# MMA: A Momentum Mamba Architecture for Human Activity Recognition with Inertial Sensors

Thai-Khanh Nguyen [§], Uyen Vo [§], Tan M. Nguyen, Thieu N. Vo, Trung-Hieu Le, Cuong Pham

*Abstract*—**Human activity recognition (HAR) from inertial sensors is essential for ubiquitous computing, mobile health, and ambient intelligence. Conventional deep models such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and transformers have advanced HAR but remain limited by vanishing or exloding gradients, high computational cost, and difficulty in capturing long-range dependencies. Structured state-space models (SSMs) like Mamba address these challenges with linear complexity and effective temporal modeling, yet they are restricted to first-order dynamics without stable long-term memory mechanisms. We introduce Momentum Mamba, a momentum-augmented SSM that incorporates second-order dynamics to improve stability of information flow across time steps, robustness, and long-sequence modeling. Two extensions further expand its capacity: Complex Momentum Mamba for frequency-selective memory scaling. Experiments on multiple HAR benchmarks demonstrate consistent gains over vanilla Mamba and Transformer baselines in accuracy, robustness, and convergence speed. With only moderate increases in training cost, momentum-augmented SSMs offer a favorable accuracy-efficiency balance, establishing them as a scalable paradigm for HAR and a promising principal framework for broader sequence modeling applications.**

*Index Terms*—**Human activity recognition, inertial sensors, state space models, Mamba architecture, momentum dynamics, vanishing gradient, complex-valued neural networks, wearable computing.**

## I. Introduction

Human activity recognition (HAR) aims to identify human actions from sensor data automatically and supports applications in healthcare [1]–[4], human-computer interaction [5], [6], lifestyle monitoring [7], [8], sports analysis, and ambient assisted living [1], [9]. Traditional HAR systems have relied on visual inputs such as RGB videos, depth maps, and skeletal data [10]–[12], but these approaches face challenges of occlusion, lighting sensitivity, computational cost, and privacy concerns [10], [13].

These limitations have driven a growing transition toward non-visual modalities, particularly inertial sensing through accelerometers and gyroscopes [14]. Due to the mass adoption of smartphones and wearable devices, IMU data have become increasingly ubiquitous and accessible for HAR applications [15]. Inertial HAR offers unique advantages: privacy preservation, low power consumption, and fine-grained motion capture suitable for continuous and real-time monitoring [16], [17]. However, modeling inertial data introduces its own challenges. Sensor signals are inherently noisy, lack spatial context, and require robust temporal modeling of multivariate time series under resource-constrained conditions [18], [19].

Deep learning has become the dominant paradigm for inertial HAR. Convolutional Neural Networks (CNN) extract local temporal patterns [20], whereas Recurrent Neural Networks (RNN) are designed to model sequential dependencies but often struggle to capture long-range temporal relationships due to issues such as vanishing gradients [21], [22]. Hybrid CNN–RNN architectures have been developed to integrate the complementary strengths of both models [23]. Transformer-based models [24] have been introduced to capture global contextual dependencies through self-attention mechanisms, with HAR-specific extensions such as GAFormer [25] and MAMC [26], but their quadratic computational and memory complexity causes the difficulty for real-time deployment.

Structured state-space models (SSMs) offer a scalable alternative to Transformers, achieving linear-time complexity in sequence length. Among these, Mamba [27] represents a significant advance, introducing an input-dependent selective mechanism that allows the model to dynamically focus on relevant information, with variants like HARMamba [28] being adapted for HAR. Despite its effectiveness, Mamba's reliance on first-order dynamics can limit its ability to maintain stable gradient flow over long sequences, a critical challenge for noisy inertial data. The principle of incorporating second-order dynamics, inspired by momentum optimization methods [29]–[31], offers a robust solution to this stability issue. This has been demonstrated in architectures such as MomentumRNN [32] and, more recently in the SSM context, by LinOSS [33], which uses oscillatory dynamics for stable long-range modeling. However, LinOSS employs time-invariant (i.e., static) state transition parameters that do not adapt to the input, thereby lacking the content-aware selectivity that is Mamba's key advantage. This exposes a clear research gap: no existing model integrates the adaptive selectivity of Mamba with the enhanced stability of second-order dynamics.

**Contribution**: We propose *Momentum Mamba (MMA)* architecture, that augments Mamba's dynamics with heavy-ball momentum. This integration enhances gradient stability, noise robustness, and temporal expressiveness while retaining linear scalability. Our contribution is three-fold:

1) We introduce *Momentum Mamba*, the first selective

§These authors contributed equally.

Thai-Khanh Nguyen and Trung-Hieu Le are with the Falcuty of Information Technology, Dainam University, Hanoi University of Science and Technology, Hanoi 10000, Vietnam.

Cuong Pham and Uyen Vo are with the Faculty of Artificial Intelligence, Posts and Telecommunications Institute of Technology, Hanoi 10000, Vietnam.

Thieu N. Vo is with the Department of Computer Science, University of Bath, Bath BA2 7AY, UK.

Tan Nguyen is with the Departments of Mathematics, National University of Singapore, Singapore 119076, Singapore.

Corresponding author: Cuong Pham (e-mail: cuongpv@ptit.edu.vn).

state-space model (SSM) that integrates momentum-driven second-order dynamics for improving gradient stability and temporal modeling in long-sequence inertial HAR.

2) We theoretically prove that Momentum Mamba achieves better-structured spectrum than the Mamba baseline Mamba, explaining to the model's stability enhancement.

3) We demonstrate that the design principles of Momentum Mamba extend naturally to other advanced momentum-based optimization methods. In particular, we introduce Complex Momentum Mamba and Adam Momentum Mamba, which achieve frequency-selective memory via complex-valued momentum and adaptive control of momentum scaling, respectively.

Extensive evaluations on a variety of inertial HAR benchmarks demonstrate that our models consistently outperform transformer, vanilla Mamba, and oscillatory SSM baselines, achieving superior trade-offs in accuracy, convergence speed, and resource efficiency.

## II. PRELIMINARIES: SELECTIVE STATE SPACE MODEL

The classical state space model describes a continuous-time system that maps an input The classical state space model describes a continuous-time system that maps an input $x(t) \in \mathbb{R}$ to an output $y(t) \in \mathbb{R}$ via a hidden state $\mathbf{h}(t) \in \mathbb{R}^{d \times 1}$, formulated as:

$$\begin{aligned} \boldsymbol{h}'(t) &= \boldsymbol{A}\boldsymbol{h}(t) + \boldsymbol{B}x(t), \\ y(t) &= \boldsymbol{C}\boldsymbol{h}(t) + Dx(t), \end{aligned} \tag{1}$$

where $\boldsymbol{A} \in \mathbb{R}^{d \times d}$, $\boldsymbol{B}, \boldsymbol{h}(t), \boldsymbol{h}'(t) \in \mathbb{R}^{d \times 1}$, $\boldsymbol{C} \in \mathbb{R}^{1 \times d}$, and $D \in \mathbb{R}$. The *structured state space sequence models (S4)* are inspired by the classical SSMs and take the following discrete form of the continuous system in Eqn. 1.

$$\boldsymbol{h}_n = \overline{\boldsymbol{A}}\boldsymbol{h}_{n-1} + \overline{\boldsymbol{B}}x_n, \tag{2a}$$

$$y_n = \boldsymbol{C}\boldsymbol{h}_n, \tag{2b}$$

where $\Delta \in \mathbb{R}$ is a timescale parameter, and $\overline{\boldsymbol{A}}, \overline{\boldsymbol{B}}$ are discretized counterparts of $\boldsymbol{A}, \boldsymbol{B}$ via zero-order hold discretization. Specifically, $\overline{\boldsymbol{A}} = \exp(\Delta \boldsymbol{A})$, and $\overline{\boldsymbol{B}} = (\Delta \boldsymbol{A})^{-1}(\exp(\Delta \boldsymbol{A}) - \boldsymbol{I}) \cdot \Delta \boldsymbol{B} \approx \Delta \boldsymbol{B}$. Structured SSMs derive their name from the requirement that the matrix $A$, which governs the temporal dynamics, must adopt a specific structure to enable efficient sequence-to-sequence transformations suitable for deep neural networks. The initial designs introduced were the diagonal plus low-rank (DPLR) structure [34] and the purely diagonal structure [35]–[37], with the latter remaining the most widely used.

*Selective State Space Models* (SSMs), such as *Mamba* [38], extend S4 by selectively attending to or disregarding inputs at each timestep. Specifically, the parameters $\boldsymbol{B}$, $\boldsymbol{C}$, and $\Delta$ are defined as functions of the input $x_n$, yielding input-dependent forms $\boldsymbol{B}_n$, $\boldsymbol{C}_n$, and $\Delta_n$. Consequently, the discretized parameters also become input-dependent, with $\overline{A}_n = \exp(\Delta_n \boldsymbol{A})$ and $\overline{\boldsymbol{B}}_n = \Delta_n \boldsymbol{B}_n$. In *Mamba*, these dependencies are instantiated as

$$\begin{aligned} \boldsymbol{B}_n &= \text{Linear}_d(x_n), \\ \boldsymbol{C}_n &= \text{Linear}_d(x_n), \\ \Delta_n &= \text{softplus}\big(\theta + \text{Broadcast}_D\big(\text{Linear}_1(x_n)\big)\big), \end{aligned}$$

where $\text{Linear}_d$ denotes a learnable projection into a $d$-dimensional space, $\text{Broadcast}_D$ expands a scalar into a $D$-dimensional vector, and $\theta \in \mathbb{R}^D$ is a learnable parameter vector.

Compared to S4, *Mamba* demonstrates better performance on information-dense data, such as language, particularly as the state dimension $d$ increases, thereby enhancing its information capacity [38], [39]. However, Mamba is still restricted to first-order dynamics, thereby still having difficulty in capture long-range dependencies in long input sequences. Our Proposition 1 proved that Mamba indeed suffers from vanishing and exploding gradient issues.

**Proposition 1** (Vanishing Gradients in Mamba). *Let $\{h_n\}$ be the hidden states of the Mamba architecture defined by*

$$h_n = \overline{A}_n h_{n-1} + \overline{B}_n x_n,$$

*where $\overline{A}_n = \exp(\Delta_n A)$ is diagonal with entries $e^{\Delta_n a_{n,i}}$ for $a_{n,i} < 0$ and $\Delta_n > 0$. Then the gradient of the loss $L$ with respect to $h_t$ satisfies*

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_T} \prod_{n=t+1}^{T} \overline{A}_n.$$

*If $\min_i a_{n,i} \ll 0$, then*

$$\left\| \prod_{n=t+1}^{T} \overline{A}_n \right\| \to 0 \quad as \quad T - t \to \infty,$$

*so gradients vanish exponentially with sequence length.*

*Proof.* By the chain rule,

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_T} \cdot \frac{\partial h_T}{\partial h_t}.$$

Unrolling the recurrence yields

$$\frac{\partial L}{\partial h_t} = \frac{\partial L}{\partial h_T} \cdot \prod_{n=t+1}^{T} \frac{\partial h_n}{\partial h_{n-1}} = \frac{\partial L}{\partial h_T} \cdot \prod_{n=t+1}^{T} \overline{A}_n,$$

since $\partial h_n / \partial h_{n-1} = \overline{A}_n$. Now, $\overline{A}_n$ is diagonal with entries $e^{\Delta_n a_{n,i}}$ where $a_{n,i} < 0$. Thus,

$$0 < e^{\Delta_n a_{n,i}} < 1 \quad \forall i,$$

meaning each factor strictly contracts the corresponding gradient component. Repeated multiplication across $T - t$ steps yields exponential decay in the gradient magnitude. This aligns with the classical description of the *vanishing gradient* [40], where Jacobians with eigenvalues strictly less than one in modulus drive gradients towards zero as the horizon grows. $\square$

We further provide empirical evidence to validate Proposition 1 in Figure 4 below.

## III. RELATED WORK

This section reviews prior studies that form the foundation of our work. We organize the discussion into three parts: (A) the evolution of deep learning models for inertial HAR, (B) structured state-space models (SSMs) with a focus on Mamba and oscillatory extensions, and (C) momentum-inspired neural dynamics in sequence modeling. Together, these perspectives contextualize our proposed integration of momentum-driven recurrences into Mamba for robust inertial HAR.

### A. Deep Learning Models for Inertial HAR

HAR using wearable inertial sensors such as accelerometers and gyroscopes is a central task in mobile health and ubiquitous computing. Traditional approaches relied on hand-crafted features and shallow classifiers such as decision trees, support vector machines, and naive Bayes [41]–[43], but deep learning has enabled end-to-end representation learning from raw multivariate time series [44].

CNN-based methods were among the earliest to show strong performance by capturing local temporal patterns directly from accelerometer and gyroscope signals [20], [45]. Lightweight variants further improved feasibility on embedded devices [46]. However, CNNs are constrained by limited receptive fields and often fail to capture long-term dependencies. To address this, RNNs such as LSTMs and GRUs were adopted to explicitly model sequential dynamics [47]. Despite their effectiveness, these models suffer from vanishing gradients, slower training, and high memory consumption. Hybrid models such as DeepConvLSTM [23] attempted to combine convolutional and recurrent layers, but the added complexity reduces suitability for real-time wearable systems.

Transformers introduced self-attention for global temporal modeling and have achieved state-of-the-art results in several HAR benchmarks [48], [49]. To further adapt Transformers to inertial sensing, GAFormer [25] integrates Gramian Angular Field representations with graph alignment to capture structural dependencies in sensor sequences, while MAMC [26] introduces modality-aware attention for heterogeneous input fusion. Other variants such as LightFormer [50] and HP-Former [51] emphasize efficiency, with HPFormer in particular reducing self-attention complexity from $O(L^2)$ to $O(L \cdot \log L)$ to improve scalability in health informatics tasks. Despite these advances, the quadratic cost of vanilla self-attention and the high energy consumption of Transformer architectures remain major obstacles for real-time, resource-constrained HAR.

Beyond inertial-only methods, multimodal fusion combining inertial data with visual or depth inputs has been explored to improve robustness in noisy or occluded environments [52]–[54]. While these approaches enhance recognition performance, they are often limited by synchronization overhead, deployment cost, and energy consumption. By contrast, inertial sensors are inexpensive, ubiquitous, and energy-efficient, making them the most practical foundation for scalable, real-time HAR.

Overall, the trajectory of deep learning for inertial HAR spans CNNs, RNNs, and transformers, with recent work emphasizing lightweight and multimodal designs. However, a persistent challenge remains: balancing temporal modeling capacity with computational efficiency for resource-constrained deployment, motivating exploration of alternative paradigms such as structured state-space models.

### B. Structured State-Space Models and Mamba

Structured state-space models (SSMs) provide an efficient alternative to attention by modeling sequence dynamics with linear-time recurrence. Early advances such as S4 [34] leveraged the HiPPO framework [55] to achieve long-horizon modeling with linear complexity, later refined by S5 and S6 for stability and precision.

Mamba (S6) [27] further improved flexibility through input-dependent selective recurrence, combining linear scalability with content-aware dynamics. Variants such as HAR-Mamba [28], ActivityMamba [56], and MHAR [57] demonstrate their effectiveness for inertial and multimodal HAR, balancing accuracy with efficiency.

Despite these advances, both Mamba and its variants remain constrained by first-order recurrences, which limit gradient regulation and robustness under noisy sensor streams. This motivates our exploration of momentum-augmented higher-order dynamics for inertial HAR.

### C. Momentum-Based Deep Learning Models

Momentum, originally introduced in convex optimization [29], accelerates convergence and stabilizes gradient trajectories, and its variants, such as Adan [58] have been proven effective in large-scale deep learning. Beyond optimization, momentum has inspired neural architectures that embed second-order recurrences directly into model design. Examples include MomentumRNN [32], which augments hidden state updates with a velocity term, and NesterovNODE [59], which reformulates neural ODEs as second-order systems for stability. Extensions such as physics-informed SSMs [60] and generative SSMs for active inference [61] further illustrate the versatility of momentum-like principles.

LinOSS [33] exemplifies this trend by augmenting first-order SSMs with oscillatory second-order dynamics, achieving stability and efficiency for long-sequence modeling. However, existing designs employ either fixed oscillatory structures or remain outside the selective SSM framework. To address this gap, we propose *Momentum Mamba*, which integrates heavy-ball momentum into Mamba's input-dependent recurrence, enhancing gradient flow, noise robustness, and temporal expressiveness while preserving linear scalability.

## IV. MOMENTUM MAMBA

### A. Incorporating Heavy-Ball Momentum into Mamba

At the core of the framework lies the *Momentum Mamba* block, which extends the Mamba (S6) architecture–originally derived from continuous-time state-space models–with a momentum mechanism inspired by iterative optimization. In addition to the hidden state $h_n$, an auxiliary momentum state $v_n$ accumulates input-driven updates through learnable parameters $(\alpha, \beta)$. This dual pathway smooths high-frequency variations,
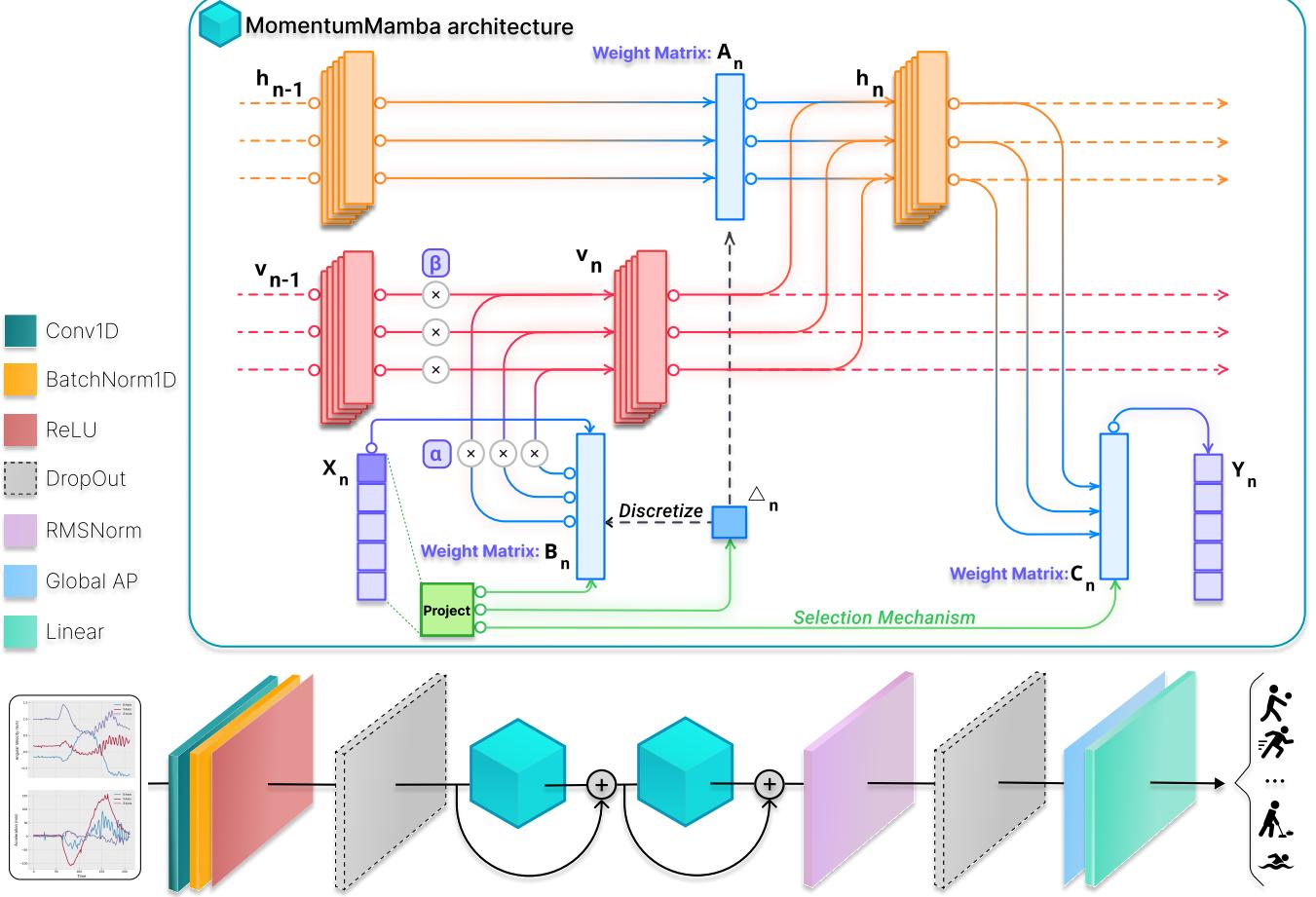
Fig. 1: Overall architecture of the proposed *Momentum Mamba* (MMA) framework for inertial HAR. The pipeline consists of three main stages: a lightweight Conv1D front-end for local feature extraction, stacked Momentum Mamba layers for temporal modeling, and a compact classification head for activity recognition. At its core, each Momentum Mamba block augments the standard Mamba recurrence with an auxiliary momentum state $v_n$ that accumulates input-driven updates through learnable parameters $(\alpha, \beta)$. This dual-state design smooths high-frequency fluctuations, stabilizes long-range dynamics, and improves robustness while preserving the linear-time scan efficiency of structured state-space models.

stabilizes long-range dynamics, and improves gradient flow while preserving the scan efficiency of SSMs.

Formally, the baseline Mamba recurrence is given by

$$h_n = \overline{\boldsymbol{A}}_n h_{n-1} + \overline{\boldsymbol{B}}_n x_n \quad (3)$$

where $\overline{\boldsymbol{A}}_n = \exp(\Delta_n A)$ and $\overline{\boldsymbol{B}}_n = (\Delta_n A)^{-1}(\overline{\boldsymbol{A}}_n - I)\Delta_n B(x_n)$ are derived from continuous-time parameters. The step size $\Delta_n = \mathrm{softplus}(\mathrm{Linear}(x_n))$, input projection $B(x_n)$, and output projection $C(x_n)$ are input-dependent, while $A$ is a fixed base matrix.

Drawing from the optimization analogy, we reinterpret $\overline{\boldsymbol{B}}_n x_n$ as a learned "gradient" and introduce a momentum accumulator to stabilize state evolution. We define an auxiliary momentum state $v_n \in \mathbb{R}^N$ and update the recurrence in three stages:

**(1) Momentum update:**

$$v_n = \beta v_{n-1} + \alpha \overline{\boldsymbol{B}}_n x_n, \quad (4)$$

where $\beta \in [0,1]$ is a learnable momentum decay coefficient and $\alpha \in \mathbb{R}$ is a learnable step size.

**(2) Hidden state update:**

$$h_n = \overline{\boldsymbol{A}}_n h_{n-1} + v_n, \quad (5)$$

so that the hidden state now evolves in response to a smoothed, temporally integrated signal rather than the raw instantaneous input.

**(3) Output projection:**

$$y_n = \boldsymbol{C}_n h_n, \quad (6)$$

with $\boldsymbol{C}_n = C(x_n)$ derived from a learned projection conditioned on input.

**Proposition 2** (Affine Recurrence Form). *The Momentum Mamba recurrence*

$$v_n = \beta v_{n-1} + \alpha \overline{\boldsymbol{B}}_n x_n, \quad (2)$$
$$h_n = \overline{\boldsymbol{A}}_n h_{n-1} + v_n, \quad (3)$$

*admits an equivalent affine formulation*

$$s_n = M'_n s_{n-1} + F'_n, \quad (7)$$

*where* $s_n = \begin{bmatrix} h_n \\ v_n \end{bmatrix} \in \mathbb{R}^{2N}$ *and*

$$M'_n = \begin{bmatrix} \overline{\boldsymbol{A}}_n & \beta I \\ 0 & \beta I \end{bmatrix}, \qquad F'_n = \begin{bmatrix} \alpha \overline{\boldsymbol{B}}_n x_n \\ \alpha \overline{\boldsymbol{B}}_n x_n \end{bmatrix}. \tag{8}$$

*Proof.* Stacking the hidden and momentum states into $s_n = [h_n^\top, v_n^\top]^\top$, the recurrence 2–3 can first be written as the constrained linear system

$$T s_n = M_n s_{n-1} + F_n, \tag{9}$$

with

$$T = \begin{bmatrix} I & -I \\ 0 & I \end{bmatrix}, \quad M_n = \begin{bmatrix} \overline{\boldsymbol{A}}_n & 0 \\ 0 & \beta I \end{bmatrix}, \quad F_n = \begin{bmatrix} 0 \\ \alpha \overline{\boldsymbol{B}}_n x_n \end{bmatrix}.$$

Multiplying both sides by $T^{-1}$, which exists and is given by

$$T^{-1} = \begin{bmatrix} I & I \\ 0 & I \end{bmatrix},$$

yields the affine recurrence equation 7. $\qquad\square$

**Remark 1** (Parallelization and Temporal Smoothing). *Proposition 2 shows that the coupled hidden–momentum dynamics admit an affine recurrence. Intuitively, this means the entire update can be expressed as repeated applications of an affine map to the augmented state. Such maps compose associatively, as captured by*

$$(a_1, a_2) \bullet (b_1, b_2) = (b_1 a_1, \; b_1 a_2 + b_2).$$

*Following the insights of LinOSS [33], this associativity ensures that the sequence $\{s_n\}_{n=1}^{L}$ can be computed in $\mathcal{O}(\log L)$ parallel time with $\mathcal{O}(L \cdot N)$ total computation, making the formulation highly compatible with GPU-based parallelization.*

*Beyond efficiency, the inclusion of the momentum state $v_n$ acts as an exponential moving average, attenuating high-frequency input fluctuations. Thus, Momentum Mamba simultaneously preserves the scan efficiency of state-space models and enhances stability through temporal smoothing, without incurring additional asymptotic cost.*

Taken together, the formal results above establish Momentum Mamba as a natural extension of the Mamba framework: it retains the algebraic structure required for efficient parallelization while enriching the dynamics with a momentum pathway. This dual perspective–optimization-inspired smoothing combined with state-space discretization–provides the foundation for several practical advantages, which we summarize below.

**Selective dynamics.** A defining feature of Mamba is its separation of dynamic modeling and input selectivity. The dynamics of the core transition is governed by the spectrum of the fixed matrix $A$, while the input dependence is injected through the learned functions $\Delta_n(x_n)$, $B_n(x_n)$, and $C_n(x_n)$. Our momentum-enhanced variant preserves this clean abstraction while inserting a smoothing buffer between the input and state transitions. The result is a more stable yet still content-aware recurrence.

**Computational efficiency.** The introduction of momentum adds moderate overhead in terms of memory footprint, primarily due to the storage of augmented momentum states; however, the complexity remains linear in sequence length and

is compatible with parallelization. The update rule remains linear in the hidden dimension and compatible with fast parallel scan operations. The parameter increase is negligible, involving only scalar momentum hyperparameters $\alpha$ and $\beta$.

**Expected benefits.** By integrating momentum into the Mamba architecture, we aim to combine the selectivity and long-range expressivity of state-space models with the smooth convergence and robustness of optimization-based dynamics. We hypothesize that this formulation leads to the following:

- More stable training across noisy or bursty input regimes;
- Improved convergence in long-horizon modeling tasks;
- Better generalization via reduced sensitivity to spurious input perturbations.

We refer to this enhanced architecture as *Momentum Mamba*.

### B. Mitigating Vanishing Gradients in Momentum Mamba

To address the inherent challenge of gradient vanishing in Mamba, we enrich the dynamics with an auxiliary momentum state $v_n$, yielding the augmented state $s_n = [h_n^\top, v_n^\top]^\top$. The momentum update

$$v_n = \beta v_{n-1} + \alpha \overline{B}_n x_n$$

acts as an exponential moving average of input signals, smoothing high-frequency variations. Embedding this update into the augmented recurrence produces the affine formulation of Equation equation 7, which fundamentally alters the Jacobian structure in backpropagation.

**Proposition 3** (Gradient Propagation in Momentum Mamba). *Let $s_n = [h_n^\top, v_n^\top]^\top$ evolve according to the affine recurrence*

$$s_n = M'_n s_{n-1} + F'_n, \qquad M'_n = \begin{bmatrix} \overline{A}_n & \beta I \\ 0 & \beta I \end{bmatrix}.$$

*Then the gradient of the loss $L$ with respect to $s_t$ is*

$$\frac{\partial L}{\partial s_t} = \frac{\partial L}{\partial s_T} \prod_{n=t+1}^{T} M'_n,$$

*which expands into the block form*

$$\frac{\partial L}{\partial s_t} = \frac{\partial L}{\partial s_T} \begin{bmatrix} \prod_{n=t+1}^{T} \overline{A}_n & \sum_{k=t+1}^{T} \Big( \prod_{n=t}^{T-k} \overline{A}_n \Big)(\beta I)^k \\ 0 & (\beta I)^{T-t+1} \end{bmatrix}.$$

*Proof.* The gradient propagation follows the chain rule:

$$\frac{\partial L}{\partial s_t} = \frac{\partial L}{\partial s_T} \cdot \frac{\partial s_T}{\partial s_t} = \frac{\partial L}{\partial s_T} \cdot \prod_{n=t+1}^{T} \frac{\partial s_n}{\partial s_{n-1}}.$$

Since $\partial s_n / \partial s_{n-1} = M'_n$, we obtain

$$\frac{\partial L}{\partial s_t} = \frac{\partial L}{\partial s_T} \prod_{n=t+1}^{T} M'_n.$$

Expanding the product of block matrices yields the closed form. The lower-right block corresponds to repeated multiplication by $(\beta I)$, while the upper-right block captures convolution-like interactions between $\overline{A}_n$ and $(\beta I)^k$. $\qquad\square$

**Remark 2** (Gradient Preservation by Momentum). *If $\beta \approx 1$, then $(\beta I)^{T-t+1} \approx I$ even for large horizons. Although the term $\prod_{n=t+1}^{T} \overline{A}_n$ may decay exponentially, the momentum pathway introduces eigenvalues close to unity. This spectral shift prevents the Jacobian product from collapsing to zero, ensuring that a non-negligible component of the gradient is preserved. As a result, Momentum Mamba maintains the gradient flow through the auxiliary state $v_t$, mitigating the vanishing gradient problem and improving the capacity to learn long-term dependencies.*

This result highlights that momentum in state-space models plays the same stabilizing role as in optimization: it accumulates and smooths gradient contributions while preventing their premature decay. Similarly to MomentumRNN [62], the velocity state $v_n$ serves as a memory buffer that allows early inputs to exert a lasting influence on the final prediction.

**Remark 3** (Exploding Gradients). *While momentum alleviates vanishing gradients, overly large $\beta$ may lead to explosive gradients. In practice, gradient clipping [63] or normalization techniques can be applied to maintain numerical stability during training.*

### C. Beyond Heavy-Ball Momentum in Mamba: Complex and Adam Extensions

*1) Complex Momentum Dynamics in Mamba:* To enrich the expressive power of state evolution, we extend Mamba with a second-order recurrence governed by complex-valued momentum:

$$v_n = \beta v_{n-1} + \alpha \overline{B}_n x_n, \tag{10}$$

$$h_n = \overline{A}_n h_{n-1} + v_n, \tag{11}$$

$$y_n = \Re(C_n h_n), \tag{12}$$

where $\beta = \rho e^{i\theta} \in \mathbb{C}$ is the complex momentum coefficient whose magnitude $\rho$ controls damping and phase $\theta$ induces oscillations, $\alpha > 0$ is the step size, $\overline{A}_n = \exp(\Delta_n A)$ the input-dependent transition, and $\overline{B}_n, C_n$ are learned projections. The real part $\Re(\cdot)$ ensures that $h_n, y_n \in \mathbb{R}^N$.

This recurrence introduces oscillatory memory traces: each past input contributes with exponential decay $|\beta|^{n-k}$ and phase rotation $e^{i(n-k)\arg(\beta)}$. Unlike classical exponential smoothing, which only damps signals, complex momentum allows constructive or destructive interference in the complex plane, enabling the model to emphasize phase-aligned components.

**Remark 4** (Frequency-Aware Filtering). *When $|\beta| \approx 1$, low damping preserves oscillations whose frequency matches $\arg(\beta)$, while out-of-phase components cancel. Thus, Complex Momentum Mamba effectively learns frequency-selective filters, useful for oscillatory or quasi-periodic signals such as inertial sensor data or wave-like physical dynamics.*

*2) Adam Momentum for Mamba:* Complex momentum equips Mamba with spectral sensitivity; in contrast, Adam-style momentum introduces variance-aware adaptivity. Classical momentum methods smooth updates uniformly, but ignore per-dimension variability. Adam [64] improves stability

by tracking both first- and second-order moments, yielding adaptive coordinate-wise learning rates.

Inspired by MomentumRNN [62], we embed this idea into Mamba:

$$v_n = \beta v_{n-1} + \alpha \overline{B}_n x_n, \tag{13}$$

$$m_n = \gamma m_{n-1} + (1-\gamma)\left(\overline{B}_n x_n\right)^2, \tag{14}$$

$$h_n = \overline{A}_n h_{n-1} + \frac{v_n}{\sqrt{m_n}+\epsilon}, \tag{15}$$

$$y_n = C_n h_n, \tag{16}$$

where $v_n$ acts as a first-order momentum state with decay $\beta$ and step size $\alpha$, while $m_n$ tracks per-coordinate variance with decay $\gamma$. The hidden state update $v_n/(\sqrt{m_n}+\epsilon)$ thus adapts the effective step size based on local variance, stabilizing learning in noisy directions and accelerating it in sparse ones.

**Proposition 4** (Adaptive Recurrence Stability). *Suppose inputs satisfy $\|\overline{B}_n x_n\| \leq B$ for all $n$. Then the normalized update $\frac{v_n}{\sqrt{m_n}+\epsilon}$ is uniformly bounded by $\frac{\alpha B}{\epsilon}$, ensuring bounded-input bounded-output stability of the hidden state $h_n$.*

*Proof.* By definition $m_n \geq 0$ component-wise. Hence $\sqrt{m_n} + \epsilon \geq \epsilon$, giving

$$\left\|\frac{v_n}{\sqrt{m_n}+\epsilon}\right\| \leq \frac{\|v_n\|}{\epsilon}.$$

Since $v_n$ is an exponentially weighted sum of $\overline{B}_k x_k$ scaled by $\alpha$, we have $\|v_n\| \leq \frac{\alpha B}{1-\beta}$, and thus the normalized term is bounded by $\alpha B/\epsilon$. $\qquad\square$

**Remark 5** (Variance-Aware Updates). *The normalization in Eq. equation 15 automatically reduces step size in high-variance directions while allowing faster adaptation in low-variance ones. Compared to standard Mamba, Adam Momentum Mamba improves stability under noisy or sparse inputs without altering the recurrence structure.*

**Discussion.** Together, Complex Momentum Mamba and Adam Momentum Mamba enrich the dynamics of state-space models in complementary ways: the former introduces phase- and frequency-selective memory through complex-valued recurrences, while the latter provides variance-aware adaptivity via per-coordinate normalization. By combining spectral sensitivity with statistical stability, these extensions enhance both the representational power and the optimization robustness of Mamba for long-sequence learning. In the context of HAR, this dual design enables the model to better capture oscillatory sensor patterns while maintaining stable and adaptive training under noisy, heterogeneous input streams.

## V. MOMENTUM MAMBA ARCHITECTURE FOR INERTIAL HUMAN ACTIVITY RECOGNITION

The overall architecture of the proposed *Momentum Mamba* framework for inertial HAR is illustrated in Fig. 1. The framework is designed to efficiently process synchronized six-axis inertial measurement unit (IMU) signals and to predict per-window activity labels through three main components: (*i*) a convolutional front-end, (*ii*) a momentum-augmented Mamba backbone, and (*iii*) a lightweight classification head.

## A. Convolutional Front-End

The first stage of the MMA employs a lightweight one-dimensional convolutional encoder that transforms raw six-axis inertial measurements (accelerometer and gyroscope) into high-dimensional temporal feature maps. A 1D convolutional layer with kernel size 3 captures local temporal dependencies within short motion windows, followed by batch normalization, ReLU activation, and dropout for regularization. This design enhances local feature extraction and mitigates sensor noise, effectively projecting raw low-dimensional signals into a richer feature space for subsequent state-space processing.

The front-end is intentionally lightweight (about 1.8K parameters for $6 \rightarrow 256$ dimensions), ensuring computational efficiency suitable for edge deployment. Its core purposes are threefold: (*i*) capture local temporal patterns characteristic of human motion dynamics such as acceleration-deceleration transitions, (*ii*) expand the representational capacity of sensor streams before entering the Momentum Mamba backbone, and (*iii*) provide initial noise filtering through learned convolutions. Consequently, this module reduces the effective sequence length while enriching feature diversity, serving as an efficient bridge between raw IMU inputs and the momentum-augmented state-space layers that model long-range temporal dependencies in human activity recognition.

## B. Momentum Mamba Backbone

The proposed framework extends the standard Mamba (S6) backbone by integrating a momentum mechanism inspired by second-order optimization dynamics. This modification, detailed in Section IV of the Momentum Mamba formulation, introduces an auxiliary velocity state alongside the hidden state, forming a dual-state recurrence that enhances long-range temporal modeling and stabilizes gradient propagation.

The momentum-augmented update incorporates a smoothed velocity term that accumulates temporal information over time. This design reshapes the Jacobian spectrum, mitigating exponential gradient decay and improving robustness against high-frequency sensor noise–all while preserving the linear-time complexity of SSMs. As a result, MMA achieves smoother hidden-state evolution, stronger resilience to transient perturbations, and more stable convergence across long sequences.

The momentum mechanism provides a physically grounded inductive bias well suited for inertial signals:

- **Inertial consistency:** The auxiliary velocity state captures the gradual, inertia-driven nature of human motion.
- **Noise robustness:** Temporal smoothing through momentum decay filters sensor noise while preserving activity-level dynamics.
- **Smooth transitions:** The model naturally represents gradual activity changes (e.g., walking $\rightarrow$ standing).
- **Extended dependencies:** The accumulated momentum state enables effective modeling of long-duration activities.

Despite the additional momentum pathway, MMA retains linear-time complexity in sequence length $L$. Parallel scan (prefix-sum) implementations allow efficient computation of both $\mathbf{v}_t$ and $\mathbf{h}_t$, preserving GPU-friendly execution. Compared to vanilla Mamba, the overhead in training time and VRAM usage is marginal ($< 5\%$ increase).

## C. Classification Head

The final stage of the architecture aggregates the temporal representations obtained from the Momentum Mamba backbone into compact embeddings for activity classification. As formulated in Section IV, the hidden representations $\mathbf{h}_t$ produced at each time step are averaged along the temporal dimension to form a global feature descriptor, which is subsequently projected through a fully connected linear layer to produce class logits for $C$ activity categories.

This simple yet effective design introduces minimal additional parameters (less than 0.5% of the total), ensuring a favorable trade-off between accuracy and efficiency. The use of global average pooling provides a holistic summary of temporal dynamics, while the linear projection enables direct mapping from latent representations to activity labels.

## D. Summary

In summary, the proposed *Momentum Mamba* framework establishes a principled and efficient approach to inertial human activity recognition by combining lightweight convolutional encoding, momentum-augmented state-space modeling, and compact classification. The integration of a momentum mechanism within the Mamba backbone provides several key advantages for HAR: it captures both transient and sustained motion patterns, preserves temporal continuity inherent in human movements, and stabilizes gradient propagation over long sequences. These properties enable the model to robustly distinguish activities that share overlapping short-term dynamics (e.g., walking vs. jogging) and to recognize complex behaviors that unfold over extended durations (e.g., cooking, exercising).

Furthermore, by maintaining linear-time complexity and a compact parameter footprint, MMA is well suited for deployment on wearable and mobile devices where computation and memory are constrained. The convolutional front-end effectively filters sensor noise and emphasizes local temporal structure, while the dual-state momentum backbone enriches temporal context through smooth, inertia-aware state transitions. Together, these components yield a model that is both physiologically interpretable and computationally efficient, advancing the robustness, accuracy, and real-time applicability of IMU-based human activity recognition systems.

## VI. EXPERIMENTS

### A. Dataset

In this work, the proposed method is evaluated on three publicly available benchmark datasets: MuWiGes [65], UESTC-MMEA-CL [66], and MMAct [67]. These datasets are designed for multimodal action recognition and provide synchronized RGB video along with inertial sensor data, including accelerometer and gyroscope signals. Fig. 2 illustrates sample data from each dataset, showcasing the integration of visual and motion modalities.
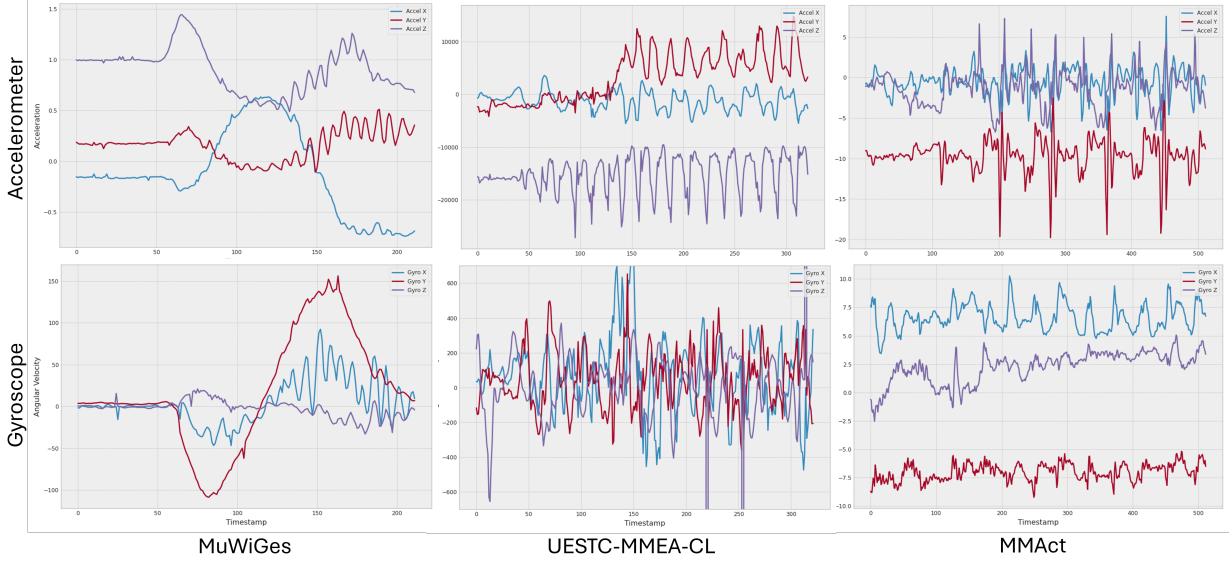
Fig. 2: Sensor signal examples from three benchmark datasets. The second and third rows display accelerometer and gyroscope signals, respectively, from the MuWiGes [65], UESTC-MMEA-CL [66], and MMAct [67] datasets. These signals reflect the temporal variations of multi-axis motion data captured during human activities.

- **MuWiGes dataset** is a carefully curated multimodal dataset acquired through a custom-designed wrist-worn device that integrates a high-resolution RGB camera and embedded IMU sensors. The camera records video streams at 1280×720 pixels with a frame rate of 30 FPS, while the inertial sensors comprising a tri-axial accelerometer and gyroscope sample motion data at 50 Hz. The dataset features gesture recordings from 50 participants (33 male and 17 female) spanning a broad age range (10–65 years), each performing 12 distinct hand gestures. Data collection was conducted across varied real-world environments, including homes and office spaces, to capture natural variations in background and user behavior. Each subject repeated each gesture spontaneously between 2 and 12 times, resulting in naturalistic gesture performance diversity. All gesture sequences were automatically segmented based on precise temporal boundaries indicating the onset and offset of each gesture. The final dataset comprises 5,408 synchronized samples that combine visual and inertial information.

- **UESTC-MMEA-CL dataset** dataset serves as a comprehensive multimodal benchmark tailored for evaluating continuous egocentric activity recognition. It comprises recordings of 32 distinct daily activities ranging from climbing stairs and drinking to shopping and playing cards performed by 10 participants across various environments, including both indoor and natural settings. Data were collected using smart glasses equipped with a first-person view camera and embedded IMU sensors. The camera records RGB video at a resolution of 640×480 pixels with a frame rate of 25 FPS, while the IMU captures motion data at 25 Hz. Each activity category includes approximately 200 synchronized samples, combining egocentric video and inertial measurements from accelerometer and gyroscope sensors.

- **MMAct dataset** is a multimodal dataset for activity recognition with 20 participants and 37 action classes. It includes seven synchronized modalities: RGB video, 2D keypoints, inertial (accelerometer, gyroscope), orientation, Wi-Fi, and barometric pressure. Data were collected using ceiling cameras (1920×1080, 30 fps), Google Glass (1280×720), and smartphones with inertial/context sensors (accelerometer 100 Hz, gyroscope 50 Hz), supplemented by smartwatches. Covering four environments, MMAct contains 35,084 labeled segments of diverse activities (e.g., typing, phone use, waving), making it a valuable benchmark for context-aware multimodal HAR.

TABLE I: Summary of three datasets used in experiments

| Datasets | MuWiGes [65] | UESTC-MMEA-CL [66] | MMAct [67] |
|---|---|---|---|
| Activity type | Hand gesture | Daily activity | Daily activity |
| Camera mounting | Wrist | Head | Ambient/Head |
| IMU mounting | Wrist | Head | Right wrist/thigh |
| Scenario | Indoor (Home, Office) | Natural | Indoor (Office) |
| Data modalities | RGB+Acc+Gyro | RGB+Acc+Gyro | RGB+Ori+Acc+Gyro |
| Total subjects | 50 | 10 | 20 |
| Number of classes | 12 | 32 | 37 |
| Total samples | 5048 | 6522 | 35084 |
| Train/test splitting | 3276/1772 | 4553/1316 | 28232/6930 |

Table I provides a summary of the three datasets tested in this study. These datasets differ in terms of the activities performed by the subjects, the placement of the cameras or IMU sensors, and the conditions under which the data was collected. In the MMAct dataset, cameras are positioned around the subjects to capture their actions from a front-facing view. Additionally, a camera mounted on Google Glass records egocentric videos, allowing a full view of the subject and their surroundings, which includes both hands and the area in front of them. In the second dataset, the camera is attached to the head, providing a view that captures both hands and the area in front, but lacks an additional observation perspective. In contrast, the final dataset, MuWiGes, has the

camera mounted on the wrist, resulting in a more limited field of view. This restricted perspective makes gesture recognition from this angle significantly more challenging.

### B. Experimental Settings

**Data preprocessing.** Raw inertial signals (tri-axial accelerometer and gyroscope) are resampled and segmented into fixed-length windows of $L = 512$ with a 50% overlap. Each channel is standardized to zero mean and unit variance using training statistics. Segments are then formatted as $X \in \mathbb{R}^{B \times L \times 6}$, where $B$ denotes the batch size.

**Model configuration.** The MMA backbone consists of a lightweight Conv1D front-end, $N = 2$ stacked Momentum Mamba layers with hidden size $d_{model} = 128$, and a compact linear classification head. Dropout with probability 0.1 is applied throughout.

**Training setup.** All models are trained end-to-end using categorical cross-entropy loss. Optimization is performed with Adam at an initial learning rate of $1 \times 10^{-3}$, weight decay of $1 \times 10^{-4}$, and cosine annealing learning rate scheduling. Mini-batches of size 16 are used for up to 50 epochs, with early stopping (patience = 10). Gradient clipping (global norm 1.0) is applied to stabilize training.

**Hardware.** Experiments are implemented in PyTorch and executed on a single NVIDIA A30 GPU with 24 GB memory.

### C. Experiment Results

*1) Overview:* We evaluated Momentum Mamba and its variants on three public benchmarks: UESTC-MMEA-CL [66], MMAct [67], and MuWiGes [65]. Comparisons are made against both transformer-based models (Transformer, GAFormer [25], MAMC [26]) and state-space approaches (Vanilla Mamba [27], LinOSS [33]). Performance is measured in terms of accuracy, precision, recall, and F1-score, as summarized in Tables II–IV.

Across all datasets, Momentum Mamba consistently outperforms strong baselines, with average accuracy gains of +2.77% over Vanilla Mamba and +1.64%-19.11% over Transformer. These improvements validate the effectiveness of embedding momentum-enhanced recurrence into the SSM backbone, yielding more stable gradient flow and stronger robustness to noisy inertial signals while preserving linear-time efficiency.

Notably, Complex Momentum Mamba achieves the highest scores across all benchmarks, indicating that frequency-sensitive dynamics provide additional memory capacity beneficial for fine-grained motion patterns. This demonstrates that our momentum-driven framework is both effective in its base form and extensible to richer recurrence mechanisms.

*2) Dataset-specific performance:*

*a) MuWiGes:* Momentum Mamba attains 98.43% accuracy, exceeding Vanilla Mamba by +1.13% and slightly surpassing GAFormer. Gesture recognition from wrist-mounted inertial sensors is inherently difficult due to subtle motion cues, subject variability, and contamination from incidental hand jitter. Momentum Mamba addresses these challenges

TABLE II: Experimental results on the MuWiGes dataset. Bold values represent the best results.

| Method | Accuracy (↑) | Precision (↑) | Recall (↑) | F1-score (↑) |
|---|---|---|---|---|
| Nguyen et al. [65] | 95.60 | - | - | - |
| GAFormer (2023) [25] | 98.33 | 98.25 | 98.33 | 98.30 |
| Vanilla Transformer | 96.79 | 97.83 | 97.75 | 97.79 |
| MAMC (2024) [26] | 96.20 | 96.56 | 95.90 | 96.21 |
| Vanilla Mamba | 97.30 | 97.46 | 97.15 | 97.30 |
| LinOSS (2025) [33] | 92.86 | 92.74 | 92.98 | 92.85 |
| Momentum Mamba (Ours) | 98.43 | 98.56 | 98.26 | 98.41 |
| **Complex Momentum Mamba (Ours)** | **98.64** | **98.67** | **98.52** | **98.59** |

TABLE III: Experimental results on the UESTC-MMEA-CL dataset. Bold values represent the best results

| Method | Accuracy (↑) | Precision (↑) | Recall (↑) | F1-score (↑) |
|---|---|---|---|---|
| Xu et al. [66] | 59.70 | - | - | - |
| GAFormer (2023) [25] | 79.00 | 78.41 | 79.94 | 79.15 |
| MAMC (2024) [26] | 68.58 | 68.20 | 69.40 | 68.81 |
| Vanilla Transformer | 73.21 | 73.50 | 72.90 | 73.17 |
| LinOSS (2025) [33] | 77.73 | 76.20 | 78.12 | 77.14 |
| Vanilla Mamba | 88.76 | 88.30 | 89.20 | 88.73 |
| Mamba-LinOSS | 87.71 | 87.80 | 87.40 | 87.57 |
| Momentum Mamba (Ours) | 92.32 | 92.10 | 92.60 | 92.39 |
| **ComplexMomentumMamba (Ours)** | **94.07** | **94.20** | **93.95** | **94.07** |

by smoothing transient fluctuations while retaining phase-consistent motion information, enabling reliable separation of gestures with highly similar dynamics (e.g., "rotate wrist" vs. "twist cap"). The Complex Momentum Mamba variant further improves accuracy to 98.64%, indicating that complex-valued recurrence enhances robustness against rapid oscillations and improves discrimination of fine-grained gesture classes.

*b) UESTC-MMEA-CL:* On this fine-grained daily activity dataset, Momentum Mamba achieves 92.32% accuracy, a +3.56% improvement over Vanilla Mamba. Recognition is difficult because many actions (e.g., "reading" or "typing") produce weak, low-frequency signals that are easily obscured by noise or incidental fluctuations. By leveraging second-order recurrence, Momentum Mamba aggregates these subtle components across long horizons while attenuating irrelevant variations, enabling the model to maintain stable and discriminative temporal patterns. In addition, the Complex Momentum Mamba variant raises performance to 94.07% accuracy and F1-score, suggesting that frequency-sensitive dynamics further sharpen the model's ability to separate overlapping or weak activity signals–an advantage particularly evident for subtle or repetitive motion classes.

*c) MMAct:* The MMAct dataset, with 37 action classes under diverse conditions, remains one of the most challenging HAR benchmarks using inertial data. Signal variations arise not only from activity type but also from device placement, environment, and synchronization quality. Momentum Mamba reaches 75.49% accuracy, improving upon Vanilla Mamba by +3.63% and Transformer by +15.58%. Its advantage lies in stabilizing hidden state evolution, preventing overreaction

TABLE IV: Experimental results on the MMAct dataset. Bold values represent the best results.

| Method | Accuracy (↑) | Precision (↑) | Recall (↑) | F1-score (↑) |
|---|---|---|---|---|
| Multi-teacher (2019) [67] | 62.67 | – | – | – |
| VLMs (2023) [68] | 64.47 | – | – | – |
| VSKD (2022) [69] | 60.14 | – | – | – |
| GAFormer (2023) [25] | 60.76 | 61.00 | 60.52 | 60.75 |
| Vanilla Transformer | 59.91 | 59.46 | 60.40 | 59.91 |
| Vanilla Mamba | 71.86 | 72.24 | 71.56 | 71.83 |
| LinOSS (2025) [33] | 72.16 | 71.28 | 72.53 | 71.90 |
| Momentum Mamba (Ours) | 75.49 | 75.86 | 75.10 | 75.45 |
| **Complex Momentum Mamba (Ours)** | **76.62** | **76.80** | **76.45** | **76.63** |

to short-lived or conflicting signals across sensor axes, and thereby capturing persistent motion trends. For this dataset, the Complex Momentum Mamba variant pushes performance to 76.62% accuracy and 76.63 F1-score, suggesting that complex-valued dynamics are especially beneficial for recognizing composite or multi-step activities involving subtle temporal transitions.

Overall, the results across MuWiGes, UESTC-MMEA-CL, and MMAct confirm that Momentum Mamba provides consistent and substantial improvements over state-of-the-art baselines by stabilizing long-range temporal modeling of inertial signals. Moreover, the Complex Momentum Mamba variant further extends these benefits with frequency-sensitive dynamics, yielding additional gains in challenging scenarios that involve subtle, noisy, or multi-step activities.

TABLE V: Comparison of Mamba, MomentumMamba, and ComplexMomentumMamba ($d_{model} = 128, n_{layers} = 2, d_{state} = 64, d_{conv} = 4, expand = 2$).

| Model | Time (s) | FLOPs (MFLOPS) | Params (K) | VRAM (MB) |
|---|---|---|---|---|
| Mamba | 0.0052 | 311.285 | 313.632 | 188.41 |
| MomentumMamba | 0.007 | 278.091 | 313.760 | 212.41 |
| CMMamba | 0.031 | 311.646 | 412.065 | 396.91 |

*3) Efficiency analysis:* In addition to accuracy, practical deployment of HAR models depends heavily on computational efficiency and memory footprint. To this end, we compare Mamba, Momentum Mamba, and Complex Momentum Mamba under identical configurations ($d_{model} = 128, n_{layers} = 2, d_{state} = 64, d_{conv} = 4, expand = 2$), as summarized in Table V.

Vanilla Mamba offers the lowest latency, with an inference time of 0.0052 seconds and a VRAM consumption of 188.41 MB. Momentum Mamba introduces only a marginal overhead, increasing runtime to 0.007 seconds and memory to 212.41 MB. Interestingly, its FLOPs (278.091 MFLOPs) are slightly lower than that of Vanilla Mamba (311.285 MFLOPs), reflecting that the integration of momentum does not impose significant computational burden. This result confirms that the proposed recurrence can be incorporated with negligible cost while delivering substantial improvements in recognition accuracy.

On the other hand, Complex Momentum Mamba achieves the strongest performance but at a notable efficiency trade-off. Its runtime (0.031 seconds) and VRAM requirement (396.91 MB) are considerably larger, and the parameter count rises to 412K compared to 314K for both Mamba and Momentum Mamba. These statistics highlight that while the complex variant is highly effective for accuracy-sensitive applications, it demands significantly greater resources.

Overall, this comparison demonstrates that Momentum Mamba provides the best balance between efficiency and performance, making it suitable for real-world HAR scenarios where computational resources are limited. In contrast, Complex Momentum Mamba represents an extensible, high-accuracy alternative tailored to scenarios where precision is prioritized over efficiency.

*4) Additional insights:*

*a) Temporal attention pattern analysis:* To gain deeper insight into the internal mechanisms of different sequence modeling architectures, we visualize and compare the temporal saliency maps produced by three models: Transformer, Vanilla Mamba, and our proposed Momentum Mamba. Figure 3, each model processes the same six-dimensional inertial input, with the raw signal overlaid in white and the background heatmap indicating time-varying saliency.

The Transformer model tends to focus its attention on isolated, high-salience peaks typically centered around abrupt transitions or local extrema in the signal. These activations are sparse and vary significantly across input channels, indicating that the model attends selectively but inconsistently across features. Furthermore, early portions of the sequence often receive negligible attention, suggesting limited capacity for long-term memory retention.

Vanilla Mamba, in contrast, demonstrates smoother and more distributed attention compared to the Transformer. It allocates saliency across longer temporal spans, with more stable patterns in the mid- and late-sequence regions. However, its focus is still less consistent than that of Momentum Mamba, particularly in terms of early sequence sensitivity and multi-channel alignment. While Mamba benefits from input-conditioned state transitions and structured recurrence, it lacks an explicit mechanism to preserve information from past states beyond first-order dynamics.

Momentum Mamba shows a clear improvement in both breadth and continuity of temporal attention. The saliency maps reveal sustained activation across wide time windows, including the initial segments of the sequence. This indicates that Momentum Mamba is capable of retaining early-stage information and using it in downstream prediction–a critical feature for HAR tasks where initial motion cues can be subtle yet discriminative. Moreover, the saliency patterns are temporally smooth and spatially aligned across channels, implying a coordinated and stable latent representation.

We attribute this improvement to the incorporation of second-order momentum into the state update mechanism. By allowing the hidden state to evolve not only based on the current input but also on the velocity of past updates, Momentum Mamba effectively propagates information over time with greater persistence. This results in richer temporal representations and a more robust ability to model long-range dependencies in multichannel sensor data.

Overall, these visualizations provide qualitative evidence that Momentum Mamba captures temporal structure more effectively than both Transformer and Vanilla Mamba, supporting its superior performance in downstream HAR tasks.

*b) Gradient flow analysis:* To better understand the optimization behavior of different recurrent state-space architectures, we examine the $\ell_2$ norm of the gradients of the loss $\mathcal{L}$ with respect to the hidden state $h_t$ across both time steps and training iterations (Fig. 4). For Vanilla Mamba (left), gradient norms diminish rapidly as the time horizon increases. This pattern is symptomatic of the vanishing gradient problem, where information from early steps in long sequences fails to propagate back effectively during training. As a result, the model has limited ability to assign credit to distant depen-
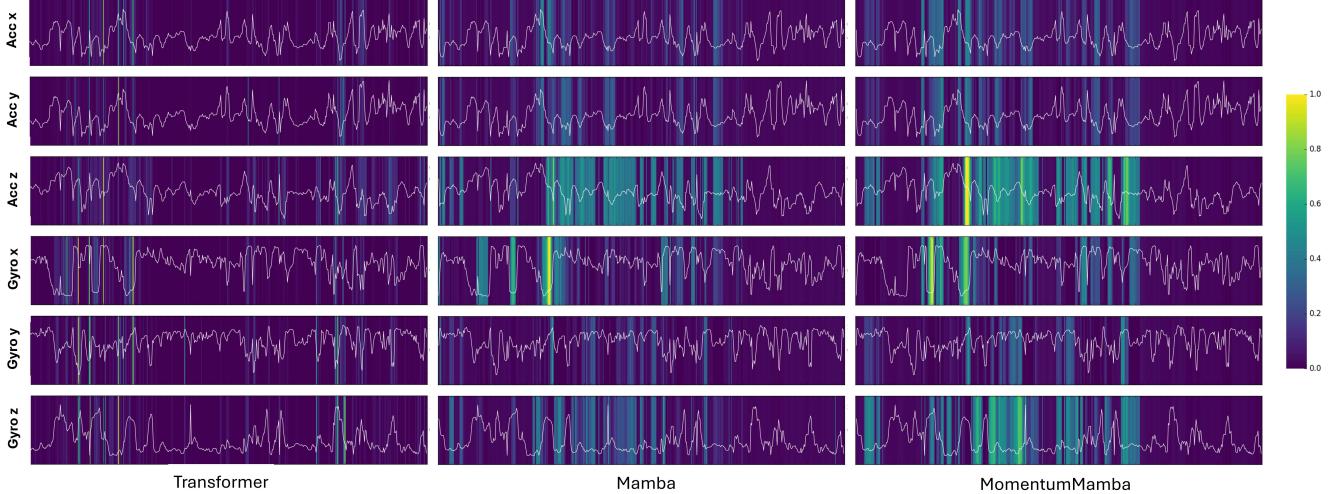
Fig. 3: Two-Step Rescaling Explanations [70]: Comparative Saliency Analysis of Transformer, Mamba, and MomentumMamba Models on UESTC Dataset, where warmer colors (yellow) indicate higher feature importance and cooler colors (purple) indicate lower importance.
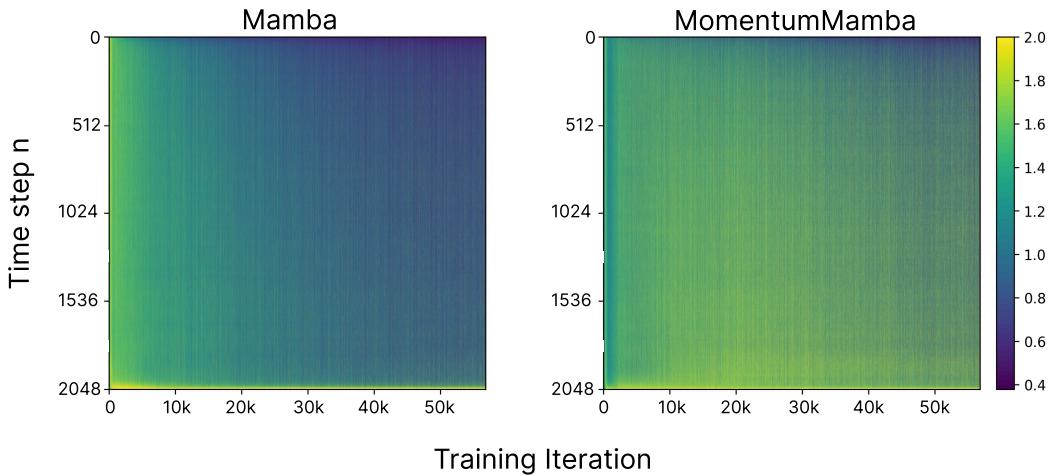


Fig. 4: $\ell_2$ norm of the gradients of the loss $\mathcal{L}$ w.r.t. the state vector $h_t$ at each time step $t$ for Vanilla Mamba (left) and MomentumMamba (right). MomentumMamba does not suffer from vanishing gradients.

dencies, which ultimately constrains its capacity to capture extended temporal structure.

In contrast, Momentum Mamba (right) exhibits substantially more stable gradient magnitudes across both shallow and deep time steps, with relatively uniform patterns persisting throughout the training process. This observation indicates that the introduction of momentum into the state update mechanism directly improves gradient flow by creating an auxiliary pathway for information propagation. The second-order recurrence not only mitigates exponential decay in gradient norms but also allows informative updates from earlier time steps to persist over long horizons.

This behavior aligns closely with findings in Momentum-RNN [62], where momentum-based recurrence was shown to counteract gradient shrinkage and maintain long-term credit assignment. By analogy, Momentum Mamba inherits these benefits within the structured state-space modeling framework, thereby providing a principled solution to one of the long-standing challenges in training deep sequential models. Crucially, the preservation of non-vanishing gradients ensures that the model can stably exploit both short-term variations and long-term trends in inertial sequences, directly contributing to the empirical performance improvements reported across MuWiGes, UESTC-MMEA-CL, and MMAct.

*c) Hyperparameter sensitivity analysis:* To evaluate the robustness of Momentum Mamba in different momentum settings, we performed a comprehensive grid search over two key hyperparameters: the momentum coefficient $\beta \in \{0.0, 0.1, 0.3, 0.6, 0.9, 0.99, 0.999\}$ and the input scaling factor $\alpha \in \{0.0, 0.1, 0.3, 0.6, 0.9, 1.0, 2.0\}$. Figure 5 shows the test accuracy achieved on three benchmark HAR datasets: MuWiGes, UESTC-MMEA-CL and MMAct.
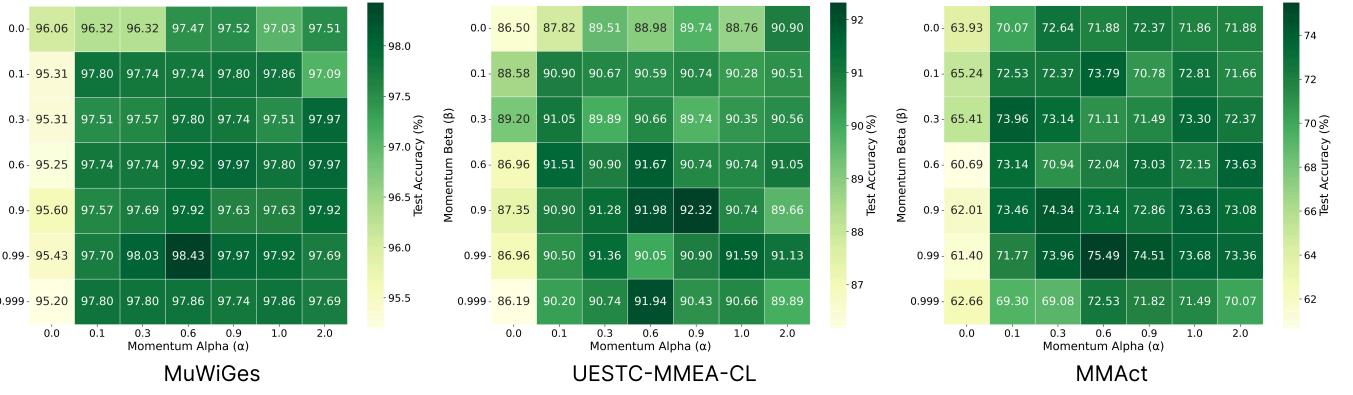
Fig. 5: Hyperparameter Grid Search Results for Momentum Mamba Model. Heatmap showing test accuracy (%) across different combinations of momentum beta ($\beta$) and momentum alpha ($\alpha$) parameters. Each cell displays the test accuracy percentage achieved with the corresponding hyperparameter combination. The color scale ranges from lower accuracy (red) to higher accuracy (green), from the MuWiGes [65], UESTC-MMEA-CL [66], and MMAct [67] datasets.

In the MuWiGes dataset, Momentum Mamba consistently achieved high accuracy across a wide range of hyperparameter settings. The highest accuracy of 98.43% was attained at ($\beta = 0.99, \alpha = 0.6$), with many other combinations achieving above 97.5%. This indicates the model's robustness to hyperparameter variations and its ability to maintain strong performance even with minimal tuning.

For UESTC-MMEA-CL, which involves more subtle and fine-grained activities, the model achieved its peak performance of 92.32% at ($\beta = 0.9, \alpha = 0.9$). While performance degraded when $\beta$ was too small (e.g., $\beta = 0.0$), moderate momentum coefficients (e.g., $\beta \in [0.6, 0.999]$) provided stable results above 90%, indicating the importance of second-order memory in modeling fine-grained temporal variations.

On the more challenging MMAct dataset, the accuracy trend revealed a stronger dependence on tuning. Performance increased steadily with $\alpha$, particularly when paired with mid-to-high momentum ($\beta \in [0.6, 0.99]$). The highest test accuracy of 75.49% was achieved at ($\beta = 0.99, \alpha = 0.6$). In particular, when $\beta = 0.0$, the accuracy plateaued below 71%, reinforcing the benefit of incorporating momentum into the state update dynamics.

In general, these results demonstrate the followings.

- Momentum Mamba benefits significantly from momentum-enhanced recurrence, particularly with $\beta$ in the range $[0.6, 0.999]$;
- The input scaling factor $\alpha$ complements $\beta$ by controlling the contribution of new input information to the momentum state;
- Properly chosen hyperparameters yield consistent improvements across datasets of varying complexity and motion dynamics.

These findings validate the utility of second-order dynamics in stabilizing the gradient flow and improving the model's ability to capture long-term dependencies in HAR tasks.

## VII. CONCLUSION

HAR with inertial sensors remains challenging due to noisy signals, limited spatial context, and strict resource constraints. To address these issues, we proposed Momentum Mamba, a momentum-augmented selective state-space model that enriches hidden-state evolution with a second-order dynamics pathway. This design mitigates vanishing gradients, enhances robustness against noise, and improves long-range temporal modeling.

We further introduced Complex Momentum Mamba for frequency-sensitive oscillatory memory and Adam Momentum Mamba for variance-aware adaptivity, extending the expressive capacity of state-space models. Comprehensive experiments on multiple HAR benchmarks confirmed that our models consistently outperform Transformer-based, RNN-based, and oscillatory SSM baselines, while preserving linear-time scalability and offering a favorable balance between accuracy and efficiency for real-time wearable deployment.

Looking forward, future work will focus on lightweight variants, adaptive trade-offs between accuracy and efficiency, and multimodal integration. Beyond HAR, momentum-augmented SSMs hold promise for a broad range of sequential domains such as biosignal analysis, robotics, speech, and multimodal temporal understanding, paving the way for scalable and interpretable sequence learning.

## REFERENCES

[1] Luis Sigcha, L. Borzì, Federica Amato, Irene Rechichi, Carlos Ramos-Romero, Andrés Cárdenas, Luis Gasco, and Gabriella Olmo. Deep learning and wearable sensors for the diagnosis and monitoring of parkinson's disease: A systematic review. *Expert Syst. Appl.*, 229:120541, 2023.

[2] Ebrahim Nemati, Shibo Zhang, Tousif Ahmed, Md. Mahbubur Rahman, Jilong Kuang, and Alex Gao. Coughbuddy: Multi-modal cough event detection using earbuds platform. *2021 IEEE 17th International Conference on Wearable and Implantable Body Sensor Networks (BSN)*, pages 1–4, 2021.

[3] Yan Gao, Yang Long, Yu Guan, Anna Purna Basu, Jessica Baggaley, and Thomas Plötz. Towards reliable, automated general movement assessment for perinatal stroke screening in infants using wearable accelerometers. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 3:1 – 22, 2019.

[4] Debarshi Bhattacharya, Deepak Kumar Sharma, Won Jae Kim, Muhammad Fazal Ijaz, and Pawan Kumar Singh. Ensem-har: An ensemble deep learning model for smartphone sensor-based human activity recognition for measurement of elderly health monitoring. *Biosensors*, 12, 2022.

[5] Cheng Zhang, Qiuyue Xue, Anandghan Waghmare, Ruichen Meng, Sumeet Jain, Yizeng Han, Xinyu Li, Kenneth A. Cunefare, Thomas Plötz, Thad Starner, Omer T. Inan, and Gregory D. Abowd. Fingerping: Recognizing fine-grained hand poses using active acoustic on-body sensing. *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, 2018.

[6] Angkoon Phinyomark and Erik J. Scheme. Emg pattern recognition in the era of big data and deep learning. *Big Data Cogn. Comput.*, 2:21, 2018.

[7] Devki Nandan Jha, Zhenghua Chen, Shudong Liu, Min Wu, Jiahan Zhang, Graham Morgan, Rajiv Ranjan, and Xiaoli Li. A hybrid accuracy- and energy-aware human activity recognition model in iot environment. *IEEE Transactions on Sustainable Computing*, 8:1–14, 2023.

[8] Kun Wang, Jun He, and L. Zhang. Attention-based convolutional neural network for weakly labeled human activities' recognition with wearable sensors. *IEEE Sensors Journal*, 19:7598–7604, 2019.

[9] Misha Karim, Shah Khalid, Aliya Aleryani, Jawad Khan, Irfan Ullah, and Z. Ali. Human action recognition systems: A review of the trends and state-of-the-art. *IEEE Access*, 12:36372–36390, 2024.

[10] Zehua Sun, Jun Liu, Qiuhong Ke, and H. Rahmani. Human action recognition from various data modalities: A review. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45:3200–3225, 2020.

[11] Cuong Pham, Linh Nguyen, Anh Nguyen, Ngon Nguyen, and Van-Toi Nguyen. Combining skeleton and accelerometer data for human fine-grained activity recognition and abnormal behaviour detection with deep temporal convolutional networks. *Multimedia Tools and Applications*, 80(19):28919–28940, 2021.

[12] Sampat Kumar Ghosh, Rashmi M, Biju R. Mohan, and Ram Mohana Reddy Guddeti. Deep learning-based multi-view 3d-human action recognition using skeleton and depth data. *Multimedia Tools and Applications*, 82:19829–19851, 2022.

[13] Yanni Yang, Pengfei Hu, Jiaxing Shen, Haiming Cheng, Zhenlin An, and Xiulong Liu. Privacy-preserving human activity sensing: A survey. *High-Confidence Computing*, 2024.

[14] W. S. Lima, E. Souto, K. El-Khatib, Roozbeh Jalali, and João Gama. Human activity recognition using inertial sensors in a smartphone: An overview. *Sensors (Basel, Switzerland)*, 19, 2019.

[15] Shibo Zhang, Yaxuan Li, Shen Zhang, Farzad Shahabi, Stephen Xia, Yuanbei Deng, and Nabil Alshurafa. Deep learning in human activity recognition with wearable sensors: A review on advances. *Sensors (Basel, Switzerland)*, 22, 2021.

[16] Andrey D. Ignatov. Real-time human activity recognition from accelerometer data using convolutional neural networks. *Appl. Soft Comput.*, 62:915–922, 2018.

[17] Wen Qi, Hang Su, and Andrea Aliverti. A smartphone-based adaptive recognition and real-time monitoring system for human activities. *IEEE Transactions on Human-Machine Systems*, 50:414–423, 2020.

[18] Md Atiqur Rahman Ahad, Anindya Das Antar, and Masud Ahmed. *Sensor-Based Human Activity Recognition: Challenges Ahead*, pages 175–189. Springer International Publishing, Cham, 2021.

[19] Kaixuan Chen, Dalin Zhang, Lina Yao, Bin Guo, Zhiwen Yu, and Yunhao Liu. Deep learning for sensor-based human activity recognition: Overview, challenges and opportunities, 2021.

[20] Ming Zeng, Le T Nguyen, Bo Yu, Ole J Mengshoel, Jiang Zhu, Pang Wu, and Joydeep Zhang. Convolutional neural networks for human activity recognition using mobile sensors. In *6th International Conference on Mobile Computing, Applications and Services*, pages 197–205. Springer, 2014.

[21] Fahmid Al Farid, Ahsanul Bari, Abu Saleh, Musa Miah, Sarina Mansor, Jia Uddin, and S. Prabha Kumaresan. A structured and methodological review on multi-view human activity recognition for ambient assisted living. *Journal of Imaging*, 11, 2025.

[22] Yoshua Bengio, Patrice Y. Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5 2:157–66, 1994.

[23] Francisco J Ordóñez and Daniel Roggen. Deep convolutional and lstm recurrent neural networks for multimodal wearable activity recognition. In *Proceedings of the ACM International Symposium on Wearable Computers*, 2016.

[24] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[25] Trung-Hieu Le, Thai-Khanh Nguyen, Trung-Kien Tran, Thanh-Hai Tran, and Cuong Pham. Gaformer: Wearable imu-based human activity recognition with gramian angular field and transformer. In *2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 297–303, 2023.

[26] Yezhuo Zhang, Zinan Zhou, Yichao Cao, Guangyu Li, and Xuanpeng Li. Mamc—optimal on accuracy and efficiency for automatic modulation classification with extended signal length. *IEEE Communications Letters*, 28(12):2864–2868, 2024.

[27] Tri Dao et al. Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2312.17143*, 2024.

[28] Shuangjian Li, Tao Zhu, Furong Duan, Liming Chen, Huansheng Ning, Christopher Nugent, and Yaping Wan. Harmamba: Efficient and lightweight wearable sensor human activity recognition based on bidirectional mamba. *IEEE Internet of Things Journal*, 2024.

[29] Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.

[30] Yurii Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. *Proceedings of the USSR Academy of Sciences*, 269:543–547, 1983.

[31] Ilya Sutskever, James Martens, George E. Dahl, and Geoffrey E. Hinton. On the importance of initialization and momentum in deep learning. In *International Conference on Machine Learning*, 2013.

[32] Tan M Nguyen, Richard Baraniuk, Andrea L Bertozzi, Stanley Osher, and Bao Wang. Momentumrnn: Integrating momentum into recurrent neural networks. In *Advances in Neural Information Processing Systems*, volume 33, pages 15014–15024, 2020.

[33] T. Konstantin Rusch and Daniela Rus. Oscillatory state-space models, 2025.

[34] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations*, 2022.

[35] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. *Advances in Neural Information Processing Systems*, 35:35971–35983, 2022.

[36] Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured state spaces. *Advances in neural information processing systems*, 35:22982–22994, 2022.

[37] Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. Simplified state space layers for sequence modeling. In *International Conference on Learning Representations*, 2023. Oral presentation.

[38] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *First Conference on Language Modeling*, 2024.

[39] Shufan Li, Harkanwar Singh, and Aditya Grover. Mamba-nd: Selective state space modeling for multi-dimensional data. In *European Conference on Computer Vision*, pages 75–92. Springer, 2024.

[40] Yoshua Bengio, Patrice Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, 1994.

[41] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, and Jorge Luis Reyes-Ortiz. Human activity recognition on smartphones using a multiclass hardware-friendly support vector machine. *International Workshop on Ambient Assisted Living*, pages 216–223, 2012.

[42] Rong Liu, Ting Chen, and Lu Huang. Research on human activity recognition based on active learning. In *2010 international conference on machine learning and cybernetics*, volume 1, pages 285–290. IEEE, 2010.

[43] Wan-Yu Deng, Qing-Hua Zheng, and Zhong-Min Wang. Cross-person activity recognition using reduced kernel extreme learning machine. *Neural Networks*, 53:1–7, 2014.

[44] Jingyuan Wang, Yiqiang Chen, Shiqiang Hao, Xiangmin Peng, and Longbing Hu. Deep learning for sensor-based activity recognition: A survey. *Pattern Recognition Letters*, 150:5–17, 2021.

[45] Jianbo Yang, Minh Nhut Nguyen, Phyo Phyo San, Xiao Li Li, and Shonali Krishnaswamy. Deep convolutional neural networks on multichannel time series for human activity recognition. In *Twenty-fourth International Joint Conference on Artificial Intelligence (IJCAI)*, 2015.

[46] Enas Abdulhay, Mohammed F Alhamid, Ahmed Ghoneim, Ghulam Muhammad, and M Shamim Hossain. Lightweight convolutional neural network for human activity recognition on wearable devices. *IEEE Access*, 9:111123–111134, 2021.

[47] Nils Y Hammerla, Shane Halloran, and Thomas Plötz. Deep, convolutional, and recurrent models for human activity recognition using

wearables. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1533–1540, 2016.

[48] Haoran Zhao, Xiangyu Zhang, Shuang Wang, Jing Shen, and Huiru Zheng. Tattn-har: Temporal self-attention network for human activity recognition using sensors. *Information Fusion*, 84:1–12, 2022.

[49] Pengfei Li, Haoran Wang, Xin Zhang, Jiantao Zhao, and Yan Li. Harformer: Transformer-based human activity recognition with hierarchical attention. *Pattern Recognition*, 138:109406, 2023.

[50] Qian Huang, Xin Zhang, Jiantao Zhao, and Yan Li. Lightformer: Lightweight transformer for human activity recognition on mobile and embedded devices. *IEEE Internet of Things Journal*, 9(20):20091–20102, 2022.

[51] Wu Lee, Yuliang Shi, Han Yu, Lin Cheng, Xinjun Wang, Zhongmin Yan, and Fanyu Kong. Hpformer: Low-parameter transformer with temporal dependency hierarchical propagation for health informatics. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.

[52] Theodoros Koutrintzes, Panagiotis Giannakeris, Anastasios Tefas, and Constantine Kotropoulos. Multimodal human action recognition via RGB and depth data fusion. *Multimedia Tools and Applications*, 82:36809–36831, 2023.

[53] Md Mahmudul Islam, Jahan Ali Nasir, and Md Mahbubur Rashid. Multimodal human activity recognition using visual and inertial sensors with CBAM-enhanced CNN and ConvLSTM. *Expert Systems with Applications*, 229:120551, 2023.

[54] Yifan Yang, Hao Zhang, and Yang Liu. Cross-modal federated human activity recognition with modality imbalance. *IEEE Transactions on Mobile Computing*, 2024.

[55] Albert Gu, Tri Dao, Stefano Ermon, and Cynthia Rudin. On the approximation power of implicit neural representations for high-frequency signals. *arXiv preprint arXiv:2006.13027*, 2020.

[56] Fei Luo, Anna Li, Bin Jiang, Salabat Khan, Kaishun Wu, and Lu Wang. Activitymamba: a cnn-mamba hybrid neural network for efficient human activity recognition. *IEEE Transactions on Mobile Computing*, 2025.

[57] Trung-Hieu Le, Thai Khanh Nguyen, Tuan-Anh Le, Mathieu Delalandre, Kien Tran Trung, Thanh-Hai Tran, and Cuong Pham. Mamba-mhar: An efficient multimodal framework for human action recognition. *Journal of Computer Science and Cybernetics*, pages 245–264, 2025.

[58] Xingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan. Adan: Adaptive nesterov momentum algorithm for faster optimizing deep models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9508–9520, 2024.

[59] Nghia H Nguyen, Tan M Nguyen, Huyen K Vo, Stanley J Osher, and Thieu N Vo. Improving neural ordinary differential equations with nesterov's accelerated gradient method. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

[60] Guangyuan Chen, Xin Huang, Zhiqiang Li, Xiaolong Zheng, Jinjun Wang, and Ji Zhang. Physics-informed state space models for dynamical system modeling and control. *arXiv preprint arXiv:2310.07981*, 2023.

[61] Ozan Çatal, Stijn Wauthier, Cedric De Boom, Tim Verbelen, and Pablo Lanillos. Learning generative state space models for active inference. *Frontiers in Computational Neuroscience*, 14:574372, 2020.

[62] Tan M Nguyen et al. Momentumrnn: Integrating momentum into recurrent neural networks. In *NeurIPS*, 2020.

[63] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318. PMLR, 2013.

[64] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. Presented at ICLR 2015.

[65] Hong-Quan Nguyen, Trung-Hieu Le, Trung-Kien Tran, Hoang-Nhat Tran, Thanh-Hai Tran, Thi-Lan Le, Hai Vu, Cuong Pham, Thanh Phuong Nguyen, and Huu Thanh Nguyen. Hand gesture recognition from wrist-worn camera for human–machine interaction. *IEEE Access*, 11:53262–53274, 2023.

[66] Linfeng Xu, Qingbo Wu, Lili Pan, Fanman Meng, Hongliang Li, Chiyuan He, Hanxin Wang, Shaoxu Cheng, and Yu Dai. Towards continual egocentric activity recognition: A multi-modal egocentric activity dataset for continual learning. *arXiv preprint arXiv:2301.10931*, 2023.

[67] Quan Kong, Ziming Wu, Ziwei Deng, Martin Klinkigt, Bin Tong, and Tomokazu Murakami. Mmact: A large-scale dataset for cross modal human action understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8658–8667, 2019.

[68] Riley Tavassoli, Mani Amani, and Reza Akhavian. Expanding frozen vision-language models without retraining: Towards improved robot perception. *arXiv preprint arXiv:2308.16493*, 2023.

[69] Jianyuan Ni, Raunak Sarbajna, Yang Liu, Anne HH Ngu, and Yan Yan. Cross-modal knowledge distillation for vision-to-sensor action recognition. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4448–4452. IEEE, 2022.

[70] Jacqueline Höllig, Cedric Kulbach, and Steffen Thoma. Tsinterpret: A python package for the interpretability of time series classification. *Journal of Open Source Software*, 8:5220, 05 2023.