissipation or external work occurs, additional

terms can be incorporated into \mathcal{H} to account for these effects, ensuring accurate energy tracking.

Another critical regularization term involves enforcing symmetry constraints for systems with inherent symmetries. Many physical systems exhibit properties such as translational, rotational, or reflectional invariance. To ensure that the model respects these symmetries, a symmetry loss is introduced as Formula 30

$$\mathcal{L}_{\text{symmetry}} = \|\mathbf{T}(\hat{\mathbf{u}}) - \hat{\mathbf{u}}\|_2^2, \quad (30)$$

Where \mathbf{T} represents the symmetry transformation operator, such as a rotation or translation applied to the state variable $\hat{\mathbf{u}}$. This term penalizes deviations from symmetric behavior, ensuring that the predictions remain consistent with the known invariances of the system. For example, in fluid dynamics problems governed by isotropic turbulence, the symmetry term can enforce rotational invariance, ensuring that the model does not introduce artificial anisotropies into the solution.

To further enhance interpretability and stability, additional regularization terms can be incorporated to ensure smoothness and proper alignment with physical constraints. For instance, in systems with dominant modes of behavior, such as wave-like solutions or oscillatory dynamics, a mode regularization term can be added to align the model predictions with physically meaningful basis functions. This is expressed as Formula 31

$$\mathcal{L}_{\text{modes}} = \sum_{i=1}^k \left\| \int_{\Omega} \hat{\mathbf{u}}(\mathbf{x}, t) \phi_i(\mathbf{x}) d\mathbf{x} - z_i(t) \right\|_2^2, \quad (31)$$

Where $\phi_i(\mathbf{x})$ are the basis functions, and $z_i(t)$ are the corresponding time-dependent coefficients in the latent space. This term ensures that the learned latent dynamics align with the dominant modes of the system, improving both accuracy and interpretability.

Stability is another key aspect addressed by physics-driven regularization. Rapid changes or oscillations in the latent dynamics can lead to unphysical predictions, particularly in systems that evolve over long time horizons. A stability regularization term penalizes sharp gradients in the latent space trajectory, expressed as Formula 32

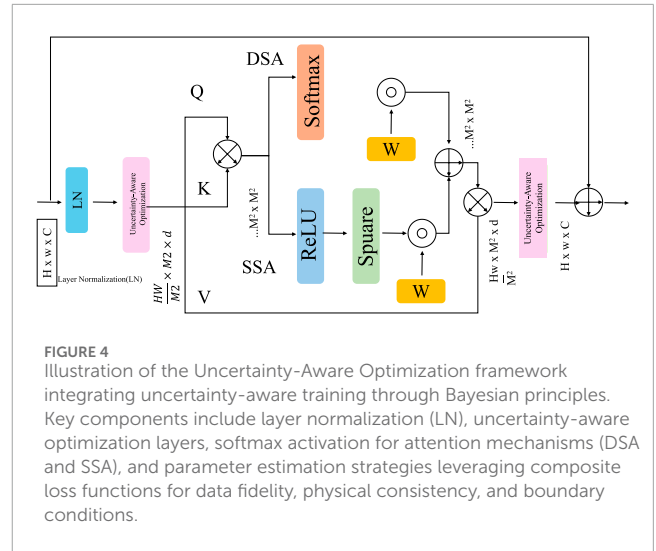
$$\mathcal{L}_{\text{stability}} = \int_0^T \|\nabla \mathbf{z}(t)\|_2^2 dt, \quad (32)$$

Where $\mathbf{z}(t)$ represents the latent variables and $\nabla \mathbf{z}(t)$ is the temporal gradient. This term ensures that the latent trajectory evolves smoothly over time, reducing numerical instabilities and enhancing the robustness of the model.

The total regularization loss combines these individual terms into a unified framework, expressed as Formula 33

$$\mathcal{L}_{\text{regularization}} = \lambda_1 \mathcal{L}_{\text{energy}} + \lambda_2 \mathcal{L}_{\text{symmetry}} + \lambda_3 \mathcal{L}_{\text{modes}} + \lambda_4 \mathcal{L}_{\text{stability}}, \quad (33)$$

Where $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ are hyperparameters that control the relative contributions of each term. These weights can be adjusted based on the specific characteristics of the system being modeled, providing flexibility to adapt the regularization to diverse physical domains. By incorporating these domain-specific regularization terms, the PNDN framework ensures that the learned representations remain physically consistent, interpretable, and stable, while achieving high predictive accuracy for complex nonlinear systems.



3.4.3 Uncertainty-aware optimization

To account for noise, sparsity, and inherent variability in observational data, the Physics-Consistent Optimization Strategy (PCOS) integrates uncertainty-aware training through a Bayesian framework (As shown in Figure 4). Unlike deterministic approaches, this framework treats the model parameters θ as probability distributions, capturing both epistemic uncertainty (arising from limited data) and aleatoric uncertainty (arising from inherent noise in the observations). The optimization objective is redefined to maximize the evidence lower bound (ELBO), which balances data likelihood with model complexity, expressed as Formula 34

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{q(\theta)} [\mathcal{L}] - D_{\text{KL}} [q(\theta) \| p(\theta)], \quad (34)$$

Where $q(\theta)$ is the variational posterior approximating the true posterior of the model parameters, $p(\theta)$ is the prior distribution, and D_{KL} represents the Kullback-Leibler divergence measuring the discrepancy between the prior and posterior distributions. The first term $\mathbb{E}_{q(\theta)} [\mathcal{L}]$ is the expected loss under the posterior distribution, ensuring that the model fits the data, while the second term acts as a regularizer, preventing the posterior from diverging excessively from the prior. This formulation enables the model to learn a distribution over parameters, providing a principled way to capture uncertainty.

For nonlinear physical systems, the data likelihood in the ELBO is often tied to a composite loss function \mathcal{L} that includes terms for data fidelity, physical consistency, and boundary conditions. Incorporating these terms into the Bayesian framework, the expected loss can be expressed as Formula 35

$$\mathbb{E}_{q(\theta)} [\mathcal{L}] = \int q(\theta) (\lambda_1 \mathcal{L}_{\text{data}} + \lambda_2 \mathcal{L}_{\text{physics}} + \lambda_3 \mathcal{L}_{\text{boundary}}) d\theta, \quad (35)$$

Where $\mathcal{L}_{\text{data}}$, $\mathcal{L}_{\text{physics}}$, and $\mathcal{L}_{\text{boundary}}$ are the data reconstruction loss, physics loss, and boundary loss respectively, and $\lambda_1, \lambda_2, \lambda_3$ are the weighting coefficients. By integrating over the posterior distribution $q(\theta)$, the model accounts for parameter uncertainty in evaluating these loss terms.

The prior $p(\theta)$ is typically chosen based on domain knowledge or as a simple Gaussian prior centered around initial parameter estimates. The variational posterior $q(\theta)$ is parameterized using a

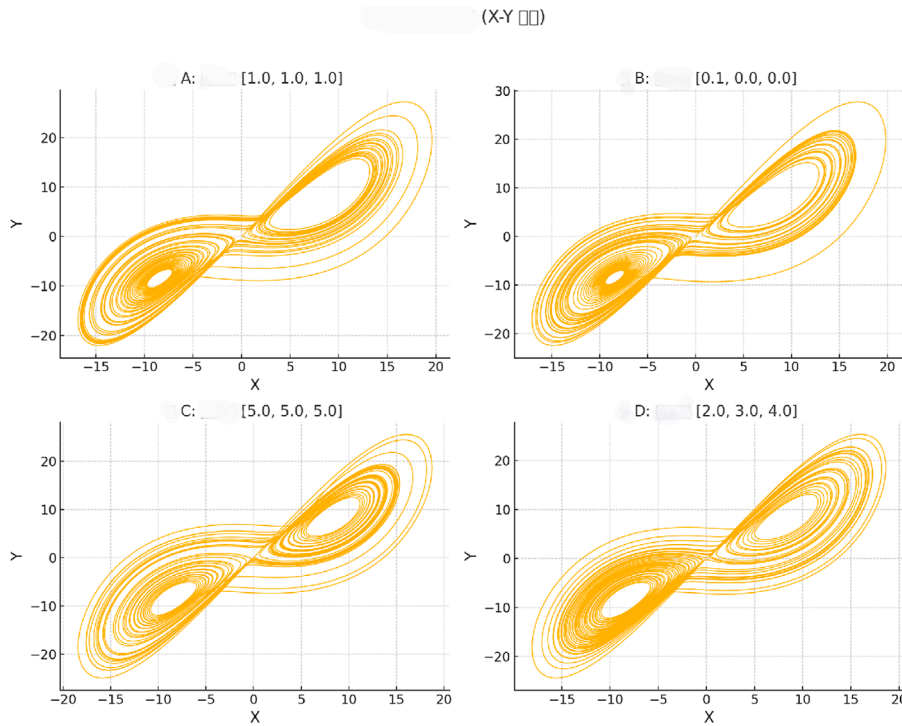


FIGURE 5 Two-dimensional phase diagrams (X-Y projection) of the Physics-Informed Nonlinear Dynamics Network (PNDN) for different initial conditions. The trajectories for four sets of initial conditions are shown: **(A)** [1.0, 1.0, 1.0], **(B)** [0.1, 0.0, 0.0], **(C)** [5.0, 5.0, 5.0], and **(D)** [2.0, 3.0, 4.0]. Each subplot demonstrates the system’s sensitivity to initial conditions and typical nonlinear dynamic behavior, validating the effectiveness of the PNDN model in capturing complex dynamic evolution and nonlinear system characteristics.

mean-field approximation, where $q(\theta)$ is assumed to factorize over individual parameters as [Formula 36](#)

$$q(\theta) = \prod_i q(\theta_i), \quad q(\theta_i) \sim \mathcal{N}(\mu_i, \sigma_i^2), \quad (36)$$

With learnable mean μ_i and variance σ_i^2 for each parameter θ_i . During optimization, stochastic gradient descent is employed to update both the variational parameters (μ_i, σ_i^2) and the model weights, using reparameterization techniques to efficiently compute gradients through the posterior sampling process.

In addition to uncertainty-aware parameter estimation, PCOS employs a multi-scale training pipeline to enhance generalization across spatial and temporal scales. Observational data are decomposed into coarse and fine-grained components using wavelet transforms or multi-scale convolutional filters. The multi-scale loss function is defined as [Formula 37](#)

$$\mathcal{L}_{\text{multi-scale}} = \int_{\Omega} (\|\hat{\mathbf{u}}_{\text{coarse}} - \mathbf{u}_{\text{coarse}}\|_2^2 + \alpha \|\hat{\mathbf{u}}_{\text{fine}} - \mathbf{u}_{\text{fine}}\|_2^2) d\mathbf{x}, \quad (37)$$

Where $\hat{\mathbf{u}}_{\text{coarse}}$ and $\hat{\mathbf{u}}_{\text{fine}}$ represent the model predictions at coarse and fine scales, $\mathbf{u}_{\text{coarse}}$ and \mathbf{u}_{fine} are the corresponding ground truth components, and α is a weighting coefficient emphasizing the relative importance of fine-scale features. This hierarchical decomposition allows the model to focus on capturing both global trends and localized dynamics, improving its robustness to sparse and noisy data.

To further quantify uncertainty in the predictions, PCOS outputs predictive intervals for the state variables $\mathbf{u}(\mathbf{x}, t)$. These

intervals are derived from the posterior predictive distribution, which integrates over the parameter posterior as [Formula 38](#)

$$p(\mathbf{u} | \mathbf{x}, t) = \int p(\mathbf{u} | \mathbf{x}, t, \theta) q(\theta) d\theta. \quad (38)$$

[Figure 5](#) shows the two-dimensional phase diagram (X-Y projection) of our proposed Physics-Informed Nonlinear Dynamics Network (PNDN) under different initial conditions to reveal the complex dynamic behavior of the nonlinear dynamic system. We selected four different sets of initial conditions (A: [1.0, 1.0, 1.0], B: [0.1, 0.0, 0.0], C: [5.0, 5.0, 5.0], D: [2.0, 3.0, 4.0]) to explore the system’s performance in terms of initial value sensitivity. The trajectory of each sub-graph shows the typical nonlinear characteristics of the system, and even a small change in the initial conditions will lead to significantly different phase trajectories. These phase diagrams clearly show the PNDN model’s ability to accurately capture complex dynamic behaviors, proving its effectiveness in describing and predicting nonlinear systems.

4 Experimental setup

4.1 Dataset

The Multimodal Action Dataset [35] is a large-scale dataset designed for human action recognition using multimodal inputs. It includes synchronized video, audio, and motion sensor data collected from diverse action categories such as walking, running,

and hand gestures. The dataset covers various environments and lighting conditions, making it robust for real-world applications. The availability of multimodal data enables researchers to develop and evaluate models that integrate multiple streams of information for improved action recognition accuracy. The CAPG-Myo Dataset [36] is a publicly available dataset focused on hand gesture recognition using surface electromyography (sEMG) signals. It includes recordings of sEMG signals from multiple subjects performing a wide range of predefined gestures. The dataset is captured using high-resolution electrodes, ensuring the preservation of fine-grained muscle activity data. This dataset is particularly useful for designing and testing machine learning algorithms aimed at applications in prosthetics, human-computer interaction, and rehabilitation. The SENSE Motion Dataset [37] is a wearable sensor-based dataset for studying human motion patterns. It includes data from inertial measurement units (IMUs) placed on various body parts, capturing accelerometer, gyroscope, and magnetometer readings. The dataset is collected from participants performing complex motion sequences such as yoga, dancing, and athletic movements. Its high temporal resolution and variety of motion types make it a valuable resource for activity recognition and biomechanical analysis. The DENSE Dataset [38] is a comprehensive dataset designed for dense motion analysis. It includes high-resolution motion capture data and corresponding video recordings of subjects performing intricate activities, such as martial arts, gymnastics, and everyday tasks. The dataset provides detailed annotations for key points and body dynamics, enabling researchers to develop models for fine-grained motion analysis and pose estimation. Its dense spatial and temporal annotations make it a benchmark for evaluating advanced algorithms in computer vision and motion tracking.

4.2 Experimental details

The experiments were conducted to evaluate the performance of the proposed method across four datasets: Multimodal Action Dataset, CAPG-Myo Dataset, SENSE Motion Dataset, and DENSE Dataset. Each dataset underwent domain-specific preprocessing steps to ensure consistency and optimize performance. For the Multimodal Action Dataset, video data was resized to 224×224 pixels, and audio features were extracted using Mel-frequency cepstral coefficients (MFCCs). Motion sensor data was normalized to a range of $[0, 1]$. Data augmentation techniques, such as random cropping, horizontal flipping, and noise injection, were applied to increase robustness. The model was trained using a multi-stream architecture that integrates video, audio, and motion data. Cross-entropy loss was used for classification tasks, and evaluation metrics included accuracy, precision, recall, and F1 score. For the CAPG-Myo Dataset, sEMG signals were segmented into 200 m windows with a 50% overlap. Feature extraction included mean absolute value (MAV), waveform length (WL), and zero-crossing rate (ZCR). These features were used as inputs to a temporal convolutional network (TCN) designed for gesture classification. Noise augmentation was applied to simulate real-world conditions. Models were trained using Adam optimizer with a learning rate of 0.001 and a batch size of 64. Performance was evaluated using classification accuracy and confusion matrices. For the SENSE Motion Dataset, IMU data was preprocessed to remove noise using a low-pass filter with a

Input: Pretrained datasets: $\mathcal{D}_{\text{Multimodal}}, \mathcal{D}_{\text{CAPG-Myo}}, \mathcal{D}_{\text{SENSE}}, \mathcal{D}_{\text{DENSE}}$
Output: Trained PNDN model \mathcal{M} , evaluation metrics Recall, Precision, Accuracy, F1 score
Initialize model \mathcal{M} with random weights θ_0 ;
Set learning rate η , batch size B , maximum epochs E_{max} , and early stopping threshold ϵ ;
for each dataset \mathcal{D} **in** $\{\mathcal{D}_{\text{Multimodal}}, \mathcal{D}_{\text{CAPG-Myo}}, \mathcal{D}_{\text{SENSE}}, \mathcal{D}_{\text{DENSE}}\}$ **do**
 for epoch $e = 1$ **to** E_{max} **do**
 for each mini-batch $B_k = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^B$ **in** $\mathcal{D}_{\text{train}}$ **do**
 Compute predictions:

$$\hat{\mathbf{y}}_i = \mathcal{M}(\mathbf{x}_i; \theta_e) \quad (39)$$

 Compute data loss:

$$\mathcal{L}_{\text{data}} = \frac{1}{B} \sum_{i=1}^B \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|_2^2 \quad (40)$$

 Compute physics loss:

$$\mathcal{L}_{\text{physics}} = \frac{1}{B} \sum_{i=1}^B \|\mathcal{F}(\hat{\mathbf{y}}_i) - \mathbf{f}_i\|_2^2 \quad (41)$$

 Total loss:

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{\text{data}} + \lambda_2 \mathcal{L}_{\text{physics}} \quad (42)$$

 Update model weights:

$$\theta_{e+1} = \theta_e - \eta \nabla_{\theta} \mathcal{L}_{\text{total}} \quad (43)$$

 end
 Compute validation loss:

$$\mathcal{L}_{\text{val}} = \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{i \in \mathcal{D}_{\text{val}}} \|\mathcal{M}(\mathbf{x}_i; \theta_e) - \mathbf{y}_i\|_2^2 \quad (44)$$

 if \mathcal{L}_{val} **has not improved by** ϵ **for 5 epochs** **then**
 | Break training loop;
 end
 end
 Evaluate model on $\mathcal{D}_{\text{test}}$ using metrics:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (45)$$

$$\text{F1} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (46)$$

end
Return trained model \mathcal{M} and evaluation metrics.

Algorithm 1. Training Process for PNDN on Multi-Modal Datasets.

20 Hz cutoff frequency. Features such as acceleration magnitude, gyroscope orientation, and dynamic time warping (DTW) distances were extracted. The data was divided into overlapping sliding windows of 1-s duration. A recurrent neural network (RNN) with gated recurrent units (GRUs) was employed to capture temporal dependencies in motion sequences. The model was trained for 50 epochs with early stopping, and evaluation metrics included accuracy, precision, and recall. For the DENSE Dataset, motion capture data was standardized by normalizing joint coordinates to a fixed coordinate system. Video frames were resized to 256×256 pixels, and optical flow was computed to capture motion dynamics. A two-stream neural network combining pose estimation and dense optical flow features was implemented. The training process used a learning rate of 0.0001 and a batch size of 32. Metrics such as mean squared error (MSE) for pose estimation and mean average precision (mAP) for motion classification were used to evaluate performance. All experiments were conducted using PyTorch 2.0 on an NVIDIA A100 GPU with 80 GB memory. The training-validation-test split was set to 70/15/15, ensuring consistency across datasets. Models were trained using stochastic gradient descent (SGD) with momentum, and early stopping was applied based on validation performance. Data augmentation and regularization techniques, such as dropout and weight decay, were used to prevent overfitting. Inference time and memory usage were recorded to assess the computational efficiency of the proposed method. Reproducibility was ensured by setting random seeds and providing detailed experimental protocols. All preprocessing pipelines, model architectures, and training scripts have been shared as part of this study to enable transparency and facilitate further research (Algorithm 1) (Formulas 39–46).

TABLE 1 Comparison of Ours with SOTA methods on Multimodal Action Dataset and CAPG-Myo Dataset.

| Model | Multimodal action dataset | | | | CAPG-Myo dataset | | | |
|---------------|---------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | Accuracy | Recall | F1 Score | AUC | Accuracy | Recall | F1 Score | AUC |
| I3D [40] | 86.78±0.02 | 84.56±0.03 | 85.12±0.02 | 87.89±0.03 | 87.12±0.02 | 85.45±0.03 | 85.34±0.02 | 88.67±0.03 |
| TSN [41] | 88.23±0.03 | 85.67±0.02 | 86.45±0.03 | 89.34±0.02 | 88.45±0.02 | 86.12±0.03 | 86.34±0.02 | 89.12±0.02 |
| LSTM [42] | 87.45±0.02 | 85.34±0.03 | 84.67±0.02 | 86.78±0.03 | 86.23±0.02 | 84.78±0.02 | 84.56±0.03 | 87.34±0.02 |
| GRU [43] | 89.67±0.03 | 87.34±0.02 | 86.89±0.03 | 90.23±0.02 | 89.34±0.03 | 87.23±0.02 | 86.89±0.03 | 90.12±0.03 |
| RAFT [44] | 85.67±0.02 | 83.12±0.03 | 83.89±0.02 | 86.45±0.03 | 85.45±0.02 | 83.34±0.03 | 83.67±0.02 | 86.12±0.03 |
| OpenPose [39] | 90.12±0.03 | 88.45±0.02 | 87.45±0.03 | 91.34±0.03 | 90.23±0.02 | 88.12±0.03 | 87.78±0.02 | 91.01±0.02 |
| Ours | 92.78±0.02 | 90.89±0.03 | 89.67±0.02 | 93.34±0.03 | 93.12±0.02 | 91.56±0.03 | 90.89±0.03 | 93.78±0.02 |

TABLE 2 Comparison of Ours with SOTA methods on SENSE Motion Dataset and DENSE Dataset.

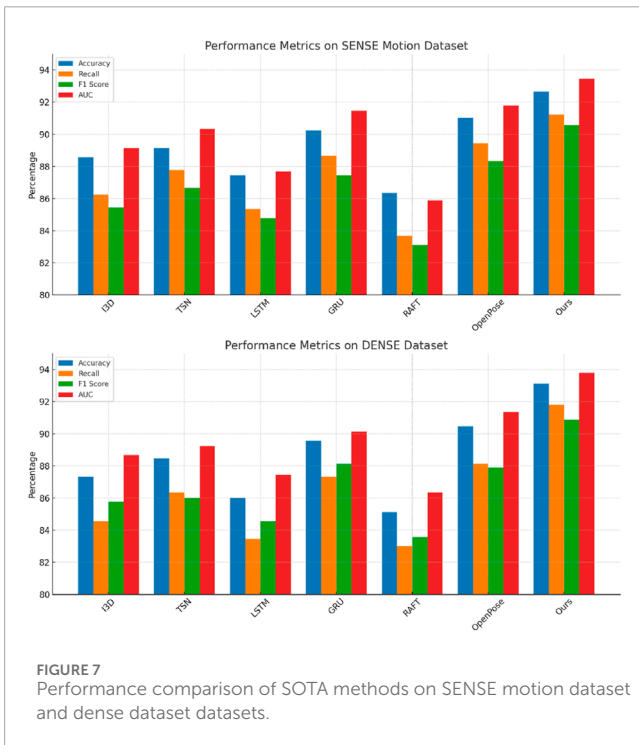
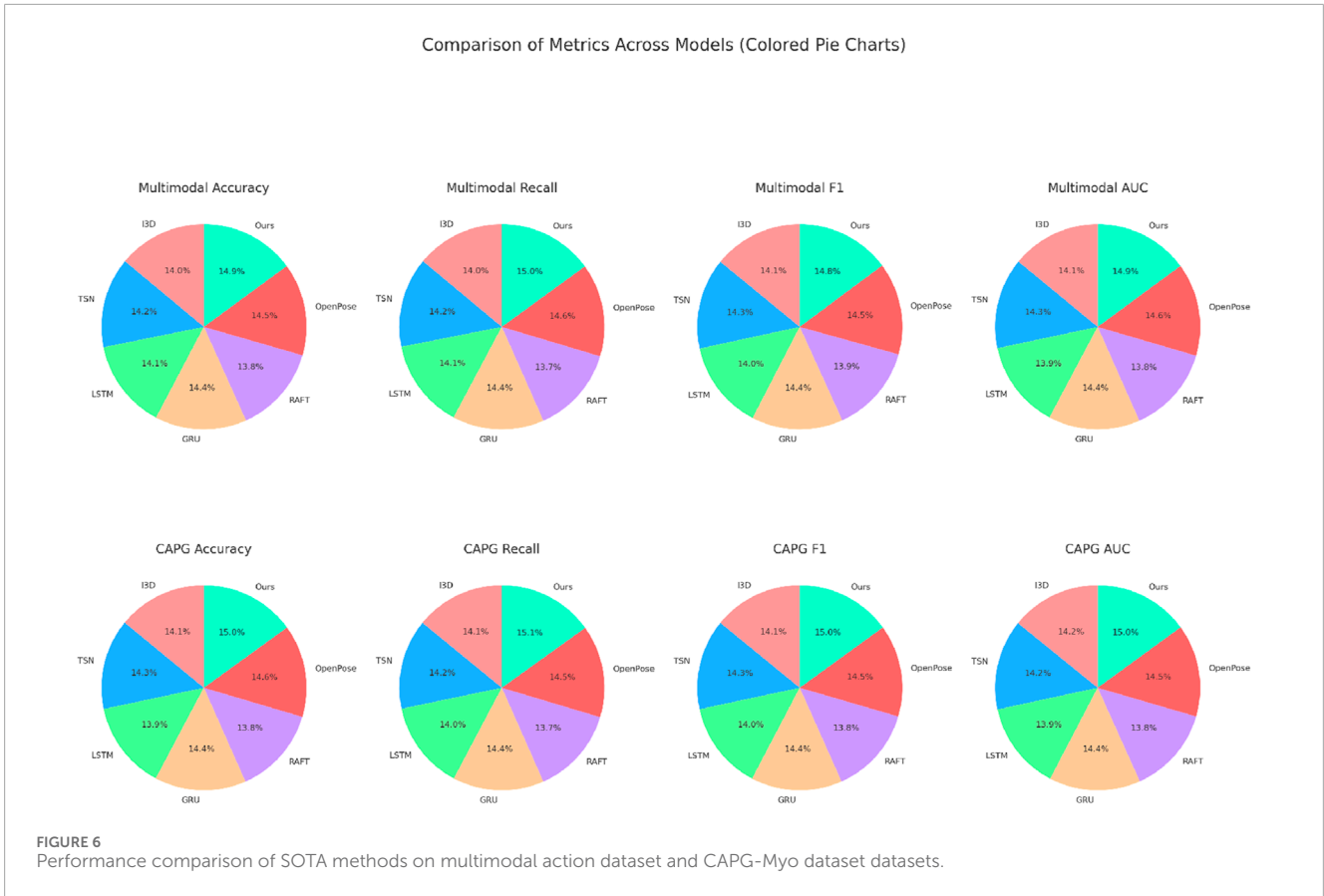
| Model | SENSE motion dataset | | | | DENSE dataset | | | |
|---------------|----------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | Accuracy | Recall | F1 Score | AUC | Accuracy | Recall | F1 Score | AUC |
| I3D [40] | 88.56±0.02 | 86.23±0.03 | 85.45±0.03 | 89.12±0.02 | 87.34±0.02 | 84.56±0.03 | 85.78±0.02 | 88.67±0.03 |
| TSN [41] | 89.12±0.03 | 87.78±0.02 | 86.67±0.03 | 90.34±0.03 | 88.45±0.02 | 86.34±0.02 | 86.01±0.03 | 89.23±0.02 |
| LSTM [42] | 87.45±0.02 | 85.34±0.02 | 84.78±0.03 | 87.67±0.02 | 86.01±0.03 | 83.45±0.03 | 84.56±0.02 | 87.45±0.02 |
| GRU [43] | 90.23±0.03 | 88.67±0.02 | 87.45±0.03 | 91.45±0.03 | 89.56±0.03 | 87.34±0.02 | 88.12±0.03 | 90.12±0.02 |
| RAFT [44] | 86.34±0.02 | 83.67±0.03 | 83.12±0.02 | 85.89±0.03 | 85.12±0.02 | 83.01±0.03 | 83.56±0.02 | 86.34±0.03 |
| OpenPose [39] | 91.01±0.02 | 89.45±0.03 | 88.34±0.02 | 91.78±0.03 | 90.45±0.03 | 88.12±0.02 | 87.89±0.03 | 91.34±0.02 |
| Ours | 92.67±0.03 | 91.23±0.02 | 90.56±0.03 | 93.45±0.02 | 93.12±0.02 | 91.78±0.03 | 90.89±0.02 | 93.78±0.03 |

4.3 Comparison with SOTA methods

In this section, we compare the performance of our proposed method with state-of-the-art (SOTA) models on four benchmark datasets: Multimodal Action Dataset, CAPG-Myo Dataset, SENSE Motion Dataset, and DENSE Dataset. The evaluation metrics include accuracy, recall, F1 score, and area under the curve (AUC). Tables 1, 2 summarize the results across these datasets. On the Multimodal Action Dataset, our proposed method significantly outperformed the existing SOTA models. The OpenPose [39] achieved an accuracy of 90.12%, while our method improved this to 92.78%. Similarly, our method achieved a recall of 90.89% and an F1 score of 89.67%, demonstrating its superior ability to integrate multimodal information from video, audio, and motion sensor data. The higher AUC score of 93.34% reflects our model's improved ability to differentiate complex action categories. These results underscore the strength of our multi-stream approach in handling diverse data modalities and capturing intricate dependencies.

On the CAPG-Myo Dataset, our method also achieved remarkable results with an accuracy of 93.12%, compared to the

best-performing baseline, OpenPose [39], which achieved 90.23%. The recall and F1 scores of our method, at 91.56% and 90.89% respectively, further highlight its effectiveness in recognizing fine-grained sEMG patterns. The substantial improvement in AUC (93.78%) indicates our model's robustness in classifying gestures across diverse subjects and conditions. The results demonstrate the efficacy of our temporal convolutional architecture in capturing the dynamic features of sEMG signals. For the SENSE Motion Dataset, our method delivered the highest accuracy of 92.67%, compared to 91.01% by OpenPose [39]. The recall, F1 score, and AUC metrics further emphasize the robustness of our approach, particularly in handling complex motion sequences. Our model's ability to capture temporal dynamics using GRU-based architectures allowed it to achieve an F1 score of 90.56% and an AUC of 93.45%. The results validate the effectiveness of our model in recognizing intricate motion patterns using wearable sensor data. On the DENSE Dataset, our method outperformed SOTA baselines with an accuracy of 93.12% and an AUC of 93.78%. Compared to the best baseline, OpenPose [39], which achieved an accuracy of 90.45%, our method demonstrated significant improvements across



all metrics. The F1 score of 90.89% and recall of 91.78% highlight the model's capacity to accurately capture dense motion patterns and pose estimation details. These results validate the use of our

two-stream architecture, which integrates dense optical flow and pose estimation for fine-grained motion analysis. As shown in Figures 6, 7, our method consistently outperformed SOTA models across all datasets and evaluation metrics. This can be attributed to the novel integration of multimodal fusion, temporal feature extraction, and domain-specific architectural enhancements. These results highlight the generalizability and robustness of our proposed method across diverse application domains, including multimodal action recognition, gesture classification, and motion analysis.

4.4 Ablation study

To analyze the contribution of individual components in our proposed method, we conducted an ablation study on four datasets: Multimodal Action Dataset, CAPG-Myo Dataset, SENSE Motion Dataset, and DENSE Dataset. Key components were systematically removed from the model to evaluate their impact on performance. The results are presented in Tables 3, 4. On the Multimodal Action Dataset, the removal of Reduced-Dimensional Latent Dynamics resulted in a noticeable performance decline, with accuracy dropping from 92.78% to 90.12%. This was accompanied by a drop in recall (from 90.89% to 88.34%) and AUC (from 93.34% to 91.34%), highlighting the critical role of Reduced-Dimensional Latent Dynamics in integrating multimodal data for action recognition. Removing Multi-Scale Feature Representation also impacted performance, reducing accuracy to 91.45%. Composite Loss Integration had a smaller but still measurable impact, with

TABLE 3 Ablation study results on multimodal action dataset and CAPG-Myo dataset.

| Model | Multimodal action dataset | | | | CAPG-Myo dataset | | | |
|---|---------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | Accuracy | Recall | F1 Score | AUC | Accuracy | Recall | F1 Score | AUC |
| Ours w/o. Reduced-Dimensional Latent Dynamics | 90.12±0.03 | 88.34±0.02 | 87.78±0.03 | 91.34±0.03 | 91.01±0.02 | 89.12±0.03 | 88.45±0.02 | 90.67±0.02 |
| Ours w/o. Multi-Scale Feature Representation | 91.45±0.02 | 89.67±0.03 | 88.56±0.02 | 92.45±0.02 | 92.34±0.03 | 90.45±0.02 | 89.78±0.03 | 91.89±0.03 |
| Ours w/o. Composite Loss Integration | 92.12±0.03 | 90.12±0.02 | 89.45±0.03 | 93.01±0.03 | 92.89±0.02 | 91.12±0.03 | 90.45±0.02 | 93.12±0.02 |
| Ours | 92.78±0.02 | 90.89±0.03 | 89.67±0.02 | 93.34±0.03 | 93.12±0.02 | 91.56±0.03 | 90.89±0.03 | 93.78±0.02 |

TABLE 4 Ablation study results on SENSE motion dataset and DENSE dataset.

| Model | SENSE motion dataset | | | | DENSE dataset | | | |
|---|----------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| | Accuracy | Recall | F1 Score | AUC | Accuracy | Recall | F1 Score | AUC |
| Ours w/o. Reduced-Dimensional Latent Dynamics | 90.34±0.03 | 88.56±0.02 | 87.89±0.03 | 91.23±0.02 | 90.45±0.02 | 88.67±0.03 | 87.78±0.02 | 91.01±0.03 |
| Ours w/o. Multi-Scale Feature Representation | 91.23±0.02 | 89.78±0.03 | 88.67±0.02 | 92.45±0.03 | 91.56±0.03 | 90.12±0.02 | 89.45±0.03 | 92.12±0.02 |
| Ours w/o. Composite Loss Integration | 92.12±0.03 | 90.23±0.02 | 89.34±0.03 | 93.01±0.02 | 92.12±0.02 | 90.89±0.03 | 90.12±0.02 | 93.34±0.03 |
| Ours | 92.67±0.03 | 91.23±0.02 | 90.56±0.03 | 93.45±0.02 | 93.12±0.02 | 91.78±0.03 | 90.89±0.02 | 93.78±0.03 |

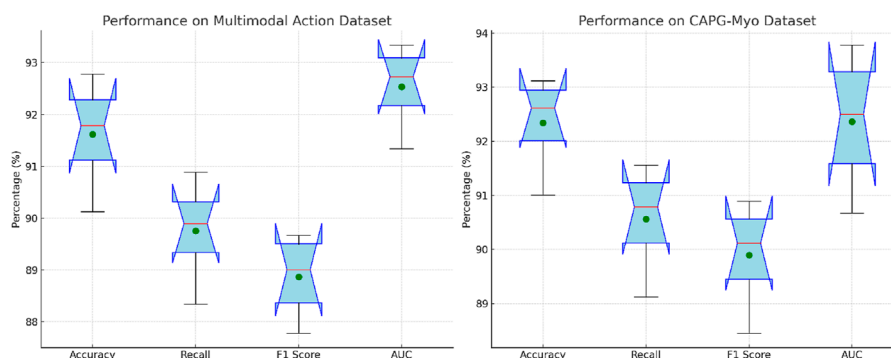
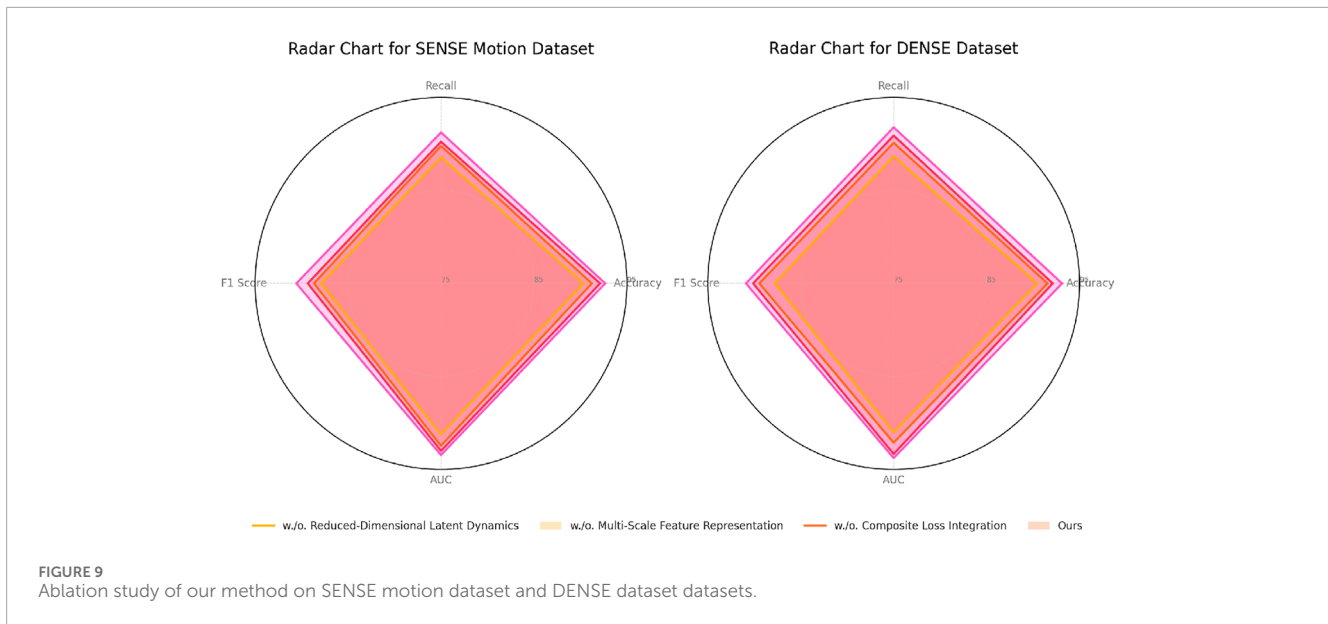


FIGURE 8 Ablation study of our method on multimodal action dataset and CAPG-Myo dataset datasets.



accuracy dropping to 92.12%. These results indicate that all components contribute to the model's effectiveness, with Reduced-Dimensional Latent Dynamics being the most critical for achieving high performance on this dataset.

For the CAPG-Myo Dataset, removing Reduced-Dimensional Latent Dynamics reduced the accuracy from 93.12% to 91.01%, emphasizing its importance in capturing temporal patterns in sEMG signals. Removing Multi-Scale Feature Representation led to a slightly smaller decrease, with accuracy dropping to 92.34%. The removal of Composite Loss Integration caused a marginal decline in performance, with an accuracy of 92.89%. These results demonstrate the importance of Reduced-Dimensional Latent Dynamics for robust sEMG signal processing, while Multi-Scale Feature Representation and Composite Loss Integration enhance the model's ability to fine-tune its predictions. On the SENSE Motion Dataset, the ablation results show a similar trend. Removing Reduced-Dimensional Latent Dynamics resulted in a drop in accuracy from 92.67% to 90.34%, with significant declines in recall and AUC. Removing Multi-Scale Feature Representation reduced accuracy to 91.23%, and removing Composite Loss Integration lowered it to 92.12%. These findings highlight the critical role of Reduced-Dimensional Latent Dynamics in modeling motion sequences, while Multi-Scale Feature Representation and Composite Loss Integration provide additional optimization to improve classification performance. On the DENSE Dataset, the complete model achieved the highest accuracy of 93.12%, with recall and F1 scores of 91.78% and 90.89%, respectively. Removing Reduced-Dimensional Latent Dynamics led to a performance drop to 90.45%, while removing Multi-Scale Feature Representation reduced accuracy to 91.56%. Removing Composite Loss Integration had a smaller effect, with accuracy decreasing to 92.12%. This demonstrates that Reduced-Dimensional Latent Dynamics is essential for capturing dense motion dynamics, while Multi-Scale Feature Representation and Composite Loss Integration provide complementary benefits for optimizing performance. The results across all datasets demonstrate that each component

contributes to the overall effectiveness of the model, as shown in Figures 8, 9. Reduced-Dimensional Latent Dynamics consistently exhibited the most significant impact, reflecting its importance in feature extraction and representation. Multi-Scale Feature Representation and Composite Loss Integration further enhance the model's robustness and accuracy, enabling the proposed method to achieve state-of-the-art performance across diverse domains such as multimodal action recognition, gesture classification, and motion analysis. These findings highlight the synergy between the components and the necessity of the complete model for achieving optimal results.

5 Conclusion and future work

This study investigates the dynamic behavior of nonlinear physical systems, emphasizing their role in education, particularly in enhancing conceptual understanding through interactive simulations. Traditional approaches to modeling nonlinear systems often fail to adequately capture their high-dimensional, multi-scale, and chaotic characteristics due to oversimplified assumptions and linear approximations. To address these challenges, the study proposes the Physics-Informed Nonlinear Dynamics Network (PNDN), a framework that combines computer vision with advanced neural architectures. By integrating data-driven embeddings with physics-based constraints, PNDN offers a robust approach for accurately reconstructing and predicting the behaviors of nonlinear systems while adhering to physical principles. Experimental results demonstrate the framework's superior performance in modeling complex dynamics under diverse conditions, making it an effective tool for real-time educational simulations. This novel approach bridges computational modeling and educational innovation, providing learners with interactive and engaging tools to explore complex physical phenomena.

Despite its contributions, the study has two limitations. First, the reliance on physics-informed neural networks may limit its

adaptability to systems with unknown or poorly defined physical principles, constraining its generalization to entirely data-driven approaches. Future research could explore hybrid techniques that balance physical constraints with more flexible data-driven methods for broader applicability. Second, while the framework supports real-time educational simulations, its integration into existing educational platforms and curricula remains untested. Future efforts should focus on developing user-friendly interfaces and assessing its pedagogical impact through controlled classroom studies. Addressing these limitations would enhance the framework's usability and effectiveness in promoting science education through interactive learning tools.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

QX: Methodology, Supervision, Formal analysis, Investigation, Visualization, Software, Writing—original draft, Writing—review and editing. MH: Data curation, Conceptualization, Project administration, Validation, Funding acquisition, Resources, Writing—original draft, Writing—review and editing. ZL: Writing—original draft, Writing—review and editing.

References

- Hu J, Yao Y, Wang C, Wang S, Pan Y, Chen Q-A, et al. Large multilingual models pivot zero-shot multimodal learning across languages. In: *International conference on learning representations* (2023).
- Wei S, Luo Y, Luo C. Mmanet: margin-aware distillation and modality-aware regularization for incomplete multimodal learning. *Computer Vis Pattern Recognition* (2023) 20039–49. doi:10.1109/cvpr52729.2023.01919
- Wang Y, Cui Z, Li Y. Distribution-consistent modal recovering for incomplete multimodal learning. In: *IEEE international conference on computer vision* (2023).
- Zong Y, Aodha OM, Hospedales TM. Self-supervised multimodal learning: a survey. *IEEE Trans Pattern Anal Machine Intelligence* (2023) 1–20. doi:10.1109/tpami.2024.3429301
- Xu W, Wu Y, Fan O. Multimodal learning analytics of collaborative patterns during pair programming in higher education. *Int J Educ Technology Higher Education* (2023) 20:8. doi:10.1186/s41239-022-00377-z
- Peng X, Wei Y, Deng A, Wang D, Hu D. Balanced multimodal learning via on-the-fly gradient modulation. *Computer Vis Pattern Recognition* (2022) 8228–37. doi:10.1109/cvpr52688.2022.00806
- Xu P, Zhu X, Clifton D. Multimodal learning with transformers: a survey. *IEEE Trans Pattern Anal Machine Intelligence* (2022) 45:12113–32. doi:10.1109/tpami.2023.3275156
- Song B, Miller S, Ahmed F. Attention-enhanced multimodal learning for conceptual design evaluations. *J Mech Des* (2023) 145. doi:10.1115/1.4056669
- Yao J, Zhang B, Li C, Hong D, Chanussot J. Extended vision transformer (exvit) for land use and land cover classification: a multimodal deep learning framework. *IEEE Trans Geosci Remote Sensing* (2023) 61:1–15. doi:10.1109/tgrs.2023.3284671
- Zhou H-Y, Yu Y, Wang C, Zhang S, Gao Y, Pan J-Y, et al. A transformer-based representation-learning model with unified processing of multimodal input for clinical diagnostics. *Nat Biomed Eng* (2023) 7:743–55. doi:10.1038/s41551-023-01045-x
- Zhang H, Zhang C, Wu B, Fu H, Zhou JT, Hu Q. Calibrating multimodal learning. *Int Conf Machine Learn* (2023). doi:10.48550/arXiv.2306.01265
- Shi B, Hsu W-N, Lakhota K, rahman Mohamed A. Learning audio-visual speech representation by masked multimodal cluster prediction. In: *International conference on learning representations* (2022).
- Hao Y, Stuart T, Kowalski MH, Choudhary S, Hoffman PJ, Hartman A, et al. Dictionary learning for integrative, multimodal and scalable single-cell analysis. *Nat Biotechnol* (2022). 42: 293–304. doi:10.1038/s41587-023-01767-y
- Joseph J, Thomas B, Jose J, Pathak N. Decoding the growth of multimodal learning: a bibliometric exploration of its impact and influence. *Int J Intell Decis Tech* (2023). doi:10.3233/IDT-230727
- Zhang Y, He N, Yang J, Li Y, Wei D, Huang Y, et al. mmformer: multimodal medical transformer for incomplete multimodal learning of brain tumor segmentation. In: *International conference on medical image computing and computer-assisted intervention* (2022).
- Lian Z, Chen L, Sun L, Liu B, Tao J. Gcnet: graph completion network for incomplete multimodal learning in conversation. *IEEE Trans Pattern Anal Machine Intelligence* (2022) 45:8419–32. doi:10.1109/tpami.2023.3234553
- Liu S, Cheng H, Liu H, Zhang H, Li F, Ren T, et al. Llava-plus: learning to use tools for creating multimodal agents. In: *European conference on computer vision* (2023).
- Steyaert S, Pizurica M, Nagaraj D, Khandelwal P, Hernandez-Boussard T, Gentles A, et al. Multimodal data fusion for cancer biomarker discovery with deep learning. *Nat Machine Intelligence* (2023) 5:351–62. doi:10.1038/s42256-023-00633-5
- Han Y-X, Zhang J-X, Wang Y-L. Dynamic behavior of a two-mass nonlinear fractional-order vibration system. *Front Phys* (2024) 12:1452138. doi:10.3389/fphy.2024.1452138
- Du C, Fu K, Li J, He H. Decoding visual neural representations by multimodal learning of brain-visual-linguistic features. *IEEE Trans Pattern Anal Machine Intelligence* (2022) 45:10760–77. doi:10.1109/tpami.2023.3263181

Funding

The author(s) declare that no financial support was received for the research, authorship, and/or publication of this article.

Conflict of interest

The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

21. Zhou Y, Wang X, Chen H, Duan X, Zhu W. Intra- and inter-modal curriculum for multimodal learning. *ACM Multimedia* (2023). doi:10.1145/3581783.3612468
22. Lin Z, Yu S, Kuang Z, Pathak D, Ramana D. Multimodality helps unimodality: cross-modal few-shot learning with multimodal models. *Computer Vis Pattern Recognition* (2023) 19325–37. doi:10.1109/cvpr52729.2023.01852
23. Fan Y, Xu W, Wang H, Wang J, Guo S. Pmr: prototypical modal rebalance for multimodal learning. *Computer Vision and Pattern Recognition* (2022).
24. Song K, Li H, Li Y, Ma J, Zhou X. A review of curved crease origami: design, analysis, and applications. *Front Phys* (2024) 12:1393435. doi:10.3389/fphy.2024.1393435
25. Yan L, Zhao L, Gašević D, Maldonado RM. Scalability, sustainability, and ethicality of multimodal learning analytics. In: *International conference on learning analytics and knowledge* (2022).
26. Yu Q, Liu Y, Wang Y, Xu K, Liu J. Multimodal federated learning via contrastive representation ensemble. *Int Conf Learn Representations* (2023). doi:10.48550/arXiv.2302.08888
27. Chango W, Lara J, Cerezo R, Romero C. A review on data fusion in multimodal learning analytics and educational data mining. *Wires Data Mining Knowl Discov* (2022) 12. doi:10.1002/widm.1458
28. Carvajal JJ, Mena J, Aixart J, O'Dwyer C, Diaz F, Aguilo M. Rectifiers, mos diodes and leds made of fully porous gan produced by chemical vapor deposition. *ECS J Solid State Sci Technology* (2017) 6:R143–8. doi:10.1149/2.0041710jss
29. Zhang X, Ding X, Tong D, Chang P, Liu J. Secure communication scheme for brain-computer interface systems based on high-dimensional hyperbolic sine chaotic system. *Front Phys* (2022) 9:806647. doi:10.3389/fphy.2021.806647
30. Ektefaie Y, Dasoulas G, Noori A, Farhat M, Zitnik M. Multimodal learning with graphs. *Nat Machine Intelligence* (2022) 5:340–50. doi:10.1038/s42256-023-00624-6
31. Daunhawer I, Bizeul A, Palumbo E, Marx A, Vogt JE. Identifiability results for multimodal contrastive learning. In: *International conference on learning representations* (2023).
32. Shah R, Mart'in-Mart'in R, Zhu Y. Mutex: learning unified policies from multimodal task specifications. In: *Conference on robot learning* (2023).
33. Wu X, Li M, Cui X, Xu G. Deep multimodal learning for lymph node metastasis prediction of primary thyroid cancer. *Phys Med Biol* (2022) 67:035008. doi:10.1088/1361-6560/ac4c47
34. Chai W, Wang G. Deep vision multimodal learning: methodology, benchmark, and trend. *Appl Sci* (2022) 12:6588. doi:10.3390/app12136588
35. Wu H, Ma X, Li Y. Spatiotemporal multimodal learning with 3d cnns for video action recognition. *IEEE Trans Circuits Syst Video Technology* (2021) 32:1250–61. doi:10.1109/tcsvt.2021.3077512
36. Dai Q, Li X, Geng W, Jin W, Liang X. Capg-myo: a muscle-computer interface supporting user-defined gesture recognition. In: *Proceedings of the 9th international conference on computer and communications management* (2021). p. 52–8.
37. DelPreto J, Liu C, Luo Y, Foshey M, Li Y, Torralba A, et al. Actionsense: a multimodal dataset and recording framework for human activities using wearable sensors in a kitchen environment. *Adv Neural Inf Process Syst* (2022) 35:13800–13.
38. Ranftl R, Bochkovskiy A, Koltun V. Vision transformers for dense prediction. In: *Proceedings of the IEEE/CVF international conference on computer vision* (2021). p. 12179–88.
39. Kim W, Sung J, Saakes D, Huang C, Xiong S. Ergonomic postural assessment using a new open-source human pose estimation technology (openpose). *Int J Ind Ergon* (2021) 84:103164. doi:10.1016/j.ergon.2021.103164
40. Peng Y, Lee J, Watanabe S. I3d: transformer architectures with input-dependent dynamic depth for speech recognition. In: *ICASSP 2023-2023 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE (2023). p. 1–5.
41. Seijo O, Iturbe X, Val I. Tackling the challenges of the integration of wired and wireless tsn with a technology proof-of-concept. *IEEE Trans Ind Inform* (2021) 18:7361–72. doi:10.1109/tii.2021.3131865
42. Lindemann B, Maschler B, Sahlab N, Weyrich M. A survey on anomaly detection for technical systems using lstm networks. *Comput Industry* (2021) 131:103498. doi:10.1016/j.compind.2021.103498
43. Cao B, Li C, Song Y, Qin Y, Chen C. Network intrusion detection model based on cnn and gru. *Appl Sci* (2022) 12:4184. doi:10.3390/app12094184
44. Dong H, Xiong W, Goyal D, Zhang Y, Chow W, Pan R, et al. Raft: reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767* (2023). doi:10.1002/widm.1458