

778 **A Supplementary material**

779 **Contents**

780	A.1 Theoretical results in continuous time	22
781	A.1.1 Definitions & assumptions	22
782	A.1.2 Proof of the continuous adjoint state method (ASM)	23
783	A.1.3 Proof of Theorem 3.1	26
784	A.1.4 Connection to Hamiltonian Echo Backpropagation (HEB)	28
785	A.2 Theoretical results in discrete time	31
786	A.2.1 Definitions & assumptions	31
787	A.2.2 Backpropagation Through Time (BPTT)	33
788	A.2.3 Proof of Theorem 3.2	35
789	A.2.4 Proof of Theorem 3.3	39
790	A.3 Models and algorithms details	41
791	A.3.1 Toy model	41
792	A.3.2 HSSM architecture	42
793	A.3.3 Linear HRU Block	43
794	A.3.4 Nonlinear HRU	45
795	A.3.5 Differentiating the time-discretization	45
796	A.3.6 Echo passes with automatic differentiation	46
797	A.4 Experimental details	48
798	A.4.1 Datasets	48
799	A.4.2 Simulation details and ressource consumption	48
800	A.4.3 Hyperparameters	48
801	A.4.4 Gradient re-scaling for dealing with numerical instabilities	49

A.1 Theoretical results in continuous time

Summary. In this section, we present all the results derived in *continuous* time. More precisely:

- We formally define our model (Def. A.1) and constrained optimization problem (Def. A.2) in *continuous time*.
- We state and prove our algorithm baseline, the *continuous adjoint state method* (ASM) for this constrained optimization problem (Theorem A.1). To ease the comparison of the continuous ASM with *Recurrent Hamiltonian Echo Learning* (RHEL) in continuous time, we state re-parametrized version of the continuous ASM where time is indexed *backwards* (Corollary A.1). Finally, we also state a variant of the continuous ASM when the loss function is only defined at the *final* timestep (Corollary A.2) to ease the comparison between the continuous ASM and *Hamiltonian Echo Backprop* (HEB, [10]).
- We introduce two technical results (Lemma A.1–A.2) which enable us to prove the *time-reversal invariance* property of our model (Lemma A.3). We then show that the direct consequence of this property is the *time-reversibility* of our model upon momentum flipping (Corollary A.3), the key mechanics which fundamentally underpins our algorithm. All these intermediate results allow us to finally introduce RHEL in continuous time and prove its equivalence with the continuous ASM (Theorem A.2).
- Lastly, we connect HEB [10] with the continuous ASM when the loss is defined only at the final time step (Theorem A.4). We highlight some key differences between HEB and RHEL.

A.1.1 Definitions & assumptions

Definition A.1 (Continuous Hamiltonian model). *Given $\theta \in \mathbb{R}^{d_\theta}$, $T \in \mathbb{R}_+^*$ and an input sequence $t \rightarrow \mathbf{u}(t) \in (\mathbb{R}^{d_u})^{[-T,0]}$, the continuous Hamiltonian model prediction $t \rightarrow \Phi(t) \in (\mathbb{R}^{d_\Phi})^{[-T,0]}$ is, by definition, implicitly given as the solution of the following ODE:*

$$\Phi(-T) = \mathbf{x}, \quad \forall t \in [-T, 0] : \partial_t \Phi(t) = \mathbf{J} \cdot \nabla_\Phi H[\Phi(t), \theta, \mathbf{u}(t)], \quad \mathbf{J} := \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ -\mathbf{I} & \mathbf{0} \end{bmatrix}.$$

We assume that:

1. H is time-reversal invariant:

$$\forall \Phi \in \mathbb{R}^{d_\Phi}, \quad \forall \theta \in \mathbb{R}^{d_\theta}, \quad \forall \mathbf{u} \in \mathbb{R}^{d_u} : \quad H[\Phi, \theta, \mathbf{u}] = H[\Sigma_z \Phi, \theta, \mathbf{u}], \quad \Sigma_z := \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & -\mathbf{I} \end{bmatrix}$$

2. $\Phi \rightarrow H[\Phi, \cdot, \cdot]$ is twice continuously differentiable,
3. $\theta \rightarrow H[\cdot, \theta, \cdot]$ is differentiable,
4. $\nabla_{1,2}^2 H$ exists and is continuous with respect to t ,
5. $\mathbf{u} \rightarrow H[\cdot, \cdot, \mathbf{u}]$ is continuous,
6. $\Phi \rightarrow \nabla_1 H[\Phi, \cdot, \cdot]$ and $\Phi \rightarrow \nabla_{2,1} H[\Phi, \cdot, \cdot]$ are Lipschitz continuous.

Remark 1. While some of these assumptions will be used explicitly in our derivations, they are all needed to guarantee the existence of partial derivatives of s as an implicit function of \mathbf{x} and θ through Eq. (1) and we refer to [28] for such claims given these assumptions.

Definition A.2 (Continuous constrained optimization problem). *Given a continuous Hamiltonian model (Def. A.1), we consider the following constrained optimization problem:*

$$\min_{\boldsymbol{\theta}} L := \int_{-T}^0 dt \ell[t, \boldsymbol{\Phi}(t), \boldsymbol{\theta}] \quad \text{s.t.} \quad \forall t \in [-T, 0] : \partial_t \boldsymbol{\Phi}(t) = \mathbf{J} \cdot \nabla_{\boldsymbol{\Phi}} H[\boldsymbol{\Phi}(t), \boldsymbol{\theta}, \mathbf{u}(t)],$$

where we assume that:

1. ℓ is time-reversal invariant:

$$\forall \boldsymbol{\Phi} \in \mathbb{R}^{d_{\boldsymbol{\Phi}}}, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^{d_{\boldsymbol{\theta}}}, \quad \forall t \in [-T, 0] : \quad \ell[t, \boldsymbol{\Phi}, \boldsymbol{\theta}] = \ell[t, \boldsymbol{\Sigma}_z \cdot \boldsymbol{\Phi}, \boldsymbol{\theta}]$$

2. $t \rightarrow \ell[t, \cdot, \cdot]$ is continuous,

3. $\boldsymbol{\theta} \rightarrow \ell[\cdot, \cdot, \boldsymbol{\theta}]$ is differentiable,

4. $t \rightarrow \nabla_{\boldsymbol{\theta}} \ell[t, \cdot, \cdot]$ is continuous.

5. $\boldsymbol{\Phi} \rightarrow \ell[\cdot, \boldsymbol{\Phi}, \cdot]$ is twice differentiable,

826

827 **Remark 2.** Note that while we did not assume in the main part of this manuscript that ℓ depended
828 on $\boldsymbol{\theta}$, we assume it in the appendix for the generality of our derivations.

829 A.1.2 Proof of the continuous adjoint state method (ASM)

Theorem A.1 (Continuous adjoint state method ([28])). *Given assumptions A.1–A.2, the gradients of L with respect to $\boldsymbol{\theta}$ and \mathbf{x} are given by:*

$$d_{\boldsymbol{\theta}} L = \int_{-T}^0 g_{\boldsymbol{\theta}}(t) dt, \quad d_{\mathbf{x}} L = \boldsymbol{\lambda}(0)$$

with $t \rightarrow g_{\boldsymbol{\theta}} \in (\mathbb{R}^{d_{\boldsymbol{\theta}}})^{[-T, 0]}$ defined as:

$$g_{\boldsymbol{\theta}}(t) := \nabla_{\boldsymbol{\theta}} \ell[t, \boldsymbol{\Phi}(t), \boldsymbol{\theta}] + \nabla_{\boldsymbol{\Phi}, \boldsymbol{\theta}}^2 H[\boldsymbol{\Phi}(t), \boldsymbol{\theta}, \mathbf{u}(t)] \cdot \mathbf{J}^{\top} \cdot \boldsymbol{\lambda}(t) \quad \forall t \in [-T, 0]$$

and $\boldsymbol{\lambda}$ solving for the **adjoint ODE**:

$$\begin{cases} \boldsymbol{\lambda}(0) &= \mathbf{0} \\ \partial_t \boldsymbol{\lambda}(t) &= -\nabla_{\boldsymbol{\Phi}}^2 H[\boldsymbol{\Phi}(t), \boldsymbol{\theta}, \mathbf{u}(t)] \cdot \mathbf{J}^{\top} \cdot \boldsymbol{\lambda}(t) - \nabla_{\boldsymbol{\Phi}} \ell[t, \boldsymbol{\Phi}(t), \boldsymbol{\theta}] \end{cases}$$

830

831 *Proof of Theorem A.1.* With slight adaptations, our proof mostly follows that of [32] and also assume
832 the existence of the partial derivatives of $\boldsymbol{\Phi} = \boldsymbol{\Phi}(t, \mathbf{x}, \boldsymbol{\theta})$ as an *implicit* function of t, \mathbf{x} and $\boldsymbol{\theta}$ – we
833 defer to [28] for the proof of this claim.

834 We start off defining the Lagrangian associated with the constraint optimization problem:

$$\mathcal{L}(\boldsymbol{\Phi}, \boldsymbol{\lambda}, \boldsymbol{\theta}, \mathbf{u}) := \int_{-T}^0 dt \left(\ell[t, \boldsymbol{\Phi}(t, \boldsymbol{\theta}), \boldsymbol{\theta}] + \boldsymbol{\lambda}^{\top}(t) \cdot (\mathbf{J} \cdot \nabla_{\boldsymbol{\Phi}} H[\boldsymbol{\Phi}(t, \boldsymbol{\theta}), \boldsymbol{\theta}, \mathbf{u}(t)] - \partial_t \boldsymbol{\Phi}(t, \boldsymbol{\theta})) \right),$$

835 where $t \rightarrow \boldsymbol{\lambda}(t) \in (\mathbb{R}^{d_{\boldsymbol{\Phi}}})^{[-T, 0]}$ denotes the Lagrangian multiplier associated to the constraint.

836 **Derivation of $d_{\boldsymbol{\theta}} L$.** For readability, we emphasize the dependence of $\boldsymbol{\Phi}$ on t and $\boldsymbol{\theta}$, as we will
837 leverage the existence of its partial derivatives with respect to these variables. Then, given Assump-
838 tions A.1, the total derivative of the Lagrangian with respect to $\boldsymbol{\theta}$ exists and reads:

$$\begin{aligned} d_{\boldsymbol{\theta}} \mathcal{L} &= \int_{-T}^0 dt \left(\partial_{\boldsymbol{\theta}} \boldsymbol{\Phi}(t, \boldsymbol{\theta})^{\top} \cdot \nabla_{\boldsymbol{\Phi}} \ell[t, \boldsymbol{\Phi}(t, \boldsymbol{\theta}), \boldsymbol{\theta}] + \nabla_{\boldsymbol{\theta}} \ell[t, \boldsymbol{\Phi}(t, \boldsymbol{\theta}), \boldsymbol{\theta}] \right) \\ &\quad + \int_{-T}^0 dt \left[\partial_{\boldsymbol{\theta}} \boldsymbol{\Phi}(t, \boldsymbol{\theta})^{\top} \cdot \nabla_{\boldsymbol{\Phi}}^2 H[\boldsymbol{\Phi}(t, \boldsymbol{\theta}), \boldsymbol{\theta}, \mathbf{u}(t)] \cdot \mathbf{J}^{\top} + \nabla_{\boldsymbol{\theta}, \boldsymbol{\Phi}}^2 H[\boldsymbol{\Phi}(t, \boldsymbol{\theta}), \boldsymbol{\theta}, \mathbf{u}(t)]^{\top} \cdot \mathbf{J}^{\top} - \partial_{\boldsymbol{\theta}, t}^2 \boldsymbol{\Phi}(t, \boldsymbol{\theta})^{\top} \right] \cdot \boldsymbol{\lambda}(t). \end{aligned}$$

839 We can transform the last term of the integrand with $\partial_{\theta,t}\Phi(t, \theta)$ by applying Schwartz's theorem and
 840 and integration by parts as:

$$\begin{aligned} -\int_{-T}^0 dt \partial_{\theta,t}\Phi(t, \theta)^\top \cdot \lambda(t) &= -\int_{-T}^0 dt \partial_{t,\theta}\Phi(t, \theta)^\top \cdot \lambda(t) \\ &= [-\partial_{\theta}\Phi(t, \theta)^\top \cdot \lambda(t)]_{-T}^0 + \int_{-T}^0 dt \partial_{\theta}\Phi(t, \theta)^\top \cdot \partial_t \lambda(t) \\ &= -\partial_{\theta}\Phi(0, \theta)^\top \cdot \lambda(0) + \int_{-T}^0 dt \partial_{\theta}\Phi(t, \theta)^\top \cdot \partial_t \lambda(t) \end{aligned}$$

841 where the contribution of the first term at $t = -T$ vanishes because:

$$\Phi(-T, \theta) = x \Rightarrow \partial_{\theta}\Phi(0, \theta) = 0$$

842 Plugging this back into $d_{\theta}\mathcal{L}$ yields:

$$\begin{aligned} d_{\theta}\mathcal{L} &= -\partial_{\theta}\Phi(0, \theta)^\top \cdot \lambda(0) \\ &\quad + \int_{-T}^0 dt \partial_{\theta}\Phi(t, x, \theta)^\top \cdot [\nabla_{\Phi}\ell[t, \Phi(t, \theta), \theta] + \nabla_{\Phi}^2 H[\Phi(t, x, \theta), \theta] \cdot J^\top \cdot \lambda(t) + \partial_t \lambda(t)] \\ &\quad + \int_{-T}^0 dt (\nabla_{\theta}\ell[t, \Phi(t, \theta), \theta] + \nabla_{\Phi, \theta}^2 H[\Phi(t, x, \theta), \theta] \cdot J^\top \cdot \lambda(t)) \end{aligned}$$

843 Denoting Φ_* the solution of the (primal) ODE and by defining λ_* as the solution of the adjoint ODE:

$$\partial_t \lambda_*(t) = -\nabla_{\Phi}^2 H[\Phi_*(t, \theta), \theta, u(t)] \cdot J^\top \cdot \lambda_*(t) - \nabla_{\Phi}\ell[t, \Phi_*(t, \theta), \theta], \quad \lambda_*(0) = 0$$

844 we have:

$$\begin{aligned} d_{\theta}L &:= d_{\theta} \int_{-T}^0 dt \ell[t, \Phi(t, \theta), \theta] \\ &= d_{\theta}\mathcal{L}(\Phi_*, \lambda_*, \theta) \\ &= \int_{-T}^0 dt (\nabla_{\theta}\ell[t, \Phi_*(t, \theta), \theta] + \nabla_{\Phi, \theta}^2 H[\Phi_*(t, x, \theta), \theta] \cdot J^\top \cdot \lambda_*(t)) \end{aligned}$$

845 **Derivation of $d_x L$.** Similarly, we can prove that:

$$\begin{aligned} d_x \mathcal{L} &= \partial_x \Phi(0, x)^\top \cdot [\nabla \ell[\Phi(0, x)] - \lambda(0)] \\ &\quad + \int_{-T}^0 dt \partial_x \Phi(t, x)^\top \cdot [\nabla_{\Phi}\ell[t, \Phi(t, \theta), \theta] + \nabla_{\Phi}^2 H[\Phi(t, x, \theta), \theta, u(t)] \cdot J^\top \cdot \lambda(t) + \partial_t \lambda(t)] \\ &\quad + \lambda(0) \end{aligned}$$

846 Using the same Φ_* and λ_* , we get $d_x L = \lambda(0)$. □

847 **Reparametrization of the continuous adjoint method.** To ease the comparison of the continuous
 848 adjoint state method with our algorithm, we slightly reparametrize the variables introduced in
 849 Theorem A.1.

Corollary A.1. *Under the same assumptions as Theorem A.1, the gradients of L with respect to θ and x are given by:*

$$d_{\theta}L = \int_0^T g_{\theta}(t) dt, \quad d_x L = \lambda(T)$$

with $t \rightarrow g_{\theta} \in (\mathbb{R}^{d_{\theta}})^{[0, T]}$ defined as:

$$g_{\theta}(t) := \nabla_{\theta}\ell[-t, \Phi(-t), \theta] + \nabla_{\Phi, \theta}^2 H[\Phi(-t), \theta, u(-t)] \cdot J^\top \cdot \lambda(t) \quad \forall t \in [0, T]$$

and λ solving for the **adjoint ODE**:

$$\begin{cases} \lambda(0) &= 0 \\ \partial_t \lambda(t) &= \nabla_{\Phi}^2 H[\Phi(-t), \theta, u(-t)] \cdot J^\top \cdot \lambda(t) + \nabla_{\Phi}\ell[-t, \Phi(-t), \theta] \end{cases}$$

850

851 *Proof of Corollary A.1.* Immediately stems from Theorem A.1. \square

852 **Edge case: loss at the final timestep only.** The backward ODE can be differently parametrized
 853 when a loss is defined only at the last time step. We need this formulation for later convenience.

Corollary A.2. *Under the same assumptions as Theorem A.1, and assuming:*

$$\ell[t, \cdot, \cdot] = 0 \quad \forall t \in [0, T), \quad \ell[T, \cdot, \cdot] := \ell_T,$$

the gradients of ℓ_T with respect to θ and x are given by:

$$d_{\theta} \ell_T[\Phi(0), \theta] = \nabla_{\theta} \ell_T[\Phi(0), \theta] + \int_0^T g_{\theta}(t) dt, \quad d_x L = \lambda(T)$$

with $t \rightarrow g_{\theta} \in (\mathbb{R}^{d_{\theta}})^{[0, T]}$ defined as:

$$g_{\theta}(t) := \nabla_{\Phi, \theta}^2 H[\Phi(-t), \theta, u(-t)] \cdot J^{\top} \cdot \lambda(t) \quad \forall t \in [0, T]$$

*and λ solving for the **adjoint** ODE:*

$$\begin{cases} \lambda(0) &= \nabla_{\Phi} \ell_T[\Phi(0), \theta] \\ \partial_t \lambda(t) &= \nabla_{\Phi}^2 H[\Phi(-t), \theta, u(-t)] \cdot J^{\top} \cdot \lambda(t) \end{cases}$$

854 *Proof of Corollary A.2.* Starting from the Lagrangian:

$$\mathcal{L}(\Phi, \lambda, \theta, u) := \ell_T[\Phi(0), \theta] + \int_{-T}^0 dt \left(\lambda^{\top}(t) \cdot (J \cdot \nabla_{\Phi} H[\Phi(t), \theta, u(t)] - \partial_t \Phi(t, \theta)) \right),$$

856 the proof reads in the exact same fashion as that of Theorem A.1. \square

857 A.1.3 Proof of Theorem 3.1

858 **Technical Lemmas.** We first introduce two technical Lemmas which will be needed for the
859 derivation of our main result.

Lemma A.1 (Block-wise Pauli matrices and associated properties). *Defining $\Sigma_x, \Sigma_y, \Sigma_z \in \mathbb{R}^{d_\Phi \times d_\Phi}$ as:*

$$\Sigma_x := \begin{bmatrix} \mathbf{0} & I_{d_\Phi/2} \\ I_{d_\Phi/2} & \mathbf{0} \end{bmatrix}, \quad \Sigma_y := \begin{bmatrix} \mathbf{0} & -iI_{d_\Phi/2} \\ iI_{d_\Phi/2} & \mathbf{0} \end{bmatrix}, \quad \Sigma_z := \begin{bmatrix} I_{d_\Phi/2} & \mathbf{0} \\ \mathbf{0} & -I_{d_\Phi/2} \end{bmatrix}$$

where i denotes the imaginary unit, the following equalities hold:

1. $J = i\Sigma_y$
2. $\Sigma_x^2 = \Sigma_y^2 = \Sigma_z^2 = I_{d_\Phi/2}$,
3. $\Sigma_x \cdot \Sigma_y = i\Sigma_z, \Sigma_y \cdot \Sigma_z = i\Sigma_x, \Sigma_z \cdot \Sigma_x = i\Sigma_y$,
4. $\Sigma_i \cdot \Sigma_j = -\Sigma_j \cdot \Sigma_i$ for any $i \neq j \in \{x, y, z\}$.

860

861 *Proof of Lemma A.1.* Because of the block-wise structure of $\Sigma_x, \Sigma_y, \Sigma_z$, these equalities can be
862 easily checked. \square

Lemma A.2. *Under the assumptions of Def. A.1, the following equalities hold for all Φ, θ, u :*

$$\begin{aligned} \nabla_\Phi H[\Phi, \theta, u] &= \Sigma_z \cdot \nabla_{\Phi^*} H[\Phi^*, \theta, u], \\ \nabla_\Phi^2 H[\Phi, \theta, u] &= \Sigma_z \cdot \nabla_{\Phi^*}^2 H[\Phi^*, \theta, u] \cdot \Sigma_z, \\ \nabla_{\Phi, \theta} H[\Phi, \theta, u] &= \nabla_{\Phi^*, \theta} H[\Phi^*, \theta, u] \cdot \Sigma_z \end{aligned}$$

863

864 *Proof of Lemma A.2.* The above equalities can be simply obtained by differentiating through the
865 time-reversal invariance hypothesis – which is possible because of the differentiability of H with
866 respect to Φ and θ – and using the chain rule. Namely, given some Φ, θ, u :

$$\begin{aligned} \nabla_\Phi H[\Phi, \theta, u] &= \nabla_\Phi (H[\Phi^*, \theta, u]) \\ &= (\partial_\Phi \Phi^*)^\top \cdot \nabla_{\Phi^*} H[\Phi^*, \theta, u] \\ &= \Sigma_z^\top \cdot \nabla_{\Phi^*} H[\Phi^*, \theta, u] = \Sigma_z \cdot \nabla_{\Phi^*} H[\Phi^*, \theta, u], \end{aligned}$$

867 since $\Phi^* := \Sigma_z \cdot \Phi$. The other equalities are derived in the same way. \square

868 **Time-reversal invariance.** We highlight here how the assumption $H[\Phi, \cdot, \cdot] = H[\Sigma_z \cdot \Phi, \cdot, \cdot]$
869 given inside Def. A.1 entails time-reversal invariance of the dynamics – up to time-reversal the input
870 sequence.

Lemma A.3 (Time-reversal invariance of the dynamics). *Under the assumptions given in Def. A.1, if Φ the solution of the ODE:*

$$\partial_t \Phi(t) = J \cdot \nabla_\Phi H[\Phi(t), \theta, u(t)],$$

the function $\tilde{\Phi}^* : t \rightarrow \Sigma_z \cdot \Phi(-t) = [\phi^\top(-t), -\pi^\top(-t)]^\top$ is solution of the same ODE with the
time-reversed input sequence $t \rightarrow \tilde{u}(t) := u(-t)$.

871

872 *Proof.* Let $t \in [-T, 0]$. We have:

$$\begin{aligned} \partial_t \tilde{\Phi}^*(t) &= -\Sigma_z \cdot \partial_t \Phi(-t) && \text{(Def. of } \tilde{\Phi}^* \text{ and chain rule)} \\ &= -\Sigma_z \cdot J \cdot \nabla_\Phi H[\Phi(-t), \theta, u(-t)] && \text{(by assumption)} \\ &= +J \cdot \Sigma_z \cdot \nabla_\Phi H[\Phi(-t), \theta, u(-t)] && \text{(Lemma A.1)} \\ &= J \cdot \Sigma_z^2 \cdot \nabla_{\Phi^*} H[\Phi^*(-t), \theta, u(-t)] && \text{(Lemma A.2)} \\ &= J \cdot \nabla_{\Phi^*} H[\tilde{\Phi}^*(t), \theta, u(-t)] && \text{(Lemma A.1)} \end{aligned}$$

873

\square

874 **Time reversibility.** A direct consequence of the time-reversal of the dynamics under consideration
 875 (Lemma A.3) is *time reversibility*: upon flipping the momentum of Φ at time $t = 0$ ($\pi(0) \leftarrow -\pi(0)$)
 876 and presenting the input sequence in reversed order, the system evolves backward to its initial state.

Corollary A.3 (Reversibility of the dynamics). *Under the same assumptions as Lemma A.3, we define Φ^e as the solution of the ODE:*

$$\Phi^e(0) = \Phi^*(0), \quad \partial_t \Phi^e(t) = \mathbf{J} \cdot \nabla_{\Phi^e} H[\Phi^e(t), \boldsymbol{\theta}, \mathbf{u}(-t)] \quad \forall t \in [0, T].$$

Then:

$$\forall t \in [0, T] : \quad \Phi^e(t) = \tilde{\Phi}^*(t)$$

877
 878 *Proof.* Φ^e and $\tilde{\Phi}^*$ satisfy: i) the same initial conditions ($\tilde{\Phi}^*(0) = \Phi^*(0) = \Phi^e(0)$), ii) the same
 879 ODE (Lemma A.3), therefore by unicity of the solution of the ODE, they are equal at all time over
 880 the domain of definition of Φ . \square

881 **Main result.** We are now ready to state and demonstrate our main result in continuous time.

Theorem A.2 (Equivalence between RHEL and the continuous ASM). *Under the assumptions of Def. A.1 and Def. A.2, let Φ be the solution of the ODE for $t \in [-T, 0]$:*

$$\Phi(-T) = \mathbf{x}, \quad \partial_t \Phi(t) = \mathbf{J} \cdot \nabla_{\Phi} H[\Phi(t), \boldsymbol{\theta}, \mathbf{u}(t)].$$

Given Φ , let Φ^e be defined as the solution of the other ODE for $t \in [0, T]$:

$$\Phi^e(0) = \Phi^*(0), \quad \partial_t \Phi^e(t, \epsilon) = \mathbf{J} \nabla_{\Phi^e} H[\Phi^e(t, \epsilon), \boldsymbol{\theta}, \mathbf{u}(-t)] - \epsilon \mathbf{J} \nabla_{\Phi^e} \ell[-t, \Phi^e(t, \epsilon), \boldsymbol{\theta}].$$

Defining:

$$\Delta_{\boldsymbol{\theta}}^{\text{RHEL}}(t, \epsilon) := \nabla_{\boldsymbol{\theta}} \ell[-t, \Phi^e(t, \epsilon), \boldsymbol{\theta}] + \frac{1}{2\epsilon} (\nabla_{\boldsymbol{\theta}} H[\Phi^e(t, \epsilon), \boldsymbol{\theta}, \mathbf{u}(-t)] - \nabla_{\boldsymbol{\theta}} H[\Phi^e(t, -\epsilon), \boldsymbol{\theta}, \mathbf{u}(-t)]),$$

$$\Delta_{\Phi}^{\text{RHEL}}(t, \epsilon) := \frac{1}{2\epsilon} \Sigma_{\mathbf{x}} \cdot (\Phi^e(t, \epsilon) - \Phi^e(t, -\epsilon)),$$

we have:

$$\forall t \in [0, T], \quad \lambda(t) = \lim_{\epsilon \rightarrow 0} \Delta_{\Phi}^{\text{RHEL}}(t, \epsilon), \quad g_{\boldsymbol{\theta}}(t) = \lim_{\epsilon \rightarrow 0} \Delta_{\boldsymbol{\theta}}^{\text{RHEL}}(t, \epsilon)$$

where λ and $g_{\boldsymbol{\theta}}$ are defined in Corollary A.1.

882
 883 *Proof.* Defining $\Delta_{\Phi}^{\text{RHEL}}(t) := \lim_{\epsilon \rightarrow 0} \Delta_{\Phi}^{\text{RHEL}}(t, \epsilon)$ and $\Delta_{\boldsymbol{\theta}}^{\text{RHEL}}(t) := \lim_{\epsilon \rightarrow 0} \Delta_{\boldsymbol{\theta}}^{\text{RHEL}}(t, \epsilon)$, note
 884 that:

$$\begin{aligned} \Delta_{\Phi}^{\text{RHEL}}(t) &= \Sigma_{\mathbf{x}} \cdot \partial_{\epsilon} \Phi^e(t, \epsilon)|_{\epsilon=0} \\ \Delta_{\boldsymbol{\theta}}^{\text{RHEL}}(t) &= \nabla_{\boldsymbol{\theta}} \ell[-t, \Phi^e(t, 0), \boldsymbol{\theta}] + \partial_{\epsilon} (\nabla_{\boldsymbol{\theta}} H[\Phi^e(t, \epsilon), \boldsymbol{\theta}, \mathbf{u}(-t)])|_{\epsilon=0} \end{aligned}$$

885 **Derivation of $\lambda(t) = \lim_{\epsilon \rightarrow 0} \Delta_{\Phi}^{\text{RHEL}}(t, \epsilon)$.** Given $t \in [0, T]$, differentiating the ODE satisfied by
 886 Φ^e with respect to ϵ at $\epsilon = 0$ yields:

$$\begin{aligned} \partial_t (\partial_{\epsilon} \Phi^e(t, \epsilon)|_{\epsilon=0}) &= \partial_{\epsilon} (\partial_t \Phi^e(t, \epsilon))|_{\epsilon=0} && \text{(Schwartz Theorem)} \\ &= \partial_{\epsilon} (\mathbf{J} \cdot \nabla_{\Phi^e} H[\Phi^e(t, \epsilon), \boldsymbol{\theta}, \mathbf{u}(-t)] - \epsilon \mathbf{J} \nabla_{\Phi^e} \ell[t, \Phi^e(t, \epsilon), \boldsymbol{\theta}])|_{\epsilon=0} \\ &= \mathbf{J} \cdot \nabla_{\Phi^e}^2 H[\Phi^e(t, 0)] \cdot \partial_{\epsilon} \Phi^e(t, \epsilon)|_{\epsilon=0} - \mathbf{J} \nabla_{\Phi^e} \ell[t, \Phi^e(t, 0), \boldsymbol{\theta}]. \end{aligned}$$

887 By Lemma A.3:

$$\Phi^e(t, 0) = \tilde{\Phi}^*(t) \quad \forall t \in [0, T],$$

888 therefore:

$$\begin{aligned} \partial_t (\partial_{\epsilon} \Phi^e(t, \epsilon)|_{\epsilon=0}) &= \mathbf{J} \cdot \nabla_{\Phi^*}^2 H[\Phi^*(-t)] \cdot \partial_{\epsilon} \Phi^e(t, \epsilon)|_{\epsilon=0} - \mathbf{J} \nabla_{\Phi^*} \ell[t, \Phi^*(-t), \boldsymbol{\theta}] \\ &= \mathbf{J} \cdot \Sigma_z \cdot \nabla_{\Phi}^2 H[\Phi(-t, 0)] \cdot \Sigma_z \cdot \partial_{\epsilon} \Phi^e(t, \epsilon)|_{\epsilon=0} - \mathbf{J} \cdot \Sigma_z \cdot \nabla_{\Phi} \ell[t, \Phi(-t), \boldsymbol{\theta}] \end{aligned} \quad \text{(Lemma A.2)}$$

889 Additionally, note that we have by Lemma A.1:

$$\mathbf{J} \cdot \Sigma_z = -\Sigma_x, \quad \Sigma_z = \mathbf{J} \cdot \Sigma_x,$$

890 so that:

$$\partial_t(\partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0}) = -\Sigma_x \cdot \nabla_{\Phi}^2 H[\Phi(-t, 0)] \cdot (\mathbf{J} \cdot \Sigma_x) \cdot \partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0} + \Sigma_x \cdot \nabla_{\Phi} \ell[t, \Phi(-t), \theta] \quad (19)$$

891 Left multiplying Eq. (19) on both sides by Σ_x yields:

$$\begin{aligned} \partial_t(\Sigma_x \cdot \partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0}) &= -\Sigma_x^2 \cdot \nabla_{\Phi}^2 H[\Phi(-t, 0)] \cdot \mathbf{J} \cdot (\Sigma_x \cdot \partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0}) + \Sigma_x^2 \cdot \nabla_{\Phi} \ell[t, \Phi(-t), \theta] \\ &= -\nabla_{\Phi}^2 H[\Phi(-t, 0)] \cdot \mathbf{J} \cdot (\Sigma_x \cdot \partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0}) + \nabla_{\Phi} \ell[t, \Phi(-t), \theta] \quad (\text{Lemma A.1}) \\ &= \nabla_{\Phi}^2 H[\Phi(-t, 0)] \cdot \mathbf{J}^\top \cdot (\Sigma_x \cdot \partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0}) + \nabla_{\Phi} \ell[t, \Phi(-t), \theta] \quad (\mathbf{J}^\top = -\mathbf{J}) \end{aligned}$$

892 Finally, note that because $\Phi^e(0) = \Phi^*$ does not depend on ϵ , we have that:

$$\partial_t(\Sigma_x \cdot \partial_\epsilon \Phi^e(0, \epsilon)|_{\epsilon=0}) = 0,$$

893 so that all in all, $\Delta_{\Phi}^{\text{RHEL}}$ satisfies:

$$\begin{cases} \Delta_{\Phi}^{\text{RHEL}}(0) &= \mathbf{0} \\ \partial_t \Delta_{\Phi}^{\text{RHEL}}(t) &= \nabla_{\Phi}^2 H[\Phi(-t), \theta, \mathbf{u}(-t)] \cdot \mathbf{J}^\top \cdot \Delta_{\Phi}^{\text{RHEL}}(t) + \nabla_{\Phi} \ell[-t, \Phi(-t), \theta] \end{cases}$$

894 Therefore $\Delta_{\Phi}^{\text{RHEL}}$ and λ (as defined in Corollary A.1) satisfy the same initial conditions and the
895 same ODE, therefore they are equal at all times.

896 **Derivation of $g_\theta(t) = \lim_{\epsilon \rightarrow 0} \Delta_\theta^{\text{RHEL}}(t, \epsilon)$.** Note that by Lemma A.3 and time-reversal invariance
897 of ℓ :

$$\begin{aligned} \Delta_\theta^{\text{RHEL}}(t) &= \nabla_\theta \ell[-t, \Phi^*(-t, 0), \theta] + \partial_\epsilon (\nabla_\theta H[\Phi^e(t, \epsilon), \theta, \mathbf{u}(-t)])|_{\epsilon=0} \\ &= \nabla_\theta \ell[-t, \Phi(-t, 0), \theta] + \partial_\epsilon (\nabla_\theta H[\Phi^e(t, \epsilon), \theta, \mathbf{u}(-t)])|_{\epsilon=0} \end{aligned}$$

898 As the first term of $\Delta_\theta^{\text{RHEL}}(t)$ and $g_\theta(t)$ coincide, the remainder of the derivation focuses on the
899 second term of $\Delta_\theta^{\text{RHEL}}(t)$. Given $t \in [0, T]$, we have:

$$\begin{aligned} \nabla_{\Phi, \theta}^2 H[\Phi(-t), \theta, \mathbf{u}(t)] \cdot \mathbf{J}^\top \cdot \lambda(t) &= \nabla_{\Phi^e, \theta}^2 H[\Phi^e(t, 0), \theta, \mathbf{u}(t)] \cdot \Sigma_z \cdot \mathbf{J}^\top \cdot \lambda(t) \quad (\text{Lemma A.2}) \\ &= -\nabla_{\Phi^e, \theta}^2 H[\Phi^e(t, 0), \theta, \mathbf{u}(t)] \cdot \Sigma_z \cdot i\Sigma_y \cdot \lambda(t) \quad (\mathbf{J} = i\Sigma_y) \\ &= +\nabla_{\Phi^e, \theta}^2 H[\Phi^e(t, 0), \theta, \mathbf{u}(t)] \cdot i\Sigma_y \cdot \Sigma_z \cdot \lambda(t) \quad (\text{Lemma A.1}) \\ &= -\nabla_{\Phi^e, \theta}^2 H[\Phi^e(t, 0), \theta, \mathbf{u}(t)] \cdot \Sigma_x \cdot \lambda(t) \quad (\text{Lemma A.1}) \\ &= -\nabla_{\Phi^e, \theta}^2 H[\Phi^e(t, 0), \theta, \mathbf{u}(t)] \cdot \Sigma_x^2 \cdot \partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0} \quad (\lambda = \Delta_{\Phi}^{\text{RHEL}}) \\ &= -\nabla_{\Phi^e, \theta}^2 H[\Phi^e(t, 0), \theta, \mathbf{u}(t)] \cdot \partial_\epsilon \Phi^e(t, \epsilon)|_{\epsilon=0} \quad (\text{Lemma A.1}) \\ &= -\partial_\epsilon (\nabla_\theta H[\Phi^e(t, \epsilon), \theta, \mathbf{u}(t)])|_{\epsilon=0}, \end{aligned}$$

900 which finishes to prove $g_\theta(t) = \Delta_\theta^{\text{RHEL}}(t)$ for $t \in [0, T]$. □

901 A.1.4 Connection to Hamiltonian Echo Backpropagation (HEB)

902 **Remark 3.** The above setup and implementation of RHEL is not exactly that of Hamiltonian Echo
903 Backprop (HEB, [10]). In particular:

- 904 • the loss function in HEB is only defined at the final time step,
- 905 • the interaction with ℓ does not happen simultaneously with H ,
- 906 • finally, we would like to recover the HEB formula giving the gradient estimate of the loss
907 with respect to the initial state of the neurons.

908 In the following corollary, we make slight algorithmic adjustments to match the seminal HEB
 909 implementation as much as possible **while preserving the generality of the sequence modelling**
 910 **setting**, i.e. dependence of H with θ and \mathbf{u} . Note that in the seminal HEB work, H does not depend
 911 on a static set of parameters θ nor on an input sequence. \mathbf{u} .

Corollary A.4. Under the assumptions of Def. A.1 and Def. A.2, and assuming additionally:

$$\ell[t, \cdot, \cdot] = 0 \quad \forall t \in [0, T], \quad \ell[T, \cdot, \cdot] := \ell_T,$$

let Φ be the solution of the ODE, for $t \in [-T, 0]$:

$$\Phi(-T) = \mathbf{x}, \quad \partial_t \Phi(t) = \mathbf{J} \cdot \nabla_{\Phi} H[\Phi(t), \theta, \mathbf{u}(t)],$$

and the solution of another ODE, for $t \in [0, \epsilon]$:

$$\partial_t \Phi(t) = \mathbf{J} \nabla_{\Phi} \ell_T[\Phi(t), \theta].$$

Let Φ^e be the solution of the following ODE, for $t \in [\epsilon, T]$:

$$\Phi^e(0, \epsilon) = \Phi(\epsilon)^*, \quad \partial_t \Phi^e(t, \epsilon) = \mathbf{J} \nabla_{\Phi^e} H[\Phi^e(t, \epsilon), \theta, \mathbf{u}(-t)],$$

Defining:

$$\begin{aligned} \Delta_{\theta}^{\text{RHEL}}(t, \epsilon) &:= \frac{1}{2\epsilon} (\nabla_{\theta} H[\Phi^e(t, \epsilon), \theta, \mathbf{u}(-t)] - \nabla_{\theta} H[\Phi^e(t, -\epsilon), \theta, \mathbf{u}(-t)]), \\ \Delta_{\Phi}^{\text{RHEL}}(t, \epsilon) &:= \frac{1}{2\epsilon} \Sigma_x \cdot (\Phi^e(t, \epsilon) - \Phi^e(t, -\epsilon)), \end{aligned}$$

we have:

$$\forall t \in [0, T], \quad \lambda(t) = \lim_{\epsilon \rightarrow 0} \Delta_{\Phi}^{\text{RHEL}}(t, \epsilon), \quad g_{\theta}(t) = \lim_{\epsilon \rightarrow 0} \Delta_{\theta}^{\text{RHEL}}(t, \epsilon)$$

where λ and g_{θ} are defined in Corollary A.2. In particular:

$$-i\epsilon d_{\mathbf{x}^*}^w \ell_T[\Phi(0), \theta] = \Phi^e(T, \epsilon)^* - \Phi(-T) + \mathcal{O}(\epsilon^2)$$

912 where $d_{\mathbf{x}^*}^w \equiv d_{\Phi}$ denotes the total Wirtinger derivative with respect to \mathbf{x}^* and i the imaginary unit.

913 *Proof of Corollary A.4.* The derivation is almost exactly similar to that of Theorem 3.1 with two key
 914 differences:

- 915 • the version of the continuous ASM against which this version of RHEL is compared is
 916 different (Corollary A.2),
- 917 • the interaction of Φ with ℓ and H do not happen simultaneously but on disjoint intervals.

918 We will simply show that the interaction with ℓ and conjugation $\Phi \rightarrow \Phi^*$ yields the correct initial
 919 conditions and defer to the proof of Theorem 3.1 for the remainder. We will also use the same
 920 notations and denote $\Delta_{\Phi}^{\text{RHEL}}(t) := \lim_{\epsilon \rightarrow 0} \Delta_{\Phi}^{\text{RHEL}}(t, \epsilon)$, $\Delta_{\theta}^{\text{RHEL}}(t) := \lim_{\epsilon \rightarrow 0} \Delta_{\theta}^{\text{RHEL}}(t, \epsilon)$.

921 **Derivation of $\lambda(t) = \lim_{\epsilon \rightarrow 0} \Delta_{\Phi}^{\text{RHEL}}(t, \epsilon)$.** Integrating the ODE satisfied by Φ between 0 and T
 922 yields:

$$\begin{aligned} \Phi(\epsilon) &= \Phi(0) + \int_0^{\epsilon} dt \mathbf{J} \cdot \nabla_{\Phi} \ell[\Phi(t), \theta] \\ &= \Phi(0) + \epsilon \mathbf{J} \cdot \nabla_{\Phi} \ell[\Phi(0), \theta] + \mathcal{O}(\epsilon^2) \end{aligned}$$

923 Therefore, the initial state of Φ^e can be written as:

$$\begin{aligned} \Phi^e(\epsilon) &= \Phi^*(0) + \epsilon \Sigma_z \cdot \mathbf{J} \cdot \nabla_{\Phi} \ell[\Phi(0), \theta] + \mathcal{O}(\epsilon^2) \\ &= \Phi^*(0) + \epsilon \Sigma_x \cdot \nabla_{\Phi} \ell[\Phi(0), \theta] + \mathcal{O}(\epsilon^2) \end{aligned} \quad (\text{Lemma A.1}).$$

924 By differentiating the last equality with respect to ϵ at $\epsilon = 0$, we obtain:

$$\Delta_{\Phi}^{\text{RHEL}}(0) = \nabla_{\Phi} \ell[\Phi(0), \theta].$$

925 Proceeding exactly as in the proof of Theorem A.2, we obtain that:

$$\begin{cases} \Delta_{\Phi}^{\text{RHEL}}(0) &= \nabla_{\Phi} \ell[\Phi(0), \theta], \\ \partial_t \Delta_{\Phi}^{\text{RHEL}}(t) &= \nabla_{\Phi}^2 H[\Phi(-t), \theta, \mathbf{u}(-t)] \cdot \mathbf{J}^{\top} \cdot \Delta_{\Phi}^{\text{RHEL}}(t) \end{cases}$$

926 Therefore $\Delta_{\Phi}^{\text{RHEL}}$ and λ (as defined in Corollary A.2) satisfy the same initial conditions and the
 927 same ODE, therefore they are equal at all times.

928 **Derivation of $g_{\theta}(t) = \lim_{\epsilon \rightarrow 0} \Delta_{\theta}^{\text{RHEL}}(t, \epsilon)$.** See proof of Theorem A.2.

929 **Connection to HEB formula.** In particular, we have:

$$\begin{aligned} d_{\mathbf{x}} \ell_T[\Phi(0), \theta] &= \lambda(T) = \Sigma_x \cdot \partial_{\epsilon} (\Phi(T, \epsilon) |_{\epsilon=0}) \\ &= \frac{1}{\epsilon} \Sigma_x (\Phi^e(T, \epsilon) - \Phi^e(T, 0)) + \mathcal{O}(\epsilon) \\ &= \frac{1}{\epsilon} \Sigma_x (\Phi^e(T, \epsilon) - \Phi^*(-T)) + \mathcal{O}(\epsilon) && \text{(Lemma A.3)} \\ &= -\frac{i}{\epsilon} \Sigma_y \cdot \Sigma_z \cdot (\Phi^e(T, \epsilon) - \Phi^*(-T)) + \mathcal{O}(\epsilon) && \text{(Lemma A.1)} \\ &= -\frac{i}{\epsilon} \Sigma_y \cdot (\Phi^e(T, \epsilon)^* - \Phi(-T)) + \mathcal{O}(\epsilon) \end{aligned}$$

930 Left multiplying on both sides by $i\epsilon \Sigma_y$ and noticing that $id_{\Phi^*}^w \equiv -i\Sigma_y \cdot d_{\Phi}$, we finally obtain:

$$-i\epsilon d_{\Phi^*}^w \ell_T[\Phi(0), \theta] = \Phi^e(T, \epsilon)^* - \Phi(-T) + \mathcal{O}(\epsilon^2)$$

931

□

A.2 Theoretical results in discrete time

Summary. In this section, we introduce all the results derived in *discrete* time. More precisely:

- We first formally define *Hamiltonian Recurrent Units* (HRUs, Definition A.3). HRUs can be regarded as the discrete-time counterpart of the continuous model introduced in the previous section (Definition A.1), namely as an *explicit* and *symplectic* integrator of the continuous Hamiltonian model which preserves the time-reversal invariance and time-reversibility properties in discrete time. We also introduce the constrained optimization problem naturally associated with HRUs (Definition A.4), which is the discrete time counterpart of the constrained continuous optimization problem introduced in the previous section (Definition A.2).
- We then formally define *Hamiltonian State Space Models* (HSSMs) as stacks of HRUs (Definition A.5) and the *multilevel* constrained optimization problem which is naturally associated to these models (Definition A.6).
- We state and prove our algorithmic baseline, *Backpropagation Through Time* (BPTT), through the lens of the *Lagrangian* formalism to establish a clear connection with the continuous ASM. We first introduce and derive BPTT in its general form (Theorem A.3) and then apply it more specifically to a HRU as defined in Definition A.3 (Corollary A.5).
- As we did in continuous time, we introduce a series of technical Lemmas needed to extend RHEL in discrete time. We first demonstrate the time-reversibility of HRUs on a single time step (Lemma A.4), which then enables us to extend the time reversibility property derived in continuous time (Corollary A.3) to *discrete* time (Corollary A.6). After introducing one last technical result (Lemma A.5), we then state and prove RHEL in discrete time when applied to HRUs (Corollary A.7). As the algorithm prescribed by Corollary A.7 includes solving an *implicit equation*, we finally introduce a slight practical (*i.e.* fully explicit) variant of RHEL in discrete time (Corollary A.8).
- Lastly, we show how to estimate gradients end-to-end in HSSMs by using *RHEL-chaining* (Theorem A.4). We also highlight that in practice, when using feedforward transformations across HRUs, the algorithm prescribed by Theorem A.4 implicitly requires to chain RHEL through HRUs and *automatic differentiation* through these feedforward transformations (Remark 8). This remark fundamentally underpins the actual algorithmic implementation of RHEL which was used throughout our experiments.

A.2.1 Definitions & assumptions

Definition A.3 (Hamiltonian Recurrent Unit). *Given $\theta \in \mathbb{R}^{d_\theta}$, $K \in \mathbb{N}^*$ and an input sequence $(\mathbf{u}_k)_{k \in [-K, 0]} \in (\mathbb{R}^{d_u})^K$, the Hamiltonian Recurrent Unit (HRU) prediction is given by:*

$$\Phi_{k+1} = \mathcal{M}_{H,\delta}[\Phi_k, \theta, \mathbf{u}_k] \quad \forall k = -K \cdots -1,$$

with $H := T + V$ and:

$$\mathcal{M}_{H,\delta} := \mathcal{M}_{T,\delta/2} \circ \mathcal{M}_{V,\delta} \circ \mathcal{M}_{T,\delta/2}, \quad \mathcal{M}_{T,\delta} := \Phi + \delta \mathbf{J} \cdot \nabla_\Phi T, \quad \mathcal{M}_{V,\delta} := \Phi + \delta \mathbf{J} \cdot \nabla_\Phi V$$

We assume that:

1. H is separable, *i.e.* V and T only depend on ϕ and π respectively:

$$V[\Phi, \theta, \mathbf{u}] = V[\phi, \theta, \mathbf{u}], \quad T[\Phi, \theta, \mathbf{u}] = T[\pi, \theta, \mathbf{u}]$$

2. T and V are time-reversal invariant:

$$\forall \Phi \in \mathbb{R}^{d_\Phi}, \forall \theta \in \mathbb{R}^{d_\theta}, \forall \mathbf{u} \in \mathbb{R}^{d_u} : \quad \begin{cases} T[\Phi, \theta, \mathbf{u}] = T[\Sigma_z \cdot \Phi, \theta, \mathbf{u}], \\ V[\Phi, \theta, \mathbf{u}] = V[\Sigma_z \cdot \Phi, \theta, \mathbf{u}] \end{cases}$$

3. T and V are twice differentiable with respect to Φ , θ and \mathbf{u} .

Remark 4. Note that $\mathcal{M}_{H,\delta}$ is simply a Leapfrog integrator associated with H . We justify each of our design choices below:

- 967 • **3 steps-parametrization.** We write the Leapfrog integrator in a three-steps fashion to yield a
968 reversible integrator.
- 969 • **Separability of the Hamiltonian.** In the case where ϕ and ϕ can be separated out in the
970 Hamiltonian function, the Leapfrog integrator becomes explicit [33].
- 971 • **T and V as functions of Φ .** Although T and V only depend on π and ϕ respectively, we
972 choose to define them as functions of Φ so that the proof of RHEL in the continuous case
973 seamlessly translates to the discrete case.

Definition A.4 (Constrained optimization problem in discrete time). *Given a continuous Hamiltonian model (Def. A.1), we consider the following constrained optimization problem:*

$$\min_{\theta} L := \sum_{k=-K+1}^0 \ell[\Phi_k, k] \quad \text{s.t.} \quad \Phi_{k+1} = \mathcal{M}_{H,\delta}[\Phi_k, \theta, \mathbf{u}_k] \quad \forall k = -K \cdots -1$$

where we assume that:

1. ℓ is time-reversal invariant:

$$\forall \Phi \in \mathbb{R}^{d_\Phi}, \forall \theta \in \mathbb{R}^{d_\theta}, \forall k = -K, \dots, 0: \quad \ell_k[\Phi, \theta] = \ell_k[\Sigma_z \cdot \Phi, \theta]$$

2. ℓ is twice differentiable with respect to Φ and θ .

Definition A.5 (Hamiltonian State Space Models). *Given $(\theta^{(1)}, \dots, \theta^{(N)}) \in (\mathbb{R}^{d_\theta})^N$, $K \in \mathbb{N}^*$ and an input sequence $(\mathbf{u}_k)_{k \in [-K, 0]} \in (\mathbb{R}^{d_u})^K$, a Hamiltonian State Space Model (HSSM) is defined as the composition of HRUs defined in Def. A.3 as:*

$$\Phi^{(0)} := \bar{\mathbf{u}}, \quad \forall \ell \in \llbracket 0, N-1 \rrbracket, \quad \forall k \in \llbracket -K, 0 \rrbracket: \quad \Phi_{k+1}^{(\ell+1)} = \mathcal{M}_{H^{(\ell)}, \delta}^{(\ell)}[\Phi_k^{(\ell+1)}, \theta^{(\ell)}, \Phi_k^{(\ell)}],$$

or in a vectorized fashion as:

$$\Phi^{(0)} := \bar{\mathbf{u}}, \quad \forall \ell \in \llbracket 0, N-1 \rrbracket: \quad \bar{\Phi}^{(\ell+1)} = \mathcal{M}_{H^{(\ell)}, \delta}^{(\ell)}[\bar{\Phi}^{(\ell+1)}, \theta^{(\ell)}, \bar{\Phi}^{(\ell)}],$$

Definition A.6 (Multilevel optimization problem in discrete time). *Given a HSSM (Def. A.5), we consider the following constrained optimization problem:*

$$\min_{\theta} L := \sum_{k=-K+1}^0 \ell[\Phi_k, k] \quad \text{s.t.} \quad \Phi^{(0)} := \bar{\mathbf{u}},$$

$$\forall \ell \in \llbracket 0, N-1 \rrbracket: \quad \bar{\Phi}^{(\ell+1)} = \mathcal{M}_{H^{(\ell)}, \delta}^{(\ell)}[\bar{\Phi}^{(\ell+1)}, \theta^{(\ell)}, \bar{\Phi}^{(\ell)}]$$

where we assume that ℓ satisfies the same assumptions as in Def. A.4.

977 A.2.2 Backpropagation Through Time (BPTT)

978 **General form.** We first state and prove Backpropagation Through Time (BPTT) for any integrator
979 $\mathcal{M}_{H,\delta}$.

Theorem A.3 (Backpropagation Through Time (BPTT)). *Given assumptions in Def. A.3–A.4, the gradients of the loss with respect to the parameters θ and the inputs \mathbf{u}_{-k} are given by:*

$$d_{\theta}L = \sum_{k=0}^{K-1} g_{\theta}(k), \quad d_{\mathbf{u}_{-(k+1)}}L = g_{\mathbf{u}}(k) \quad \forall k \in \llbracket 0, K-1 \rrbracket,$$

with:

$$\begin{cases} g_{\theta}(k) &= \nabla_2 \ell[\Phi_{-k}, \theta, -k] + \partial_2 \mathcal{M}_{H,\delta}[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}]^{\top} \cdot \lambda_k \\ g_{\mathbf{u}}(k) &= \partial_3 \mathcal{M}_{H,\delta}[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}]^{\top} \cdot \lambda_k, \end{cases}$$

and where (λ_k) satisfy the following recursion relationship:

$$\begin{cases} \lambda_0 = \nabla_1 \ell[\Phi_0, 0], \\ \lambda_{k+1} = \partial_1 \mathcal{M}_{H,\delta}(\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)})^{\top} \cdot \lambda_k + \nabla_1 \ell[\Phi_{-(k+1)}, -(k+1)] \quad \forall k = 0, \dots, K-1 \end{cases}$$

981 *Proof of Theorem A.3.* BPTT can classically be derived through the application of the “chain rule”
982 backward through the inference computational graph defined in Def. A.3. Another useful viewpoint
983 though, which directly connects BPTT as the discrete counterpart of the continuous ASM and will be
984 useful later in the appendix, is to derive it through *the method of Lagrangian multipliers*. Namely, the
985 Lagrangian associated to the constrained optimization problem in Def. A.4 reads as:

$$\mathcal{L}(\Phi, \lambda, \theta, \mathbf{u}) = \sum_{k=0}^{K-1} \ell[\Phi_{-k}, \theta, -k] + \lambda_k^{\top} \cdot (\mathcal{M}_{H,\delta}[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}] - \Phi_{-k})$$

986 Extremizing \mathcal{L} with respect to Φ and λ yield $\Phi_{k,*}$ and $\lambda_{k,*}$:

$$\forall k = 0, \dots, K-1 : \quad \partial_{\lambda_k} \mathcal{L}(\Phi_*, \lambda_*, \theta, \mathbf{u}) = \mathcal{M}_{H,\delta}[\Phi_{-(k+1),*}, \theta, \mathbf{u}_{-(k+1)}] - \Phi_{-k,*} = 0,$$

$$\partial_{\Phi_0} \mathcal{L}(\Phi_*, \lambda_*, \theta, \mathbf{u}) = \nabla_1 \ell[\Phi_{0,*}, \theta, 0] - \lambda_{0,*} = 0$$

$$\forall k = 1, \dots, K-1 : \quad \partial_{\Phi_{-k}} \mathcal{L}(\Phi_*, \lambda_*, \theta, \mathbf{u}) = \nabla_1 \ell[\Phi_{-k}, \theta, -k] + \partial_1 \mathcal{M}_{H,\delta}[\Phi_{-k}, \theta, \mathbf{u}_{-k}]^{\top} \cdot \lambda_{k-1} - \lambda_k = 0,$$

987 Finally, the total derivative of L with respect to θ reads as:

$$\begin{aligned} d_{\theta}L &= d_{\theta} \mathcal{L}(\Phi_*, \lambda_*, \theta, \mathbf{u}) \\ &= \partial_{\theta} \mathcal{L}(\Phi_*, \lambda_*, \theta, \mathbf{u}) + \underbrace{\partial_{\theta} \Phi_*^{\top} \cdot \partial_{\Phi} \mathcal{L}(\Phi_*, \lambda_*, \theta, \mathbf{u})}_{=0} + \underbrace{\partial_{\theta} \lambda_*^{\top} \cdot \partial_{\lambda} \mathcal{L}(\Phi_*, \lambda_*, \theta, \mathbf{u})}_{=0} \\ &= \sum_{k=0}^{K-1} \nabla_2 \ell[\Phi_{-k}, \theta, -k] + \partial_2 \mathcal{M}_{H,\delta}[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}]^{\top} \cdot \lambda_k \end{aligned}$$

988 The total derivative of L with respect to \mathbf{u}_{-k} is derived in the exact same fashion. \square

989 **Remark 5.** Note that using the vectorized notations introduced in subsection 3.3, the Lagrangian of
990 the constrained optimization problem defined in Def. A.4 re-writes:

$$\mathcal{L} = \mathbf{1}^{\top} \cdot \ell[\tilde{\Phi}, \theta] + \text{Tr} \left[\left(\mathcal{M}_{H,\delta}[\tilde{\Phi}, \theta, \tilde{\mathbf{u}}] - \tilde{\Phi} \right) \cdot \tilde{\lambda}^{\top} \right]$$

991 with Tr denoting the trace matrix operator and:

$$\begin{aligned} \mathbf{1} &:= \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^{K \times 1}, \quad \ell[\tilde{\Phi}, \theta] := \begin{bmatrix} \ell[\Phi_0, \theta, 0] \\ \vdots \\ \ell[\Phi_{-(K-1)}, \theta, -(K-1)] \end{bmatrix} \in \mathbb{R}^{K \times 1}, \\ \mathcal{M}_{H,\delta}[\tilde{\Phi}, \theta, \tilde{\mathbf{u}}] &:= \begin{bmatrix} \mathcal{M}_{H,\delta}[\Phi_0, \theta, \mathbf{u}] \\ \vdots \\ \mathcal{M}_{H,\delta}[\Phi_{-(K-1)}, \theta, \mathbf{u}] \end{bmatrix} \in \mathbb{R}^{K \times d_{\Phi}} \end{aligned}$$

992 **Detailed BPTT.** For the needs of our derivation of RHEL in discrete time, we now introduce a
 993 finer-grained version of BPTT given model assumptions given in Def. A.3.

Corollary A.5 (Detailed BPTT). *Given assumptions in Def. A.3–A.4, the gradients of the loss with respect to the parameters θ and the inputs \mathbf{u}_{-k} are given by:*

$$d_{\theta}L = \sum_{k=0}^{K-1} g_{\theta}(k), \quad d_{\mathbf{u}_{-k}}L = g_{\mathbf{u}}(k),$$

with:

$$\begin{aligned} g_{\theta}(k) &:= \nabla_2 \ell[\Phi_{-k}, \theta, -k] + \frac{\delta}{2} \nabla_{1,2}^2 T[\Phi_{-(k+1/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_k \\ &+ \delta \nabla_{1,2}^2 V[\Phi_{-(k+2/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_{k+1/3} + \frac{\delta}{2} \nabla_{1,2}^2 T[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_{k+2/3}, \end{aligned}$$

$$\begin{aligned} g_{\mathbf{u}}(k) &:= \frac{\delta}{2} \nabla_{1,3}^2 T[\Phi_{-(k+1/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_k \\ &+ \delta \nabla_{1,3}^2 V[\Phi_{-(k+2/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_{k+1/3} + \frac{\delta}{2} \nabla_{1,3}^2 T[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_{k+2/3}, \end{aligned}$$

and where (λ_k) satisfy the following recursion relationship, with $\lambda_0 = \nabla_{\Phi} \ell[\Phi_0]$ and $\forall k \in [0, K-1]$:

$$\begin{cases} \lambda_{k+1/3} &= \lambda_k + \frac{\delta}{2} \nabla_1^2 T[\Phi_{-(k+1/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_k \\ \lambda_{k+2/3} &= \lambda_{k+1/3} + \delta \nabla_1^2 V[\Phi_{-(k+2/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_{k+1/3} \\ \lambda_{k+1} &= \lambda_{k+2/3} + \frac{\delta}{2} \nabla_1^2 T[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \lambda_{k+2/3} + \nabla_1 \ell[\Phi_{-(k+1)}, \theta] \end{cases}$$

994

995 *Proof of Corollary A.5.* Direct application of Theorem A.3 with the inference computational graph
 996 details defined inside Def. A.3. \square

997 **A.2.3 Proof of Theorem 3.2**

998 **Time reversibility in discrete time.** We first derive the discrete counterpart of Lemma A.3 as a
 999 technical pre-requisite for the extension of RHEL to the discrete-time setting.

Lemma A.4 (Reversibility of $\mathcal{M}_{H,\delta}$). *Under the assumptions of Def. A.3–A.4:*

$$\forall \Phi_k \in \mathbb{R}^{d_\Phi}, \forall \theta \in \mathbb{R}^{d_\theta}, \forall u \in \mathbb{R}^{d_u} : \quad \Phi_{k+1} = \mathcal{M}_{H,\delta}[\Phi_k, \theta, u] \Rightarrow \mathcal{M}_{H,\delta}[\Phi_{k+1}^*, \theta, u] = \Phi_k^*$$

1001 *Proof of Lemma A.4.* Let $\Phi_k, \Phi_{k+1} \in \mathbb{R}^{d_\Phi}$ be such that:

$$\Phi_{k+1} = \mathcal{M}_{H,\delta}[\Phi_k, \theta, u]$$

1002 which rewrites, given Def. A.3:

$$\mathcal{M}_{H,\delta} : \begin{cases} \phi_{k+1/3} &= \phi_k + \frac{\delta}{2} \nabla_\pi T[\pi_k, \theta, u] \\ \pi_{k+1/3} &= \pi_k \\ \phi_{k+2/3} &= \phi_{k+1/3} \\ \pi_{k+2/3} &= \pi_{k+1/3} - \delta \nabla_\phi V[\phi_{k+1/3}, \theta, u] \\ \phi_{k+1} &= \phi_{k+2/3} + \frac{\delta}{2} \nabla_\pi T[\pi_{k+2/3}, \theta, u] \\ \pi_{k+1} &= \pi_{k+2/3} \end{cases} \quad (20)$$

1003 It becomes apparent from Eq. (20) that $\mathcal{M}_{H,\delta}$ is invertible with respect to its first argument and that
 1004 inverting $\mathcal{M}_{H,\delta}$ amounts to change δ to $-\delta$:

$$\mathcal{M}_{H,\delta}^{-1} : \begin{cases} \phi_{k+2/3} &= \phi_{k+1} - \frac{\delta}{2} \nabla_\pi T[\pi_{k+1}, \theta, u] \\ \pi_{k+2/3} &= \pi_{k+1} \\ \phi_{k+1/3} &= \phi_{k+2/3} \\ \pi_{k+1/3} &= \pi_{k+2/3} + \delta \nabla_\phi V[\phi_{k+2/3}, \theta, u] \\ \phi_k &= \phi_{k+1/3} - \frac{\delta}{2} \nabla_\pi T[\pi_{k+1/3}, \theta, u] \\ \pi_k &= \pi_{k+1/3} \end{cases} \quad (21)$$

1005 and therefore:

$$\Phi_k = \mathcal{M}_{H,\delta}^{-1}[\Phi_{k+1}, \theta, u] = \mathcal{M}_{H,-\delta}[\Phi_{k+1}, \theta, u],$$

1006 Denoting $\pi^* := -\pi$, note that by time-reversal invariance hypothesis in Def. A.3 and Lemma A.2,
 1007 we have that $\nabla_\pi T[\pi, \theta, u] = -\nabla_{\pi^*} T[\pi^*, \theta, u]$. Therefore, Eq. (21) rewrites as:

$$\begin{cases} \phi_{k+2/3} &= \phi_{k+1} + \frac{\delta}{2} \nabla_{\pi^*} T[\pi_{k+1}^*, \theta, u] \\ \pi_{k+2/3}^* &= \pi_{k+1}^* \\ \phi_{k+1/3} &= \phi_{k+2/3} \\ \pi_{k+1/3}^* &= \pi_{k+2/3}^* - \delta \nabla_\phi V[\phi_{k+2/3}, \theta, u] \\ \phi_k &= \phi_{k+1/3} + \frac{\delta}{2} \nabla_{\pi^*} T[\pi_{k+1/3}^*, \theta, u] \\ \pi_k^* &= \pi_{k+1/3}^* \end{cases}, \quad (22)$$

1008 where equations bearing on π have been multiplied on both sides by -1 . Finally note that Eq. (22)
 1009 simply rewrites as:

$$\mathcal{M}_{H,\delta}[\Phi_{k+1}^*, \theta, u] = \Phi_k^*$$

Corollary A.6. *Under the assumptions of Def. A.3, if Φ satisfies the following recursive equations:*

$$\Phi_{-K} = x, \quad \forall k = -K, \dots, -1 : \quad \Phi_{k+1} = \mathcal{M}_{H,\delta}[\Phi_k, \theta, u_k]$$

and Φ^e is subsequently defined as:

$$\Phi_0^e = \Phi_0^*, \quad \forall t = 0, \dots, K-1 : \quad \Phi_{k+1}^e = \mathcal{M}_{H,\delta}[\Phi_k^e, \theta, u_{-(k+1)}]$$

Then:

$$\forall t = 0, \dots, K-1 : \quad \Phi_k^e = \Phi_{-k}^*$$

1012 *Proof of Corollary A.6.* This result is immediately obtained by iterating Lemma A.4 over the whole
 1013 trajectory. \square

1014 **A technical pre-requisite.** Finally, we need one last technical Lemma to handle subtleties pertaining
 1015 to Jacobian evaluation which only occur in discrete time.

Lemma A.5. *Under the assumptions of Def. A.3, if we have, for some $\theta \in \mathbb{R}^{d_\theta}$ and $u \in \mathbb{R}^{d_u}$:*

$$\Phi_{k+1} = \mathcal{M}_{H,\delta}[\Phi_k, \theta, u],$$

then:

$$\begin{cases} T[\Phi_{k+1/3}, \theta, u] = T[\Phi_k, \theta, u], \\ V[\Phi_{k+2/3}, \theta, u] = V[\Phi_{k+1/3}, \theta, u] \\ T[\Phi_{k+1}, \theta, u] = T[\Phi_{k+2/3}, \theta, u] \end{cases}$$

1016

1017 *Proof.* This can be seen by simply writing $\mathcal{M}_{H,\delta}$ explicitly for ϕ and π :

$$\mathcal{M}_{H,\delta} : \begin{cases} \phi_{k+1/3} &= \phi_k + \frac{\delta}{2} \nabla_\pi T[\pi_k, \theta, u] \\ \pi_{k+1/3} &= \pi_k \\ \phi_{k+2/3} &= \phi_{k+1/3} \\ \pi_{k+2/3} &= \pi_{k+1/3} - \delta \nabla_\phi V[\phi_{k+1/3}, \theta, u] \\ \phi_{k+1} &= \phi_{k+2/3} + \frac{\delta}{2} \nabla_\pi T[\pi_{k+2/3}, \theta, u] \\ \pi_{k+1} &= \pi_{k+2/3} \end{cases}$$

1018

□

1019 **Discrete-time RHEL.** We are now ready to state the main result of this section.

Corollary A.7. *Under the assumptions of Def. A.3–A.4, let $(\Phi_k)_k$ satisfy the recursive equation:*

$$\Phi_{-K} = x, \quad \Phi_{k+1} = \mathcal{M}_{H,\delta}[\Phi_k, \theta, u_k] \quad \forall k = -K, \dots, -1,$$

and let Φ^e satisfy:

$$\begin{cases} \Phi_0^e(\epsilon) = \Phi_0^* + \epsilon \Sigma_x \cdot \nabla_\Phi \ell[\Phi_0, \theta, 0], \\ \forall k = 0, \dots, K-1 : \\ \quad \Phi_{k+1/3}^e(\epsilon) = \mathcal{M}_{T,\delta/2}[\Phi_k^e(\epsilon), \theta, u_{-(k+1)}] \\ \quad \Phi_{k+2/3}^e(\epsilon) = \mathcal{M}_{V,\delta}[\Phi_{k+1/3}^e(\epsilon), \theta, u_{-(k+1)}] \\ \quad \Phi_{k+1}^e(\epsilon) = \mathcal{M}_{T,\delta/2}[\Phi_{k+2/3}^e(\epsilon), \theta, u_{-(k+1)}] - \epsilon J \cdot \nabla_{\Phi^e} \ell[\Phi_{k+1}^e, \theta, -(k+1)], \end{cases}$$

Then defining:

$$\begin{aligned} H^{1/2}[\Phi_k^e(\epsilon), \theta, u_{-(k+1)}] &:= \frac{1}{2} \left(H[\Phi_{k+1/3}^e(\epsilon), \theta, u_{-(k+1)}] + H[\Phi_{k+2/3}^e(\epsilon), \theta, u_{-(k+1)}] \right), \\ \Delta_\theta^{\text{RHEL}}(k, \epsilon) &:= \nabla_2 \ell[\Phi_k^e(\epsilon), \theta, -k] \\ &\quad - \frac{\delta}{2\epsilon} \left(\nabla_2 H^{1/2}[\Phi_k^e(\epsilon), \theta, u_{-(k+1)}] - \nabla_2 H^{1/2}[\Phi_k^e(-\epsilon), \theta, u_{-(k+1)}] \right), \\ \Delta_u^{\text{RHEL}}(k, \epsilon) &:= -\frac{\delta}{2\epsilon} \left(\nabla_3 H^{1/2}[\Phi_k^e(\epsilon), \theta, u_{-(k+1)}] - \nabla_3 H^{1/2}[\Phi_k^e(-\epsilon), \theta, u_{-(k+1)}] \right), \\ \Delta_\Phi^{\text{RHEL}}(k, \epsilon) &:= \frac{1}{2\epsilon} \Sigma_x \cdot (\Phi_k^e(\epsilon) - \Phi_k^e(-\epsilon)), \end{aligned}$$

we have:

$$\begin{aligned} \forall k = 0, \dots, K-1 : \quad \lambda_k &= \lim_{\epsilon \rightarrow 0} \Delta_\theta^{\text{RHEL}}(k, \epsilon), \quad g_\theta(k) = \lim_{\epsilon \rightarrow 0} \Delta_\theta^{\text{RHEL}}(k, \epsilon), \\ g_u(k) &= \lim_{\epsilon \rightarrow 0} \Delta_u^{\text{RHEL}}(k, \epsilon), \end{aligned}$$

where $(\lambda_k)_{k \in \llbracket 0, K \rrbracket}$, $(g_\theta(k))_{k \in \llbracket 0, K-1 \rrbracket}$ and $(g_u(k))_{k \in \llbracket 0, K-1 \rrbracket}$ are defined inside Corollary (A.5).

1020

1021 *Proof of Corollary A.7.* Let $k \in [0, K - 1]$. We define:

$$\Delta_{\Phi}^{\text{RHEL}}(k) := \lim_{\epsilon \rightarrow 0} \Delta_{\Phi}^{\text{RHEL}}(k, \epsilon) = \Sigma_x \cdot \partial_{\epsilon} \Phi(k, \epsilon)|_{\epsilon=0}$$

$$\Delta_{\theta}^{\text{RHEL}}(k) := \lim_{\epsilon \rightarrow 0} \Delta_{\theta}^{\text{RHEL}}(k, \epsilon) = \nabla_2 \ell[\Phi_k^e(0), \theta, -k] - \delta \partial_{\epsilon} \left(\nabla_2 H^{1/2}[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \right) \Big|_{\epsilon=0}$$

$$\Delta_{\mathbf{u}}^{\text{RHEL}}(k) := \lim_{\epsilon \rightarrow 0} \Delta_{\mathbf{u}}^{\text{RHEL}}(k, \epsilon) = -\delta \partial_{\epsilon} \left(\nabla_3 H^{1/2}[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \right) \Big|_{\epsilon=0}$$

1022 **Derivation of $\lambda_k = \lim_{\epsilon \rightarrow 0} \Delta_{\theta}^{\text{RHEL}}(k, \epsilon)$.** We proceed exactly as in Theorem A.2 with some subtle
 1023 adaptations which we highlight. Differentiating the dynamics of Φ^e between k and $k + 1/3$ and
 1024 proceeding as in the proof of Theorem A.2 using Lemma A.1, Lemma A.2 and Corollary A.6 (as the
 1025 discrete counterpart of Lemma A.3 which was used for Theorem A.2), we obtain:

$$\Delta_{\Phi}^{\text{RHEL}}(k + 1/3) = \Delta_{\Phi}^{\text{RHEL}}(k) + \frac{\delta}{2} \nabla_{\Phi}^2 T[\Phi_{-k}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \Delta_{\Phi}^{\text{RHEL}}(k)$$

1026 However note that this does not correctly match the dynamics satisfied by λ inside Corollary A.5
 1027 between k and $k + 1/3$: the Hessian $\nabla_{\Phi}^2 T$ should instead be evaluated at $\Phi_{-(k+1/3)}$. Fortunately,
 1028 using Lemma A.5:

$$\nabla_{\Phi}^2 T[\Phi_{-k}, \theta, \mathbf{u}_{-(k+1)}] = \nabla_{\Phi}^2 T[\Phi_{-(k+1/3)}, \theta, \mathbf{u}_{-(k+1)}]$$

1029 Therefore we get:

$$\Delta_{\Phi}^{\text{RHEL}}(k + 1/3) = \Delta_{\Phi}^{\text{RHEL}}(k) + \frac{\delta}{2} \nabla_{\Phi}^2 T[\Phi_{-(k+1/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \Delta_{\Phi}^{\text{RHEL}}(k)$$

1030 Proceeding the same way on $\Phi_{k+2/3}^e$ and Φ_{k+1}^e , we get altogether:

$$\begin{cases} \Delta_{\Phi}^{\text{RHEL}}(k + 1/3) = \Delta_{\Phi}^{\text{RHEL}}(k) + \frac{\delta}{2} \nabla_{\Phi}^2 T[\Phi_{-(k+1/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \Delta_{\Phi}^{\text{RHEL}}(k) \\ \Delta_{\Phi}^{\text{RHEL}}(k + 2/3) = \Delta_{\Phi}^{\text{RHEL}}(k + 1/3) + \delta \nabla_{\Phi}^2 V[\Phi_{-(k+2/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \Delta_{\Phi}^{\text{RHEL}}(k + 1/3) \\ \Delta_{\Phi}^{\text{RHEL}}(k + 1) = \Delta_{\Phi}^{\text{RHEL}}(k + 2/3) + \frac{\delta}{2} \nabla_{\Phi}^2 T[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^{\top} \cdot \Delta_{\Phi}^{\text{RHEL}}(k + 2/3) \\ \quad + \nabla_{\Phi} \ell[\Phi_{-(k+1)}, \theta, -(k+1)] \end{cases}$$

1031 $\Delta_{\Phi}^{\text{RHEL}}$ satisfying the same equations as $(\lambda_k)_k$ given by BPTT, together with same initial conditions:

$$\Delta_{\Phi}^{\text{RHEL}}(0) = \nabla_{\Phi} \ell[\Phi_0, \theta, 0]$$

1032 yields the desired equality.

1033 **Derivation of $g_{\theta}(k) = \lim_{\epsilon \rightarrow 0} \Delta_{\theta}^{\text{RHEL}}(k, \epsilon)$.** Proceeding in the same way as in the derivation of
 1034 Theorem A.2, starting from the expression of $g_{\theta}(k)$ derived in Corollary A.5, using $\Delta_{\Phi}^{\text{RHEL}}(k) = \lambda_k$
 1035 and paying attention to evaluating Jacobian at the right places using Lemma A.5, we obtain $\forall k =$
 1036 $0, \dots, K - 1$:

$$g_{\theta}(k) = \nabla_2 \ell[\Phi_k^e(0), \theta, -k] - \frac{\delta}{2} \partial_{\epsilon} \left\{ \nabla_{\theta} T[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] + 2 \nabla_{\theta} V[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] + \nabla_{\theta} T[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \right\} \Big|_{\epsilon=0}.$$

1037 There again, using Lemma A.5:

$$\nabla_{\theta} T[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] = \nabla_{\theta} T[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}], \quad \nabla_{\theta} V[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] = \nabla_{\theta} V[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}],$$

1038 therefore:

$$\begin{aligned} g_{\theta}(k) &= \nabla_2 \ell[\Phi_k^e(0), \theta, -k] - \frac{\delta}{2} \partial_{\epsilon} \left\{ \nabla_{\theta} T[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] + \nabla_{\theta} V[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \right. \\ &\quad \left. + \nabla_{\theta} V[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] + \nabla_{\theta} T[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \right\} \Big|_{\epsilon=0} \\ &= \nabla_2 \ell[\Phi_k^e(0), \theta, -k] - \frac{\delta}{2} \partial_{\epsilon} \left(H[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] + H[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \right) \Big|_{\epsilon=0} \\ &= \nabla_2 \ell[\Phi_k^e(0), \theta, -k] - \delta \partial_{\epsilon} \left(H^{1/2}[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \right) \Big|_{\epsilon=0} = \lim_{\epsilon \rightarrow 0} \Delta_{\theta}^{\text{RHEL}}(k, \epsilon) \end{aligned}$$

1039 **Derivation of** $g_u(k) = \lim_{\epsilon \rightarrow 0} \Delta_u^{\text{RHEL}}(k, \epsilon)$. Strictly identical to the above paragraph. \square

1040 **Remark 6.** Note that the echo dynamics prescribed by Corollary A.7 are *implicit*:

$$\begin{cases} \Phi_0^e(\epsilon) = \Phi_0^* + \epsilon \Sigma_x \cdot \nabla_{\Phi} \ell[\Phi_0, \theta, 0], \\ \forall k = 0, \dots, K-1 : \\ \quad \Phi_{k+1/3}^e(\epsilon) = \mathcal{M}_{T, \delta/2}[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \\ \quad \Phi_{k+2/3}^e(\epsilon) = \mathcal{M}_{V, \delta}[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \\ \quad \Phi_{k+1}^e(\epsilon) = \mathcal{M}_{T, \delta/2}[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] - \epsilon \mathbf{J} \cdot \nabla_{\Phi^e} \ell[\Phi_{k+1}^e, \theta, -(k+1)], \end{cases}$$

1041 where Φ_{k+1}^e appears, as highlighted in red, on both sides of the last step. A straightforward way to
1042 make the echo dynamics explicit while preserving the theoretical guarantees of Corollary A.7 is to
1043 linearize the nudging signal, namely using instead the following set of equations:

$$\begin{cases} \Phi_0^e(\epsilon) = \Phi_0^* + \epsilon \Sigma_x \cdot \nabla_{\Phi} \ell[\Phi_0, \theta, 0], \\ \forall k = 0, \dots, K-1 : \\ \quad \Phi_{k+1/3}^e(\epsilon) = \mathcal{M}_{T, \delta/2}[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \\ \quad \Phi_{k+2/3}^e(\epsilon) = \mathcal{M}_{V, \delta}[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \\ \quad \Phi_{k+1}^e(\epsilon) = \mathcal{M}_{T, \delta/2}[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] + \epsilon \Sigma_x \cdot \nabla_{\Phi^e} \ell[\Phi_{-(k+1)}, \theta, -(k+1)], \end{cases}$$

1044 Note that the $\epsilon \mathbf{J}$ of the original implicit equation becomes $-\epsilon \Sigma_x$ in its linearized counterpart.

1045 This remark leads us to a slight variant of Corollary A.7.

Corollary A.8. Under the assumptions of Def. A.3–A.4, let $(\Phi_k)_k$ satisfy the recursive equation:

$$\Phi_{-K} = \mathbf{x}, \quad \Phi_{k+1} = \mathcal{M}_{H, \delta}[\Phi_k, \theta, \mathbf{u}_k] \quad \forall k = -K, \dots, -1,$$

and let Φ^e satisfy:

$$\begin{cases} \Phi_0^e(\epsilon) = \Phi_0^* + \epsilon \Sigma_x \cdot \mathbf{y}_0, \\ \forall k = 0, \dots, K-1 : \\ \quad \Phi_{k+1/3}^e(\epsilon) = \mathcal{M}_{T, \delta/2}[\Phi_k^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \\ \quad \Phi_{k+2/3}^e(\epsilon) = \mathcal{M}_{V, \delta}[\Phi_{k+1/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] \\ \quad \Phi_{k+1}^e(\epsilon) = \mathcal{M}_{T, \delta/2}[\Phi_{k+2/3}^e(\epsilon), \theta, \mathbf{u}_{-(k+1)}] + \epsilon \Sigma_x \cdot \mathbf{y}_{-(k+1)}, \end{cases}$$

where $\mathbf{y} \in \mathbb{R}^{K \times d_{\Phi}}$ does not depend on Φ . Then the same conclusions as Corollary A.7 hold, with $(\lambda_k)_{k \in \llbracket 0, K-1 \rrbracket}$ satisfying $\lambda_0 = \mathbf{y}_0$ and $\forall k \in \llbracket 0, K-1 \rrbracket$:

$$\begin{cases} \lambda_{k+1/3} &= \lambda_k + \frac{\delta}{2} \nabla_1^2 T[\Phi_{-(k+1/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^\top \cdot \lambda_k \\ \lambda_{k+2/3} &= \lambda_{k+1/3} + \delta \nabla_1^2 V[\Phi_{-(k+2/3)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^\top \cdot \lambda_{k+1/3} \\ \lambda_{k+1} &= \lambda_{k+2/3} + \frac{\delta}{2} \nabla_1^2 T[\Phi_{-(k+1)}, \theta, \mathbf{u}_{-(k+1)}] \cdot \mathbf{J}^\top \cdot \lambda_{k+2/3} + \mathbf{y}_{-(k+1)} \end{cases}$$

1046

1047 *Proof Corollary A.8.* Identical to the proof of Corollary A.7. \square

1048 **A.2.4 Proof of Theorem 3.3**

Theorem A.4. Assuming a HSSM (Def. A.5) and the optimization problem depicted in Def. A.6, we have:

$$\forall \ell = 0, \dots, N-1 : \quad d_{\theta^{(\ell)}} L = \lim_{\epsilon \rightarrow 0} \Delta \theta^{(\ell)}(\epsilon),$$

where $\Delta \theta^{(\ell)}(\epsilon)$ can be recursively computed backwards, starting from the **top-most block** as:

$$\begin{aligned} \bar{\Phi}^{e,(N)}(\epsilon) &= \mathcal{M}_{H^{(N-1)}, \delta}[\bar{\Phi}^{e,(N)}(\epsilon), \theta^{(N-1)}, \tilde{\Phi}^{(N-1)}] + \epsilon \Sigma_x \cdot \nabla_{\Phi^{(N)}} L \\ \Delta \theta^{(N-1)}(\epsilon) &= \sum_{k=0}^{K-1} \nabla_2 \ell[\Phi_k^e, \theta^{(N-1)}, -k] \\ &\quad - \frac{\delta}{2\epsilon} \sum_{k=0}^{K-1} \left(\nabla_2 H^{1/2}[\Phi_k^{e,(L)}(\epsilon), \theta^{(N-1)}, \Phi_{-k}^{(N-1)}] - \nabla_2 H^{1/2}[\Phi_k^{e,(N)}(-\epsilon), \theta^{(N-1)}, \Phi_{-k}^{(N-1)}] \right), \\ \bar{\Delta}_{\Phi^{(N-1)}}(\epsilon) &= -\frac{\delta}{2\epsilon} \left(\nabla_3 H^{1/2}[\bar{\Phi}^{e,(N)}(\epsilon), \theta^{(N-1)}, \tilde{\Phi}^{(N-1)}] - \nabla_3 H^{1/2}[\bar{\Phi}^{e,(L)}(-\epsilon), \theta^{(N-1)}, \tilde{\Phi}^{(N-1)}] \right) \end{aligned}$$

and subsequently for **upstream blocks**, i.e. $\forall \ell = N-2, \dots, 0$:

$$\begin{aligned} \bar{\Phi}^{e,(\ell+1)}(\epsilon) &= \mathcal{M}_{H^{(\ell)}, \delta}[\bar{\Phi}^{e,(\ell+1)}(\epsilon), \theta^{(\ell)}, \tilde{\Phi}^{(\ell)}] + \epsilon \Sigma_x \cdot \tilde{\Delta}_{\Phi^{(\ell+1)}}(\epsilon) \\ \Delta \theta^{(\ell)}(\epsilon) &= -\frac{\delta}{2\epsilon} \sum_{k=0}^{K-1} \left(\nabla_2 H^{1/2}[\Phi_k^{e,(\ell+1)}(\epsilon), \theta^{(\ell)}, \Phi_{-k}^{(\ell)}] - \nabla_2 H^{1/2}[\Phi_k^{e,(\ell+1)}(-\epsilon), \theta^{(\ell)}, \Phi_{-k}^{(\ell)}] \right) \\ \Delta_{\Phi^{(\ell)}}(\epsilon) &= -\frac{\delta}{2\epsilon} \left(\nabla_3 H^{1/2}[\bar{\Phi}^{e,(\ell+1)}(\epsilon), \theta^{(\ell)}, \tilde{\Phi}^{(\ell)}] - \nabla_3 H^{1/2}[\bar{\Phi}^{e,(\ell+1)}(-\epsilon), \theta^{(\ell)}, \tilde{\Phi}^{(\ell)}] \right) \end{aligned}$$

1049

1050 *Proof.* Re-using the vectorized notations introduced in subsection 3.3 and used in Remark 5, the
1051 Lagrangian associated to the optimization problem depicted in Def. A.6 reads:

$$\begin{aligned} \mathcal{L} &= \mathbf{1}^\top \cdot \ell[\tilde{\Phi}^{(L)}, \theta^{(L-1)}] + \sum_{\ell=0}^{L-1} \text{Tr} \left[\left(\mathcal{M}_{H^{(\ell)}, \delta}[\tilde{\Phi}^{(\ell+1)}, \theta^{(\ell)}, \tilde{\mathbf{u}}^{(\ell)}] - \tilde{\Phi}^{(\ell+1)} \right) \cdot \left(\bar{\lambda}^{(\ell+1)} \right)^\top \right] \\ &= \sum_{k=0}^{K-1} \ell[\Phi_{-k}^{(L)}, \theta^{(L)}, -k] + \sum_{\ell=0}^{L-1} \left(\lambda_k^{(\ell+1)} \right)^\top \cdot \left(\mathcal{M}_{H^{(\ell)}, \delta}[\Phi_{-(k+1)}^{(\ell+1)}, \theta^{(\ell)}, \Phi_{-(k+1)}^{(\ell)}] - \Phi_{-k}^{(\ell+1)} \right) \end{aligned}$$

1052 We proceed by induction on the block index starting from $\ell = L$.

1053 **Initialization** ($\ell = L$). Let $\Phi_k^{(N)}$ and $\lambda_k^{(N)}$ for $k \in \llbracket 0, K-1 \rrbracket$ be the critical points of \mathcal{L} . By
1054 Theorem A.3:

$$d_{\theta^{(N-1)}} L = \sum_{k=0}^{K-1} \nabla_2 \ell[\Phi_{-k}^{(N)}, \theta^{(N-1)}, -k] + \partial_2 \mathcal{M}_{H^{(N-1)}, \delta}[\Phi_{-(k+1)}^{(N)}, \theta^{(N-1)}, \Phi_{-(k+1)}^{(N-1)}]^\top \cdot \lambda_k^{(N)},$$

1055 with $(\lambda_k^{(N)})_{k \in \llbracket 0, K-1 \rrbracket}$ satisfying the following recursion relationship:

$$\begin{cases} \lambda_0^{(N)} = \nabla_1 \ell[\Phi_0^{(N)}, \theta^{(N-1)}, 0], \\ \forall k = 0, \dots, K-1 : \\ \lambda_{k+1}^{(N)} = \partial_1 \mathcal{M}_{H^{(N-1)}, \delta}[\Phi_{-(k+1)}^{(N)}, \theta^{(N-1)}, \Phi_{-(k+1)}^{(N-1)}]^\top \cdot \lambda_k^{(N)} + \nabla_1 \ell[\Phi_{-(k+1)}^{(N)}, \theta, -(k+1)] \end{cases}$$

1056 Given the definition of the dynamics of $\Phi^{e,(N)}$ by hypothesis, we can directly apply Corollary A.7 to
1057 obtain:

$$\begin{aligned} d_{\theta^{(N-1)}} L &= \sum_{k=0}^{K-1} \nabla_2 \ell[\Phi_{-k}^{(N)}, \theta^{(N-1)}, -k] - \delta \partial_\epsilon \left(\nabla_2 H^{(N-1), 1/2}[\Phi_k^{e,(N)}(\epsilon), \theta^{(L-1)}, \Phi_{-(k+1)}^{(N-1)}] \right) \Big|_{\epsilon=0}, \\ d_{\Phi_{k-1}^{(N-1)}} L &= \partial_3 \mathcal{M}_{H^{(N-1)}, \delta}[\Phi_{-(k+1)}^{(N)}, \theta^{(N-1)}, \Phi_{-(k+1)}^{(N-1)}]^\top \cdot \lambda_k^{(N)} \\ &= -\delta \partial_\epsilon \left(\nabla_3 H^{(N-1), 1/2}[\Phi_k^{e,(N)}(\epsilon), \theta^{(L-1)}, \Phi_{-(k+1)}^{(N-1)}] \right) \Big|_{\epsilon=0} \end{aligned}$$

1058 **Induction** ($\ell + 1 \rightarrow \ell$). Let us assume that the desired property is satisfied at layer $\ell + 1$. We denote
 1059 again $\Phi_k^{(\ell)}$ and $\lambda_k^{(\ell)}$ for $k \in \llbracket 0, K - 1 \rrbracket$ the critical point of \mathcal{L} . We have, for $k \in \llbracket 0, K - 1 \rrbracket$:

$$\begin{cases} \lambda_0^{(\ell)} = \partial_3 \mathcal{M}_{H^{(\ell)}, \delta} \left[\Phi_0^{(\ell+1)}, \theta^{(\ell)}, \Phi_0^{(\ell)} \right]^\top \cdot \lambda_0^{(\ell+1)} \\ \forall k = 0, \dots, K - 1 : \\ \lambda_{k+1}^{(\ell)} = \partial_1 \mathcal{M}_{H^{(\ell-1)}, \delta} \left[\Phi_{-(k+1)}^{(\ell)}, \theta^{(\ell-1)}, \Phi_{-(k+1)}^{(\ell-1)} \right]^\top \cdot \lambda_k^{(\ell)} + \partial_3 \mathcal{M}_{H^{(\ell)}, \delta} \left[\Phi_{-(k+1)}^{(\ell+1)}, \theta^{(\ell)}, \Phi_{-(k+1)}^{(\ell)} \right]^\top \cdot \lambda_k^{(\ell+1)} \end{cases}$$

1060 Using the induction hypothesis at layer $(\ell + 1)$:

$$\begin{aligned} \partial_3 \mathcal{M}_{H^{(\ell)}, \delta} \left[\Phi_{-(k+1)}^{(\ell+1)}, \theta^{(\ell)}, \Phi_{-(k+1)}^{(\ell)} \right]^\top \cdot \lambda_k^{(\ell+1)} &= -\delta \partial_\epsilon \left(\nabla_3 H^{(\ell), 1/2} \left[\Phi_k^{e, (\ell+1)}(\epsilon), \theta^{(\ell)}, \Phi_{-(k+1)}^{(\ell)} \right] \right) \Big|_{\epsilon=0} \\ &= \lim_{\epsilon \rightarrow 0} \Delta_{\Phi_{(k+1)}^{(\ell)}}(\epsilon) \end{aligned}$$

1061 Therefore on the one hand, denoting $\Delta_{\Phi^{(\ell)}} := \lim_{\epsilon \rightarrow 0} \Delta_{\Phi^{(\ell)}}(\epsilon) \in \mathbb{R}^{K \times d_\Phi}$, the dynamics on λ
 1062 rewrite:

$$\begin{cases} \lambda_0^{(\ell)} = \Delta_{\Phi_0^{(\ell)}} \\ \forall k = 0, \dots, K - 1 : \\ \lambda_{k+1}^{(\ell)} = \partial_1 \mathcal{M}_{H^{(\ell-1)}, \delta} \left[\Phi_{-(k+1)}^{(\ell)}, \theta^{(\ell-1)}, \Phi_{-(k+1)}^{(\ell-1)} \right]^\top \cdot \lambda_k^{(\ell)} + \Delta_{\Phi_{(k+1)}^{(\ell)}} \end{cases}$$

1063 On the other hand, the dynamics of $\Phi^{e, (\ell)}$ read by hypothesis:

$$\begin{cases} \Phi_0^{(\ell)} = \left(\Phi_0^{(\ell)} \right)^\star + \epsilon \Sigma_x \cdot \Delta_{\Phi_0^{(\ell)}}(\epsilon) \\ \forall k = 0, \dots, K - 1 : \\ \Phi_{k+1}^{e, (\ell)} = \mathcal{M}_{H^{(\ell-1)}, \delta} \left[\Phi_{k+1}^{e, (\ell)}, \theta^{(\ell-1)}, \Phi_{-(k+1)}^{(\ell-1)} \right] + \epsilon \Sigma_x \cdot \Delta_{\Phi_{(k+1)}^{(\ell)}} \end{cases}$$

1064 using Corollary A.8 with $y = \Delta_{\Phi^{(\ell)}}$, we conclude that:

$$d_{\theta^{(\ell)}} L = \lim_{\epsilon \rightarrow 0} \Delta \theta^{(\ell)}(\epsilon)$$

1065 □

1066 **Remark 7.** Theorem A.4, and more generally our definition of HSSMs (Def. A.5) assume that the
 1067 connectivity pattern of the HRU units is a linear chain. Note that while we chose this hypothesis for
 1068 the sake of clarity of our results and their derivations, Theorem A.4 could be seamlessly extended to
 1069 any Directed Acyclic Graph (DAG) of HRUs. This allows, for instance as a simple and realistic case,
 1070 to use skip connections across HRUs within HSSMs.

1071 **Remark 8.** Note that RHEL chaining as prescribed by Theorem A.4 **implicitly chains RHEL and**
 1072 **automatic differentiation.** Indeed, if H explicitly parametrizes feedforward mappings across HRUs
 1073 as:

$$H^{(\ell)} \left[\Phi_k^{e, (\ell+1)}(\epsilon), \theta^{(\ell)}, \Phi_{-(k+1)}^{(\ell)} \right] = H^{(\ell)} \left[\Phi_k^{e, (\ell+1)}(\epsilon), \theta_\alpha^{(\ell)}, F \left[\Phi_{-(k+1)}^{(\ell)}, \theta_\beta^{(\ell)} \right] \right], \quad (23)$$

1074 then, denoting $u^{(\ell)} := F \left[\Phi_{-(k+1)}^{(\ell)}, \theta_\beta^{(\ell)} \right]$, we have:

$$\begin{aligned} &\partial_\epsilon \left(\nabla_3 H^{(\ell), 1/2} \left[\Phi_k^{e, (\ell+1)}(\epsilon), \theta^{(\ell)}, \Phi_{-(k+1)}^{(\ell)} \right] \right) \Big|_{\epsilon=0} \\ &= \partial_1 F \left[\Phi_{-(k+1)}^{(\ell)}, \theta_\beta^{(\ell)} \right]^\top \cdot \partial_\epsilon \left(\nabla_3 H^{(\ell), 1/2} \left[\Phi_k^{e, (\ell+1)}(\epsilon), \theta_\alpha^{(\ell)}, u^{(\ell)} \right] \right) \Big|_{\epsilon=0} \\ &\approx \partial_1 F \left[\Phi_{-(k+1)}^{(\ell)}, \theta_\beta^{(\ell)} \right]^\top \cdot \frac{1}{2\epsilon} \left(\nabla_3 H^{(\ell), 1/2} \left[\Phi_k^{e, (\ell+1)}(\epsilon), \theta_\alpha^{(\ell)}, u^{(\ell)} \right] - \nabla_3 H^{(\ell), 1/2} \left[\Phi_k^{e, (\ell+1)}(-\epsilon), \theta_\alpha^{(\ell)}, u^{(\ell)} \right] \right) \end{aligned}$$

1075 The red part is done by automatic differentiation and the blue part by RHEL. This underpins the
 1076 implementation of RHEL chaining we used in our own code.

1077 A.3 Models and algorithms details

1078 **Summary.** In this section, we provide details about our models and algorithms. More precisely:

- 1079 • In section A.3.1, we first describe our *toy model* used inside Fig. 2 in terms of its Hamilto-
1080 nian, resulting continuous-time dynamics and RHEL gradient estimators as prescribed by
1081 Theorem 3.1.
- 1082 • In section A.3.2, we provide details about the HSSMs which we used in our experiments.
1083 We describe in greater details HSSMs made up of *linear* HRU blocks (section A.3.3). We
1084 describe their Hamiltonian, their resulting dynamics, the associated gradient estimators
1085 prescribed by RHEL and explain how *parallel scan* can be used on these models, especially
1086 when applying RHEL. Similarly, we describe HSSMs made up of *nonlinear* HRU blocks
1087 (section A.3.4).
- 1088 • In section A.3.5, we show how the time discretization δ can itself be trained by absorbing it
1089 into the definition of the Hamiltonian function.
- 1090 • Finally, in the light of Remark 8, we highlight in section A.3.6 how our implementation
1091 hybridizes *Automatic Differentiation* (AD) and RHEL using Algorithms 3–4.

1092 A.3.1 Toy model

1093 In this section, we provide the gradient estimators for the parameters of the toy model. The toy model
1094 is a simple network of six mechanically coupled oscillators. Each oscillator is described by a state
1095 $\Phi^i = (\phi^i, \pi^i)$ where $\phi^i \in \mathbb{R}$ is the position of the oscillator and $\pi^i \in \mathbb{R}$ is its momentum. It has a
1096 mass parameter m_i and spring parameter k_i . Any pair of oscillators (i, j) is coupled via the spring
1097 parameter k_{ij} . The input to the models is a time-varying external force $\mathbf{u}(t) \in \mathbb{R}$ coupled to oscillator
1098 1. During the echo passes, the nudging force is modelled by a spring coupling with parameter $\epsilon \in \mathbb{R}$
1099 to an external force $y(t) \in \mathbb{R}$. The Hamiltonian of the system is given by:

$$H[\Phi, \theta, \mathbf{u}] = \sum_i \frac{(\pi^i)^2}{2m_i} + \frac{1}{2} \sum_i k_i (\phi^i)^2 + \frac{1}{2} \sum_i \sum_{j>i} k_{ij} (\phi^j - \phi^i)^2 + \mathbf{u} \phi^1 \quad (24)$$

1100 Which gives the following equations of motion:

$$\begin{cases} \dot{\phi}^i = \frac{\pi^i}{m_i} & \text{for all } i \in \{1, 6\} \\ \dot{\pi}^i = -k_i \phi^i + \sum_{j, j \neq i} k_{ij} (\phi^j - \phi^i), & i \in \{2, 3, 5, 6\} \\ \dot{\pi}^1 = -k_1 \phi^1 + \sum_{j, j \neq 1} k_{1j} (\phi^j - \phi^1) + \mathbf{u} \\ \dot{\pi}^4 = -k_4 \phi^4 + \sum_{j, j \neq 4} k_{4j} (\phi^j - \phi^4) - \delta_e \epsilon (\phi^4 - y) \end{cases}$$

1101 where δ_e is the indicator function of the echo pass, it's equal to 1 during the echo pass and 0 otherwise.

1102 **RHEL gradient estimators of the model parameters.** For the mass m_i , we have:

$$\begin{aligned} \Delta_{m_i}^{RHEL} &= -\frac{1}{2\epsilon} (\nabla_{m_i} H[\Phi^e(t, \epsilon), \theta, \mathbf{u}] - \nabla_{m_i} H[\Phi^e(t, -\epsilon), \theta, \mathbf{u}]) \\ &= \frac{1}{2\epsilon} \left(\frac{(\pi^i(t, \epsilon))^2}{2m_i^2} - \frac{(\pi^i(t, -\epsilon))^2}{2m_i^2} \right) \end{aligned} \quad (25)$$

1103 For the spring parameters k_i , we have:

$$\begin{aligned} \Delta_{k_i}^{RHEL} &= -\frac{1}{2\epsilon} (\nabla_{k_i} H[\Phi^e(t, \epsilon), \theta, \mathbf{u}] - \nabla_{k_i} H[\Phi^e(t, -\epsilon), \theta, \mathbf{u}]) \\ &= -\frac{1}{2\epsilon} ((\phi^i(t, \epsilon))^2 - (\phi^i(t, -\epsilon))^2) \end{aligned} \quad (26)$$

$$(27)$$

1104 For the coupling parameters k_{ij} , we have:

$$\begin{aligned}\Delta_{k_{ij}}^{RHEL} &= -\frac{1}{2\epsilon} (\nabla_{k_{ij}} H[\Phi^e(t, \epsilon), \theta, \mathbf{u}] - \nabla_{k_{ij}} H[\Phi^e(t, -\epsilon), \theta, \mathbf{u}]) \\ &= -\frac{1}{2\epsilon} ((\phi^j(t, \epsilon) - \phi^i(t, \epsilon))^2 - (\phi^j(t, -\epsilon) - \phi^i(t, -\epsilon))^2)\end{aligned}\quad (28)$$

1105 **RHEL gradient estimators of the state sensitivities** For the state position sensitivities, we have:

$$\begin{aligned}\Delta_{\phi^i}^{RHEL} &= -\frac{1}{2\epsilon} \Sigma_x (\Phi^e(t, \epsilon) - \Phi^e(t, -\epsilon)) \\ &= -\frac{1}{2\epsilon} (\pi^i(t, \epsilon) - \pi^i(t, -\epsilon))\end{aligned}\quad (29)$$

1106 For the state momentum sensitivities, we have:

$$\begin{aligned}\Delta_{\pi^i}^{RHEL} &= -\frac{1}{2\epsilon} \Sigma_x (\Phi^e(t, \epsilon) - \Phi^e(t, -\epsilon)) \\ &= -\frac{1}{2\epsilon} (\phi^i(t, \epsilon) - \phi^i(t, -\epsilon))\end{aligned}\quad (30)$$

1107 A.3.2 HSSM architecture

1108 In this section, we outline the detailed architecture of a full multi-layer HSSM architecture. For the
1109 inference, we re-use the same stacking architecture of recurrent blocks and feedforward elements as
1110 the LinOSS model [29]. The model starts by encoding an input sequence $\bar{\mathbf{u}} \in \mathbb{R}^{d_u \times K}$ via an affine
1111 transformation. The transformed sequence then progresses through multiple HSSM blocks, linear (see
1112 Appendix A.3.3), or nonlinear (see Appendix A.3.4), directly followed by nonlinear transformations.
1113 These transformations include the Gaussian error linear unit (GELU) [55] and the Gated Linear Unit
1114 (GLU) [56], defined as $\text{GLU}(\mathbf{x}) = \text{sigmoid}(\mathbf{W}_1 \mathbf{x}) \circ \mathbf{W}_2 \mathbf{x}$ where $\mathbf{W}_{1,2}$ represent trainable weight
1115 matrices, accompanied by a residual connection. The sequence output from the final block undergoes
1116 a second affine transformation to produce the model output.

1117 The full linear and nonlinear HSSM is further presented in Algorithm 2. For clarity, when applying
1118 operations to sequence elements denoted with an overline (e.g., $\bar{\mathbf{u}}$), these operations are implicitly
1119 broadcast across the time dimension. Specifically, for any function f applied to $\bar{\mathbf{u}}$, we have $f(\bar{\mathbf{u}})_t =$
1120 $f(\mathbf{u}_t)$ for all time steps $t \in \{1, 2, \dots, K\}$.

1121 For the inference of the linear HSSMs, we keep the same recurrent block as LinOSS with a slight
1122 change in the integrator (see A.3.3). For the nonlinear HSSM, we replace the recurrent block by a
1123 UniCORRN recurrent block [30] for which we use the same integrator as for the linear HSSM (see
1124 A.3.4).

Algorithm 2 HSSM model

```

1: Input: Input sequence  $\bar{\mathbf{u}}$ , model type type  $\in \{\text{linear}, \text{nonlinear}\}$ 
2: Output: HSSM output sequence  $\bar{\mathbf{u}}$ 
3:  $\bar{\mathbf{u}}^{(0)} \leftarrow \mathbf{W}_{enc} \bar{\mathbf{u}} + \mathbf{b}_{enc}$ 
4: for  $\ell = 1, \dots, N$  do
5:   if type = linear then
6:      $\bar{\Phi}^{(\ell)}, \_, \_ \leftarrow \text{LINEARHRU}(\bar{\mathbf{u}}^{(\ell-1)}, 0)$  ▷ Via parallel scan
7:   else
8:      $\bar{\Phi}^{(\ell)}, \_, \_ \leftarrow \text{NONLINEARHRU}(\bar{\mathbf{u}}^{(\ell-1)}, 0)$  ▷ Sequentially
9:   end if
10:   $\bar{\mathbf{x}}^{(\ell)} \leftarrow \mathbf{C} \bar{\Phi}^{(\ell)} + \mathbf{D} \bar{\mathbf{u}}^{\ell-1}$ 
11:   $\bar{\mathbf{x}}_g^{(\ell)} \leftarrow \text{GELU}(\bar{\mathbf{x}}^{(\ell)})$ 
12:   $\bar{\mathbf{u}}^{(\ell)} \leftarrow \text{GLU}(\bar{\mathbf{x}}_g^{(\ell)} + \bar{\mathbf{u}}^{(\ell-1)})$ 
13: end for
14:  $\bar{\mathbf{o}} \leftarrow \mathbf{W}_{dec} \bar{\mathbf{u}}^{(N)} + \mathbf{b}_{dec}$ 

```

1125 A.3.3 Linear HRU Block

1126 **Hamiltonian of the recurrence.** The linear HRU block is the composition of a nonlinear spatial
 1127 transformation and a linear recurrent transformation (see Eq. 17). Here we provide more details
 1128 about the linear recurrence that is computed with the RHEL gradient estimator. The linear recurrence
 1129 is defined by the following Hamiltonian:

$$\begin{aligned} H[\Phi, \theta, u] &= T[\pi, \theta, u] + V[\phi, \theta, u] \\ &= \frac{1}{2} \|\pi\|^2 + \left(\frac{1}{2} \phi^\top A \phi - \phi^\top B u \right) \end{aligned} \quad (31)$$

1130 **Dynamics.** The dynamics of the linear HRU block are defined by the following equations:

$$\begin{cases} \dot{\phi} = \pi \\ \dot{\pi} = -A\phi + Bu \end{cases} \quad (32)$$

1131 Which, after time-discretization with the integrator defined in A.2.1, gives the following equations:

$$\begin{cases} \phi_{k+1/3} = \phi_k + \frac{\delta}{2} \nabla_{\pi} T[\pi_k, \theta, u_k] \\ \quad = \phi_k + \frac{\delta}{2} \pi_k \\ \pi_{k+1/3} = \pi_k \\ \phi_{k+2/3} = \phi_{k+1/3} \\ \pi_{k+2/3} = \pi_{k+1/3} + \delta \nabla_{\phi} V[\phi_{k+1/3}, \theta, u_k] \\ \quad = \pi_{k+1/3} - \delta A \phi_{k+1/3} + \delta B u_k \\ \phi_{k+1} = \phi_{k+2/3} + \frac{\delta}{2} \nabla_{\pi} T[\pi_{k+2/3}, \theta, u_k] \\ \quad = \phi_{k+2/3} + \frac{\delta}{2} \pi_{k+2/3} \\ \pi_{k+1} = \pi_{k+2/3} \end{cases} \quad (33)$$

1132 with the initial condition $\Phi_{-K} = (\phi_{-K}^\top, \pi_{-K}^\top)^\top = x$.

1133 For the echo passes, the initial condition are $\Phi_0^e = \Phi_0^* \pm \epsilon \Sigma_x \cdot \Delta_{\Phi} \ell[\Phi_0, 0]$. The dynamics equations
 1134 follow below, with modifications from equation 33 highlighted in blue::

$$\begin{cases} \phi_{k+1/3}^e = \phi_k^e + \frac{\delta}{2} \pi_k^e \\ \pi_{k+1/3}^e = \pi_k^e \\ \phi_{k+2/3}^e = \phi_{k+1/3}^e \\ \pi_{k+2/3}^e = \pi_{k+1/3}^e - \delta A \phi_{k+1/3}^e + \delta B u_{-(k+1)} \\ \phi_{k+1}^e = \phi_{k+2/3}^e + \frac{\delta}{2} \pi_{k+2/3}^e + \epsilon \nabla_{\pi} \ell[\Phi_{-(k+1)}] \\ \pi_{k+1}^e = \pi_{k+2/3}^e + \epsilon \nabla_{\phi} \ell[\Phi_{-(k+1)}] \end{cases} \quad (34)$$

1135 **RHEL gradient estimators.** The gradient estimators of the parameters of the linear HRU are:

$$\begin{aligned}\Delta_{\mathbf{A}}^{RHEL}(k, \epsilon) &= -\frac{\delta}{2\epsilon} \left(\nabla_{\mathbf{A}} H^{1/2}[\Phi_k^e(\epsilon), \boldsymbol{\theta}, \mathbf{u}_{-(k+1)}] - \nabla_{\mathbf{A}} H^{1/2}[\Phi_k^e(-\epsilon), \boldsymbol{\theta}, \mathbf{u}_{-(k+1)}] \right) \\ &= -\frac{\delta}{4\epsilon} \left[\left(\phi_{k+1/3}^e(\epsilon)^\top \phi_{k+1/3}^e(\epsilon) + \phi_{k+2/3}^e(\epsilon)^\top \phi_{k+2/3}^e(\epsilon) \right) \right. \\ &\quad \left. - \left(\phi_{k+1/3}^e(-\epsilon)^\top \phi_{k+1/3}^e(-\epsilon) + \phi_{k+2/3}^e(-\epsilon)^\top \phi_{k+2/3}^e(-\epsilon) \right) \right] \\ &= -\frac{\delta}{2\epsilon} \left[\left(\phi_{k+1/3}^e(\epsilon) \right)^\top \left(\phi_{k+1/3}^e(\epsilon) \right) - \left(\phi_{k+1/3}^e(-\epsilon) \right)^\top \left(\phi_{k+1/3}^e(-\epsilon) \right) \right] \quad (35)\end{aligned}$$

$$\begin{aligned}\Delta_{\mathbf{B}}^{RHEL}(k, \epsilon) &= -\frac{\delta}{2\epsilon} \left(\nabla_{\mathbf{B}} H^{1/2}[\Phi_k^e(\epsilon), \boldsymbol{\theta}, \mathbf{u}_{-(k+1)}] - \nabla_{\mathbf{B}} H^{1/2}[\Phi_k^e(-\epsilon), \boldsymbol{\theta}, \mathbf{u}_{-(k+1)}] \right) \\ &= \frac{\delta}{4\epsilon} \left[\left(\phi_{k+1/3}^e(\epsilon) + \phi_{k+2/3}^e(\epsilon) \right) \mathbf{u}_{-(k+1)}^\top \right. \\ &\quad \left. - \left(\phi_{k+1/3}^e(-\epsilon) + \phi_{k+2/3}^e(-\epsilon) \right) \mathbf{u}_{-(k+1)}^\top \right] \\ &= \frac{\delta}{2\epsilon} \left[\phi_{k+1/3}^e(\epsilon) \mathbf{u}_{-(k+1)}^\top - \phi_{k+1/3}^e(-\epsilon) \mathbf{u}_{-(k+1)}^\top \right] \quad (36)\end{aligned}$$

1136 The gradient estimator with respect to the input of the recurrent transformation is:

$$\begin{aligned}\Delta_{\mathbf{u}}^{RHEL}(k, \epsilon) &= -\frac{\delta}{2\epsilon} \left(\nabla_{\mathbf{u}} H^{1/2}[\Phi_k^e(\epsilon), \boldsymbol{\theta}, \mathbf{u}_{-(k+1)}] - \nabla_{\mathbf{u}} H^{1/2}[\Phi_k^e(-\epsilon), \boldsymbol{\theta}, \mathbf{u}_{-(k+1)}] \right) \\ &= \frac{\delta}{4\epsilon} \left[\mathbf{B}^\top \left(\phi_{k+1/3}^e(\epsilon) + \phi_{k+2/3}^e(\epsilon) \right) \right. \\ &\quad \left. - \mathbf{B}^\top \left(\phi_{k+1/3}^e(-\epsilon) + \phi_{k+2/3}^e(-\epsilon) \right) \right] \\ &= \frac{\delta}{2\epsilon} \left[\mathbf{B}^\top \phi_{k+1/3}^e(\epsilon) - \mathbf{B}^\top \phi_{k+1/3}^e(-\epsilon) \right] \quad (37)\end{aligned}$$

1137 **Parallel Scan.** Similarly to the LinOSS model [29], we can compute the recurrence of Linear HRU
1138 with a parallel scan [5] to reduce the computational time. To implement the parallel scan, we need to
1139 put the discretized dynamics in a form that can be computed in parallel. For this we vectorize the
1140 Equation 33:

$$\begin{aligned}\Phi_{k+1} &= \begin{bmatrix} \mathbf{I} & \frac{\delta}{2}\mathbf{I} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \cdot \Phi_{k+2/3} \\ &= \begin{bmatrix} \mathbf{I} & \frac{\delta}{2}\mathbf{I} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \cdot \left(\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\delta\mathbf{A} & \mathbf{I} \end{bmatrix} \cdot \Phi_{k+1/3} + \begin{bmatrix} \mathbf{0} \\ \delta\mathbf{B} \cdot \mathbf{u}_k \end{bmatrix} \right) \\ &= \begin{bmatrix} \mathbf{I} & \frac{\delta}{2}\mathbf{I} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\delta\mathbf{A} & \mathbf{I} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{I} & \frac{\delta}{2}\mathbf{I} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \cdot \Phi_k + \begin{bmatrix} \frac{\delta^2}{2}\mathbf{B} \cdot \mathbf{u}_k \\ \delta\mathbf{B} \cdot \mathbf{u}_k \end{bmatrix}\end{aligned}$$

1141 which gives us the discrete dynamics in matrix form:

$$\Phi_{k+1} = M\Phi_k + F_k \quad (38)$$

1142 with:

$$M = \begin{bmatrix} \mathbf{I} - \frac{\delta^2}{2}\mathbf{A} & \delta[\mathbf{I} - \frac{\delta^2}{4}\mathbf{A}] \\ -\delta\mathbf{A} & \mathbf{I} - \frac{\delta^2}{2}\mathbf{A} \end{bmatrix}, \quad F_k = \begin{bmatrix} \frac{\delta^2}{2}\mathbf{B} \cdot \mathbf{u}_k \\ \delta\mathbf{B} \cdot \mathbf{u}_k \end{bmatrix} \quad (39)$$

1143 To compute the echo passes $\Phi_k^e, k \in \{1, \dots, K\}$, we only need to adapt F_k to get for the positive
1144 nudging $+\epsilon$:

$$F_k^e = \begin{bmatrix} \frac{\delta^2}{2}\mathbf{B} \cdot \mathbf{u}_{-(k+1)} + \epsilon \nabla_{\boldsymbol{\pi}} \ell[\Phi_{-(k+1)}] \\ \delta\mathbf{B} \cdot \mathbf{u}_{-(k+1)} + \epsilon \nabla_{\boldsymbol{\phi}} \ell[\Phi_{-(k+1)}] \end{bmatrix} \quad (40)$$

1145 A.3.4 Nonlinear HRU

1146 **Hamiltonian of the recurrence.** The nonlinear HRU block is the composition of a nonlinear
 1147 spatial transformation and a nonlinear recurrent transformation (see Eq. 18). The Hamiltonian of the
 1148 nonlinear HRU block is defined by the following equations:

$$\begin{aligned} H[\Phi, \theta, u] &= T[\pi, \theta, u] + V[\phi, \theta, u] \\ &= \frac{1}{2}\|\pi\|^2 + \frac{\alpha}{2}\|\phi\|^2 + (\mathbf{A}^{-\top} \cdot \log(\cosh(\mathbf{A} \cdot \phi + \mathbf{B} \cdot u + \mathbf{b}))) \end{aligned} \quad (41)$$

1149 **Dynamics.** The dynamics of the nonlinear HRU block are defined by the following equations:

$$\begin{cases} \dot{\phi} = \pi \\ \dot{\pi} = -(\tanh(\mathbf{A}\phi + \mathbf{B}u + \mathbf{b}) + \alpha\phi) \end{cases} \quad (42)$$

1150 Which, after time-discretization with the integrator defined in A.2.1, gives the following equations:

$$\begin{cases} \phi_{k+1/3} = \phi_k + \frac{\delta}{2} \nabla_{\pi} T[\pi_k, \theta, u_k] \\ \quad = \phi_k + \frac{\delta}{2} \pi_k \\ \pi_{k+1/3} = \pi_k \\ \phi_{k+2/3} = \phi_{k+1/3} \\ \pi_{k+2/3} = \pi_{k+1/3} + \delta \nabla_{\phi} V[\phi_{k+1/3}, \theta, u_k] \\ \quad = \pi_{k+1/3} - \delta(\tanh(\mathbf{A}\phi_{k+1/3} + \mathbf{B}u_k + \mathbf{b}) + \alpha\phi_{k+1/3}) \\ \phi_{k+1} = \phi_{k+2/3} + \frac{\delta}{2} \nabla_{\pi} T[\pi_{k+2/3}, \theta, u_k] \\ \quad = \phi_{k+2/3} + \frac{\delta}{2} \pi_{k+2/3} \\ \pi_{k+1} = \pi_{k+2/3} \end{cases} \quad (43)$$

1151 A.3.5 Differentiating the time-discretization

1152 In training HRUs, one can also train the multidimensional time-discretization $\delta \in \mathbb{R}^{d_{\Phi} \times d_{\Phi}}$ that is
 1153 assumed to be diagonal. For this, it suffices to reparametrize the integrator and the Hamiltonian, so
 1154 that the time discretization becomes a parameter of the Hamiltonian. Hence the HRU equation 9
 1155 becomes:

$$\Phi_{k+1} = \mathcal{M}_{\hat{H},1}[\Phi_k, \theta, u_k, \delta] \quad \forall k = -K \dots -1, \quad (44)$$

1156 with \hat{H} is a reparametrization of the Hamiltonian H such that $\mathcal{M}_{\hat{H},1}[\Phi_k, \theta, u_k, \delta] =$
 1157 $\mathcal{M}_{H,\delta}[\Phi_k, \theta, u_k]$.

1158 For instance, for the linear HRU, we can reparametrize the Hamiltonian as:

$$\hat{H}[\Phi, \theta, u, \delta] = \frac{1}{2} \pi^{\top} \delta \pi + \left(\frac{1}{2} \phi^{\top} (\delta \mathbf{A}) \phi - \phi^{\top} (\delta \mathbf{B}) u \right) \quad (45)$$

1159 And for the nonlinear HRU, we can reparametrize the Hamiltonian as:

$$\hat{H}[\Phi, \theta, u, \delta] = \frac{1}{2} \pi^{\top} \delta \pi + \frac{\alpha}{2} \phi^{\top} \delta \phi + \frac{1}{2} (\delta \mathbf{A}^{-\top} \cdot \log(\cosh(\mathbf{A} \cdot \phi + \mathbf{B} \cdot u + \mathbf{b}))) \quad (46)$$

(47)

1160 Note that to match with previous implementations of the nonlinear HRU [30], we used the following
 1161 parametrization for the discretization step:

$$\hat{H}[\Phi, \theta, u, \delta] = \frac{1}{2} \pi^{\top} \sigma(\delta) \pi + \frac{\alpha}{2} \phi^{\top} \sigma(\delta) \phi + \frac{1}{2} (\sigma(\delta) \mathbf{A}^{-\top} \cdot \log(\cosh(\mathbf{A} \cdot \phi + \mathbf{B} \cdot u + \mathbf{b}))) \quad (48)$$

1162 where $\sigma(\delta)$ is a diagonal matrix with $\sigma(\delta)_{ii} = 0.5 + 0.5 \tanh(\delta_{ii}/2)$.

Algorithm 3 Custom Reverse Mode Automatic Differentiation for an arbitrary HRU

```
1: @custom_autodiff
2: function HRU( $\theta_{HRU}, \bar{u}$ )
3:    $\bar{\Phi} \leftarrow \text{HRU-HELPER}(\mathbf{0}, \theta_{HRU}, \bar{u}, \mathbf{0})$ 
4:   return  $\bar{\Phi}$ 
5: end function

6: function FORWARDHRU( $\theta_{HRU}, \bar{u}$ )
7:    $\bar{\Phi} \leftarrow \text{HRU-HELPER}(\mathbf{0}, \theta_{HRU}, \bar{u}, \mathbf{0})$ 
8:   return  $\bar{\Phi}, \Phi_K$ 
9: end function

10: function BACKWARDHRU( $r, \bar{g}$ , input_forward)
11:    $\theta_{HRU}, \bar{u} \leftarrow \text{input\_forward}$ 
12:    $\tilde{u} \leftarrow \text{REVERSE}(\bar{u})$  ▷ Reverse the input sequence
13:    $\Phi_K \leftarrow r$ 
14:    $\Phi_{0, \pm\epsilon}^e \leftarrow \Sigma_x \Phi_K + \epsilon \Sigma_x g_K$ 
15:    $\bar{\Phi}_{\pm\epsilon}^e \leftarrow \text{HRU-HELPER}(\Phi_{0, \pm\epsilon}^e, \theta_{HRU}, \bar{u}, \bar{g})$ 
16:    $\Delta_{\theta_{HRU}}, \Delta_{\bar{u}} \leftarrow \text{LEARNINGHRU}(\bar{\Phi}_{\pm\epsilon}^e)$ 
17:   return  $\Delta_{\theta_{HRU}}, \Delta_{\bar{u}}$  ▷ Loss gradient with respect to the input of FORWARDHRU( $\cdot, \cdot$ )
18: end function

19: HRU.define_custom_autodiff(ForwardHRU, BackwardHRU)
```

Algorithm 4 Helper functions for custom Automatic Differentiation

```
1: function HRU-HELPER( $\Phi_0, \theta_{HRU}, \bar{u}, \bar{n}$ ) ▷ See Dynamics. in App. A.3.3 and A.3.4 for
   concrete examples
2:   Input: Initial State  $\Phi_0$ , Parameters of the recurrence  $\theta_{HRU}$ , Inputs  $\bar{u}$ , Nudging  $\bar{n}$ 
3:   Output: State sequence  $\bar{\Phi}$ 
4: end function

5: function LEARNINGHRU( $\bar{\Phi}_{\pm\epsilon}^e, \bar{u}$ ) ▷ See RHEL gradient estimator in App. A.3.3 for a concrete
   example
6:   Input: Echo passes  $\bar{\Phi}_{\pm\epsilon}^e$ , input sequence time-reversed  $\tilde{u}$ 
7:   Output: Loss gradient w.r.t to the input of FORWARDHRU( $\cdot, \cdot$ ):  $\Delta_{\theta_{HRU}}, \Delta_{\bar{u}}$ 
8: end function
```

1163 A.3.6 Echo passes with automatic differentiation

1164 As mentioned in Remark 8, learning in HSSMs with RHEL involves chaining RHEL gradient
1165 estimators with Automatic Differentiation (AD). This design makes RHEL implementation readily
1166 compatible with modern automatic differentiation frameworks such as PyTorch or JAX. These
1167 libraries provide the capability to create custom functions that, when invoked within a computational
1168 graph, are automatically differentiated. Algorithm 3 demonstrates how to implement RHEL/AD
1169 chaining for the recurrent component of an HRU block.

1170 To implement a custom autodiff function HRU, we must define two additional functions required by
1171 AD: ForwardHRU and BackwardHRU. The ForwardHRU function is invoked when the HRU function
1172 is called within a computational graph to be differentiated. The ForwardHRU function returns two
1173 elements: its output for computing the remainder of the computational graph $\bar{\Phi}$, and a *residual* to
1174 be stored for the backward pass. For the HRU, the residual consists only of the final state Φ_K from
1175 the forward pass. The BackwardHRU function is called during graph backpropagation. It receives as
1176 input the loss gradient from the upstream portion of the computational graph \bar{g} , the input from the
1177 forward pass \bar{u} , and the residual Φ_K . The BackwardHRU function computes the gradient of the loss
1178 with respect to the HRU block parameters, which is subsequently used to update the HRU parameters.

1179 It also computes the gradient of the loss with respect to the forward pass input \bar{u} , which is then used
1180 to propagate the loss gradient backward through the computational graph.

1181 These three functions are then registered with the autodiff library using the
1182 `HRU.define_custom_autodiff` function, enabling the library to automatically differenti-
1183 ate the HRU function when it is called within a computational graph.

1184 The three functions utilize two helper functions: `HRU-HELPER` and `LearningHRU`. The `HRU-HELPER`
1185 function implements the HRU dynamics and is employed in both the forward and backward passes.
1186 The `LearningHRU` function implements the RHEL gradient estimator and is used in the backward
1187 pass to compute the gradient of the loss with respect to the HRU block parameters.s

1188 A.4 Experimental details

1189 **Summary.** This last section provides additional experimental details. More precisely:

- 1190 • We provide details about the datasets at use, the devices we used, as well as detailed
1191 hyperparameters.
- 1192 • Importantly, we detail **how RHEL can be subject to numerical *underflow* and how we**
1193 **mitigated this issue** – see Fig. 5–6 for details.

1194 A.4.1 Datasets

1195 This series of tasks was recently introduced as a subset of the University of East Anglia (UEA)
1196 datasets with the longest sequences for increased difficulty and recently used to benchmark the linear
1197 HSSM previously introduced [29]

1198 The classification datasets are drawn from a recently introduced benchmark [34] that selects a
1199 subset of the University of East Anglia (UEA) datasets [35], specifically choosing those with
1200 the longest sequences to increase difficulty, and which has been recently employed to evaluate
1201 the linear HSSM model [29]. These datasets include EigenWorms (17,984 sequence length, 5
1202 classes), SelfRegulationSCP1 (896 length, 2 classes), SelfRegulationSCP2 (1,152 length, 2 classes),
1203 EthanolConcentration (1,751 length, 4 classes), Heartbeat (405 length, 2 classes), and MotorImagery
1204 (3,000 length, 2 classes).

1205 Additionally, we evaluate our HSSMs on the PPG-DaLiA dataset [36], a multivariate time series
1206 regression dataset designed for heart rate prediction using data collected from a wrist-worn device. It
1207 includes recordings from fifteen individuals, each with approximately 150 minutes of data sampled at
1208 a maximum rate of 128 Hz. The dataset consists of six channels: blood volume pulse, electrodermal
1209 activity, body temperature, and three-axis acceleration. After splitting the data, a sliding window of
1210 length 49920 and step size 4992 is applied, representing a challenging very long-range interaction
1211 task.

1212 A.4.2 Simulation details and ressource consumption

1213 **Simulation details** The code to run the experiments is implemented using the JAX auto-
1214 differentiation framework [57]. All experiments were run on Nvidia V100 GPUs, except for the PPG
1215 experiments, which were run on Nvidia Tesla A100 GPUs due to larger memory demands.

1216 **Ressource consumption.** For the linear HSSM, training time is approximately 30 minutes for both
1217 classification (SCP1 dataset) and regression (PPG dataset) tasks across all training algorithms (RHEL
1218 and BPPT). The nonlinear HSSM requires approximately 1.5 hours for classification tasks with both
1219 algorithms. For the regression tasks, both BPTT and RHEL require approximately between 10 and
1220 20 hours.

1221 RAM consumption remained consistent across all conditions (learning algorithms, datasets, and
1222 models), reaching 30GB for regression tasks and 24.15GB for classification (SCP1 dataset). This
1223 represents 75% pre-allocation of GPU memory by JAX. While the regression task exceeded the
1224 V100’s 32GB memory capacity, it fit within the A100’s 40GB RAM.

1225 A.4.3 Hyperparameters

1226 We adopted the hyperparameters of [29] without modification, as their experiments utilized the IMEX
1227 integrator, which, like our approach, derives from the LeapFrog integrator. The hyperparameters
1228 are: learning rate (lr), number of layers (#blocks), number of hidden neurons (hidden dim), state-
1229 space dimension (state dim), and whether the time dimension is sent as input (include time). These
1230 hyperparameters were found by grid search and are presented in Table 3.

1231 We note that the implementation of the linear HSSM models is based on the code of [29] and uses
1232 complex numbers for implementation, which corresponds to doubling the state-space dimension
1233 mentioned in Table 3. We found that this dimensional doubling was necessary to recover the
1234 performance reported in the original paper, so we maintained this approach. We did not need to apply
1235 the same doubling for the nonlinear HSSM covering the performance reported in the paper, so we
1236 kept

Table 3: Hyperparameters for the linear and nonlinear HSSM model

Dataset	lr	hidden dim	state dim	#blocks	include time
Worms	0.0001	64	16	2	False
SCP1	0.0001	64	256	6	False
SCP2	0.00001	64	256	6	True
Ethanol	0.00001	16	256	4	False
Heartbeat	0.00001	64	16	2	True
Motor	0.0001	16	256	6	True
PPG	0.0001	64	16	2	True

For the *RHEL* algorithm we have two additional hyperparameters: the nudging strength ϵ and the scaling factor γ (see Appx. A.4.4). The nudging strength ϵ was set to 10^{-1} without prior tuning. For the scaling factor γ we did a grid search over the values $\{10^0, 10^4, 10^8, 10^{12}\}$ for the regression task (PPG-DaLiA) and found that the best performing parameter was 10^4 . For the classification tasks, we did a grid search over the values $\{10^0, 10^1, 10^2, 10^4\}$ and found that the best performing scaling was 10^4 for the averaged score.

We employ different initialization schemes for linear and nonlinear HSSM variants. For the linear HSSM model, we follow the same initialization scheme as [29]:

- $\mathbf{A} \sim \text{Uniform}(0, 1)$
- $\mathbf{B} \sim \text{Uniform}\left(-\frac{1}{\text{hidden_dim}}, \frac{1}{\text{hidden_dim}}\right)$
- $\mathbf{C} \sim \text{Uniform}\left(-\frac{1}{\text{state_dim}}, \frac{1}{\text{state_dim}}\right)$
- $\mathbf{D} \sim \mathcal{N}(0, 1)$
- $\delta \sim \text{Uniform}(0, 1)$

For the nonlinear HSSM model, we adopt the initialization scheme from [30]:

- $\mathbf{A} \sim \text{Uniform}(0.5, 1)$
- $\mathbf{B} \sim \text{Uniform}\left(-\frac{1}{\text{hidden_dim}}, \frac{1}{\text{hidden_dim}}\right)$
- $\mathbf{C} = 0$
- $\mathbf{D} \sim \mathcal{N}(0, 1)$
- $\mathbf{b} \sim \mathcal{N}(0, 1)$
- $\alpha \sim \text{Uniform}(0.1, 1.0)$
- $\delta \sim \text{Uniform}(-1, 1)$

A.4.4 Gradient re-scaling for dealing with numerical instabilities

Analysis From the theoretical results on RHEL, the nudging strength ϵ should be chosen as small as possible to accurately estimate the gradient of the loss function. However, naive numerical implementation can encounter underflow issues. In RHEL, gradient information is encoded in small perturbations to the state Φ that generate the echo trajectory Φ^e . This presents numerical challenges when the perturbations and state values differ by several orders of magnitude.

In practice, when using finite precision representations (such as floating-point numbers), the perturbations may be lost due to discretization error, where small values cannot be accurately represented or distinguished from zero in finite precision arithmetic.

To address this numerical challenge, we employ a simple gradient rescaling method. We multiplicatively scale the output loss $\ell[\cdot, \cdot]$ by a constant $\gamma > 1$ during the echo dynamics computation. This amplification ensures that the perturbation-induced changes in the loss remain within the representable range of floating-point arithmetic. Subsequently, we divide the RHEL parameter gradient estimate $(\Delta_{\theta}^{\text{RHEL}}(k, \epsilon))$ by the same scaling factor γ to recover the unbiased gradient.

1272 We now demonstrate the effectiveness of this gradient rescaling method on both Linear and Nonlinear
 1273 HSSM. To evaluate the effect of gradient scaling, we compare the gradients of RHEL and BPTT.
 1274 Given a HSSM model, we randomly sample a tuple $(\mathbf{u}_i, \mathbf{y}_i) \sim \mathcal{D}$ from the SCP1 dataset (see
 1275 App. A.4.1) and compute the gradients of the loss function with respect to the parameters of the
 1276 recurrence of the HRUs in Fig. 5A for the linear HSSM and Fig. 6A for the nonlinear HSSM.

1277 To gain a more fine-grained understanding, we also compute the gradients with respect to the inputs
 1278 of the third layer of the HSSM (see Eq. 17 and 18) in Fig. 5B and Fig. 6B. Compared to the
 1279 parameter gradients, these input gradients are not time-averaged and hence reveal more detail about
 1280 the underflow issues.

1281 **Gradient rescaling recovers the underflow issues.** As a general pattern across both architectures,
 1282 we observe that the unscaled ($\gamma = 10^0$) parameter and input gradients tend to be biased: they have
 1283 a norm ratio above 1.0 and cosine similarity below 1.0. We also note that for both the linear and
 1284 nonlinear cases, there exists a scaling factor (10^6 for linear and 10^4 for nonlinear) that recovers a
 1285 nearly perfect match between RHEL and BPTT.

1286 Additionally, for the input gradients that are not time-averaged (first column of Fig. 5B and Fig. 6B),
 1287 there is a large amount of noise in the unscaled gradients, which is eliminated for the best-performing
 1288 scaling factor (10^6 for linear and 10^4 for nonlinear). We conjecture that this noise is due to the
 1289 underflow effects described above.

1290 **Linear HSSM: scaling improves gradient matching** For the linear analysis, we observe a general
 1291 pattern: the more we scale the RHEL gradient, the better it matches BPTT, both for input gradients
 1292 and parameter gradients (Fig. 5A and B).

1293 **Nonlinear HSSM: optimal scaling balances underflow and linearization errors** For the nonlin-
 1294 ear analysis, we also observe that scaling the RHEL gradient improves the match between RHEL
 1295 and BPTT gradients. However, for higher values of the scaling factor ($\gamma = 10^6$), we observe that
 1296 both input and parameter gradients show a drop in matching quality. We conjecture that this is due to
 1297 linearization errors. If the loss gradient is too large, the nudging it drives will be large, and the finite
 1298 difference method used for the learning rule will no longer be valid.

1299 **Solution Implemented.** For the experiments, we conducted a grid search over the scaling parameter
 1300 γ . In our initial implementation, we applied gradient scaling without the corresponding downscaling
 1301 step, effectively amplifying the gradient magnitude throughout training. This provided valuable
 1302 insights into the sensitivity of RHEL dynamics to gradient scaling and improved performance. The
 1303 complete rescaling procedure (including both upscaling and downscaling) will be implemented in
 1304 the camera-ready version to provide a more comprehensive analysis of the numerical stabilization
 1305 approach.

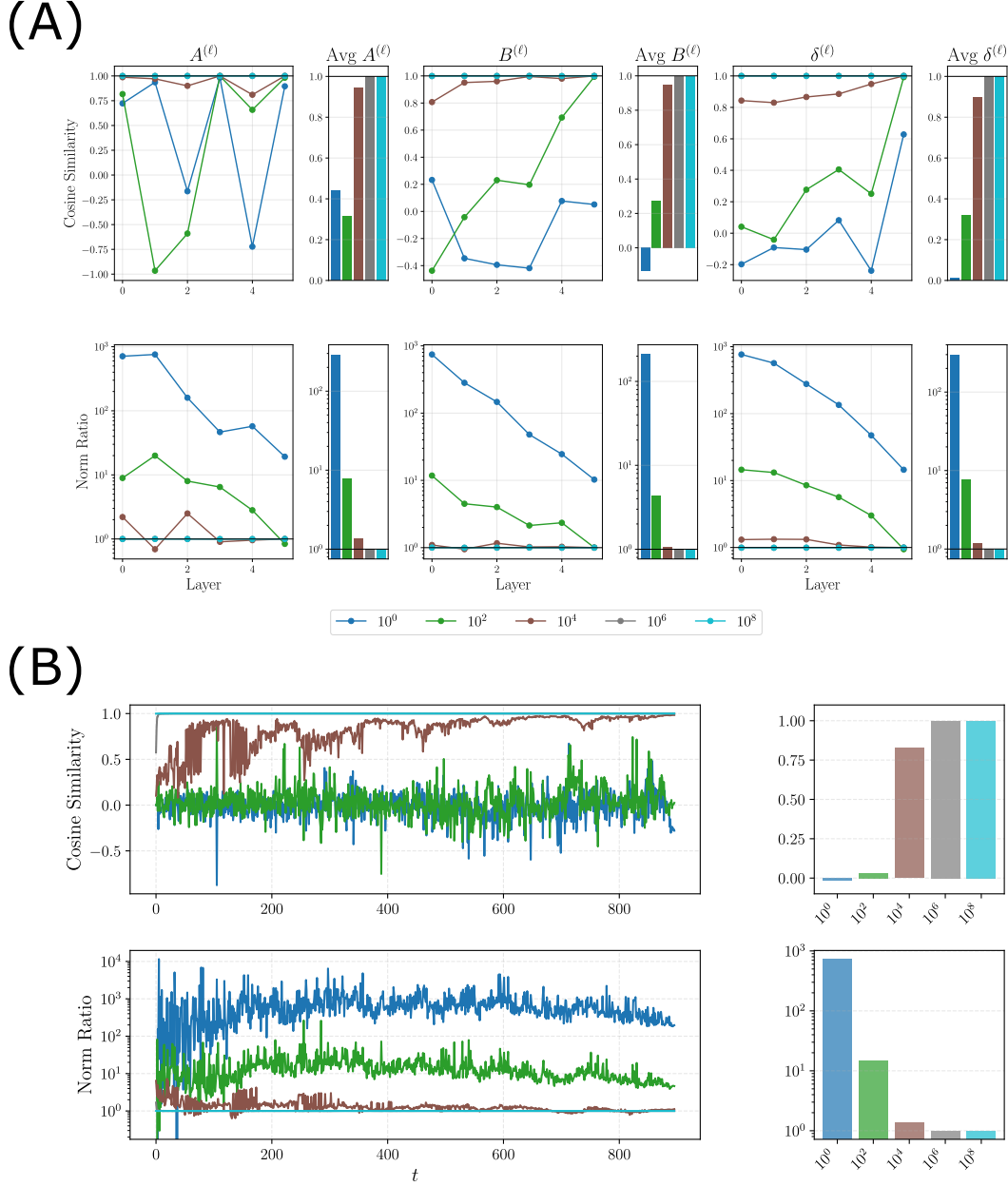
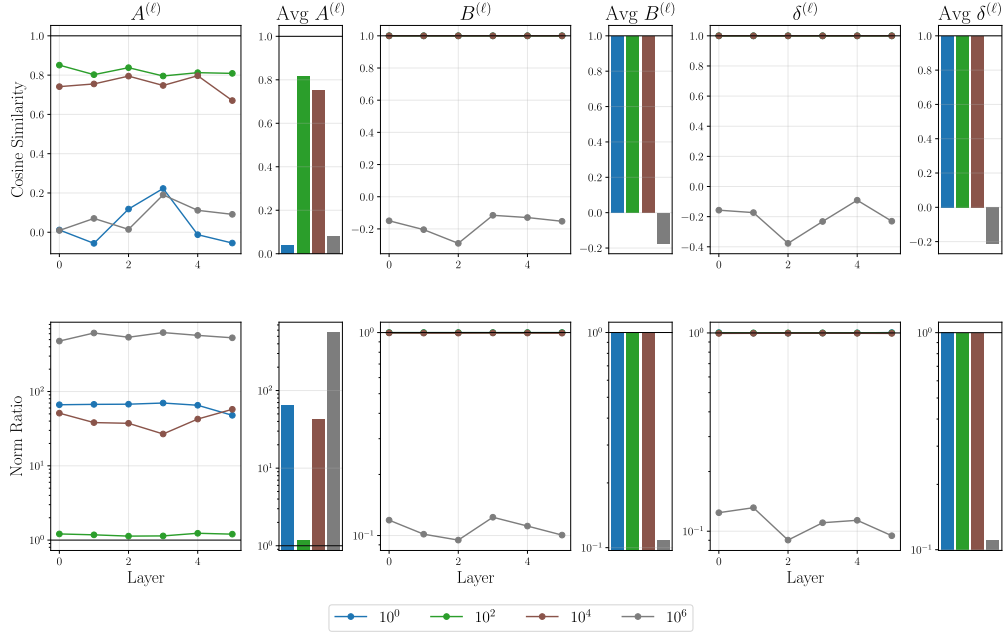


Figure 5: **(A): Parameters gradients comparison between RHEL and BPTT for HSSM.** Given some $(u_i, y_i) \sim \mathcal{D}$, we perform BPTT and RHEL on six blocks-deep linear HSSMs for different scaling factor γ of RHEL (different colors). We measure, per layer (line plot) and when averaged across layers (bar-plot), the cosine similarity (top panels) and norm ratio (bottom panels) between RHEL and BPTT parameters gradients of a linear HSSM (Eq. (17)). **(B): Inputs gradient comparison between RHEL and BPTT for HSSM.** Same setting as (A) but we focus on the gradients with respect to the inputs of the third layer ($u^{(3)}$, see Eq. 17). We measure, per time steps (line plot) and when averaged across time (bar-plot), the cosine similarity (top panels) and norm ratio (bottom panels) between RHEL and BPTT inputs gradients.

(A)



(B)

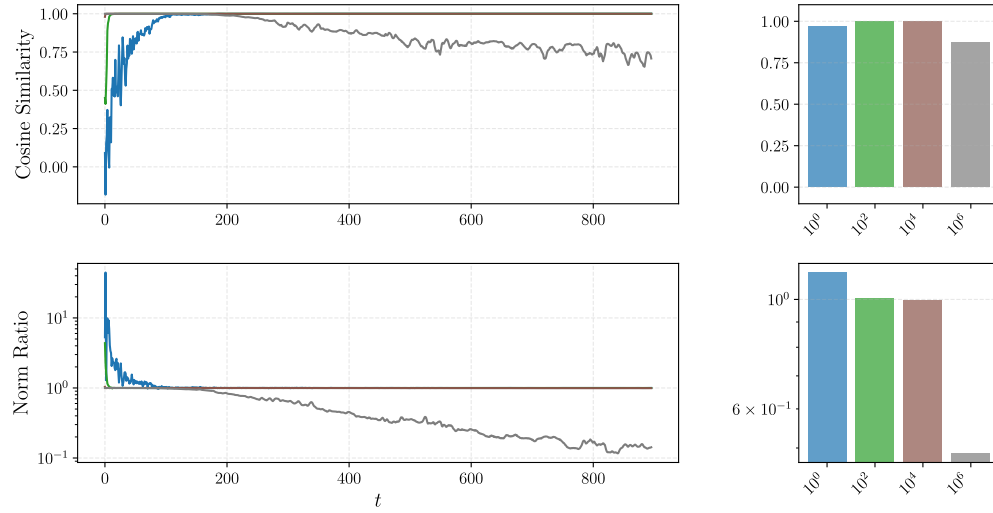


Figure 6: **(A): Parameters gradients comparison between RHEL and BPTT for HSSM.** Given some $(\mathbf{u}_i, \mathbf{y}_i) \sim \mathcal{D}$, we perform BPTT and RHEL on six blocks-deep nonlinear HSSMs for different scaling factor γ of RHEL (different colors). We measure, per layer (line plot) and when averaged across layers (bar-plot), the cosine similarity (top panels) and norm ratio (bottom panels) between RHEL and BPTT parameters gradients of a nonlinear HSSM (Eq. (18)). **(B): Inputs gradient comparison between RHEL and BPTT for HSSM.** Same setting as (A) but we focus on the gradients with respect to the inputs of the third layer ($\mathbf{u}^{(3)}$, see Eq. 18). We measure, per time steps (line plot) and when averaged across time (bar-plot), the cosine similarity (top panels) and norm ratio (bottom panels) between RHEL and BPTT inputs gradients. For (A) and (B), $\gamma = 10^8$ was also tested but produced numerical instabilities leading to NaN values and is therefore omitted from the results.