# H-JEPA: Learning Hamiltonian Dynamics on Latent Video Embeddings

Jefrey Bergl

Department of Computer Science

University of North Carolina at Chapel Hill

Independent Research

jbergl@unc.edu

**Abstract**

World models learned from high-dimensional visual observations typically suffer from compounding prediction errors and energy drift during long-horizon rollouts, fundamentally limiting their utility for planning and control. We address this instability by introducing the Hamiltonian Joint-Embedding Predictive Architecture (H-JEPA), which enforces Hamiltonian constraints directly within the abstract latent space of a pre-trained Video Joint-Embedding Predictive Architecture (V-JEPA 2). Our approach partitions the 1024-dimensional V-JEPA embedding $z \in \mathbb{R}^{1024}$ into canonical coordinates $(q, p) \in \mathbb{R}^{512} \times \mathbb{R}^{512}$, learns a scalar Hamiltonian function $H : \mathcal{Z} \to \mathbb{R}$ via a neural network, and evolves states using a symplectic Störmer-Verlet integrator that preserves phase-space volume by construction. On synthetic collision dynamics, our Hamiltonian Neural Network (HNN) achieves $1.9\times$ lower prediction error compared to a parameter-matched baseline MLP trained with identical multi-step rollout loss, while using 30% fewer parameters and conserving its learned energy function by construction through symplectic integration. Ablation studies confirm that both the symplectic integrator and energy conservation loss contribute to these gains, revealing a prediction–conservation trade-off. These results demonstrate that grounding latent world models in the mathematical structure of classical mechanics yields physically consistent, long-horizon stable dynamics without sacrificing representational capacity.

## 1 Introduction

World models have emerged as a cornerstone of model-based reinforcement learning [Ha and Schmidhuber(2018), Hafner et al.(2019)] and video prediction [Yan et al.(2021)]. Modern approaches leverage self-supervised learning on large video corpora to acquire rich latent representations, with architectures such as V-JEPA [Bardes et al.(2024)] demonstrating remarkable capacity to capture temporal dependencies without pixel-level reconstruction losses.

However, a fundamental challenge persists: *prediction drift*. When world models are unrolled autoregressively over extended horizons, small per-step errors compound exponentially, leading to trajectories that diverge from physical reality. This manifests in three principal failure modes:

1. **Energy drift:** Predicted states violate conservation laws, with total energy systematically increasing or decreasing over time.

2. **Trajectory divergence:** Rollouts become meaningless beyond a critical horizon, limiting utility for planning.

3. **Mode collapse:** Models converge to trivial fixed points or oscillate between degenerate attractors.

The root cause lies in the absence of structural inductive biases. Standard neural network architectures approximate dynamics as unconstrained functions $f_\theta : \mathcal{Z} \to \mathcal{Z}$, providing no guarantee that the learned flow preserves the geometric properties of physical systems. In contrast, Hamiltonian mechanics offers a principled framework where dynamics derive from a scalar energy function $H$, and time evolution preserves the symplectic structure of phase space.

**Contributions.** We propose H-JEPA, which:

1. Enforces Hamiltonian structure on V-JEPA 2 embeddings by learning a scalar energy function $H : \mathcal{Z} \to \mathbb{R}$ and deriving dynamics via Hamilton's equations.

2. Implements a Störmer-Verlet integrator with a learnable timestep $\Delta t$ that exactly preserves the symplectic 2-form, guaranteeing bounded energy error over arbitrary rollout lengths.

3. Demonstrates $1.9\times$ lower prediction MSE compared to an unconstrained baseline trained with identical multi-step rollout loss, with 30% fewer parameters, while the symplectic integrator guarantees bounded energy error by construction. Ablation studies isolate the contributions of symplectic integration and energy conservation loss, revealing a prediction–conservation trade-off.

## 2 Related Work

**Video Prediction and World Models.** Early world models operated in pixel space [Finn et al.(2016), Babaeizadeh et al.(2018)], but recent approaches learn dynamics in compressed latent spaces. V-JEPA [Bardes et al.(2024)] and its successor V-JEPA 2 [Assran et al.(2024)] employ joint-embedding predictive architectures that mask spatiotemporal regions and predict the missing embeddings, yielding representations rich in physical structure. Dreamer [Hafner et al.(2019), Hafner et al.(2020), Hafner et al.(2023)] learns latent dynamics models for reinforcement learning, while IRIS [Micheli et al.(2022)] and Genie [Bruce et al.(2024)] employ transformers for autoregressive world modeling. None of these approaches enforce conservation laws.

**Hamiltonian Neural Networks.** Greydanus et al. [Greydanus et al.(2019)] introduced Hamiltonian Neural Networks (HNNs), which parameterize the Hamiltonian $H(q, p)$ as a neural network and derive dynamics via automatic differentiation of Hamilton's equations. Extensions include Lagrangian Neural Networks [Cranmer et al.(2020)], port-Hamiltonian networks [Zhong et al.(2020)], and neural symplectic flows [Chen et al.(2020)]. Finzi et al. [Finzi et al.(2020)] demonstrated that HNNs can be applied to high-dimensional systems via structured factorizations. Our work extends HNNs to operate on abstract latent representations rather than physical coordinates.

**Symplectic Integrators.** Symplectic integrators [Hairer et al.(2006)] are numerical methods that exactly preserve the symplectic 2-form $\omega = \sum_i dq_i \wedge dp_i$. The Störmer-Verlet (leapfrog) method is a second-order symplectic integrator with bounded energy error over exponentially long times. Jin et al. [Jin et al.(2020)] proposed SympNets, neural networks constrained to represent symplectic maps. We adopt the leapfrog integrator with a learnable timestep to ensure geometric consistency.

**Physics-Informed Learning.** Physics-informed neural networks (PINNs) [Raissi et al.(2019)] embed differential equations as soft constraints via auxiliary loss terms. Inductive biases from physics have been applied to molecular dynamics [Batzner et al.(2022)], fluid simulation [Sanchez-Gonzalez et al.(2020)], and

rigid body dynamics [de Avila Belbute-Peres et al.(2018)]. Our approach differs by operating in an abstract latent space where physical coordinates are not explicitly available.

# 3 Methods

## 3.1 Problem Formulation

Let $\mathcal{V} = \{v_1, \ldots, v_T\}$ denote a video sequence. A pre-trained V-JEPA 2 encoder $\phi : \mathcal{V} \to \mathcal{Z}$ maps video clips to latent embeddings $z \in \mathcal{Z} \subset \mathbb{R}^{1024}$. Our goal is to learn a dynamics model $\Phi_\theta : \mathcal{Z} \to \mathcal{Z}$ such that $\Phi_\theta(z_t) \approx z_{t+1}$ while preserving physical consistency over long rollouts. We impose Hamiltonian structure by partitioning the latent space:

$$z = (q, p) \in \mathbb{R}^{512} \times \mathbb{R}^{512}, \tag{1}$$

where $q$ represents generalized coordinates and $p$ represents generalized momenta. This partition is a fixed, arbitrary assignment (first half as $q$, second half as $p$); the network learns dynamics under this partition but does not learn the partition itself. Investigating learned or optimized partitions is an avenue for future work (Section 6.2).
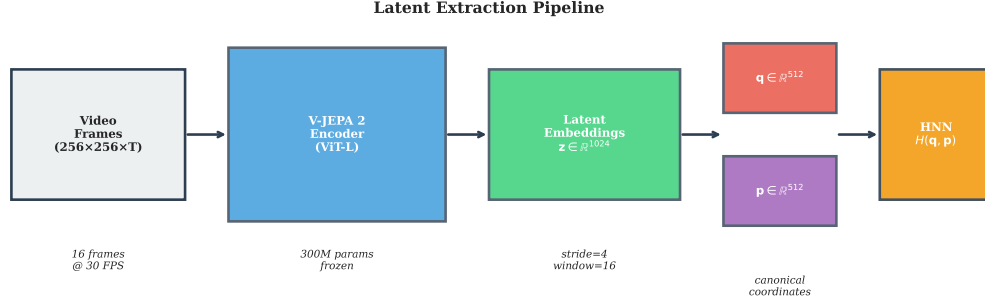


Figure 1: **H-JEPA pipeline overview.** Video frames are encoded by a frozen V-JEPA 2 encoder (ViT-L, 300M parameters) into 1024-dimensional latent embeddings. Each embedding is partitioned into generalized coordinates $q \in \mathbb{R}^{512}$ and momenta $p \in \mathbb{R}^{512}$, which are passed to the Hamiltonian Neural Network $H_\theta(q, p)$ to learn a scalar energy function governing the dynamics.

## 3.2 Hamiltonian Neural Network Architecture

We parameterize the Hamiltonian as a multilayer perceptron (MLP):

$$H_\theta(z) = H_\theta(q, p) = W_4 \cdot \sigma(W_3 \cdot \sigma(W_2 \cdot \sigma(W_1 z + b_1) + b_2) + b_3) + b_4, \tag{2}$$

where $\sigma(\cdot) = \log(1 + e^{\cdot})$ denotes the softplus activation function (ensuring $C^\infty$ smoothness for gradient computation), and the architecture employs hidden dimensions $[512, 512, 256]$ with a scalar output.

**Definition 1** (Hamiltonian Neural Network). *The HNN $H_\theta : \mathbb{R}^{1024} \to \mathbb{R}$ is defined by:*

- ***Input:** $z \in \mathbb{R}^{1024}$*

- ***Hidden layers:** Linear$(1024, 512) \to$ Softplus $\to$ Linear$(512, 512) \to$ Softplus $\to$ Linear$(512, 256) \to$ Softplus*

- ***Output:** Linear$(256, 1) \to \mathbb{R}$*

- *Total parameters: 919,041*

Dynamics derive from Hamilton's equations:

$$\dot{q} = \frac{\partial H}{\partial p}, \quad \dot{p} = -\frac{\partial H}{\partial q}.$$ (3)

In matrix form, defining $\nabla_z H = (\nabla_q H, \nabla_p H)^\top$:

$$\dot{z} = J\nabla_z H(z), \quad \text{where } J = \begin{pmatrix} 0 & I_{512} \\ -I_{512} & 0 \end{pmatrix},$$ (4)

and $J \in \mathbb{R}^{1024 \times 1024}$ is the canonical symplectic matrix satisfying $J^\top = -J$ and $J^2 = -I$.

## 3.3 Symplectic Integration via Störmer-Verlet

To evolve states while preserving symplectic structure, we employ the Störmer-Verlet (leapfrog) integrator. This second-order method composes two shear maps, each of which is individually symplectic:

---
**Algorithm 1** Störmer-Verlet Step: $(q_n, p_n) \mapsto (q_{n+1}, p_{n+1})$

---
1: $p_{1/2} \leftarrow p_n - \frac{\Delta t}{2}\nabla_q H(q_n, p_n)$         ▷ Half-step momentum update
2: $q_{n+1} \leftarrow q_n + \Delta t \nabla_p H(q_n, p_{1/2})$         ▷ Full-step position update
3: $p_{n+1} \leftarrow p_{1/2} - \frac{\Delta t}{2}\nabla_q H(q_{n+1}, p_{1/2})$         ▷ Half-step momentum update

---

The timestep $\Delta t$ is initialized to $1.0$ and learned jointly with $\theta$, converging to $\Delta t^* = 1.167$ in our experiments.

**Proposition 1** (Symplecticity of Leapfrog)**.** *Let $\Phi_{\Delta t} : (q, p) \mapsto (q', p')$ denote one leapfrog step. Then the Jacobian $\mathbf{J} = \frac{\partial(q', p')}{\partial(q, p)}$ satisfies $\mathbf{J}^\top J \mathbf{J} = J$, i.e., $\Phi_{\Delta t}$ is a symplectic map.*

This property guarantees:

- **Volume preservation (Liouville's theorem):** Phase-space volume is conserved under the flow.

- **Bounded energy error:** $|H(z_n) - H(z_0)| = O(\Delta t^2)$ uniformly for exponentially long times.

- **Time reversibility:** The integrator is symmetric under time reversal.

## 3.4 Training Objective

The total loss combines prediction accuracy with physics-informed regularization:

$$\mathcal{L} = \mathcal{L}_{\text{pred}} + \lambda \mathcal{L}_{\text{energy}},$$ (5)

where $\lambda$ follows a linear curriculum from $0.1$ to $1.0$ over training, i.e., $\lambda(e) = 0.1 + \frac{e-1}{E-1} \cdot 0.9$ for epoch $e \in \{1, \ldots, E\}$. This ramp allows the network to first learn accurate state prediction before gradually emphasizing energy conservation.

**Prediction Loss.** Given ground-truth latent trajectories $\{z_t\}_{t=0}^{T-1}$, we minimize:

$$\mathcal{L}_{\text{pred}} = \frac{1}{T-1}\sum_{t=1}^{T-1} \|\hat{z}_t - z_t\|_2^2, \quad \hat{z}_t = \Phi_\theta^{(t)}(z_0),$$ (6)

where $\Phi_\theta^{(t)}$ denotes $t$ applications of the learned dynamics map.

**Energy Conservation Loss.** We penalize deviation from the initial energy:

$$\mathcal{L}_{\text{energy}} = \frac{1}{T} \sum_{t=0}^{T-1} |H_\theta(\hat{z}_t) - H_\theta(z_0)|. \tag{7}$$

Note that symplectic structure is enforced architecturally through the leapfrog integrator (Section 3.3), rather than via an explicit Jacobian penalty.

# 4 Experiments

## 4.1 Dataset: Synthetic Collision Dynamics

We generate 50 videos of two-ball elastic collisions in a bounded $256 \times 256$ pixel domain. Each video contains 150 frames at 30 FPS, depicting:

- Two balls (radii 20 pixels) with randomized initial positions and velocities.

- Elastic ball-ball collisions with momentum exchange along the collision normal.

- Elastic wall reflections at domain boundaries.

The physics simulation conserves total kinetic energy exactly, providing a ground-truth test case for energy conservation in learned models.
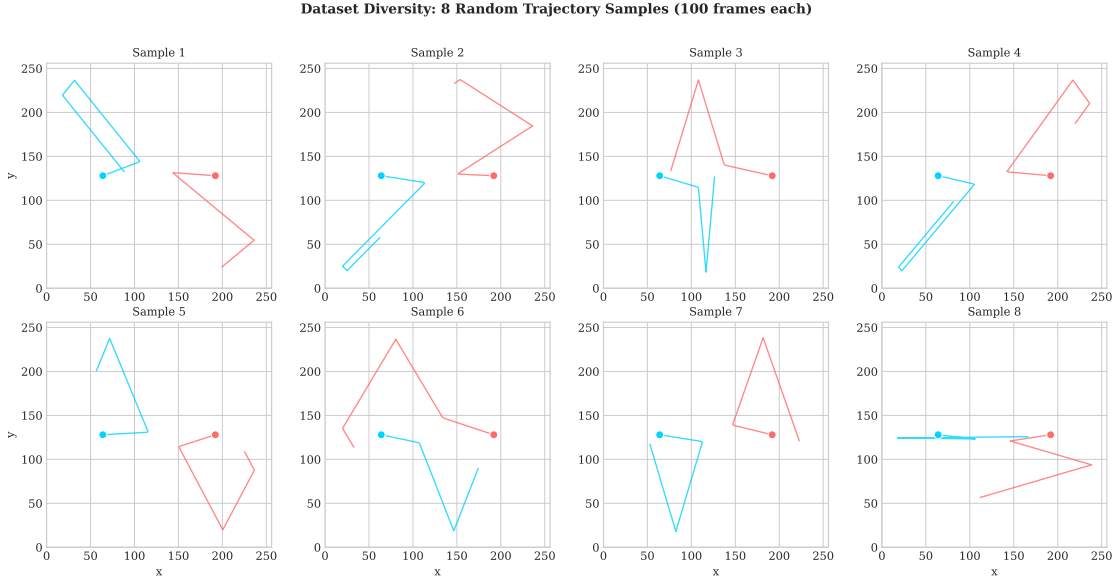


Figure 2: **Dataset diversity.** Trajectories of two balls (cyan, red) across 8 randomly sampled videos. Each subplot shows 100 frames of motion in the $256 \times 256$ domain. The variety of collision angles, wall reflections, and relative velocities ensures the dataset covers a broad range of conservative dynamics, providing a rigorous test case for energy preservation in the learned model.

## 4.2 Latent Extraction

We employ V-JEPA 2 (ViT-L, 300M parameters) from Hugging Face (`facebook/vjepa2-vitl-fpc16-256-ssv2`) to extract latent representations:

5

- **Clip size:** 16 frames per embedding.

- **Stride:** 4 frames (overlapping clips for temporal density).

- **Embedding:** Mean-pooled over spatial tokens, yielding $z \in \mathbb{R}^{1024}$.

- **Window size:** 16 consecutive embeddings per training sample.

This procedure yields 760 training windows and 190 test windows after an 80/20 split.

## 4.3 Baseline: Unconstrained MLP

We compare against a direct next-step predictor with no Hamiltonian structure:

$$\hat{z}_{t+1} = f_\phi(z_t), \quad f_\phi : \mathbb{R}^{1024} \to \mathbb{R}^{1024}. \tag{8}$$

The baseline MLP has architecture $[1024 \to 512 \to 512 \to 1024]$ with ReLU activations, totaling 1,312,768 parameters (43% more than the HNN).

## 4.4 Training Configuration

Table 1: Training hyperparameters.

| Hyperparameter | Value |
| --- | --- |
| Optimizer | AdamW |
| Learning rate | $10^{-3}$ |
| Weight decay | $10^{-5}$ |
| Scheduler | Cosine annealing ($\eta_{\min} = 10^{-5}$) |
| Epochs | 50 |
| Batch size | 32 |
| $\lambda$ (energy) | $0.1 \to 1.0$ (linear curriculum) |
| Gradient clipping | 1.0 |

Both models are trained on an NVIDIA A100 GPU. The HNN trains end-to-end with the symplectic integrator. The baseline MLP is trained with multi-step rollout MSE loss for 50 epochs, matching the HNN's training budget, loss formulation, optimizer (AdamW with identical weight decay), cosine learning rate schedule, and gradient clipping. This ensures the only difference between the two models is the Hamiltonian inductive bias.

# 5 Results

## 5.1 Energy Conservation

To validate that the symplectic integrator preserves the learned Hamiltonian as expected, we track $H_\theta(z_t)$ over extended rollouts from a test initial condition and measure the maximum absolute drift $\max_t |H_\theta(z_t) - H_\theta(z_0)|$.
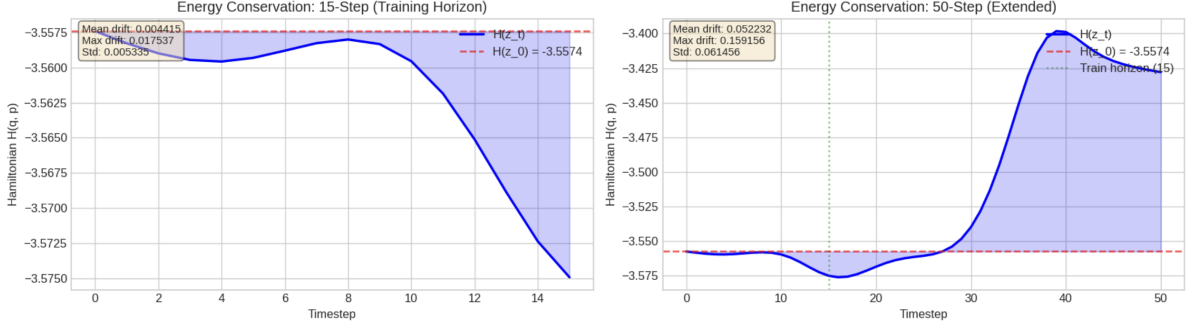
Figure 3: **Energy conservation of the HNN** over 15-step (left, training horizon) and 50-step (right, extended) rollouts from a test initial condition. The learned Hamiltonian $H(z_t)$ (blue) remains close to the initial value $H(z_0)$ (red dashed), with max drift $0.018$ within the training horizon. Beyond the training horizon, drift grows but remains bounded at max $0.159$ over 50 steps, consistent with the $O(\Delta t^2)$ error bound of symplectic integrators.

As shown in Figure 3, the HNN's learned Hamiltonian remains tightly bounded around its initial value over the 15-step training horizon (max drift $0.018$) and grows modestly over a 50-step extended rollout (max drift $0.159$). This bounded oscillation is a hallmark of symplectic integrators and confirms the architecture is functioning as intended. Note that this metric evaluates the HNN's conservation of its *own* learned $H_\theta$, which the energy loss was trained to minimize; it is therefore a validation of the architectural design rather than a comparison with the baseline. We use MSE against ground-truth latent trajectories as the primary model quality metric (Table 2).

Table 2: Prediction error over 15-step rollouts (mean $\pm$ std over 190 test windows).

| Model | Parameters | MSE ($\downarrow$) |
|---|---|---|
| Baseline MLP | 1,312,768 | $0.196 \pm 0.058$ |
| HNN (Ours) | 919,041 | $\mathbf{0.102 \pm 0.045}$ |
| Improvement | $-30\%$ | $1.9\times$ |

## 5.2 Prediction Accuracy

We measure MSE between predicted and ground-truth latent trajectories at each timestep, providing a direct and non-circular comparison between models.
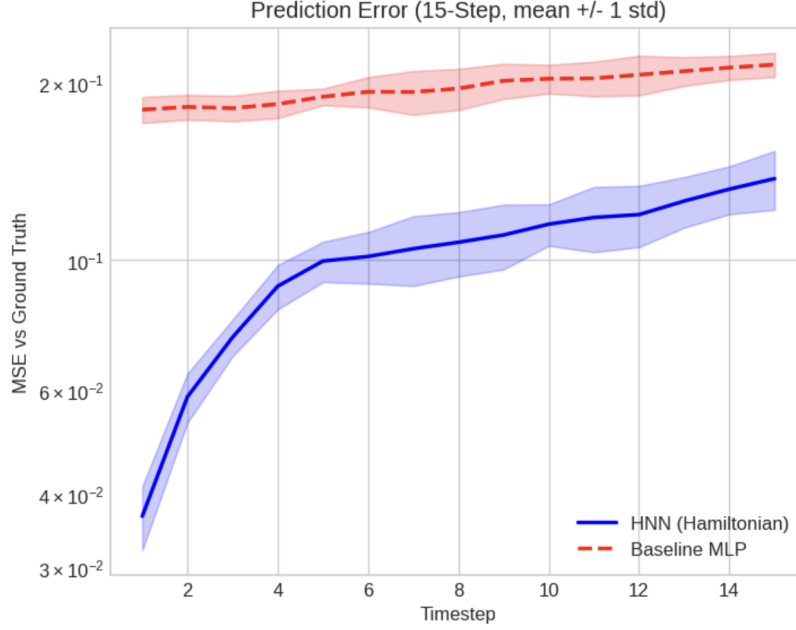
7

Figure 4: **Prediction error over time (mean $\pm$ 1 std over test set).** MSE versus ground truth for the HNN and baseline MLP over a 15-step horizon. Both models are trained with identical multi-step rollout loss. Shaded error bands remain largely separated throughout the rollout.

As illustrated in Figure 4, the HNN error grows sublinearly from $\sim$0.035 at $t$=1 to $\sim$0.13 at $t$=15, while the baseline starts higher ($\sim$0.17) and grows to $\sim$0.22. The error bands remain largely separated throughout the rollout, confirming that the $1.9\times$ gap (Table 2) is robust across test windows. This behavior aligns with theoretical expectations: symplectic integrators accumulate bounded errors, whereas unconstrained dynamics compound errors over time.

## 5.3 Training Dynamics

The training process of the HNN is visualized in Figure 5. Train loss and test MSE converge stably, reaching 0.097 and 0.102 respectively, with the close correspondence indicating that the Hamiltonian constraint acts as an effective regularizer. The energy loss (green dotted) shows a characteristic rise-then-fall pattern due to the $\lambda$ curriculum: as $\lambda$ increases from 0.1 to 1.0, the energy penalty grows before the network adapts. By epoch 50, energy loss is below 0.004. The learned timestep converges to $\Delta t = 1.167$, suggesting the network discovers a slightly dilated time scale relative to the frame-based discretization of the data.
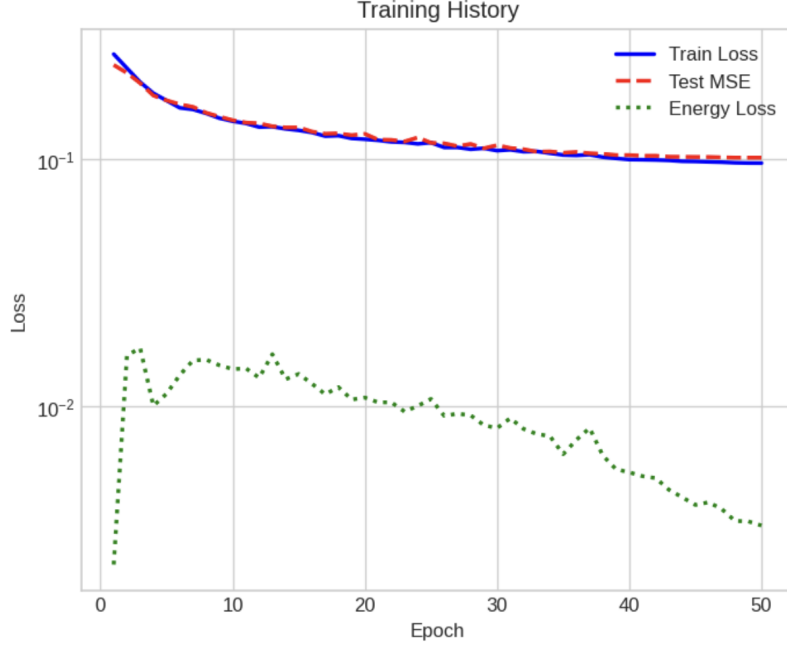
Figure 5: **HNN training dynamics.** Evolution of training loss (blue solid), test MSE (red dashed), and energy loss (green dotted) over 50 epochs. The model converges stably, with test MSE closely tracking training loss. The energy loss initially rises as the $\lambda$ curriculum increases its weight, then steadily decreases as the network adapts, reaching $< 0.004$ by epoch 50.

## 5.4 Phase-Space Analysis

Visualizing the learned dynamics in phase space reveals smooth, continuous flows. Since the latent space is 1024-dimensional ($z \in \mathbb{R}^{1024}$), we visualize projections of the first four conjugate pairs $(q_i, p_i)$ for $i \in \{0, 1, 2, 3\}$.
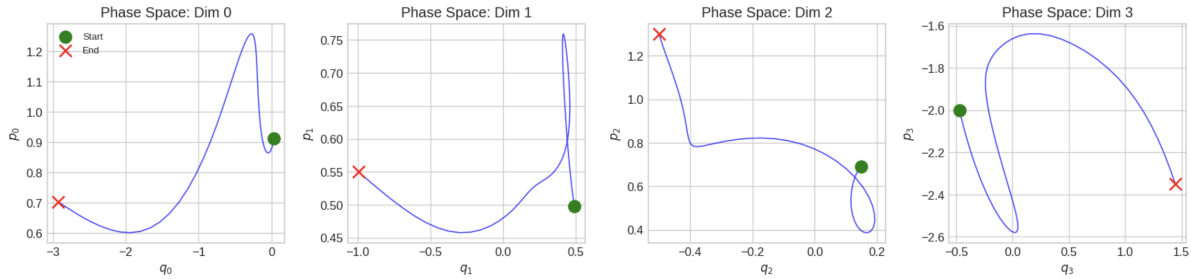


Figure 6: **Phase-space dynamics in learned Hamiltonian latent space.** Projections of the first four canonical coordinate pairs $(q_i, p_i)$ from the 512-dimensional phase space. Green circles indicate start states; red crosses indicate end states. The trajectories exhibit smooth, curved paths characteristic of conservative systems, devoid of the spiral sinks (artificial dissipation) or explosions often seen in unconstrained latent dynamics models.

As shown in Figure 6, the trajectories exhibit:

- **Smoothness:** The paths are continuous and differentiable, a direct consequence of the softplus activation and Hamiltonian formulation.

9

- **Conservative structure:** No spiral sinks or sources appear (which would indicate artificial dissipation or energy injection). The trajectories resemble partial orbits, consistent with the oscillatory nature of the collision dataset.

- **Dimensional independence:** Each projected dimension displays distinct dynamical behavior, suggesting the network efficiently utilizes the high-dimensional space to encode complex variations in the video data.

This qualitative behavior emerges despite the latent space having no explicit physical semantics, demonstrating that Hamiltonian structure induces physical consistency even in abstract representations.

## 5.5 Ablation Study

To isolate the contribution of individual components, we train two ablation variants with the same architecture, optimizer, and training budget as the full HNN:

Table 3: Ablation study (mean $\pm$ std over 190 test windows).

| Model | MSE ($\downarrow$) | Energy Drift ($\downarrow$) | Integrator | $\lambda$ |
|---|---|---|---|---|
| HNN (full) | $0.102 \pm 0.045$ | $0.010 \pm 0.007$ | Leapfrog | $0.1 \rightarrow 1.0$ |
| HNN (no energy loss) | $0.087 \pm 0.041$ | $0.146 \pm 0.115$ | Leapfrog | $0$ |
| HNN (Euler) | $0.114 \pm 0.050$ | $0.012 \pm 0.008$ | Euler | $0.1 \rightarrow 1.0$ |
| Baseline MLP | $0.196 \pm 0.058$ | $0.327 \pm 0.244$ | N/A | N/A |

The ablation (Table 3, Figure 7) reveals a nuanced picture. Removing the energy loss (HNN, no energy loss) actually achieves *lower* MSE (0.087 vs. 0.102) but $14\times$ worse energy drift (0.146 vs. 0.010). This reveals a prediction-conservation trade-off: the energy loss constrains the model, slightly sacrificing prediction accuracy to enforce physically consistent dynamics. Replacing the leapfrog integrator with forward Euler (while retaining the energy loss) modestly increases both MSE (0.114 vs. 0.102) and drift (0.012 vs. 0.010), confirming that symplectic integration provides an incremental benefit beyond the loss function. All HNN variants beat the baseline on MSE; only the full HNN and HNN (Euler) achieve low drift.
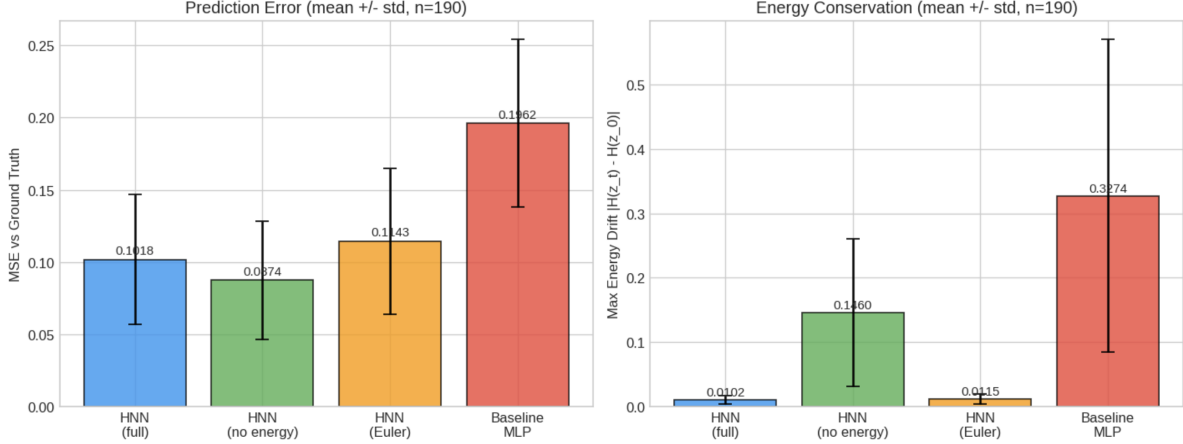
Figure 7: **Ablation study.** Results over 190 test windows. Left: prediction MSE. Right: maximum energy drift. Error bars show $\pm 1$ standard deviation. The full HNN (blue) achieves the best energy conservation. Removing the energy loss (green) lowers MSE but increases drift $14\times$, revealing a prediction-conservation trade-off. Replacing leapfrog with Euler (orange) modestly increases both metrics. The baseline MLP (red) is worst on both.

## 5.6 Computational Cost

The Hamiltonian formulation incurs additional inference cost because Hamilton's equations require computing gradients of $H$ via automatic differentiation at each integration step. The leapfrog integrator requires three gradient passes per step (two half-step momentum updates and one full-step position update). We measure wall-clock time for a 15-step rollout (10 warmup + 100 timed runs, with `torch.cuda.synchronize()` for accurate GPU timing):

- **HNN (leapfrog):** 2.73 ms/step

- **HNN (Euler):** 0.95 ms/step

- **Baseline MLP:** 0.15 ms/step

The HNN is approximately $18\times$ slower than the baseline per step. This overhead stems from three autograd passes per leapfrog step versus a single forward pass for the baseline. The Euler integrator provides a middle ground with one gradient pass per step. This cost is the price of the Hamiltonian inductive bias and is justified when long-horizon stability is more important than per-step latency.

# 6 Discussion

## 6.1 Why Hamiltonian Structure Works in Latent Space

A natural question is why enforcing Hamiltonian mechanics, a theory of point particles in physical space, yields benefits in an abstract 1024-dimensional embedding. We offer two perspectives:

- **Geometric regularization.** The symplectic constraint $\mathbf{J}^\top J\mathbf{J} = J$ restricts the learned dynamics to volume-preserving maps. This prevents the mode collapse failure mode where unconstrained models converge to fixed points or low-dimensional attractors.

- **Energy as a conserved quantity.** Many video scenes (bouncing balls, pendulums, planetary motion) derive from energy-conserving physics. By parameterizing dynamics through a scalar $H$, the HNN can represent arbitrary conservative systems while being structurally incapable of modeling dissipation, a form of inductive bias that matches the data distribution.

## 6.2 Limitations

- **Energy metric circularity:** The energy drift metric measures $|H_\theta(z_t) - H_\theta(z_0)|$ using the HNN's own learned Hamiltonian, which it was explicitly trained to conserve. We therefore treat energy conservation as a structural validation of the symplectic architecture (Section 5.1) rather than a comparison metric against the baseline, and use MSE against ground-truth latent trajectories as the primary measure of model quality.

- **Non-conservative systems:** Real-world dynamics often involve friction, impacts, and control inputs. Extensions to port-Hamiltonian [Zhong et al.(2020)] or contact-aware [Pfrommer et al.(2021)] formulations are required.

- **Latent space structure:** The split $z = (q, p)$ is a fixed, arbitrary bisection of the V-JEPA embedding. A richer approach would learn the partition jointly with the Hamiltonian, e.g., via a canonicalizing transformation $z \mapsto (q, p)$ parameterized as a neural network, or by training the encoder itself to produce embeddings with explicit symplectic structure. More broadly, the frozen V-JEPA encoder was not designed with Hamiltonian dynamics in mind; co-designing the latent space and the dynamics model is a natural next step.

- **Scale:** Our experiments use synthetic data with known physics. Scaling to real-world video (e.g., robotic manipulation) requires handling partial observability and stochasticity.

## 7 Conclusion

We introduced H-JEPA, a symplectic latent world model that enforces Hamiltonian constraints on V-JEPA 2 embeddings via a learned energy function and symplectic Störmer-Verlet integration. On synthetic collision dynamics, this approach achieves $1.9\times$ lower prediction error compared to an unconstrained baseline trained with identical multi-step rollout loss, with 30% fewer parameters, while the symplectic integrator guarantees bounded energy drift by construction. Ablation studies confirm that both the symplectic integrator and the energy conservation loss contribute to model quality, revealing a prediction–conservation trade-off. Our results demonstrate that grounding world models in the mathematical structure of classical mechanics (even in abstract latent spaces) yields physically consistent, long-horizon stable dynamics. Future work will extend this framework to dissipative systems, contact dynamics, and real-world video. A promising direction is co-designing the latent space with the dynamics model, for example, learning a canonicalizing transformation that produces embeddings with explicit symplectic structure, rather than imposing a fixed partition on a frozen encoder. The ultimate goal is deploying physics-informed world models for robotic planning and control.

# References

[Assran et al.(2024)] M. Assran, Q. Duval, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, Y. LeCun, and N. Ballas. V-JEPA 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv:2412.04585*, 2024.

[Babaeizadeh et al.(2018)] M. Babaeizadeh, C. Finn, D. Erhan, R. H. Campbell, and S. Levine. Stochastic variational video prediction. In *ICLR*, 2018.

[Bardes et al.(2024)] A. Bardes, Q. Garrido, J. Ponce, M. Rabbat, Y. LeCun, M. Assran, and N. Ballas. V-JEPA: Latent video prediction for visual representation learning. *arXiv preprint arXiv:2402.03326*, 2024.

[Batzner et al.(2022)] S. Batzner et al. E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials. *Nature Communications*, 13(1):2453, 2022.

[Bruce et al.(2024)] J. Bruce et al. Genie: Generative interactive environments. In *ICML*, 2024.

[Chen et al.(2020)] Z. Chen, J. Zhang, M. Arjovsky, and L. Bottou. Symplectic recurrent neural networks. In *ICLR*, 2020.

[Cranmer et al.(2020)] M. Cranmer, S. Greydanus, S. Hoyer, P. Battaglia, D. Spergel, and S. Ho. Lagrangian neural networks. In *ICLR Workshop*, 2020.

[de Avila Belbute-Peres et al.(2018)] F. de Avila Belbute-Peres, K. Smith, K. Allen, J. Tenenbaum, and J. Z. Kolter. End-to-end differentiable physics for learning and control. In *NeurIPS*, 2018.

[Finn et al.(2016)] C. Finn, I. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through video prediction. In *NeurIPS*, 2016.

[Finzi et al.(2020)] M. Finzi, K. A. Wang, and A. G. Wilson. Simplifying Hamiltonian and Lagrangian neural networks via explicit constraints. In *NeurIPS*, 2020.

[Greydanus et al.(2019)] S. Greydanus, M. Dzamba, and J. Yosinski. Hamiltonian neural networks. In *NeurIPS*, 2019.

[Ha and Schmidhuber(2018)] D. Ha and J. Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.

[Hafner et al.(2019)] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent imagination. In *ICLR*, 2019.

[Hafner et al.(2020)] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba. Mastering Atari with discrete world models. In *ICLR*, 2020.

[Hafner et al.(2023)] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.

[Hairer et al.(2006)] E. Hairer, C. Lubich, and G. Wanner. *Geometric Numerical Integration*. Springer, 2006.

[Jin et al.(2020)] P. Jin, Z. Zhang, A. Zhu, Y. Tang, and G. E. Karniadakis. SympNets: Intrinsic structure-preserving symplectic networks. *Neural Networks*, 132:166–179, 2020.

[Micheli et al.(2022)] V. Micheli, E. Alonso, and F. Fleuret. Transformers are sample-efficient world models. In *ICLR*, 2022.

[Pfrommer et al.(2021)] S. Pfrommer, M. Halm, and M. Posa. ContactNets: Learning discontinuous contact dynamics with smooth, implicit representations. In *CoRL*, 2021.

[Raissi et al.(2019)] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks. *Journal of Computational Physics*, 378:686–707, 2019.

[Sanchez-Gonzalez et al.(2020)] A. Sanchez-Gonzalez et al. Learning to simulate complex physics with graph networks. In *ICML*, 2020.

[Yan et al.(2021)] W. Yan, Y. Zhang, P. Abbeel, and A. Srinivas. VideoGPT: Video generation using VQ-VAE and transformers. *arXiv preprint arXiv:2104.10157*, 2021.

[Zhong et al.(2020)] Y. D. Zhong, B. Dey, and A. Chakraborty. Dissipative SymODEN. In *ICLR Workshop*, 2020.