

1 We thank the four reviewers for their excitement about our work and detailed feedback! Overall, reviewers enjoyed
 2 the problem and approach: **"potential to impact the wide array of applications"** **"appealingly simple,"** **"novel ap-**
 3 **proach"**, **"very extensive and promising"** **"better than existing techniques"**; criticism centered around description
 4 of prior work, and clarity of problem formulation and details. We address these concerns here

5 **(a) Regularizer design and d_E (R2, R1).** We thank R2, and we hope that they will reconsider. We've revised the text
 6 in order to clarify the problem framing and fix notational ambiguities and some typos, which we believe clarify R2's
 7 major concerns (and a related point about d_E by R1). d is the number of independent dynamical variables (number of
 8 ODEs); $d_F \leq d$ is attractor manifold dimension (non-integer for chaotic systems with fractal attractors; $d_F = 1$ for
 9 limit cycles); d_E is embedding dimension. The hyperparameter L sets the maximum d_E expressible by the autoencoder
 10 (AE). However, the AE does not pick an integer d_E ; rather, d_E is estimated continuously post-training via the relative
 11 variance of each latent variable averaged across train data (similar to PCA weights). Dimensionality score thus compares
 12 explained variance (a function of latent index) between reconstruction and full dimensional system. Our experiments on
 13 systems with known d seek and find that $d_E \approx d$. Our regularizer takes a batch of latent activations, and estimates \bar{F}_m ,
 14 the proportion of new false neighbors indexed by number of latent dimensions. Since \bar{F}_m is intensive, it only weakly
 15 depends on batch size (we have added new physics references related to this observation).

16 **(b) Dependence on L (R2, R1).** We have added new experiments showing that L does not affect learned d_E of latent
 17 embedding as long as L is larger than d (see Fig H1); thus for unknown systems any large L can be used.

18 **(c) Prior Work (R1, R3).** We performed new experiments and extended the main text discussion of Ref. 35, which
 19 (to our knowledge) is the primary prior application of AE to attractor reconstruction. Ref. 35 embeds datasets via a
 20 one-layer AE with tanh activation and MSE loss, similar to our existing baseline unregularized MLP (see appendix).
 21 We performed a new set of experiments with a model exactly matching Ref. 35 (Fig H1), and find it is comparable to
 22 our other baselines for the noisy prediction task (same task as ref. 35). We have added these results to the paper.

23 **(d) Higher dimensions (R4)** Thank you! We have clarified that the ecosystem results are 10D; pendulum is 4D. We've
 24 added discussion and references to physics papers about mathematical limitations of embedding in the high- d limit.

25 **(e) Existing work on state space modelling (R3):** Thank you, and we hope that you will reconsider. Our paper
 26 doesn't claim to be the first AE applied to embedding (see (c) above). Indeed, 30% of our submission is a literature
 27 review of state space modelling (SSM); there we demonstrate several clear areas of novelty: (1) Our paper focuses on a
 28 **fully novel loss function and regularizer** rooted in the classical theory of dynamical systems, and we shows that this
 29 regularizer strongly constrains and improves AE representations, in contrast to prior AEs used for SSM. (2) Our paper
 30 uses a variety of **novel measures of attractor fidelity**—e.g. topology, neighbor coverage, fractal dimension—which go
 31 beyond few-timestep RMSE forecasting errors (the primary metric in previous works). (3) We show strong results for
 32 embedding **consistency across replicates, and robustness to Brownian stochasticity** (a more complex noise source
 33 than the measurement errors studied in prior works), two desirable embedding properties not explored previously.

34 **R1 & R2 additional comments:** Thank you so much for
 35 detailed feedback; we've fixed all wording, framing, and
 36 added suggested references; we regret that space limits us
 37 to major concerns not covered above: **R1 Misc: 8.1.2a,b**
 38 We revised hyperparameter discussion to add nuance: we
 39 mean that our experiments achieve strong results only by
 40 varying learning rate and regularizer strength, and the for-
 41 mer is only tuned to ensure that train loss plateaus. Rather
 42 than pre-select embedding timescale, we favor fixing large
 43 T and batch size, and letting the AE learn how much to
 44 weigh different timepoints. **8.1.2c,d,e** See (a) above. **8.3**
 45 We've moved details from appendix 5 into main text. There

46 are few widely-used definitions of attractor similarity (many SSM papers from ML authors focus on prediction, not
 47 verisimilitude, and many physics papers are qualitative), and so we report multiple established and novel metrics in
 48 order to give a holistic view. We've added the caveat that FNN-AE is more expensive than ETD/ICA, especially on
 49 small datasets, but only marginally more expensive than unregularized AE. **8.5** We revised to clarify that time series are
 50 Fourier-resampled only to ensure consistency across datasets; preprocessing/filtering otherwise has little effect (hence
 51 noise results) **8.6** We re-ran experiments and confirmed with pendulum data; we will add this to the appendix.

52 **R2 Misc: 3.2, 8.8, 8.12, 8.14-16** See (a, b); we've also moved Appendix 5 details to main text to clarify scoring
 53 metrics. We always refer to size L latent space as embedding space: T time delays only serve as a featurization of input
 54 to the AE, which seeks (and we find achieves) $d_E \approx d$ (not T). $d_E \leq L$ is computed *post-hoc* from relative latent
 55 activations by finding variance of the L latent variables across train set, giving continuous measure of dimensionality
 56 (thus unaffected by zero-padding). Thus d_E is neither a hyperparameter nor a direct AE output. **R2.8.2** The method
 57 and the code we're releasing now works for multivariate time series; we will highlight this. **R2.8.11** No leakage; when
 58 available, we use 2 different datasets or initial conditions; otherwise we use first N and last N points of a time series
 59 with length $\gg 2N$. **R2.8.16** We scale lower bound to mean, not theoretical min (see Nassar et al ICLR 2019, Eq. 25).

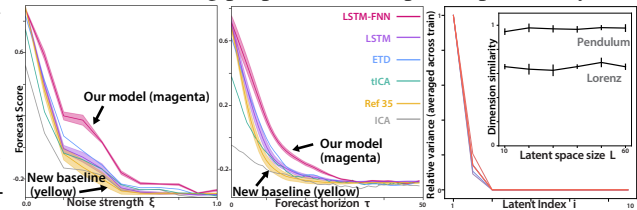


Figure H1: (A, B) Updated Fig 4 with ref. 35 baseline (yellow). (C) Similar activity patterns in first 10 latent units for fnn-AEs trained with $L = 10, 20, 30, 40, 50$ (blue to red). (Inset) dimension accuracy vs L also shows no dependence.