

---

# ON THE DYNAMICS OF COHERENT MEMORY STRUCTURES IN NEURAL FIELDS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Memory in biological neural networks is often supported by coherent spatiotemporal patterns, such as traveling waves and [neural activity confined to low dimensional manifolds, which are captured at mesoscopic scales by continuum neural field models](#). Substantial progress has been made in [mechanistically analyzing both biological and artificial neural network architectures](#). Recent works obtain interpretable latent states by imposing traveling waves or low-dimensional invariant manifolds, but typically do not provide data-driven explanations for when and why such structures emerge during task training. We develop a theoretical framework for studying [latent dynamics based on the Mori-Zwanzig projection-operator formalism](#). Our approach casts memory as a family of time-dependent projections that reveal how coupled dynamics support memory encoding and decoding. We instantiate a neural-field-inspired architecture, and evaluate it on both long-range benchmarks and [neuroscience applications involving EEG and ECoG prediction](#). Across these tasks, we observe robust long-range accuracy and interpretable memory modes in the learned latent dynamics.

## 1 INTRODUCTION

Biological neural networks exhibit a range of coherent dynamical phenomena—such as stable attractor states and traveling waves (Engel et al., 2001; Wang, 1999; Engel & Steinmetz, 2019)—that are increasingly implicated in working memory and large-scale coordination across cortical regions. These phenomena have been studied using a variety of dynamical modeling frameworks, including neural field models (Ermentrout, 1998; Coombes, 2005a)—continuous, spatially extended dynamical systems that describe mesoscopic activity of densely interconnected neuron populations—as well as methods that reduce the neural dynamics to low-dimensional manifolds (Marrouch et al., 2020). These models provide mechanistic insight into how coherent population activity supports cognition, and have been used to explain core functions including the integration of sensory input, consolidation of long-term memories, and the organization of decisions, motor actions, and temporal sequences (Muller et al., 2014; Massimini et al., 2004; Rubino et al., 2006; Wimmer et al., 2014). These same attractor and wave phenomena are emergent in artificial neural networks (ANNs) (Rajan et al., 2016; Karuvally et al., 2024), and have been explicitly manipulated through architectural design to improve memory retention, sequence processing, and structured computation (Hopfield, 1982; Rusch & Mishra, 2021; Keller et al., 2024). Many of these architectures draw direct inspiration from biological neural networks, using insights from cortical dynamics to guide the design of memory and sequence-processing mechanisms in artificial systems. Some studies have attempted to bridge these motifs, demonstrating waves emerging from attractor instabilities (Coombes, 2005a) or attractor basins organizing wave propagation (Laing & Chow, 2002). However, they are typically studied in isolation (Ságodi et al., 2024; Karuvally et al., 2024; Keller et al., 2024), and a unifying theoretical framework to explain the flow of information in neural systems remains absent (Liu et al., 2025; Lei et al., 2024). Developing a framework that reconciles stability, propagation, and recall in biological and artificial networks remains a key challenge to explaining the dynamical behavior of intelligence (Alamia et al., 2025; Keller, 2025).

A central goal in neuroscience is to understand how cognitive functions emerge from complex neural activity, and dynamical models, particularly those grounded in attractor dynamics and traveling waves have proven essential for this purpose. Fixed-point and stationary bump attractors have been shown to stabilize persistent activity patterns in working memory (Wang, 1999; Wimmer et al.,

---

054 2014). Sequential and metastable attractors govern transitions between internal states, modeling  
055 decisions, motor plans, and temporal sequences (Friston, 1997; Kelso, 2012). Stochastic attractor  
056 exploration during rest and sleep facilitates internal simulation, cognitive flexibility, and memory  
057 consolidation through spontaneous traversal of neural state space (Deco et al., 2009; Chaudhuri  
058 et al., 2019). Traveling waves appear to serve complementary roles in coordinating activity across  
059 space and time, and in some instances emerge as instabilities of attractor states. Wavefronts, of-  
060 ten stimulus-evoked (Muller et al., 2014) or multistability-driven (Laing & Chow, 2002), support  
061 sensory integration (Muller et al.), attentional shifts (Maris et al., 2013), and large-scale coordina-  
062 tion (Takahashi et al., 2011) by propagating sustained activation through cortical maps. Wave pulses  
063 are localized, transient bursts, shaped by excitation-inhibition balance (Brunel, 2000) or excitabil-  
064 ity thresholds (Douglas et al., 1995), and are implicated in timing, signal relay, and motor plan-  
065 ning (Rubino et al., 2006; Latash et al., 2010). Spontaneous waves emerge endogenously during  
066 anesthesia (Townsend et al., 2015), sleep (Massimini et al., 2004), or perception (Davis et al., 2020),  
067 traverse cortical attractor landscapes to support memory consolidation (Lee & Wilson, 2002), inter-  
068 nal simulation, and synaptic refinement (Feller, 1999). [Attractor dynamics and traveling waves are  
069 tightly linked in pattern-formation and continuum neural field models \(Coombes, 2005b; Bressloff,  
070 2013\).](#) Here we instead consider high-dimensional neural systems trained for specific computa-  
071 tions, where attractor-based working-memory models and traveling-wave models are typically stud-  
072 ied within separate frameworks. For such learned systems, it is still unclear how changes in the  
underlying dynamics reorganize memory representations over time (Liu et al., 2025).

073 Recurrent neural networks (RNNs) often struggle to retain information over long timescales  
074 due to vanishing/exploding gradients (EVGP) and the compression of long histories into finite-  
075 dimensional hidden states Bengio et al. (1994); Pascanu et al. (2013); attention mitigates the fixed-  
076 context bottleneck Bahdanau et al. (2014), even as EVGP remains an issue (Zucchet & Orvieto,  
077 2024). Recent successes in deep state-space models (Gu & Dao, 2024) and transformer (Vaswani  
078 et al., 2017) architectures have overcome these challenges through structured state updates and  
079 self-attention mechanisms, respectively. However, even state-of-the-art transformers and deep  
080 state-space models can struggle with long-range dependencies and structured sequence tasks (Je-  
081 lassi et al., 2024), highlighting the importance of understanding memory mechanisms. Inspired by  
082 biological systems, many RNNs have been imbued with (stable) attractor-like (Rusch & Mishra,  
083 2021; Keller & Welling, 2023; Ságodi et al., 2024) or wave-like (Keller et al., 2024; Keller, 2025;  
084 Liu et al., 2025) structures to bolster memory retention and sequence processing. Intriguingly, even  
085 standard RNNs trained on history-dependent dynamical systems reveal latent waves under coordi-  
086 nate transforms (Karuvally et al., 2024). Inspired by neural fields, researchers have extended these  
087 ideas to practical applications, emulating cortical wave propagation for image segmentation (Li-  
088 boni et al., 2025), modeling spatially working memory geometries (Lei et al., 2024) and sensory  
input (Xie et al., 2022).

089 Mori-Zwanzig (MZ) formalism offers an exact decomposition of a dynamical system into an equa-  
090 tion over chosen variables, that explicitly accounts for the memory effects that shape their future  
091 behavior (Mori, 1965b; Zwanzig, 1961; Nakajima, 1958). Classical MZ is a technique developed  
092 for statistical mechanics that has been used to study molecular dynamics (Meyer et al., 2017), vis-  
093 cous Burgers flows (Stinis, 2012), and the Euler equations (Stinis, 2007). Data-driven machine  
094 learning approaches using MZ (Chorin et al., 2002; Lin et al., 2021; 2023) are a bottom-up ap-  
095 proach to reduced-order modeling (Givon et al., 2004; Gupta et al., 2024) similar to time-delay  
096 embeddings (Woodward et al., 2025), which have shown recent success in modeling isotropic tur-  
097 bulence (Tian et al., 2021) and hypersonic boundary layer transitions (Woodward, 2023). More re-  
098 cently, MZ has been used as a framework for deep learning (Venturi & Li, 2023), where it has been  
099 used to inform the latent state of LSTMS (Maulik et al., 2020), as an effective auto-encoder (Gupta  
100 et al., 2024), to predict time-dependent PDEs using neural operators (Buitrago et al., 2025), and  
101 to enhance the explainability of neural networks (Menier et al., 2023). However, two assumptions  
102 made by MZ inspired deep learning architectures oversimplify the latent dynamics. First, many  
103 MZ architectures formulate memory using a time-delay of the latent state, neglecting the inclusion  
104 of the generalized fluctuation-dissipation relation (GFDR) in the memory kernel (Lin et al., 2023).  
105 Second, many MZ inspired architectures assume an at equilibrium state for the latent dynamics ne-  
106 glecting the effects of time-dependent memory kernels (Grabert, 2006; Héry & Netz, 2024; Netz,  
107 2024; Venturi & Li, 2023). Moreover, the approaches that properly assume the structure of the  
latent dynamics neglect to account for the additional degrees of freedom often introduced during the  
encoding of information into the latent state. This approach is critical for learning time varying be-

havior of information in the latent state, where the latent state itself contains an over-representation of information. Recent work has linked these dynamics to the ability of data-driven MZ to discover emergent organization (Rupe & Crutchfield, 2024). To our knowledge, prior MZ-based approaches do not explicitly track how neural dynamics reorganize in latent space over the course of learning.

## 1.1 OUR CONTRIBUTION.

We present a novel theoretical framework for modeling the time-dependent dynamics of latent representations of an ANN during sequence learning. In particular:

1. We derive a generalized Langevin equation that accounts for intrinsic degrees of freedom using a family of time-dependent projections.
2. We provide practical guidance by implementing a biologically-inspired Neural Wave Field architecture equipped with MZ dynamics. By considering wave and oscillatory dynamics we are able to study information encoding and retrieval most naturally tied to the brain.
3. We empirically validate our approach by evaluating it on several long-range learning benchmarks and real-world neuroscience applications. We observe robust long-range recall, minimal memory dimension, and interpretable latent modes.

## 1.2 RELATED WORK

Several deep learning architectures using MZ-inspired time-delay memory, e.g. in neural operators and autoencoders (Buitrago et al., 2025; Gupta et al., 2024), do not explicitly enforce GFDR consistency as discussed in (Lin et al., 2023). Some data-driven techniques do enforce GFDR, e.g. through iterative regression (Lin et al., 2023; 2021). However, these approaches do not focus on coherent behaviors. Neural oscillators and traveling-wave networks directly encode coherent latent dynamics, but lack theoretical justifications for their dynamic behavior (Rusch & Rus, 2025; Keller et al., 2024). Transformers and structured state-space-models are often deployed as high-capacity predictors optimized for task performance (Gu et al., 2022; Fu et al., 2023; Gu & Dao, 2024).

By contrast, our approach aims to leverage coherent latent dynamics and enhance their mechanistic interpretability. This approach enables the model to suppress uninformative latents, elevate coherent structure, and shorten effective memory. See Appendix A for additional related works.

## 2 BACKGROUND

We first motivate our mathematical approach using a traveling wave example. We then formalize the background projection operator theory. Finally, we return to the traveling wave example to introduce notions of invariant-trivialization and projection-induced coherence.

### 2.1 MOTIVATION: COMPRESSING TRAVELING WAVE INFORMATION

To motivate our work, we consider a traveling wave over a *neural field*, which is a coarse-grained representation of cortex. Consider the field  $u : [0, L] \times [0, T] \mapsto \mathbb{R}$  with traveling wave dynamics

$$\partial_t u(x, t) = -\nu \partial_x u(x, t) = \mathcal{S}u(x, t), \quad \nu > 0, \quad 0 \leq x \leq L, \quad u(0, t) = f(t),$$

where  $\mathcal{S}$  is an advection (shift) operator implementing left-to-right transport at rate  $\nu$  and the right boundary is an outflow. Information is injected at  $x = 0$  by  $f(t)$  and then travels across the field without re-entering or reflecting. We have chosen this input/free-flow configuration in contrast to prior works Keller et al. (2024) that impose periodic boundary conditions; see Appendix B.1 for a detailed discussion. This choice is also biologically motivated. For example, visual cortical input is injected at specific input layers and then flows feedforward through downstream populations.

Our aim is to study the compression of information in the latent state  $u(x, \cdot)$  by restricting our view to a subset of the domain  $z \subsetneq [0, L]$ . After spatial discretization, we write the full state as  $\mathbf{u}(t) \in \mathbb{R}^N$ , let  $\mathbf{v}(t) \in \mathbb{R}^m$  denote the restriction of  $\mathbf{u}(t)$  to the *resolved* region  $z$ , and let  $\mathbf{w}(t) \in \mathbb{R}^{N-m}$  denote the remaining *unresolved* components. We define a projection  $P$  onto the resolved subspace by  $P\mathbf{u} = (\mathbf{v}, 0)$  and  $Q\mathbf{u} = (I - P)\mathbf{u} = (0, \mathbf{w})$ . The dynamics of  $\mathbf{u}(t)$  can be exactly decomposed as

$$\frac{d}{dt} \mathbf{u}(t) = \begin{pmatrix} PSP & PSQ \\ QSP & QSQ \end{pmatrix} \mathbf{u}(x, t) + \begin{pmatrix} Pe_0 \\ Qe_0 \end{pmatrix} f(t) = \begin{pmatrix} S_{vv} & S_{vw} \\ S_{wv} & S_{ww} \end{pmatrix} \begin{pmatrix} \mathbf{v}(t) \\ \mathbf{w}(t) \end{pmatrix} + \begin{pmatrix} \mathbf{b}_v \\ \mathbf{b}_w \end{pmatrix} f(t) = \frac{d}{dt} \begin{pmatrix} \mathbf{v}(t) \\ \mathbf{w}(t) \end{pmatrix},$$

where the blocks  $S_{vv}$ ,  $S_{vw}$ ,  $S_{wv}$ ,  $S_{ww}$  are induced by the restriction of  $S$  to the resolved and unresolved coordinates.

The dynamics of  $v(t)$  are then described by the following non-Markovian system

$$\frac{d}{dt}v(t) = S_{vv}v(t) + \int_0^t S_{vw}e^{(t-s)S_{ww}}S_{wv}v(s)ds + S_{vw}e^{tS_{ww}}w(0) + \int_0^t S_{vw}e^{(t-s)S_{ww}}b_w f(s)ds + b_v f(t),$$

see Appendix C.5 for a derivation. Thus, even though the full system  $u(t)$  evolves Markovianly under  $S$ , the compressed latent  $v(t)$  is non-Markovian.

**Lifting Information** To remove the dependence on the boundary we will introduce a lifting operator that allows us to simplify the dynamics of  $v$ . In particular, we will assume (A1) that we trade temporal encoding of information on the boundary for spatially-resolved information at some time  $k$ . We do this by introducing  $\mathcal{L}_k : (u(0), f_{[0,k]}) \rightarrow u(k)$ . Then starting from time  $k$  for  $\tau = t + k$

$$\frac{d}{dt}v(\tau) = S_{vv}v(\tau) + \int_k^\tau S_{vw}e^{(\tau-s)S_{ww}}S_{wv}v(s)ds + S_{vw}e^{\tau S_{ww}}w(k).$$

With this non-Markovian system in hand, we reflect on the generalization of this technique before revisiting this example.

## 2.2 PROJECTION OPERATOR FORMALISM

This technique is formalized for a fixed choice of resolved *measurement* (i.e., representation of information) by near-equilibrium MZ (NE-MZ) (Mori, 1965a; Zwanzig, 2001). An extension to dynamic representations is given by the far-from-equilibrium MZ (FFE-MZ) (Grabert, 2006).

**Near-equilibrium** The NE-MZ formalism provides an exact decomposition for the evolution of a measurement  $\partial_t g = \mathcal{L}g$  with *resolved*  $\hat{g} = Pg$  and *unresolved*  $\tilde{g} = Qg$  components. The resolved evolution is given by the generalized Langevin equation (GLE)

$$\frac{\partial}{\partial t}\hat{g}(t) = \underbrace{P\mathcal{L}\hat{g}(t)}_{\text{Markov}} + \underbrace{\int_0^t P\mathcal{L}e^{(t-s)Q\mathcal{L}}Q\mathcal{L}\hat{g}(s)ds}_{\text{Memory}} + \underbrace{P\mathcal{L}e^{tQ\mathcal{L}}Qg(0)}_{\text{Fluctuating Force}}. \quad (1)$$

Equation 1 consists of three distinct terms (underscored). The Markov term represents the instantaneous drift from the resolved dynamics. The Memory term re-introduces the influence of dynamics previously *forgotten*, i.e. prior resolved information that has been projected into the unresolved subspace. The Fluctuating Force term<sup>1</sup> captures the residual influence of the unresolved initial state.

**Far-from-equilibrium** For a neural network architecture, it may not be possible, and potentially unreasonable to ascribe to each element of the latent state a static representation of *what* information is encoded. FFE-MZ (Grabert, 2006) is a prior approach to handling time-dependent measurements by introducing time-dependent  $P(t)$ . Suppose  $P(t)$  is differentiable (A2), then the resulting GLE is

$$\frac{\partial}{\partial t}\hat{g}(t) = P(t)\mathcal{L}\hat{g}(t) + \dot{P}(t)g(t) + \int_0^t P(t)\mathcal{L}G(t,s)Q(s)\mathcal{L}\hat{g}(s)ds + P(t)\mathcal{L}G(t,0)Q(0)g(0). \quad (2)$$

The two-time memory kernel  $G(t,s) = \mathcal{T}_- \exp\left(\int_s^t Q(u)\mathcal{L}du\right)$  is the negatively time-ordered exponential operator that captures the *extrinsic* influence from the evolution of the subspaces. The Kinematic term  $\dot{P}(t)g(t)$  and captures the intrinsic evolution of the resolved subspace (Meyer et al., 2017). The take-away is that the time-dependent projection operator acts as a moving frame of reference tied to the desired measurement.

Our approach is distinguished from FFE-MZ in that we do not obtain a two-time memory kernel; although, we derive a similar drift term. This is because our lift operator assumption allows us to suppose that the information of interest is intrinsically in the latent state. As a result, we do not recover extrinsic influences from the evolution of the subspaces. We leave extensions to partially observed systems and lift-free derivations to future work.

<sup>1</sup>The third term referred to by (Mori, 1965a) as a random force and by (Zwanzig, 2001) as a fluctuating force, is frequently called the noise term in data-driven and stochastic applications.

### 2.3 COHERENT DYNAMICS IN A COMPRESSED TRAVELING-WAVE

We return to the compression of the traveling-wave to illustrate two regimes of coherent behavior that arise when the latent is an advecting wave but the readout only observes a subset of coordinates.

Consider the long-range copy task, a benchmark designed to test long-range information retention (Graves et al., 2014; Arjovsky et al., 2016; Keller et al., 2024). The task consists of an input sequence of  $N$  random scalar integers in  $\{1, \dots, 9\}$ , followed by  $T + N$  count of 0's. The target for this task is a sequence of the same length of all 0's except the last  $N$  elements that are set to the initial sequence. The duration is  $T + N$  and the resolved space is  $\hat{y} \in \mathbb{R}^1$ .

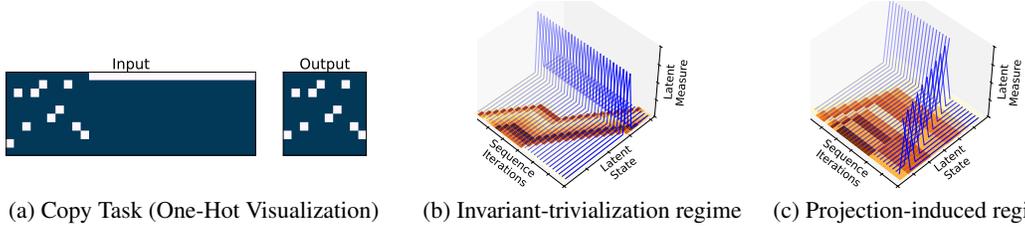


Figure 1: (a) The copy task. (b) Invariant trivialization. The readout  $P$  is fixed (blue), while the latent wave (orange) transports symbols to the readout location. (c) Projection-induced coherence. With constrained latent width, the latent saturates; a time-varying readout  $\{P_k\}$  sweeps the compressed latent to reconstruct the sequence at the output.

Let  $f_k \in \{0, \dots, 9\}$  denote the input at time  $k$ . Let  $S$  be a strictly subdiagonal shift, we introduce the softmax  $\sigma$  and parameter  $\theta$  to transition between  $S$  and a parameterized matrix  $W$

$$\mathbf{u}_{k+1} = \sigma(\theta)S\mathbf{u}_k + (1 - \sigma(\theta))W\mathbf{u}_k + e_0 f_k, \quad \hat{y}_k = P_k \mathbf{u}_k,$$

where  $e_0$  injects the input at the left boundary, and  $P_k$  restricts the resolved coordinates/readout. When the latent has width  $d$  and the task duration is  $T + N$ , two qualitatively distinct regimes appear as illustrated in Figure 1. In Figure 1(b) recall is carried entirely by the traveling wave with a fixed readout, and in Figure 1(c) the latent compresses and a time-varying readout sweeps over it. Both are examples of what we will refer to as coherent dynamics.

**Coherent dynamics** We call the dynamics *coherent* when either the latent wave or the readout subspace remains simple and persistent over extended intervals—e.g., a steady traveling front in the latent or a slowly varying/constant readout that consistently extracts stored content.

**Invariant trivialization (time-independent readout)** If the latent is wide enough to carry the entire history ( $d \geq T + N + 1$ ), we can choose a time-independent projector ( $\forall k, P = P_k$ ) that simply selects the coordinate where the advected information arrives at the end of the delay window (e.g., index  $T + N + 1$ ). In this case  $\hat{y}_k = f_{k-(T+N+1)}$  exactly and all time-variation needed for recall is supplied by the latent advection; the readout is stationary. We refer to this as *invariant trivialization* because, after the (co-moving) lift that centers the wave, the projection subspace does not change over time. This is illustrated in Figure 1b where the blue readout remains flat while the orange latent front carries the symbols downstream.

**Projection-induced emergence (time-dependent readout)** When the latent is narrow ( $d \leq T + N + 1$ ), the wave cannot store the entire history at distinct spatial positions. The learned evolution operator may deviate from the dynamics of the lift operator so that after the inputs cease (at step  $N$ ), the latent often settles into a compact, stable pattern that retains only a compressed trace of the sequence. Successful recall then requires the projection itself to evolve in time—a family  $\{P_k\}$  that sweeps across, or re-weights, the compressed latent so that the appropriate symbol is routed to the output at each step. Coherence therefore emerges from the projection dynamics: even though the latent has become stationary or a slow-moving manifold, the time-varying readout produces a systematic, wave-like replay at the output. This is illustrated in Figure 1c, and we term this projection-induced emergence.

These two coherent regimes in the simple copy serve as our motivating example. We formalize them using time-dependent projection operators and an intrinsic generalized Langevin equation.

---

### 3 MORI-ZWANZIG FORMALISM FOR COHERENCE AND EMERGENCE

We first formalize the two assumptions made, then our intrinsic GLE, and finally, coherence.

#### 3.1 AN INTRINSIC TIME-DEPENDENT GENERALIZED LANGEVIN EQUATION

We model a fixed resolved subspace of information  $\mathcal{V}_*$  which receives time-varying input from an unresolved space  $g(t) \in \mathcal{W}_t$ , and the aim is to output information  $h(t) \in \mathcal{V}_t \subset \mathcal{V}_*$  from a time-varying subspace of the resolved space. That is, we assume  $\mathcal{V}_*$  covers the entire family of incoming information  $\{\mathcal{W}_t\}$ , while the outputs evolve in the (possibly smaller) time-dependent subspaces  $\{\mathcal{V}_t\}$ . We formalize our lifting assumption as follows.

**Assumption 3.1.** (*Encoding Time-Dimension Tradeoff*) *The embedding of a generic input  $g_t \in \mathcal{W}_t$  is a transport map  $T_t : \mathcal{W}_t \rightarrow \mathcal{V}_*$  to a time-invariant subspace.*

For example, a traveling-wave is a co-moving shift, i.e., a linear lifting operator that increases dimensionality and centers moving patterns in the latent state. Kuramoto models are non-linear lifts that trade time-dependence in the signal for dimensions in phase coordinates.

We will learn a family of projection operators  $\{P_{\mu_t}\}$  by parameterizing their measures  $\{\mu_t\}$  over time. Consider a family of measures  $\{\mu_t\}_{t \in [0, T]}$  with  $\mu_t \ll \mu_*$  for all  $t$ , i.e.  $\mu_t$  is absolutely continuous with respect to  $\mu_*$ . Our time-dependent projection operators are defined by

$$P_{\mu_t} : \mathcal{V} \rightarrow \mathcal{V}_t, \quad P_{\mu_t} g = \mathbb{E}_{\mu_t} [g | \mathcal{G}], \quad \mathbb{E}_{\mu_t} [f | \mathcal{G}] = \frac{\mathbb{E}_{\mu_*} [\rho_t f | \mathcal{G}]}{\mathbb{E}_{\mu_*} [\rho_t | \mathcal{G}]},$$

which is the conditional-expectation onto the  $\sigma$ -algebra  $\mathcal{G}$  of the fixed resolved space but with weights  $\mu_t$ . Note that  $\rho_t = \frac{d\mu_t}{d\mu_{d0}}$  is the Radon-Nikodym derivative further discussed in Appendix C. We formalize our assumption of the derivative of  $P_{\mu_t}$  as follows.

**Assumption 3.2.** (*Differentiability of  $P_{\mu_t}$* ) *Suppose the time-dependent conditional expectation operator  $P_{\mu_t} : L^2(\mu_*) \rightarrow L^2(\mu_t)$  is Fréchet-differentiable with derivative  $\dot{P}_{\mu_t}$ .*

#### 3.2 INTRINSIC TIME-DEPENDENT GLE

Here we present our GLE. Appendix C provides a formal derivation, with proofs in Appendix D.

**Proposition 3.1.** (*Intrinsic Time-Dependent GLE*) *Let  $g(t)$  evolve under the operator  $\mathcal{L}$  on a fixed Hilbert space  $\mathcal{H} = L^2(\mathcal{M}, \mathcal{F}, \mu_*)$ . Let  $P_{\mu_*} : \mathcal{H} \rightarrow \mathcal{V} \subset \mathcal{H}$  be an orthogonal projection onto  $\mathcal{V} = L^2(\mathcal{M}, \mathcal{G}, \mu_*)$  with  $\mathcal{G} \subset \mathcal{F}$ . For a family of  $C^1$  measures  $\{\mu_t\}_{t \in [0, T]}$  let  $P_{\mu_t} : \mathcal{V} \rightarrow \mathcal{V}_t$  be the corresponding family of projections defining a Hilbert bundle  $\{\mathcal{V}_t\}_{t \in [0, t]}$  with  $\mathcal{V}_t = L^2(\mathcal{M}, \mathcal{G}, \mu_t)$ . The evolution of the resolved variable  $P_{\mu_t} g(t)$  satisfies the following GLE*

$$\begin{aligned} \frac{d}{dt} (P_{\mu_t} P_{\mu_*} g(t)) &= P_{\mu_t} \dot{P}_{\mu_t} Q_{\mu_t} P_{\mu_*} g(t) + P_{\mu_t} \mathcal{L} P_{\mu_*} g(t) \\ &+ \int_0^t P_{\mu_t} P_{\mu_*} \mathcal{L} e^{(t-s)Q_{\mu_*}} \mathcal{L} P_{\mu_*} g(s) ds + P_{\mu_t} \mathcal{L} e^{tQ_{\mu_*}} \mathcal{L} Q_{\mu_*} g(0). \end{aligned} \quad (3)$$

The additional term  $P_{\mu_t} \dot{P}_{\mu_t} Q_{\mu_t} P_{\mu_*} g(t)$  captures the instantaneous drift of the resolved state caused by the time-dependent rotation of the projection subspace, i.e., the transfer of latent information. This additional term is similar to the FFE-MZ. However, our approach does not result in a two-time memory kernel, and our drift depends only the dynamics of intrinsic subspaces  $P_{\mu_t}$  and  $Q_{\mu_t}$ .

#### 3.3 COHERENCE AND EMERGENCE

We now use the kinematic drift term from Proposition 3.1 to characterize how time-dependent projections give rise to coherent memory dynamics.

**Invariant trivialization** Our projections  $P_{\mu_t}$  live on a family of  $L^2(\mu_t)$  spaces whose inner products change with time. A trivialization is a way of re-expressing every  $L^2(\mu_t)$  inside a single reference space  $L^2(\mu_0)$ . In general, as  $t$  varies this change-of-measure will warp basic functions and

cause the projection subspace to rotate. Using the Radon-Nikodym densities (Appendix C), we obtain a spatially uniform change-of-measure. The trivialization becomes a global rescaling that does not distort the directions in  $L^2$ . In this case, there exists a basis that is preserved for all  $t$ , where the projection is time-invariant. We call this situation an *invariant trivialization*, and it corresponds to the lift-driven coherence regime from the copy example. We formalize this as follows.

**Proposition 3.2.** (Coherence Under Invariant Trivialization) *If the densities  $\rho_t(x)$  are spatially constant,  $\rho_t(x) = \alpha(t)$ , then the family of subspaces  $\{\mathcal{V}_t\}$  is unitarily equivalent to the fixed subspace  $\mathcal{V}_0$ . Then  $\{P_{\mu_t}\}$  is coherent under the invariant trivialization  $T_t$  where  $T_t P_{\mu_t} T_t^{-1} = P_{\mu_t} = P_{\mu_0}$ .*

In this case, the drift term will vanish, providing an **explainable and controllable mechanism**.

**Corollary 3.1.** (Vanishing Drift Under an Invariant Trivialization) *Suppose the Radon-Nikodym densities satisfy  $\rho_t(x) = \alpha(t)$ , and  $\alpha > 0$  independent of  $x$ . Then  $P_{\mu_t} = P_{\mu_0}$ , hence  $\dot{P}_{\mu_t} = 0$ .*

**Projection-induced coherence** In the intrinsic GLE of Proposition 3.1, the only place where time-dependence of the projection enters is through the kinematic drift term  $P_{\mu_t} \dot{P}_{\mu_t} Q_{\mu_t}$ . Intuitively, once we compress the latent dynamics onto an  $r$ -dimensional resolved subspace  $\mathcal{V}_t$ , any extra motion coming from  $\dot{P}_{\mu_t}$  must be funneled through a small set of directions inside  $\mathcal{V}_t$ . We refer to this situation—where nonzero  $D_t$  organizes the resolved dynamics along a few stable directions in  $\mathcal{V}_t$ —as projection-induced coherence. We make this statement precise in the following proposition.

**Proposition 3.3.** (Low-rank drift under latent compression) *Let  $\mathcal{V}_t$  be the resolved subspace at time  $t$ , with dimension  $\dim \mathcal{V}_t = r$ , and define the kinematic drift*

$$D_t := P_{\mu_t} \dot{P}_{\mu_t} Q_{\mu_t} : \mathcal{H} \rightarrow \mathcal{V}_t.$$

Then  $\text{rank}(D_t) \leq r$ . In particular, all additional drift induced by the time-dependence of the projection is confined to at most  $r$  independent directions in  $\mathcal{V}_t$ .

### 3.4 A NEURAL WAVE FIELD ARCHITECTURE

We instantiate the intrinsic GLE in a Neural Wave Field architecture that factorizes the latent evolution into a fixed lift and boundary update, an MZ-driven latent dynamics module, and a projection module illustrated in Figure 2. At each time step  $t$ , the input  $x_t$  is embedded as a *ghost boundary* into a 1-D latent field  $h_t \in \mathbb{R}^n$ . A fixed shift operator  $S$  implements the traveling-wave lift (co-moving frame), while a small gating network latent update mixes three simple behaviors, identity, pure shift and direct boundary injection. This mixture controls how much new input overwrites or augments the boundary versus how much of the existing wave is transported downstream.

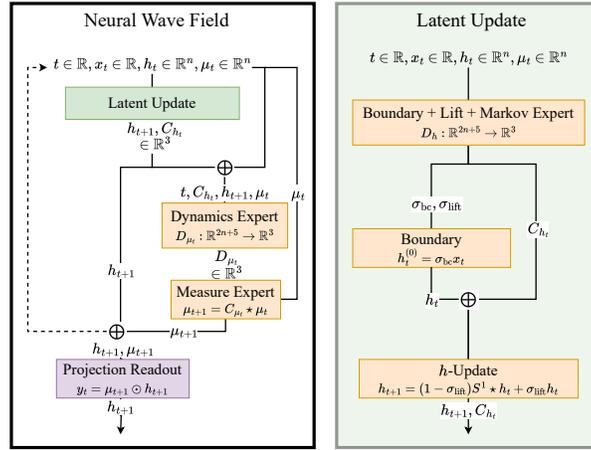


Figure 2: A diagram of the Neural Wave Field architecture.

The MZ drift and the projection dynamics are then modeled by two additional conditional experts. A dynamics expert  $D_{h_t}$  takes as context the current latent state, the boundary state, and the lift output, and returns coefficients that define a local convolution kernel  $C_h(h_t, x_t, t)$ . This kernel parameterizes the instantaneous drift term  $P_{\mu_t} \mathcal{L} P_{\mu_t}$  in the intrinsic GLE. In parallel, a measure expert  $D_{\mu_t}$  operates on the current measure  $\mu_t$ , the latent features, and time to produce a normalized convolution kernel  $C_{\mu_t}(h_t, \mu_t, t)$  that updates the measure  $\mu_{t+1}$  and thus the projection  $P_{\mu_t}$ . This yields a learned, context-dependent kinematic drift  $C_{\mu_t}$ . Under invariant trivialization,  $C_{\mu_t}$  vanishes, while under projection-induced coherence it remains nonzero but low-rank. The final output is the element-wise product  $y_t = h_t \odot \mu_t$ , so that the decoding is explicitly tied to the current projection.

## 4 EXPERIMENTAL RESULTS

To empirically evaluate our theoretical framework, we test our architectures ability to learn coherent dynamics for traveling-wave and non-linear oscillatory models. These results further support Proposition 3.1 across a range of task including long-range benchmarks and real-world EEGs. We find that the derived GFDR provides enhanced the robustness of the coherence across all tasks further supporting the use of MZ formalism. Furthermore, we demonstrate how these emergent behaviors can help characterize memory encoding, retention and retrieval similar to biological neural networks.

For a comparison on long-range benchmarks, we consider WaveRNN (Keller et al., 2024), Mamba (Gu & Dao, 2024), Alibi (Press et al., 2021), NoPe (Kazemnejad et al., 2023), and RoPe (Su et al., 2024). For a comparison on real-world data, we consider several baseline models including ShallowFBCSPNet Ang et al. (2008), Deep4 Schirmeister et al. (2017), EEGNet Lawhern et al. (2018) and TIDNet Kostas & Rudzicz (2020). Additional details including hyperparameters and optimization procedures and can be found in Appendix F.

### 4.1 COPY TASK

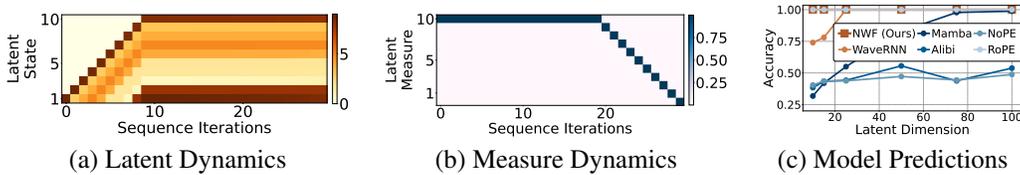


Figure 3: Explainability of the Neural Wave Field equipped with a traveling-wave lift operator in (a) the latent state (b) the measure and (c) a model comparison.

For comparison, we evaluate the accuracy of each architecture as the size of the latent stat is reduced. In particular, we systematically constrain the latent dimension from 100 to 10, the minimal size needed to represent the information in the long-range copy task for  $T = 10$ . This forces each model to rely on its latent dynamics rather than excess dimensionality, and allows us to test whether the core mechanism can efficiently encode and preserve information. In Figure 1(a), we see that our architecture maintains high accuracy even at the minimal latent dimension, where the information content fully saturates the latent state.

### 4.2 SELECTIVE COPY TASK

The selective copy task (Jing et al., 2019; Gu & Dao, 2024) modifies the copy task by randomizing the spacing of the  $N$  tokens over the first  $N + T$  inputs. The target is the same as the copy task. Due to this randomization, it requires more data-dependent reasoning to solve the task. From the FFE–MZ perspective, the task highlights how the lifting operation is tied to the time-dependent projections. Specifically, when the projection operator evolves in time, the lifting operator is static.

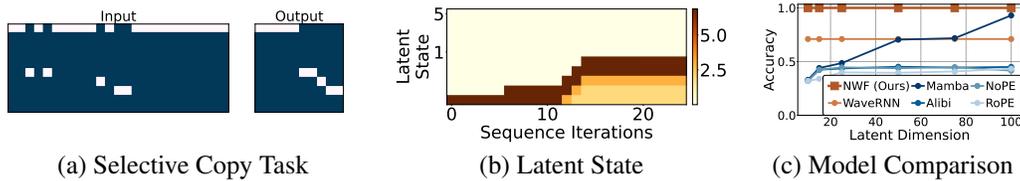


Figure 4: Results for the selective copy task (a) the latent state (b) the model predictions and (c) a memory capacity comparison against baselines.

In Figure 4 we illustrate the (a) latent state (b) model predictions and (c) a comparison of architectures as the memory capacity is reduced. As a result, of the tie between the lifting operator and the measure dynamics, we observe that the latent state encodes only the relevant information. Here it achieves a minimal yet sufficient representation in the latent state. As a result, it not only preserves accuracy but also outperforms all baseline architectures under identical constraints.

### 4.3 EEG DATASET

Using the braindecode library, we further benchmark our method on two neuroscience datasets from the BNCI IV competition 1) The BNCI IV-2a dataset, which contains EEG recordings from 9 subjects performing four motor imagery tasks: left hand, right hand, foot, and tongue movements. Each subject completed two sessions across different days, with 288 trials per session (12 per class per run, 6 runs per session). 2) The BNCI IV-4 dataset, which contains recordings and simultaneous finger-flexion measurements from three epilepsy patients at Harborview Hospital, Seattle. Each subject wore a subdural platinum-electrode grid (62, 48, and 64 channels for Subjects 1–3) sampled at 1000 Hz (0.15–200 Hz band-pass) and referenced to a common average; finger movements of all five digits were captured via a 5-sensor data glove (25 Hz, up-sampled to 1 kHz). This benchmark allows us to assess whether the traveling-wave inductive bias meaningfully improves decoding under practical EEG conditions.

To test whether latent propagating dynamics can serve as an effective inductive bias for EEG/ECoG decoding, we insert our Neural Wave Field module at the front of the network, directly operating on raw EEG/ECoG signals. This module compresses the raw sequence into a dynamic latent state by simulating learned traveling waves in feature space using gated, memory-aware updates derived from the Mori–Zwanzig formalism. The output is a sequence representation that is then passed into a standard CNN-based classification pipeline, similar to ShallowFBCSPNet. In this way, we can assess if the latent traveling wave representation enhances the expressivity of the models.

By placing the NeuralField before conventional spatial-temporal filtering, we evaluate whether traveling-wave dynamics can serve as an effective neural preprocessor, enhancing downstream performance. This setting allows us to test the expressiveness and utility of our proposed inductive bias in a realistic, cue-based EEG classification task.

Model	BNCI IV-2a (Accuracy $\uparrow$ )	BNCI IV-4 ( $r$ -value $\uparrow$ )
ShallowFBCSPNet	72.9	0.311
Deep4	56.25	<b>0.653</b>
EEGNet	<b>77.08</b>	0.354
TIDNet	40.97	0.356
NWF (Ours)	<b>74.31</b>	<b>0.375</b>

Table 1 presents the accuracy on the BNCI IV-2a dataset and the Pearson  $r$ -score on the BNCI IV-4 dataset. The Neural Wave Field is the second best performer with a single channel latent state size of 30, which maintains a compressed traveling wave representation of the full 22 channel input. Again the Neural Wave Field is the second best performer with a single channel latent state size of 20, obtaining a compressed traveling wave representation of the full 62 channel input. Moreover, it showed strong improvement over the direct baseline ShallowFBCSPNet. It also suggest the potential to include alternative lifting operators that align better with the underlying dynamics.

Table 1: Accuracy on the BNCI IV-2a dataset and Pearson’s  $r$  on the BNCI IV-4 dataset. The Neural Wave Field including a traveling wave lifting operator ranks as the second-best performer.

## 5 CONCLUSION

We introduced an intrinsic time-dependent framework for the Mori-Zwanzig formalism and used it to derive a structured model of latent memory dynamics. In particular we observed how a lifting operator and the latent drift were coupled. Building on this, we proposed the Neural Wave Field architecture, which utilizes traveling wave lifting operations to learn both drift and memory closure end-to-end. Empirically, we validated our theoretical observations about the expressivity of the architecture, and showed that it reliably discovers coherent memory structures, achieves minimal latent representations and outperforms baselines on long-range sequence tasks.

**Limitations and Future Work** While our Neural Wave Field provides a clear proof of concept, it is only one instantiation of a much richer framework to be explored in future works. In particular, our preliminary insights into EEG and ECoG datasets warrant further exploration of oscillatory models as lifting mechanisms. We made two assumptions regarding continuity and support of the measure  $\mu_t$  in our framework. Empirically the first assumption stabilizes training as shown in the copy task of Section 4. The second assumption on the differentiability of  $\mu_t$  may not always be assumed, e.g. for the ordered recall task a variant of the copy task in which numbers are recalled in order. A framework that handles discontinuous  $\mu_t$  is non-trivial and would provide additional insights into higher level cognitive capabilities and we leave this to future work.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

---

## REFERENCES

- Andrea Alamia, Antoine Grimaldi, Frederic Chavane, and Martin Vinck. What do neural travelling waves tell us about information flow?, February 2025.
- Kai Keng Ang, Zheng Yang Chin, Haihong Zhang, and Cuntai Guan. Filter Bank Common Spatial Pattern (FBCSP) in Brain-Computer Interface. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pp. 2390–2397, June 2008. doi: 10.1109/IJCNN.2008.4634130.
- Martin Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1120–1128, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/arjovsky16.html>.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- 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.
- P.C. Bressloff. *Waves in Neural Media: From Single Neurons to Neural Fields*. Lecture Notes on Mathematical Modelling in the Life Sciences. Springer New York, 2013. ISBN 978-1-4614-8866-8.
- Nicolas Brunel. Dynamics of sparsely connected networks of excitatory and inhibitory spiking neurons. *Journal of computational neuroscience*, 8(3):183–208, 2000.
- Steven L. Brunton, Bingni W. Brunton, Joshua L. Proctor, Eurika Kaiser, and J. Nathan Kutz. Chaos as an intermittently forced linear system. *Nature Communications*, 8(1):19, May 2017. ISSN 2041-1723. doi: 10.1038/s41467-017-00030-8.
- Ricardo Buitrago, Tanya Marwah, Albert Gu, and Andrej Risteski. On the benefits of memory for modeling time-dependent PDEs. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=o9kqa5K3tB>.
- Herbert B. Callen and Theodore A. Welton. Irreversibility and Generalized Noise. *Physical Review*, 83(1):34–40, 1951. doi: 10.1103/PhysRev.83.34.
- Rishidev Chaudhuri, Bülent Gerçek, Biraj Pandey, Adrien Peyrache, and Ila Fiete. The intrinsic attractor manifold and population dynamics of a canonical cognitive circuit across waking and sleep. *Nature Neuroscience*, 22(9):1512–1520, 2019. doi: 10.1038/s41593-019-0460-x.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- Alexandre J. Chorin, Ole H. Hald, and Raz Kupferman. Optimal prediction with memory. *Physica D: Nonlinear Phenomena*, 166(3-4):239–257, June 2002. ISSN 01672789. doi: 10.1016/S0167-2789(02)00446-3.
- Norman. Coddington, Earl A.; Levinson. *Statistical Mechanics of Nonequilibrium Liquids*. McGraw-Hill Book Company, New York, 1955. ISBN 9780070992566.
- S. Coombes. Waves, bumps, and patterns in neural field theories. *Biological cybernetics*, 93(2): 91–108, August 2005a. ISSN 0340-1200. doi: 10.1007/s00422-005-0574-y.
- Stephen Coombes. Waves, bumps, and patterns in neural field theories. *Biological cybernetics*, 93(2):91–108, 2005b.
- Zachary W. Davis, Lyle Muller, Julio Martinez-Trujillo, Terrence Sejnowski, and John H. Reynolds. Spontaneous travelling cortical waves gate perception in behaving primates. *Nature*, 587(7834): 432–436, November 2020. ISSN 1476-4687. doi: 10.1038/s41586-020-2802-y.

---

540 Gustavo Deco, Edmund T. Rolls, and Ranulfo Romo. Stochastic dynamics as a principle of brain  
541 function. *Progress in Neurobiology*, 88(1):1–16, 2009. doi: 10.1016/j.pneurobio.2009.01.006.  
542

543 R. J. Douglas, C. Koch, M. Mahowald, K. A. Martin, and H. H. Suarez. Recurrent excitation in  
544 neocortical circuits. *Science (New York, N.Y.)*, 269(5226):981–985, August 1995. ISSN 0036-  
545 8075. doi: 10.1126/science.7638624.

546 Andreas K. Engel, Pascal Fries, and Wolf Singer. Dynamic predictions: Oscillations and synchrony  
547 in top–down processing. *Nature Reviews Neuroscience*, 2(10):704–716, October 2001. ISSN  
548 1471-0048. doi: 10.1038/35094565.

549 Tatiana A Engel and Nicholas A Steinmetz. New perspectives on dimensionality and variability  
550 from large-scale cortical dynamics. *Current Opinion in Neurobiology*, 58:181–190, October 2019.  
551 ISSN 0959-4388. doi: 10.1016/j.conb.2019.09.003.  
552

553 Bard Ermentrout. Neural networks as spatio-temporal pattern-forming systems. *Reports on Progress  
554 in Physics*, 61(4):353–430, April 1998. ISSN 0034-4885, 1361-6633. doi: 10.1088/0034-4885/  
555 61/4/002.

556 M. B. Feller. Spontaneous correlated activity in developing neural circuits. *Neuron*, 22(4):653–656,  
557 April 1999. ISSN 0896-6273. doi: 10.1016/s0896-6273(00)80724-2.  
558

559 Karl J. Friston. Transients, Metastability, and Neuronal Dynamics. *NeuroImage*, 5(2):164–171,  
560 February 1997. ISSN 1053-8119. doi: 10.1006/nimg.1997.0259.

561 Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, and Christopher Re.  
562 Hungry hungry hippos: Towards language modeling with state space models. In *The Eleventh  
563 International Conference on Learning Representations*, 2023. URL [https://openreview.  
564 net/forum?id=COZDy0WYGg](https://openreview.net/forum?id=COZDy0WYGg).  
565

566 Dror Givon, Raz Kupferman, and Andrew Stuart. Extracting macroscopic dynamics: Model prob-  
567 lems and algorithms. *Nonlinearity*, 17(6):R55–R127, November 2004. ISSN 0951-7715, 1361-  
568 6544. doi: 10.1088/0951-7715/17/6/R01.

569 Hermann Grabert. *Projection operator techniques in nonequilibrium statistical mechanics*, vol-  
570 ume 95. Springer, 2006.  
571

572 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint  
573 arXiv:1410.5401*, 2014.

574 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *First  
575 Conference on Language Modeling*, 2024. URL [https://openreview.net/forum?id=  
576 tEYskw1VY2](https://openreview.net/forum?id=tEYskw1VY2).

577 Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured  
578 state spaces. In *International Conference on Learning Representations*, 2022. URL [https://  
579 openreview.net/forum?id=uYLFoz1v1AC](https://openreview.net/forum?id=uYLFoz1v1AC).  
580

581 Varun Gumma, Pranjal A Chitale, and Kalika Bali. On the interchangeability of positional embed-  
582 dings in multilingual neural machine translation models. *arXiv e-prints*, pp. arXiv–2408, 2024.  
583

584 Priyam Gupta, Peter J. Schmid, Denis Sipp, Taraneh Sayadi, and Georgios Rigas. Mori-zwanzig  
585 latent space koopman closure for nonlinear autoencoder, 2024. URL [https://arxiv.org/  
586 abs/2310.10745](https://arxiv.org/abs/2310.10745).

587 Benjamin J A Héry and Roland R Netz. Derivation of a generalized Langevin equation from a  
588 generic time-dependent Hamiltonian. *Journal of Physics A: Mathematical and Theoretical*, 57  
589 (50):505003, November 2024. ISSN 1751-8121. doi: 10.1088/1751-8121/ad91ff.

590 J J Hopfield. Neural networks and physical systems with emergent collective computational abilities.  
591 *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, April 1982. doi: 10.1073/  
592 pnas.79.8.2554.  
593

Herbert Jaeger. Echo state network. *scholarpedia*, 2(9):2330, 2007.

- 
- 594 Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and eran malach. Repeat after me: Trans-  
595 formers are better than state space models at copying. In *Forty-first International Conference on*  
596 *Machine Learning*, 2024. URL <https://openreview.net/forum?id=duRRoGeoQT>.  
597
- 598 Li Jing, Caglar Gulcehre, John Peurifoy, Yichen Shen, Max Tegmark, Marin Soljagic, and Yoshua  
599 Bengio. Gated orthogonal recurrent units: On learning to forget. *Neural computation*, 31(4):  
600 765–783, 2019.
- 601 Arjun Karuvally, Terrence J. Sejnowski, and Hava T. Siegelmann. Hidden traveling waves bind  
602 working memory variables in recurrent neural networks. In *Proceedings of the 41st International*  
603 *Conference on Machine Learning*, ICML’24. JMLR.org, 2024.  
604
- 605 Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva  
606 Reddy. The impact of positional encoding on length generalization in transformers. *Advances*  
607 *in Neural Information Processing Systems*, 36:24892–24928, 2023.
- 608 T Anderson Keller. Nu-wave state space models: Traveling waves as a bio-  
609 logically plausible context. *Science Communications Worldwide*, 2025. doi:  
610 10.57736/b30b-8eed. URL [https://www.world-wide.org/cosyne-25/  
611 nu-wave-state-space-models-traveling-3803805f](https://www.world-wide.org/cosyne-25/nu-wave-state-space-models-traveling-3803805f).  
612
- 613 T. Anderson Keller and Max Welling. Neural wave machines: Learning spatiotemporally struc-  
614 tured representations with locally coupled oscillatory recurrent neural networks. In Andreas  
615 Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scar-  
616 lett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202  
617 of *Proceedings of Machine Learning Research*, pp. 16168–16189. PMLR, 23–29 Jul 2023. URL  
618 <https://proceedings.mlr.press/v202/keller23a.html>.
- 619 T. Anderson Keller, Lyle Muller, Terrence Sejnowski, and Max Welling. Traveling waves encode the  
620 recent past and enhance sequence learning. In *The Twelfth International Conference on Learning*  
621 *Representations*, 2024. URL <https://openreview.net/forum?id=p4S5Z6Sah4>.  
622
- 623 J.A.S. Kelso. Multistability and metastability: Understanding dynamic coordination in the brain.  
624 *Philosophical Transactions of the Royal Society B: Biological Sciences*, 367(1591):906–918,  
625 2012. ISSN 0962-8436. doi: 10.1098/rstb.2011.0351.
- 626 Demetres Kostas and Frank Rudzicz. Thinker invariance: enabling deep neural networks for bci  
627 across more people. *Journal of Neural Engineering*, 17(5):056008, 2020.  
628
- 629 Carlo R. Laing and Carson C. Chow. A Spiking Neuron Model for Binocular Rivalry. *Journal of*  
630 *Computational Neuroscience*, 12(1):39–53, January 2002. ISSN 1573-6873. doi: 10.1023/A:  
631 1014942129705.
- 632 Samuel Lanthaler, T. Konstantin Rusch, and Siddhartha Mishra. Neural oscillators are universal.  
633 In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=QGQsOZcQ2H>.  
634
- 635 Mark L. Latash, Mindy F. Levin, John P. Scholz, and Gregor Schöner. Motor control theories and  
636 their applications. *Medicina (Kaunas, Lithuania)*, 46(6):382–392, 2010. ISSN 1648-9144.  
637
- 638 Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and  
639 Brent J Lance. Eegnet: a compact convolutional neural network for eeg-based brain–computer  
640 interfaces. *Journal of neural engineering*, 15(5):056013, 2018.
- 641 Albert K. Lee and Matthew A. Wilson. Memory of sequential experience in the hippocampus during  
642 slow wave sleep. *Neuron*, 36(6):1183–1194, December 2002. ISSN 0896-6273. doi: 10.1016/  
643 s0896-6273(02)01096-6.  
644
- 645 Xiaoxuan Lei, Takuya Ito, and Pouya Bashivan. Geometry of naturalistic object representations in  
646 recurrent neural network models of working memory. In *The Thirty-eighth Annual Conference on*  
647 *Neural Information Processing Systems*, 2024. URL [https://openreview.net/forum?  
id=N2RaC7LO6k](https://openreview.net/forum?id=N2RaC7LO6k).

---

648 Luisa H. B. Liboni, Roberto C. Budzinski, Alexandra N. Busch, Sindy Löwe, Thomas A. Keller,  
649 Max Welling, and Lyle E. Muller. Image segmentation with traveling waves in an exactly  
650 solvable recurrent neural network. *Proceedings of the National Academy of Sciences*, 122(1):  
651 e2321319121, 2025. doi: 10.1073/pnas.2321319121.

652 Yen Ting Lin, Yifeng Tian, Daniel Livescu, and Marian Anghel. Data-driven learning for the mori–  
653 zwanzig formalism: A generalization of the koopman learning framework. *SIAM Journal on*  
654 *Applied Dynamical Systems*, 20(4):2558–2601, 2021. doi: 10.1137/21M1401759.

655 Yen Ting Lin, Yifeng Tian, Danny Perez, and Daniel Livescu. Regression-based projection for  
656 learning mori–zwanzig operators. *SIAM Journal on Applied Dynamical Systems*, 22(4):2890–  
657 2926, 2023.

658 Chenghao Liu, Shuncheng Jia, Hongxing Liu, Xuanle Zhao, Chengyu T. Li, Bo Xu, and Tielin  
659 Zhang. Recurrent neural networks with transient trajectory explain working memory encoding  
660 mechanisms. *Communications Biology*, 8(1):1–13, January 2025. ISSN 2399-3642. doi: 10.  
661 1038/s42003-024-07282-3.

662 Mantas Lukoševičius and Herbert Jaeger. Reservoir computing approaches to recurrent neural net-  
663 work training. *Computer science review*, 3(3):127–149, 2009.

664 Eric Maris, Thilo Womelsdorf, Robert Desimone, and Pascal Fries. Rhythmic neuronal synchron-  
665 ization in visual cortex entails spatial phase relation diversity that is modulated by stimulation  
666 and attention. *NeuroImage*, 74:99–116, July 2013. ISSN 1095-9572. doi: 10.1016/j.neuroimage.  
667 2013.02.007.

668 Natasza Marrouch, Joanna Slawinska, Dimitrios Giannakis, and Heather L. Read. Data-driven  
669 Koopman operator approach for computational neuroscience. *Annals of Mathematics and Ar-  
670 tificial Intelligence*, 88(11):1155–1173, December 2020. ISSN 1573-7470. doi: 10.1007/  
671 s10472-019-09666-2.

672 Marcello Massimini, Reto Huber, Fabio Ferrarelli, Sean Hill, and Giulio Tononi. The sleep  
673 slow oscillation as a traveling wave. *The Journal of Neuroscience: The Official Journal*  
674 *of the Society for Neuroscience*, 24(31):6862–6870, August 2004. ISSN 1529-2401. doi:  
675 10.1523/JNEUROSCI.1318-04.2004.

676 Romit Maulik, Arvind Mohan, Bethany Lusch, Sandeep Madireddy, Prasanna Balaprakash, and  
677 Daniel Livescu. Time-series learning of latent-space dynamics for reduced-order model closure.  
678 *Physica D: Nonlinear Phenomena*, 405:132368, 2020.

679 Emmanuel Menier, Sebastian Kaltenbach, Mouadh Yagoubi, Marc Schoenauer, and Petros  
680 Koumoutsakos. Interpretable learning of effective dynamics for multiscale systems. *CoRR*,  
681 abs/2309.05812, 2023. URL <https://doi.org/10.48550/arXiv.2309.05812>.

682 Hugues Meyer, Thomas Voigtmann, and Tanja Schilling. On the non-stationary generalized langevin  
683 equation. *The Journal of chemical physics*, 147(21), 2017.

684 Hugues Meyer, Thomas Voigtmann, and Tanja Schilling. On the dynamics of reaction coordinates  
685 in classical, time-dependent, many-body processes. *Journal of Chemical Physics*, 150:174118,  
686 May 2019. ISSN 0021-9606.

687 Hazime Mori. Transport, collective motion, and brownian motion. *Progress of theoretical physics*,  
688 33(3):423–455, 1965a.

689 Hazime Mori. Transport, Collective Motion, and Brownian Motion\*). *Progress of Theoretical*  
690 *Physics*, 33(3):423–455, March 1965b. ISSN 0033-068X. doi: 10.1143/PTP.33.423.

691 Lyle Muller, Frédéric Chavane, John Reynolds, and Terrence J. Sejnowski. Cortical travelling waves:  
692 Mechanisms and computational principles. 19(5):255–268. ISSN 1471-0048. doi: 10.1038/nrn.  
693 2018.20. URL <https://www.nature.com/articles/nrn.2018.20>.

694 Lyle Muller, Alexandre Reynaud, Frédéric Chavane, and Alain Destexhe. The stimulus-evoked  
695 population response in visual cortex of awake monkey is a propagating wave. *Nature Communi-  
696 cations*, 5(1):3675, April 2014. ISSN 2041-1723. doi: 10.1038/ncomms4675.

---

702 Sadao Nakajima. On Quantum Theory of Transport Phenomena: Steady Diffusion. *Progress of*  
703 *Theoretical Physics*, 20(6):948–959, December 1958. ISSN 0033-068X. doi: 10.1143/PTP.20.  
704 948.

705 Roland R. Netz. Derivation of the nonequilibrium generalized langevin equation from a time-  
706 dependent many-body hamiltonian. *Phys. Rev. E*, 110:014123, Jul 2024. doi: 10.1103/  
707 PhysRevE.110.014123. URL [https://link.aps.org/doi/10.1103/PhysRevE.](https://link.aps.org/doi/10.1103/PhysRevE.110.014123)  
708 [110.014123](https://link.aps.org/doi/10.1103/PhysRevE.110.014123).

709 Mitchell Ostrow, Adam Eisen, and Ila Fiete. Delay embedding theory of neural sequence models,  
710 2024. URL <https://arxiv.org/abs/2406.11993>.

711 Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural  
712 networks. In *International conference on machine learning*, pp. 1310–1318. Pmlr, 2013.

713 Ofir Press, Noah Smith, and Mike Lewis. Train Short, Test Long: Attention with Linear Biases  
714 Enables Input Length Extrapolation. In *International Conference on Learning Representations*,  
715 October 2021.

716 Kanaka Rajan, Christopher D. Harvey, and David W. Tank. Recurrent Network Models of Sequence  
717 Generation and Memory. *Neuron*, 90(1):128–142, April 2016. ISSN 0896-6273. doi: 10.1016/j.  
718 neuron.2016.02.009.

719 Yulia Rubanova, Ricky TQ Chen, and David K Duvenaud. Latent ordinary differential equations for  
720 irregularly-sampled time series. *Advances in neural information processing systems*, 32, 2019.

721 Doug Rubino, Kay A. Robbins, and Nicholas G. Hatsopoulos. Propagating waves mediate informa-  
722 tion transfer in the motor cortex. *Nature Neuroscience*, 9(12):1549–1557, December 2006. ISSN  
723 1546-1726. doi: 10.1038/nn1802.

724 Adam Rupe and James P. Crutchfield. On principles of emergent organization. *Physics Reports*,  
725 1071:1–47, 2024. ISSN 0370-1573. doi: 10.1016/j.physrep.2024.04.001.

726 T. Konstantin Rusch and Siddhartha Mishra. Coupled oscillatory recurrent neural network  
727 (co{rnn}): An accurate and (gradient) stable architecture for learning long time dependencies. In  
728 *International Conference on Learning Representations*, 2021. URL [https://openreview.](https://openreview.net/forum?id=F3s69XzWOia)  
729 [net/forum?id=F3s69XzWOia](https://openreview.net/forum?id=F3s69XzWOia).

730 T. Konstantin Rusch and Daniela Rus. Oscillatory state-space models. In *The Thirteenth Interna-*  
731 *tional Conference on Learning Representations*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=GRMfXcAAfH)  
732 [forum?id=GRMfXcAAfH](https://openreview.net/forum?id=GRMfXcAAfH).

733 Ábel Ságodi, Guillermo Martín-Sánchez, Piotr A Sokol, and Il Memming Park. Back to the continu-  
734 ous attractor. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,  
735 2024. URL <https://openreview.net/forum?id=fvG6ZHrH0B>.

736 Robin Tibor Schirrmester, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin  
737 Glasstetter, Katharina Eggenberger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and  
738 Tonio Ball. Deep learning with convolutional neural networks for eeg decoding and visualization.  
739 *Human brain mapping*, 38(11):5391–5420, 2017.

740 Panagiotis Stinis. Higher Order Mori–Zwanzig Models for the Euler Equations. *Multiscale Model-*  
741 *ing & Simulation*, 6(3):741–760, January 2007. ISSN 1540-3459. doi: 10.1137/06066504X.

742 Panagiotis Stinis. Mori-Zwanzig reduced models for uncertainty quantification I: Parametric uncer-  
743 tainty, November 2012.

744 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. RoFormer: En-  
745 hanced transformer with Rotary Position Embedding. *Neurocomputing*, 568:127063, February  
746 2024. ISSN 0925-2312. doi: 10.1016/j.neucom.2023.127063.

747 Kazutaka Takahashi, Maryam Saleh, Richard D. Penn, and Nicholas G. Hatsopoulos. Propagating  
748 waves in human motor cortex. *Frontiers in Human Neuroscience*, 5:40, 2011. ISSN 1662-5161.  
749 doi: 10.3389/fnhum.2011.00040.

---

756 Yifeng Tian, Yen Ting Lin, Marian Anghel, and Daniel Livescu. Data-driven learning of mori-  
757 zwanzig operators for isotropic turbulence. *Physics of Fluids*, 33(12), 2021.  
758

759 Rory G. Townsend, Selina S. Solomon, Spencer C. Chen, Alexander N. J. Pietersen, Paul R. Martin,  
760 Samuel G. Solomon, and Pulin Gong. Emergence of complex wave patterns in primate cerebral  
761 cortex. *The Journal of Neuroscience: The Official Journal of the Society for Neuroscience*, 35  
762 (11):4657–4662, March 2015. ISSN 1529-2401. doi: 10.1523/JNEUROSCI.4509-14.2015.

763 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
764 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-  
765 tion processing systems*, 30, 2017.

766 Daniele Venturi and Xiantao Li. The Mori–Zwanzig formulation of deep learning. *Research  
767 in the Mathematical Sciences*, 10(2):23, May 2023. ISSN 2197-9847. doi: 10.1007/  
768 s40687-023-00390-2.

770 Xiao-Jing Wang. Synaptic Basis of Cortical Persistent Activity: The Importance of NMDA Re-  
771 ceptors to Working Memory. *The Journal of Neuroscience*, 19(21):9587–9603, November 1999.  
772 ISSN 0270-6474. doi: 10.1523/JNEUROSCI.19-21-09587.1999.

773 Klaus Wimmer, Duane Q Nykamp, Christos Constantinidis, and Albert Compte. Bump attractor  
774 dynamics in prefrontal cortex explains behavioral precision in spatial working memory. *Nature  
775 Neuroscience*, 17(3):431–439, March 2014. ISSN 1546-1726. doi: 10.1038/nn.3645.

777 Michael Woodward. *Reduced Lagrangian and Mori-Zwanzig Models: Applications To Turbulent  
778 Flows*. The University of Arizona, 2023.

779 Michael Woodward, Yen Ting Lin, Yifeng Tian, Christoph Hader, Hermann Fasel, and Daniel  
780 Livescu. Mori-Zwanzig mode decomposition: Comparison with time-delay embeddings, May  
781 2025.

782

783 Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico  
784 Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural Fields in Visual Com-  
785 puting and Beyond, April 2022.

786 Qunxi Zhu, Yao Guo, and Wei Lin. Neural delay differential equations. In *International Confer-  
787 ence on Learning Representations*, 2021. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=Q1jmmQz72M2)  
788 [Q1jmmQz72M2](https://openreview.net/forum?id=Q1jmmQz72M2).

789

790 Nicolas Zucchet and Antonio Orvieto. Recurrent neural networks: vanishing and exploding gradi-  
791 ents are not the end of the story. In *The Thirty-eighth Annual Conference on Neural Information  
792 Processing Systems*, 2024. URL <https://openreview.net/forum?id=46Jr4sgTwa>.

793 Robert Zwanzig. Memory Effects in Irreversible Thermodynamics. *Physical Review*, 124(4):983–  
794 992, November 1961. doi: 10.1103/PhysRev.124.983.

795

796 Robert Zwanzig. *Nonequilibrium statistical mechanics*. Oxford university press, 2001.  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

---

## 810 A RELATED WORKS

811  
812 Data-driven MZ (Chorin et al., 2002; Lin et al., 2021) and time-delay embeddings (Brunton et al.,  
813 2017; Woodward et al., 2025; Zhu et al., 2021; Ostrow et al., 2024) fix a static projection (e.g.,  
814 EDMD modes or a stack of delays) to learn stationary memory kernels. Deep learning exten-  
815 sions of these approaches use sequential inputs to learn memory kernels (Lin et al., 2023) in auto-  
816 encoders (Gupta et al., 2024) and neural operators (Buitrago et al., 2025). Time-dependent MZ  
817 formalism has been used to characterized deep learning (Venturi & Li, 2023).

818  
819 **Structured models.** Neural oscillators (Lanthaler et al., 2023; Rusch & Mishra, 2021; Keller &  
820 Welling, 2023), traveling-wave networks (Keller et al., 2024; Liboni et al., 2025; Keller, 2025), at-  
821 tractor embeddings (Ságodi et al., 2024), reservoir computing (Jaeger, 2007; Lukoševičius & Jaeger,  
822 2009) and neural delay-difference equations (Zhu et al., 2021) bake in known dynamical motifs to  
823 encode memory explicitly. WaveRNN (Keller et al., 2024) structures its latent updates as linear  
824 advection, yielding transparent memory dynamics as traveling waves.

825  
826 **Models.** GRUs, LSTMs, residual and deep equilibrium models are well established recurrent and  
827 feed-forward NNs with augmented memory-mechanisms. Continuous-time models such as neural  
828 ODEs (Chen et al., 2018) and ODE-RNNs (Rubanova et al., 2019) encode history through flows  
829 in state space. Recent state-of-the-art performance has been achieved by structured state-space-  
830 models (SSMs) (Gu et al., 2022; Fu et al., 2023; Gu & Dao, 2024), notably the oscillatory SSM-  
831 LinOSS (Rusch & Rus, 2025).

832  
833 **Global context embeddings.** Transformers (Press et al., 2021; Su et al., 2024; Gumma et al.,  
834 2024) use full self-attention for global sequence dependencies. Recent work investigating the per-  
835 formance of positional encodings (Jelassi et al., 2024) has demonstrated that various positional en-  
836 coding strategies (Press et al., 2021; Kazemnejad et al., 2023; Su et al., 2024) outperform SSMs on  
837 copying tasks.

838  
839 **Memory Neural Operator (MemNO)** MemNO (Buitrago et al., 2025) interleaves a memory op-  
840 erator (sequential model) into the layer updates of a neural operator, in order to capture memory  
841 effects in a GLE inspired manner. The goal of the memory operator is to re-introduce projected  
842 variables, and is theoretically motivated by a theorem demonstrating the divergence of solutions  
843 with and without memory for a second-order elliptic PDE. The approach is empirically validated by  
844 testing super resolution capacity of architectures, i.e. reducing the input resolution and maintaining  
845 the output resolution during the training of an encoder-decoder framework. A further ablation study  
846 is performed on the window size of the memory operator, where the performance improves as the  
847 window size increases, i.e. in a time-delay embedding fashion.

848  
849 At a high level, both works aim to resolve a memory closure using sequential linear layers (S4 in the  
850 case of MemNO). MemNO uses a multi-layer FNO as an embedding and read out of the latent state,  
851 whereas NWF uses linear layers based on projection operators. Additionally, NWF directly induces  
852 wave like phenomena into the latent state, and studies the rise of phenomena like coherence and  
853 emergence. An interesting future direction would be to characterize MemNO’s memory operator  
854 using the theory developed here-in.

855  
856 **Time-dependent GLE** This time-dependent relevant ensemble  $\rho(t)$  has been extended to a bundle  
857 of trajectories, i.e. measurements for a distribution of moving points in the phase space (Meyer et al.,  
858 2017). The resolved subspace is changing in time and the two-time memory kernel appears again.  
859 The time-dependent projection operator is an average over all possible trajectories.

860  
861 A discrete analogue of the time-dependent GLE has been proposed in the context of deep residual  
862 neural networks (Venturi & Li, 2023). In this formulation, each layer  $n$  is associated with a pro-  
863 jection operator  $P_n$  and the hidden state evolves with a Markov term, two-time memory kernel and  
864 layer-wise fluctuating force. Although their streaming term does not explicitly include a kinematic  
865 component, it implicitly accounts for the evolution of the projection subspace across layers through  
866 the residual propagator. While (Venturi & Li, 2023) notes that MZ formalism can be used to reduce  
867 the total number of degrees of freedom in the neural network, in practice their approach does not  
868 provide a mechanism by which they may go about reducing the number of variables.

---

## 864 B MOTIVATING EXAMPLE

### 866 B.1 BOUNDARY CONDITIONS

867  
868 **Corollary B.1.** (*Toroidal Latent Manifold*) Suppose we constrain each latent coordinate  $h_i(t)$  to live  
869 on a circle of period  $L_i$  and we enforce that both the learned drift and memory-kernel parameters  
870 depend on  $h$  only through these periodic coordinates. Then the entire latent trajectory  $h(t)$  evolves  
871 on the  $m$ -dimensional torus  $\mathbb{T}^m$ . As a result, the network can only represent—and learn—functions  
872 defined on this compact, boundary-free manifold.

873  
874 *Proof.* By the assumption of periodicity then each of the MZ terms descent to well-defined maps on  
875 the quotient  $\mathbb{R}^m / (L_1\mathbb{Z} \times \dots \times L_m\mathbb{Z})$ , and the initial condition  $h(0) \in S_{L_1}^1 \times \dots \times S_{L_m}^m$  uniquely  
876 determines a solution  $h(t)$  that never leaves the torus.

877 Therefore any decoder  $F : \mathbb{R}^m \rightarrow Y$  must descent to a well-defined map  $\hat{F} : \mathbb{T}^m \rightarrow Y$ , i.e., those  
878 maps that are periodic in each coordinate.  $\square$   
879

## 881 C MORI-ZWANZIG FORMALISM

882  
883 In this section, we present the preliminary background. We treat the latent state of the neural network  
884 as observations of an underlying dynamical system. The near-equilibrium MZ formalism (NE-MZ)  
885 describes the evolution of a time-invariant subset of observations. Time invariance may be overly  
886 restrictive for the latent states of a neural network, in which case we employ time-dependent operator  
887 formalism for far-from-equilibrium systems (FFE-MZ). Finally, we recall the important distinction  
888 of MZ-type memory, the generalized fluctuation-dissipation relation (GFDR).

### 889 C.1 PRELIMINARIES

890 Suppose the underlying system evolves dynamically on a smooth manifold  $\mathcal{M} \subset \mathbb{R}^n$ , called the  
891 phase-state, described by the following (ergodic and possibly nonlinear) autonomous ODE  
892

$$893 \frac{d\Phi(t)}{dt} = S(\Phi(t)), \quad \Phi(0) = x_0, \quad (4)$$

894 where  $S : \mathcal{M} \rightarrow \mathbb{R}^n$  is  $C^1$ . By the Picard-Lindelöf (Coddington, 1955) theorem, Equation 4 admits  
895 a unique solution  $\Phi_t(x_0) = \Phi(t)$  for all  $t$  in  $\mathbb{T} \subseteq \mathbb{R}$ , inducing the flow  $\Phi_t : \mathcal{M} \rightarrow \mathcal{M}$ .

896 Let the collection  $(\mathcal{M}, \mathcal{F}, \mu)$  be the phase-state manifold  $\mathcal{M}$  equipped with a  $\sigma$ -algebra  $\mathcal{F}$  and a  
897 finite, flow-invariant probability measure  $\mu$ . A system *observation*  $g : \mathcal{M} \rightarrow \mathbb{R}$  is a real-valued  
898 square-integrable function, i.e.  $g \in \mathcal{H} := L^2(\mathcal{M}, \mu)$  where  $\mathcal{H}$  is a *separable* Hilbert space.

899 **Definition C.1.** (*Liouville Operator*) The Liouville operator  $\mathcal{L} : \mathcal{H} \rightarrow \mathcal{H}$  describes the infinitesimal  
900 evolution of an observable  $g \in \mathcal{H}$  along the flow  $\Phi_t$ . In general we will take it to be  $\frac{d}{dt}g(t) = \mathcal{L}g(t)$ .

901 Remarkably, the evolution of the observations can be expressed in terms of linear operators on  $\mathcal{H}$ ,  
902 despite the underlying system being possibly nonlinear and mildly complex (ergodic). However,  
903 this is a linear operator that acts on the space of all observables, which may be infinite dimensional.

### 904 C.2 NEAR-EQUILIBRIUM MORI-ZWANZIG FORMALISM (NE-MZ)

905 Using the separability of  $\mathcal{H}$ , the space of observations can be separated into a set of *resolved* ob-  
906 servables and complementary *unresolved* observables. In particular, for any closed subspace  $\mathcal{V} \subset \mathcal{H}$   
907 there is a decomposition  $\mathcal{H} = \mathcal{V} \oplus \mathcal{V}^\perp$  realized by the unique orthogonal projection  
908

$$909 P : \mathcal{H} \rightarrow \mathcal{V}, \quad Q = I - P : \mathcal{H} \rightarrow \mathcal{V}^\perp \quad (P^2 = P, Q^2 = Q, P = P^*, Q = Q^*, PQ = 0).$$

910 These projections can be linear operators (Mori, 1965a), or as we adopt, (non)-linear opera-  
911 tors (Zwanzig, 2001) realized as conditional expectations  $P = \mathbb{E}[\cdot | \mathcal{G}]$  on the sub- $\sigma$ -algebra  $\mathcal{G} \subset \mathcal{F}$ .  
912 In NE-ZM formalism the subset of observables—and therefore the projection  $P$ —is *time-invariant*.  
913

**The generalized Langevin equation.** NE-MZ formalism describes the exact evolution of a time-invariant subset of observables by decomposing  $\mathcal{H}$  into *resolved*  $\hat{g} \in \mathcal{V}$  and *unresolved*  $\tilde{g} \in \mathcal{V}^\top$  observables. The result is the generalized Langevin equation (GLE)

$$\frac{\partial}{\partial t} \hat{g}(t) = \underbrace{P\mathcal{L}\hat{g}(t)}_{\text{Markov}} + \underbrace{\int_0^t P\mathcal{L}e^{(t-s)Q\mathcal{L}}Q\mathcal{L}\hat{g}(s)ds}_{\text{Memory}} + \underbrace{P\mathcal{L}e^{tQ\mathcal{L}}Qg(0)}_{\text{Fluctuating Force}}. \quad (2)$$

Equation 1 consists of three distinct terms (underscored). The Markov term represents the instantaneous drift from the resolved dynamics. The Memory term re-introduces the influence of dynamics previously *forgotten*, i.e. prior resolved information that has been projected into the unresolved subspace. The Fluctuating Force term<sup>2</sup> captures the residual influence of the unresolved initial state.

### C.3 FAR-FROM-EQUILIBRIUM MORI-ZWANZIG FORMALISM (FFE-MZ)

For a neural network architecture, it may not be possible—and potentially unreasonable—to ascribe to each element of its latent state a static representation. Our approach models the latent state by using FFE-MZ (Grabert, 2006) which allows  $P(t)$  to evolve.

The resulting GLE is given by

$$\frac{\partial}{\partial t} \hat{g}(t) = P(t)\mathcal{L}\hat{g}(t) + \dot{P}(t)g(t) + \int_0^t P(t)\mathcal{L}G(t,s)Q(s)\mathcal{L}\hat{g}(s)ds + P(t)\mathcal{L}G(t,0)Q(0)g(0). \quad (5)$$

The two-time memory kernel  $G(t,s) = \mathcal{T}_- \exp\left(\int_s^t Q(u)\mathcal{L}du\right)$  is the negatively time-ordered exponential operator that captures the extrinsic influence from the evolution of the subspaces. The Kinematic term  $\dot{P}(t)g(t)$  captures the intrinsic evolution of the resolved subspace (Meyer et al., 2017).

Critically, the time-dependent projection operator acts as moving frame of reference that is tied to the relevant ensemble. The source of the time dependence is *extrinsic* to the resolved observables i.e. it is driven. As a result, this non-stationarity cannot be removed by a simple change of coordinates.

### C.4 GENERALIZED FLUCTUATION DISSIPATION RELATION (GFDR)

We now observe the critical distinction between MZ memory and auto-regressive or time-delay mechanisms, that MZ assumes an underlying principle of detailed balance. The principle of detailed balance states that at equilibrium, each process is in equilibrium with its reverse process. For NE-MZ this is formalized via the fluctuation-dissipation theorem (Callen & Welton, 1951) directly. The *generalized* fluctuation-dissipation relation (GFDR) is the extension to FFE-MZ (Meyer et al., 2019)

$$K(t,s) = \langle F(t|s), F(s) \rangle C(s)^{-1} \quad (6)$$

$K(t,s) = P(t)\mathcal{L}G(t,s)Q(s)\mathcal{L}$ ,  $C(s) = \langle \hat{g}(s), \hat{g}(s) \rangle$ ,  $F(s) = Q(s)\mathcal{L}\hat{g}(s)$ ,  $F(t|s) = G(t,s)F(s)$  which relates the memory kernel  $K(t,s)$  to the level of noise  $\langle F(t|s), F(s) \rangle$  relative to the covariance of the resolved observable  $C$ . Instead of treating noise in the black-box model as a limitation of explainability, MZ formalism allows us to model noise predictably from the memory kernel itself.

For more details on the distinction of architectures, we refer the reader to (Lin et al., 2023).

### C.5 INFORMAL DERIVATIONS

**The derivation of the GLE.** The instantaneous evolution of  $g$  is given by

$$\frac{d}{dt} e^{t\mathcal{L}}g(0) = \mathcal{L}e^{t\mathcal{L}}g(0),$$

which can be decomposed into its two projected dynamics yielding two coupled equations

$$\begin{aligned} \frac{d}{dt} P e^{t\mathcal{L}}g(0) &= P\mathcal{L}P e^{t\mathcal{L}}g(0) + P\mathcal{L}Q e^{t\mathcal{L}}g(0), \\ \frac{d}{dt} Q e^{t\mathcal{L}}g(0) &= Q\mathcal{L}Q e^{t\mathcal{L}}g(0) + Q\mathcal{L}P e^{t\mathcal{L}}g(0). \end{aligned}$$

<sup>2</sup>The third term referred to by (Mori, 1965a) as a random force and by (Zwanzig, 2001) as a fluctuating force, is frequently called the noise term in data-driven and stochastic applications.

We rewrite the second equation for  $v(t) = Qe^{t\mathcal{L}}g(0)$  where  $A(t) = Qe^{t\mathcal{L}}g(0)$  and  $F(t) = Q\mathcal{L}Pe^{t\mathcal{L}}g(0)$ ,

$$\frac{d}{dt}v(t) = A(t)v(t) + F(t).$$

The solution is given by Dyson's identity

$$v(t) = e^{tA}v(0) + \int_0^t e^{(t-s)A}F(s)ds.$$

Notice that  $v(0) = Qg(0)$ . Substituting for  $v, A, F$ , we have

$$Qe^{t\mathcal{L}}g(0) = e^{tQ\mathcal{L}}Qg(0) + \int_0^t e^{(t-s)Q\mathcal{L}}Pe^{s\mathcal{L}}g(0)ds = e^{tQ\mathcal{L}}g(0) + \int_0^t e^{(t-s)Q\mathcal{L}}Pg(s)ds.$$

The GLE results from substituting the prior result into the dynamics for  $\frac{d}{dt}Pg(t)$

$$\frac{\partial}{\partial t}Pg(t) = P\mathcal{L}Pg(t) + \int_0^t P\mathcal{L}e^{(t-s)Q\mathcal{L}}Q\mathcal{L}Pg(s)ds + P\mathcal{L}e^{tQ\mathcal{L}}Qg(0).$$

**The connection to Koopman operator theory.** The Koopman operator  $\mathcal{K}^t : \mathcal{H} \rightarrow \mathcal{H}$  is a bounded linear operator that evolves any observable  $g \in \mathcal{H}$  along the flow  $T \subset \mathbb{R}$  on the phase manifold

$$\mathcal{K}^t g(x_0) = g(T(x_0, t)).$$

Because  $\mathcal{H}$  is infinite dimensional, in practice one often restricts attention to a finite resolved subspace  $\mathcal{V} = \text{Span}\{g^{(1)}, \dots, g^{(r)}\} \subset \mathcal{H}$  with orthogonal complement  $\mathcal{V}^\top$ .

The evolution of  $\hat{g} \in \mathcal{V}$  in this reduced subspace, with restricted evolution operator  $\hat{\mathcal{K}}$ , accumulates an error term

$$\hat{g} \circ T = \hat{\mathcal{K}}^t \hat{g} + r, \\ r \in \mathcal{V}^\top.$$

where  $\hat{g} \in \mathcal{V}$ . The residual  $r$  is the closure problem, which is addressed via the Mori–Zwanzig formalism by projecting onto  $\mathcal{V}$  while accounting for the influence of  $\mathcal{V}^\top$ .

## D THEORETICAL DETAILS

In this section, we provide proofs of the corresponding propositions from Section 3.

### D.1 ASSUMPTIONS

For completeness, we restate our assumptions below. In addition, we will provide some more context to the significance of these assumptions.

**Assumption 3.2.** (*Differentiability of  $P_{\mu_t}$* ) Suppose the time-dependent conditional expectation operator  $P_{\mu_t} : L^2(\mu_*) \rightarrow L^2(\mu_t)$  is Fréchet-differentiable with derivative  $\dot{P}_{\mu_t}$ .

This assumption is critical to ensuring that the GLE is well-posed. In practice it forces us to choose a feature-map basis whose dependence on  $t$  makes  $P_{\mu_t}$  a smooth function of time—only then can the model reliably learn the evolving dynamics.

**Assumption D.1.** (*Support Coverage Assumption*) Let  $\tilde{\mu} = \sum_{t=1}^T \mu_t$ . We require  $\mu_* \ll \tilde{\mu}$  or equivalently  $\text{supp}(\mu_*) \subseteq \bigcup_{t=1}^T \text{supp}(\mu_t)$ .

Assumption D.1 ensures that every region with positive mass under  $\mu_*$  is observed at some time  $t$ , so that all potential degrees of freedom in the reference measure are, in principle, observable. This condition ensures that the projected dynamics  $P_{\mu_t}$  can act on the entire latent state: there are no hidden modes in  $\mu_*$  that fall completely outside the supports of the training measures. Equivalently, it removes any degrees of freedom from the latent state, so that our GLE truly governs all of the relevant latent dynamics.

D.2 PROOFS

**Proposition 3.1.** (*Intrinsic Time-Dependent GLE*) Let  $g(t)$  evolve under the operator  $\mathcal{L}$  on a fixed Hilbert space  $\mathcal{H} = L^2(\mathcal{M}, \mathcal{F}, \mu_*)$ . Let  $P_{\mu_*} : \mathcal{H} \rightarrow \mathcal{V} \subset \mathcal{H}$  be an orthogonal projection onto  $\mathcal{V} = L^2(\mathcal{M}, \mathcal{G}, \mu_*)$  with  $\mathcal{G} \subset \mathcal{F}$ . For a family of  $C^1$  measures  $\{\mu_t\}_{t \in [0, T]}$  let  $P_{\mu_t} : \mathcal{V} \rightarrow \mathcal{V}_t$  be the corresponding family of projections defining a Hilbert bundle  $\{\mathcal{V}_t\}_{t \in [0, t]}$  with  $\mathcal{V}_t = L^2(\mathcal{M}, \mathcal{G}, \mu_t)$ . The evolution of the resolved variable  $P_{\mu_t}g(t)$  satisfies the following GLE

$$\begin{aligned} \frac{d}{dt}(P_{\mu_t}P_{\mu_*}g(t)) &= P_{\mu_t}\dot{P}_{\mu_t}Q_{\mu_t}P_{\mu_*}g(t) + P_{\mu_t}\mathcal{L}P_{\mu_*}g(t) \\ &+ \int_0^t P_{\mu_t}P_{\mu_*}\mathcal{L}e^{(t-s)Q_{\mu_*}}\mathcal{L}P_{\mu_*}g(s)ds + P_{\mu_t}\mathcal{L}e^{tQ_{\mu_*}}\mathcal{L}Q_{\mu_*}g(0). \end{aligned} \quad (3)$$

*Proof.* By Assumption 3.2  $P_{\mu_t}$  is differentiable, so that the GLE is given by chain rule as

$$\frac{d}{dt}(P_{\mu_t}P_{\mu_*}g(t)) = \dot{P}_{\mu_t}P_{\mu_*}g(t) + P_{\mu_t}P_{\mu_*}\frac{d}{dt}g(t) = \dot{P}_{\mu_t}P_{\mu_*}g(t) + P_{\mu_t}P_{\mu_*}\mathcal{L}g(t).$$

Let  $\mathcal{H}$  and  $\mathcal{V}$  be decomposed as  $\mathcal{H} = \mathcal{V} \oplus \mathcal{V}^\perp = \text{ran}(P_{\mu_*}) \oplus \text{ran}(Q_{\mu_*})$ , and  $\mathcal{V} = \text{ran}(P_{\mu_t}) \oplus \text{ran}(Q_{\mu_t})$  for all  $t$ . First, using the decomposition of  $\mathcal{V}$ , we rewrite

$$\dot{P}_{\mu_t}P_{\mu_*}g(t) = \dot{P}_{\mu_t}P_{\mu_t}P_{\mu_*}g(t) + \dot{P}_{\mu_t}Q_{\mu_t}P_{\mu_*}g(t)$$

using the identities in Section 2. Multiplying on the left by  $P_{\mu_t}$  yields  $P_{\mu_t}\dot{P}_{\mu_t}P_{\mu_t} = 0$ , so  $P_{\mu_t}\dot{P}_{\mu_t}\hat{g}(t) = P_{\mu_t}\dot{P}_{\mu_t}Q_{\mu_t}\hat{g}(t)$ .

Inserting the fixed-time decomposition for  $\mathcal{L}g(t)$ , we see  $P_{\mu_t}\mathcal{L}g(t) = P_{\mu_t}\mathcal{L}(P_{\mu_*} + Q_{\mu_*})g(t)$  hence

$$\frac{d}{dt}(P_{\mu_t}P_{\mu_*}g(t)) = P_{\mu_t}\dot{P}_{\mu_t}P_{\mu_*}g(t) + P_{\mu_t}\mathcal{L}P_{\mu_*}g(t) + P_{\mu_t}\mathcal{L}Q_{\mu_*}g(t)$$

Finally, using Dyson's identity to solve for  $v(t) = Q_{\mu_*}g(t)$  as in the standard MZ formalism, we find

$$\begin{aligned} \frac{d}{dt}(P_{\mu_t}P_{\mu_*}g(t)) &= P_{\mu_t}\dot{P}_{\mu_t}Q_{\mu_t}P_{\mu_*}g(t) + P_{\mu_t}\mathcal{L}P_{\mu_*}g(t) \\ &+ \int_0^t P_{\mu_t}P_{\mu_*}\mathcal{L}e^{(t-s)Q_{\mu_*}}\mathcal{L}P_{\mu_*}g(s)ds + P_{\mu_t}\mathcal{L}e^{tQ_{\mu_*}}\mathcal{L}Q_{\mu_*}g(0). \end{aligned}$$

□

**Corollary 3.1.** (*Vanishing Drift Under an Invariant Trivialization*) Suppose the Radon-Nikodym densities satisfy  $\rho_t(x) = \alpha(t)$ , and  $\alpha > 0$  independent of  $x$ . Then  $P_{\mu_t} = P_{\mu_0}$ , hence  $\dot{P}_{\mu_t} = 0$ .

*Proof.* For any  $g \in L^2(\mathcal{M}, \mu_t)$ ,  $P_{\mu_t}$  is defined by the requirement

$$\int_G f d\mu_t = \int_G (P_{\mu_t}f) d\mu_t \quad \text{for all measurable } G.$$

Since  $\mu_t = \alpha(t)\mu_0$

$$\int_G f d\mu_t = \alpha(t) \int_G f d\mu_0, \quad \int_G (P_{\mu_t}f) d\mu_t = \alpha(t) \int_G (P_{\mu_t}f) d\mu_0$$

Therefore

$$\int_G f d\mu_0 = \int_G (P_{\mu_t}f) d\mu_0 \quad \text{for all measurable } G,$$

which by the uniqueness of the conditional-expectation operator in  $L^2$  characterizes  $P_{\mu_0}$ . We thus conclude that  $P_{\mu_t} = P_{\mu_0}$  for all  $t$ , and as a result the time-derivative vanishes, i.e.,  $\dot{P}_{\mu_t} = 0$ . □

**Note on trivialization and isometries.** Collectively, the family  $\{\mathcal{V}_t\}_{t \in [0, T]}$  together with the projection map

$$\pi = \bigsqcup_t \mathcal{V}_t \rightarrow [0, T], \quad \pi(v) = t$$

constitutes a Hilbert bundle over the interval  $[0, T]$ . In this bundle picture, fibers are the individual  $\mathcal{V}_t$ , a section is a time-indexed observable  $g(t) \in \mathcal{V}_t$ . Here we describe trivialization with respect to a fixed reference within the bundle, which is given by the Radon-Nikodym isometry

$$\mathcal{T}_t : \mathcal{V}_0 \rightarrow \mathcal{V}_t, \quad \mathcal{T}_t(g) = \sqrt{\frac{d\mu_t}{d\mu_0}}(x)g(x) = \rho_t^{\frac{1}{2}}g.$$

## E INTRINSIC GENERALIZED FLUCTUATION DISSIPATION RELATION

**Assumption E.1.** (*Intrinsic Time-Dependent GFDR*) For the fixed reference measure  $\mu^*$  and projection  $P_{\mu^*}$ , the standard Mori–Zwanzig memory kernel and fluctuating force satisfy the generalized fluctuation–dissipation relation (5). In our intrinsic setting we enforce that the learned memory kernel  $K_{int}$  and noise  $F_{int}$  are compatible with this structure in the sense that

$$K_{int}(t-s) = P_{\mu_t} K_{\mu^*}(t-s), \quad F_{int}(t) = P_{\mu_t} F_{\mu^*}(t),$$

where  $K_{\mu^*}, F_{\mu^*}$  satisfy the GFDR under  $\mu^*$ .

## F METHODOLOGICAL DETAILS

### F.1 ARCHITECTURAL DETAILS

**Neural Wave Field** The Neural Wave Field maintains two coupled latent state  $h_t \in \mathbb{R}^n$  and  $\mu_t \in \mathbb{R}^n$ , which evolve under a Mori–Zwanzig inspired network and an accompanying measure-update expert. At each time step  $t$  the raw input  $x_t$  is first embedded into the feature space as a ghost boundary point. That is, it is available to be uptaken by the memory kernel provided the gating mechanism allows it.

For this reason, the MZ-NET  $\sigma_{mem}$  and  $\sigma_{force}$  are critical for determining the amount of long history information to retain, and the amount of new information to incorporate into the memory state. Whether the information is ultimately taken into the latent state is governed by  $\sigma_{closure}$ . These signals jointly determine a convolutional kernel  $C_{h_t}$  and padded hidden state  $\tilde{h}_t$  for updating  $h_{t+1} = C_{h_t} \star \tilde{h}_t$ .

A measure-dynamics expert network  $D_\mu$  determines the update for the measure between two time periods. This module enforces that  $\mu_t$  remains a valid probability density via softmax with a large temperature of 100.

Given our assumptions on the conditional-expectation projections of  $P_{\mu_t}$ , we train using the MSE loss across all tasks.

**WaveRNN** The WaveRNN architecture is most similar to the Neural Wave Field in its construction of a latent state. There are two particular differences in the approaches. First, the WaveRNN utilizes periodic boundary conditions which are a limiting factor as described by Corollary B.1. Moreover, the architecture relies on a static decoder and encoder which forces the projection dynamics to be invariant. As a result, the architecture will be unable to achieve a minimal latent state representation. Furthermore, it will be prohibited from accurately learning the selective copy task.

**Mamba** The Mamba architecture is a state-of-the-art structured state-space model. It has achieved particular success in modeling long-range tasks. It has done so by balancing long-range and short range updates to the latent state.

**Transformers** The positional encoding-based (or replacement) transformers aim to use various methods to replace fixed positional encoding mechanisms with relative positional encoding mechanisms. These have shown strong results in memory tasks such as the copy task.

## F.2 ADDITIONAL EXPERIMENTS

### F.2.1 CHAOTIC DYNAMICAL SYSTEMS

We evaluate how well our architecture can learn a highly non-periodic, chaotic manifold in accordance with Corollary B.1. For this reason, we compare against the WaveRNN baseline, which uses periodic boundary conditions in its latent state. We train both models to reconstruct the full phase-state from only its  $x$ -coordinate, using 300-step input sequences ( $\Delta t = 0.01$ ), and a latent dimension of 3.

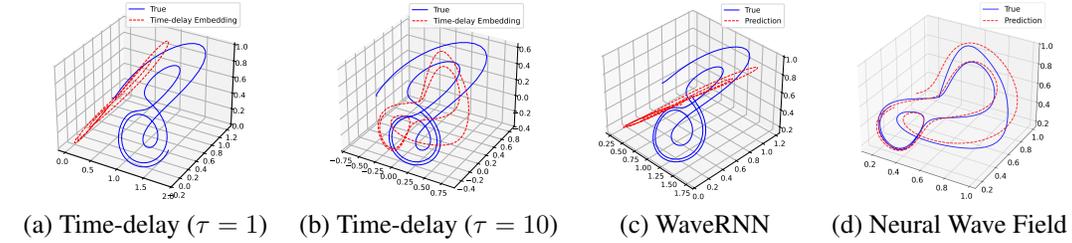


Figure 5: Time-delay and predicted trajectories of the Lorenz attractor using the time delays of  $\tau = 1$  and 10, and WaveRNN and Neural Wave Field models. We observe that the WaveRNN performs comparably to the under resolved  $\tau = 1$  time-delay embedding. In contrast, the Neural Wave Field achieves strong trajectory matching that degrades over time as errors slowly accumulate.

Figure 5 presents the time-delay and predicted latent trajectories of the Lorenz attractor using two classical delay embeddings ( $\tau = 1$  and  $\tau = 10$ ), as well as the learned embeddings from the WaveRNN and Neural Wave Field models. In our Neural Wave Field model, the latent trajectory forms smooth, closed loops that align with the true attractor and only gradually diverge as errors accumulate. Although this is slightly relaxed behavior from Proposition ??, it is attributable to approximation errors in the memory kernel, the drift dynamics, and the dynamics of  $\mu_t$ . By contrast, the WaveRNN fits the dynamics into a toroidal manifold introducing distortion and misalignment, especially over long time horizons, coinciding with Corollary B.1.

### F.3 TASK DETAILS

**Lorenz Attractor** We simulate the Lorenz system

$$\dot{x} = \sigma(y - x), \quad \dot{y} = x(\rho - z) - y, \quad \dot{z} = xy - \beta z$$

with standard parameters  $(\sigma, \rho, \beta) = (10, 28, 8/3)$  using a fourth-order Runge-Kutta integrator at step size  $\Delta t = 0.01$ . At each time step only the  $x$ -coordinate is provided as input; the models must reconstruct the full state  $(x_t, y_t, z_t)$ .

For all experiments, we use a training batch size of 128 and test using a batch size of 32. All batches are generated randomly to obtain the trajectory of 300 time-steps. The loss is only computed on the last 280 time-steps. For all models we use the Adam optimizer with a learning rate of 0.001 for 1000 batches.

For our comparisons, we use the following configurations. For WaveRNN (Keller et al., 2024), we use one channel, an identity activation, and a hidden dimension of 20 to have a more direct comparison to our model. The loss is mean squared error (MSE).

**Copy** For all experiments, we use a training batch size of 128 and test using a batch size of 50. All batches are generated randomly to obtain the sequence of 10 tokens to be memorized. We use  $T = 20$ , so the total sequence length is 30. The loss is only computed on the last 10 tokens; the intermediate outputs are not considered. That is, we only care about the model’s ability to reproduce the sequence of 10 tokens at the final 10 timesteps. For all models we use the Adam optimizer with a learning rate of 0.001 for 1000 batches.

For our comparisons, we use the following configurations. For WaveRNN (Keller et al., 2024), we use one channel and an identity activation to have a more direct comparison to our model. The loss

---

1188 is mean squared error (MSE). For Mamba and the transformer models, we use cross entropy loss,  
1189 as they naturally output logits over the vocabulary size. We found that these models needed at least  
1190 2 layers to perform on the task, which we use in our experiments. For the transformers, we use a  
1191 single attention head.

1192  
1193 **Selective Copy** By randomizing token positions and focusing evaluation solely on the terminal  
1194 outputs, this task highlights each model’s ability to selectively attend to and retain the correct in-  
1195 formation. Our architecture’s time-dependent projection and delay-coordinate closure enable it to  
1196 isolate the  $N$  informative tokens with minimal overhead, even as memory capacity is constrained.

1197  
1198 F.4 ASSUMPTIONS NOTE

1199  
1200 As a note on the practical implications of the assumptions made. When the size of the latent state is  
1201 larger than the minimal representation but not large enough to trivialize the dynamics of the measure,  
1202 then the additional degrees of freedom provide many non-unique and non-trivial solutions. In this  
1203 case, we experience large standard deviations in the training loss between runs with differing initial  
1204 conditions. In the case where memory is sufficiently large to trivialize the measure dynamics, the  
1205 learning became significantly more consistent.

1206  
1207 In addition, the continuity assumptions on the measure make it impossible to use the current frame-  
1208 work to effectively learn a version of the copy task where the predicted output is required to be  
1209 placed in order. However, on this task, we observe that the Mamba and transformer architectures  
1210 perform exceptionally well.

1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241