
Effective Latent Differential Equation Models via Attention and Multiple Shooting

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The GOKU-net is a continuous-time generative model that allows leveraging prior
2 knowledge in the form of differential equations. We present GOKU-UI, an evolu-
3 tion of the GOKU-nets, which integrates attention mechanisms and a novel
4 multiple shooting training strategy in the latent space. On simulated data, GOKU-UI
5 significantly improves performance in reconstruction and forecasting, outperform-
6 ing baselines even with 16 times less training data. Applied to empirical human
7 brain data, using stochastic Stuart-Landau oscillators, it is able to effectively cap-
8 ture complex brain dynamics, surpassing baselines in reconstruction and better
9 predicting future brain activity up to 15 seconds ahead. Ultimately, our research
10 provides further evidence on the fruitful symbiosis given by the combination of
11 established scientific insights and modern machine learning.

12 1 Introduction

13 Scientific Machine Learning (SciML) is an emerging field combining scientific models with mod-
14 ern data-driven techniques, often yielding increased interpretability, generalizability, and data effi-
15 ciency. (Baker et al., 2019; von Rueden et al., 2023; Shen et al., 2023). Latent Ordinary Differential
16 Equations (Latent ODEs) (Chen et al., 2018; Rubanova et al., 2019) are VAE-like generative models
17 that encode time series data into a latent space ruled by a differential equation which is parametrized
18 by a neural network. Building on Latent ODEs, Linial et al. (2021) introduced GOKU-nets (Gener-
19 ative ODE Modeling with Known Unknowns), which fundamental difference with the former is the
20 inclusion of a predefined differential equation structure as a prior for the latent dynamics. Compared
21 to LSTM and Latent-ODE on pendulum videos and cardiovascular system modeling, GOKU-net ex-
22 celled in reconstruction, forecasting, reduced training data needs, and offered better interpretability.

23 We propose an enhancement to the original GOKU-net architecture which adds attention mecha-
24 nisms to the main part of the model that infers the parameters of the differential equations. Moreover,
25 to overcome the inherent difficulties of training, we developed a novel strategy to train the GOKU-
26 net based on the multiple shooting technique (Bock & Plitt, 1984; Ribeiro et al., 2020; Turan &
27 Jäschke, 2021) in the latent space. Testing on simulated stochastic oscillators and empirical brain
28 data derived from resting state human functional Magnetic Resonance Imaging (fMRI), our *GOKU-*
29 *nets with Ubiquitous Inference* (GOKU-UI) surpassed both the original GOKU-net and baselines in
30 accuracy and data efficiency. GOKU-UI exemplifies the potential of melding traditional scientific
31 insights with modern machine learning.

32 2 Methods

33 2.1 Basic GOKU-nets

34 A general model class that we denominate Latent Differential Equation model (Latent DE), illus-
35 trated in Figure 1, begins by independently processing each temporal frame x_i with a *Feature Ex-*

36 *tractor*. The sequence then passes through a *Pattern Extractor* which aims to learn the distribution
 37 of the initial conditions and possibly of the parameters for the DE that will be subsequently inte-
 38 grated. Lastly, a *Reconstructor* transforms the solution back to the input space. Training follows the
 39 standard VAE approach, optimizing the evidence lower bound (ELBO) (Kingma & Welling, 2013).

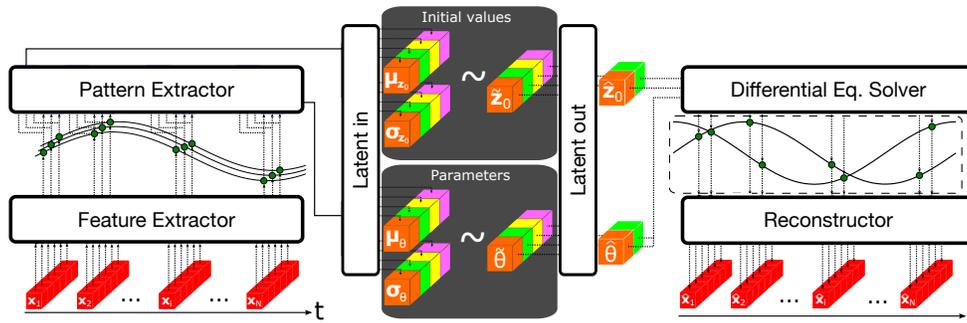


Figure 1: Diagram of a Latent Differential Equation model.

40 The original GOKU-net by Linial et al. (2021), understood as a Latent DE model, uses a ResNet
 41 with dense NNs as the Feature Extractor. Its Pattern Extractor employs an RNN for the initial condi-
 42 tions and a bidirectional LSTM for the ODE parameters. Unlike Latent ODEs, which parameterize
 43 the differential equation using another NN, in this case, the differential equation is predefined, incor-
 44 porating prior system knowledge. GOKU-net employs fully connected layers for the Latent in/out,
 45 and a ResNet similar to the Feature Extractor for the Reconstructor.

46 2.2 GOKU-UI

47 **Attention mechanism** In our GOKU-UI model, we incorporate a basic attention mecha-
 48 nism (Vaswani et al., 2017) into the Pattern Extractor, specifically when learning the differential
 49 equation parameters. Namely, instead of keeping the last element of the bidirectional LSTM (BiL-
 50 STM) used in the original GOKU-net model, all of its sequential outputs pass through a dense layer
 51 with softmax activation to calculate the attentional scores that weight the sum of all the BiLSTM
 52 outputs in order to obtain its final output.

53 **Multiple Shooting** Gradients calculations through differential equations often lead to highly com-
 54 plex loss landscapes (Ribeiro et al., 2020; Metz et al., 2021). Turan & Jäschke (2021) showed that
 55 training Neural ODEs even on basic oscillatory data could be problematic, resulting in trajectories
 56 resembling moving averages. To address this, they utilized *multiple shooting* techniques (Bock
 57 & Plitt, 1984; Diehl et al., 2006; Baake et al., 1992; Ribeiro et al., 2020). This method divides
 58 the differential equation’s time span into segments, independently inferring each segment’s initial
 59 conditions. These segments are then merged, enforcing continuity during the optimization.

60 However, applying the multiple shooting method to GOKU-nets is not straightforward. Firstly, in
 61 most cases that use this method, such as in Turan & Jäschke (2021), the differential equations are
 62 typically directly modeling the observable data, having direct access to the true initial conditions for
 63 each window. In the case of GOKU-nets, the dynamics modeled by differential equations occur in
 64 the latent space, which is being learned simultaneously; as a result, such true initial conditions are
 65 not available. Secondly, it is necessary to determine how the method will behave in relation to the
 66 parameters of the differential equation, which in the case of Neural ODEs are implicitly learned as
 67 part of their parameterization through the neural network.

68 Our proposal for extending the multiple shooting method to GOKU-nets is as follows. After passing
 69 through the Feature Extractor, we divide the temporal interval in the latent space in such a way that
 70 the Pattern Extractor generates in parallel different initial conditions for each temporal window, but
 71 provides a single set of parameters for the differential equations that will be shared by all windows.
 72 By this strategy, we maintain the potential benefits inherent to the multiple shooting method while
 73 leveraging the information available in a wider temporal range for the task of parameter inference,
 74 which is generally more challenging than estimating initial conditions. As mentioned before, we
 75 do not have access to target true initial conditions, however, what we can strive to achieve is the
 76 continuity of trajectories across different windows. To this end, these intervals are defined by over-
 77 lapping the last temporal point of each window with the first one of the following and the goal is

78 to minimize the distance between these points. Specifically, we employ regularization in the cost
 79 function when training the model, quadratically penalizing the discrepancy in the latent space of the
 80 overlapping points, that is, between the initial condition of each window and the end point of its
 81 preceding segment.

82 Non-variational GOKU-nets outperformed their variational counterparts in our experiments (see
 83 Supplementary Information B.3). Consequently, we utilized non-variational GOKU-nets for subse-
 84 quent results. Rather than sampling from normal distributions in the latent space (Figure 1), we used
 85 mean values μ_{z_0} and μ_θ . The resulting model’s cost function excludes the KL divergence term from
 86 the ELBO but includes the reconstruction term, computed as the normalized mean squared error
 87 between the model’s outputs and inputs. Continuity regularization from multiple shooting training
 88 is also incorporated.

89 2.3 Experiments

90 We assessed our attention and multiple shooting enhancements on synthetic data based stochastic
 91 Stuart-Landau oscillators and empirical human brain data. We compared various GOKU-model
 92 variations (basic, attention-enhanced) using either the original single shooting or the new multiple
 93 shooting method. Baseline models included LSTM, Latent ODE, and a naïve model. For fairness,
 94 the LSTM and Latent ODE maintained the GOKU-net’s architecture, differing only at the differ-
 95 ential equation layer. Here, the Neural ODE replaced the differential equation layer for the Latent
 96 ODE model, while an LSTM did for the other. The Latent ODE and LSTM sizes were adjusted to
 97 match GOKU-UI’s total number parameters. Naïve predictors utilized time-averaged input values
 98 for constant predictions. Detailed training procedures, models architectures, and hyperparameters
 99 are available in the Supplementary Information.

100 **Simulated data** Our simulated data derives from networks of stochastic Stuart-Landau (SL) oscil-
 101 lators, a standard model for resting state fMRI brain dynamics (Jobst et al., 2017; Deco et al., 2017).
 102 The dynamics for an oscillator node in an N node network is:

$$\begin{aligned} \dot{x}_j &= \text{Re}(\dot{z}_j) = [a_j - x_j^2 - y_j^2]x_j - \omega_j y_j + G \sum_{i=1}^N C_{ij}(x_i - x_j) + \beta \eta_j(t) \\ \dot{y}_j &= \text{Im}(\dot{z}_j) = [a_j - x_j^2 - y_j^2]y_j + \omega_j x_j + G \sum_{i=1}^N C_{ij}(y_i - y_j) + \beta \eta_j(t) \end{aligned} \quad (1)$$

103 where C_{ij} is the network’s connectivity matrix, G is a global coupling factor, and η_j is Gaussian
 104 noise. Each node has distinct bifurcation parameters a_j and frequencies ω_j . During the construction
 105 of our dataset, we perform a dimensionality augmentation on the network of oscillators, which are
 106 utilized as latent dynamics. Specifically, we apply a fixed random linear transformation, $f: \mathbb{R}^{2N} \rightarrow$
 107 \mathbb{R}^D , to the latent trajectories of each sample, with $D = 784$. Each sample corresponds to a unique
 108 random set of initial conditions and parameters for the $N = 3$ coupled oscillators.

109 **Empiric data** We assessed our models on resting-state fMRI data from the human brain, sourced
 110 from 153 subjects in the Track-On HD study (Klöppel et al., 2015). After preprocessing as outlined
 111 in Polosecki et al. (2020), we applied a 20-component Canonical ICA (Varoquaux et al., 2010),
 112 retaining 11 components post artifact removal. This resulted in 306 samples, each with 160 time
 113 points collected every 3 seconds. We reserved 20% of the data (n=60) for testing, ensuring balanced
 114 representation, and used the remaining (n = 246) for training and validation. The first 114 time
 115 points from each of these samples were used for model training, with the remainder reserved for
 116 validation. The GOKU-UI model employed 20 Stuart-Landau oscillators (Eq. 1) in its latent space.

117 3 Results

118 We assessed four GOKU-net variants for both reconstruction and forecast tasks, considering sin-
 119 gle/multiple shooting methods and the presence or absence of attention. We compared them to three
 120 baselines: LSTM, Latent ODEs, and a naïve predictor. Although models were exclusively trained
 121 for reconstruction, we evaluated their forecasting during testing. The normalized root mean square
 122 error (NRMSE) measures the prediction error against the target ground truth.

123 **Simulated data** Figure 2a depicts GOKU-net variants with multiple shooting having significantly
 124 reduced errors on a synthetic dataset of three stochastic Stuart-Landau oscillators. Attention mech-
 125 anism enhanced performance, especially in single shooting. GOKU-UI, combining attention and

126 multiple shooting, remained the top performer, with Wilcoxon tests confirming significance at p-
 127 values < 0.02 after Holm correction. Latent ODEs underperformed, resembling the naïve predictor,
 128 while LSTMs surpassed basic GOKU-nets. Notably, when trained on just 150 unique samples,
 129 GOKU-UI outperformed all other single shooting models, even those trained on datasets 32 times
 130 larger. In forecasting, as shown in Figure 2b, attention-aided GOKU-nets excelled over LSTMs, with
 131 GOKU-UI standing out up to 150 samples (p-values < 0.02 , Wilcoxon signed-rank, Holm corrected),
 132 beyond which its performance was statistically indistinguishable from that of the basic GOKU model
 133 with multiple shooting (p-values > 0.05 , Wilcoxon signed-rank tests, Holm corrected).

134 **Empirical data** Figure 2c demonstrates that the attention mechanism didn't boost single shooting
 135 GOKU-net performance. Yet, multiple shooting training significantly improved results, with the
 136 GOKU-UI model, merging both techniques, reducing the median reconstruction error by five times
 137 compared to single shooting baselines. Furthermore, GOKU-UI had a significantly lower recon-
 138 struction NRMSE than the multiple shooting GOKU basic model (p < 0.04 , Wilcoxon signed-rank
 139 test). In forecasting, GOKU-UI outperformed other models, achieving lower forecast errors for up
 140 to 15 seconds of brain activity.

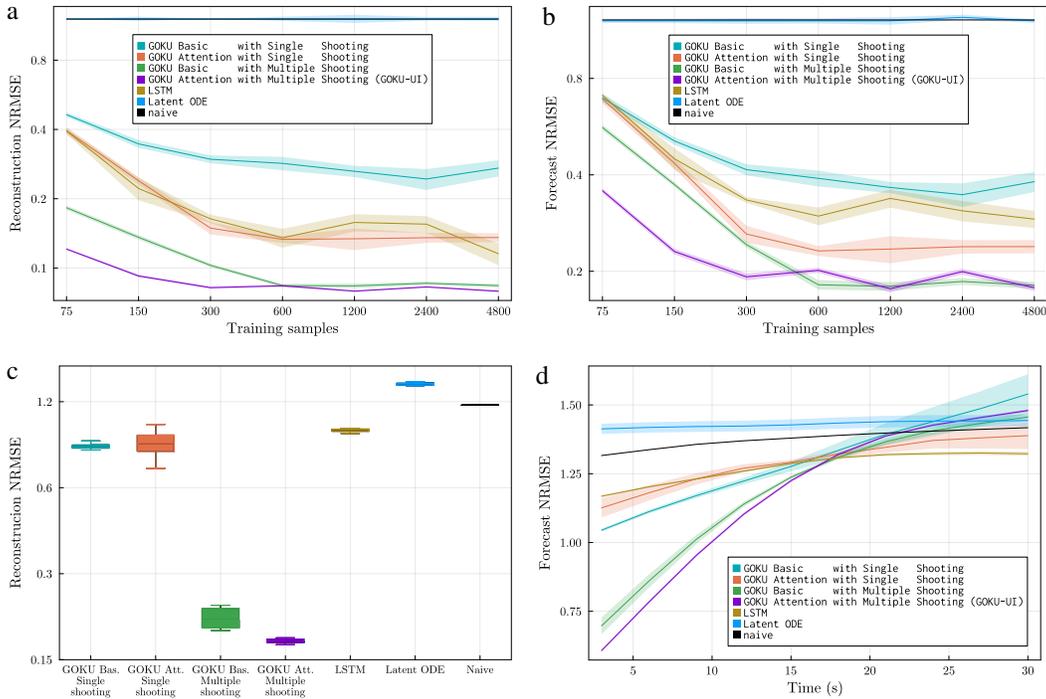


Figure 2: Comparison of reconstruction (left panels) and forecast (right panels) performances on the synthetic Stuart-Landau (top panels) and fMRI (bottom panels) test data sets, using normalized RMSE. Averages are taken across input dimensions and time span, with shaded areas indicating standard errors from multiple training runs with varied random seeds. Forecasts in panel b are assessed over a 20 time-step horizon.

141 4 Conclusion

142 We enhanced the GOKU-nets with attention and multiple shooting, with the latter yielding the most
 143 impact. The resulting model, GOKU-UI, showed improved performance and data efficiency on
 144 both synthetic and empirical brain data. By leveraging established scientific insights into modern
 145 machine learning, GOKU-UI was able to encode whole-brain dynamics into a latent representation,
 146 learning a low-dimensional interpretable dynamical system model that could offer insights into brain
 147 functionality and open avenues for multiple practical applications.

148 Applying GOKU-UI to new problems might not be as straightforward as a general-purpose black-
 149 box neural network model due to the need for a specific differential equation. Still, with guidance
 150 from dynamical systems theory, it is not only feasible but also beneficial.

151 References

- 152 Baake, E., Baake, M., Bock, H., and Briggs, K. Fitting ordinary differential equations to chaotic
153 data. *Physical Review A*, 45(8):5524, 1992.
- 154 Baker, N., Alexander, F., Bremer, T., Hagberg, A., Kevrekidis, Y., Najm, H., Parashar, M., Patra,
155 A., Sethian, J., Wild, S., Willcox, K., and Lee, S. Workshop report on basic research needs for
156 scientific machine learning: Core technologies for artificial intelligence, 2 2019. URL <https://www.osti.gov/biblio/1478744>.
- 157
- 158 Bock, H. G. and Plitt, K.-J. A multiple shooting algorithm for direct solution of optimal control
159 problems. *IFAC Proceedings Volumes*, 17(2):1603–1608, 1984.
- 160 Chen, R. T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. Neural ordinary differential equa-
161 tions, 2018. URL <https://arxiv.org/abs/1806.07366>.
- 162 Datsersis, G., Isensee, J., Pech, S., and Gál, T. Drwatson: the perfect sidekick for your scientific
163 inquiries. *Journal of Open Source Software*, 5(54):2673, 2020. doi: 10.21105/joss.02673. URL
164 <https://doi.org/10.21105/joss.02673>.
- 165 Deco, G., Kringelbach, M. L., Jirsa, V. K., and Ritter, P. The dynamics of resting fluctuations in the
166 brain: metastability and its dynamical cortical core. *Scientific reports*, 7(1):3095, 2017.
- 167 Diehl, M., Bock, H. G., Diedam, H., and Wieber, P.-B. Fast direct multiple shooting algorithms
168 for optimal robot control. *Fast motions in biomechanics and robotics: optimization and feedback*
169 *control*, pp. 65–93, 2006.
- 170 Innes, M., Saba, E., Fischer, K., Gandhi, D., Rudilosso, M. C., Joy, N. M., Karmali, T., Pal, A.,
171 and Shah, V. Fashionable modelling with flux. *CoRR*, abs/1811.01457, 2018. URL <https://arxiv.org/abs/1811.01457>.
- 172
- 173 Jobst, B. M., Hindriks, R., Laufs, H., Tagliazucchi, E., Hahn, G., Ponce-Alvarez, A., Stevner, A.,
174 Kringelbach, M. L., and Deco, G. Increased stability and breakdown of brain effective connec-
175 tivity during slow-wave sleep: mechanistic insights from whole-brain computational modelling.
176 *Scientific reports*, 7(1):1–16, 2017.
- 177 Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*,
178 2013.
- 179 Klöppel, S., Gregory, S., Scheller, E., Minkova, L., Razi, A., Durr, A., Roos, R. A., Leavitt, B. R.,
180 Papoutsi, M., Landwehrmeyer, G. B., et al. Compensation in preclinical huntington’s disease:
181 evidence from the track-on hd study. *EBioMedicine*, 2(10):1420–1429, 2015.
- 182 Linial, O., Ravid, N., Eytan, D., and Shalit, U. Generative ODE modeling with known unknowns.
183 In *Proceedings of the Conference on Health, Inference, and Learning*. ACM, apr 2021. doi:
184 10.1145/3450439.3451866. URL <https://doi.org/10.1145/3450439.3451866>.
- 185 Ma, Y., Gowda, S., Anantharaman, R., Laughman, C., Shah, V., and Rackauckas, C. Modeling-
186 toolkit: A composable graph transformation system for equation-based modeling, 2021.
- 187 Metz, L., Freeman, C. D., Schoenholz, S. S., and Kachman, T. Gradients are not all you need. *arXiv*
188 *preprint arXiv:2111.05803*, 2021.
- 189 Misra, D. Mish: A self regularized non-monotonic activation function. *arXiv preprint*
190 *arXiv:1908.08681*, 2019.
- 191 Polosecki, P., Castro, E., Rish, I., Pustina, D., Warner, J. H., Wood, A., Sampaio, C., and Cec-
192 chi, G. A. Resting-state connectivity stratifies premanifest huntingtons disease by longitudinal
193 cognitive decline rate. *Scientific reports*, 10(1):1–15, 2020.
- 194 Rackauckas, C. and Nie, Q. Differentialequations.jl—a performant and feature-rich ecosystem for
195 solving differential equations in julia. *Journal of Open Research Software*, 5(1):15, 2017.

- 196 Rackauckas, C., Ma, Y., Martensen, J., Warner, C., Zubov, K., Supekar, R., Skinner, D., Ramadhan,
197 A., and Edelman, A. Universal differential equations for scientific machine learning, 2020. URL
198 <https://arxiv.org/abs/2001.04385>.
- 199 Ribeiro, A. H., Tiels, K., Umenberger, J., Schön, T. B., and Aguirre, L. A. On the smoothness of
200 nonlinear system identification. *Automatica*, 121:109158, 2020.
- 201 Rubanova, Y., Chen, R. T., and Duvenaud, D. K. Latent ordinary differential equations for
202 irregularly-sampled time series. *Advances in neural information processing systems*, 32, 2019.
- 203 Shen, C., Appling, A. P., Gentine, P., Bandai, T., Gupta, H., Tartakovsky, A., Baity-Jesi, M., Fencia,
204 F., Kifer, D., Li, L., et al. Differentiable modelling to unify machine learning and physical models
205 for geosciences. *Nature Reviews Earth & Environment*, pp. 1–16, 2023.
- 206 Turan, E. M. and Jäschke, J. Multiple shooting for training neural differential equations on time
207 series. *IEEE Control Systems Letters*, 6:1897–1902, 2021.
- 208 Varoquaux, G., Sadaghiani, S., Pinel, P., Kleinschmidt, A., Poline, J.-B., and Thirion, B. A group
209 model for stable multi-subject ica on fmri datasets. *Neuroimage*, 51(1):288–299, 2010.
- 210 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polo-
211 sukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30,
212 2017.
- 213 von Rueden, L., Mayer, S., Beckh, K., Georgiev, B., Giesselbach, S., Heese, R., Kirsch, B., Pfrom-
214 mer, J., Pick, A., Ramamurthy, R., Walczak, M., Garcke, J., Bauckhage, C., and Schuecker, J.
215 Informed machine learning a taxonomy and survey of integrating prior knowledge into learning
216 systems. *IEEE Transactions on Knowledge and Data Engineering*, 35(1):614–633, 2023. doi:
217 10.1109/TKDE.2021.3079836.

218 **Supplementary Information**

219 **A Models architectures**

220 Referring to the diagram in Figure 1, the specific architecture used for the different models, for both
221 simulated and empirical data experiments, is as follows:

222 **A.1 Basic GOKU-nets**

223 **Feature Extractor**

224 ResNet with 4 fully-connected layers, each with 200 neurons and using mish activation func-
225 tions (Misra, 2019). Input dimension = number of dimensions in the input data. Output dimension
226 = 128.

227 **Pattern Extractor**

228 Initial values path: an RNN with 2 layers and 64 neurons in each with ReLU activations. Input dim
229 = 128. Output dim = 64.

230 Parameters path: Bidirectional LSTM with 2 layers and 64 neurons in each. Input dimension =
231 128. Output dimension = 128. Note that the dimension of the output of the forward LSTM and the
232 backward LSTM are 64 but when concatenating them, the resulting output dimension is the given
233 one.

234 **Latent in**

235 Initial values path: single-layered fully connected NN. Input dim = 64. Output dim = 64.

236 Parameters path: fully connected NN with 1 layer. Input dim = 128. Output dim = 128.

237 **Latent out**

238 Initial values path: fully connected NN with 2 layers and 200 neurons in the hidden layer, using
239 no activation function (identity). Input dim = 64. Output dim = number of state variables of the
240 differential equation.

241 Parameters path: fully connected NN with 2 layers and 200 neurons in the hidden layer, using sig-
242 moid activation function. The parameters are projected from the interval [0, 1] to the desired range
243 when integrating the differential equation. Input dim = 128. Output dim = number of parameters of
244 the differential equation.

245 **Differential Equation layer**

246 The predefined differential equation is solved numerically for each of the sets of parameters and ini-
247 tial conditions provided by the previous layer. The output is the trajectories at time points equivalent
248 to the input data.

249 **Reconstructor**

250 ResNet is similar to the one in the Feature Extractor, except that in this case the input dimension
251 is the number of state variables of the differential equation and the output dimension is the one
252 corresponding to the input data.

253 **A.2 GOKU-nets with attention**

254 With the exception of the Pattern Extractor, the rest of the layers in the GOKU-nets with attention
255 model remain identical to those in the basic GOKU-nets.

256 **Pattern Extractor**

257 Initial values path: LSTM with 1 layer. Input dimension = 128. Output dimension = 128.

258 Parameters path: Bidirectional LSTM (BiLSTM) with 1 layer. Input dim = 128. Output dim = 128.
259 A fully connected NN with input and output dimensions of 128 is used for the attention mechanism.
260 This attention NN processes all the output sequences of the BiLSTM, after which a softmax is
261 applied across the time dimension in order to obtain the attentional scores that will be used in the
262 weighted sum of all the time steps returned by the BiLSTM.

263 **A.3 LSTM baseline model**

264 The whole architecture is the same as in the basic GOKU-net, except for the Differential Equation
265 layer, which is replaced by an LSTM:

266 **LSTM layer**

267 We used a single-layered LSTM with input and output dimensions set to z_dim . This value is
268 determined in each experiment to ensure that the total number of parameters in the LSTM model
269 closely matches that of the corresponding GOKU-UI. For the simulated dataset experiments, we set
270 $z_dim = 42$. In the case of the empirical dataset experiments, $z_dim = 105$. The LSTM operates
271 recursively. It takes as its first input the value equivalent to the initial condition in differential
272 equations. Subsequently, the model feeds back its last output as the new input, continuing this
273 process until the number of time steps matches that of the model's input.

274 **A.4 Latent ODE baseline model**

275 The whole architecture is the same as in the basic GOKU-net, except for the Differential Equation
276 layer, which is replaced by a Neural ODE:

277 **Neural ODE layer**

278 Neural ODE is parametrized by a fully connected NN with 3 layers and $node_hidden_dim$ neurons
279 in each. The input and output dimensions are given by z_dim , which is the number of state variables.
280 In the case of the simulated dataset experiments, the number of state variables was selected to match

281 the true latent dimension $z_dim = 6$ and the number of neurons in each layer was adjusted so that
282 the total number of parameters in the model matched as closely as possible that of the corresponding
283 GOKU-UI, resulting in $node_hidden_dim = 137$. On the other hand, in the case of the fMRI
284 experiments, the number of state variables was set to $z_dim = 20$ and $node_hidden_dim = 317$,
285 also matching the total number of parameters of the corresponding GOKU-UI model.

286 B Comprehensive description of experiments

287 B.1 Simulated dataset generation

288 The high-dimensional simulated dataset used for training the model was constructed based on the
289 simulations of 3 coupled Stuart-Landau oscillators (Eqs. 1) with different random sets of parameters.
290 Each set of parameters corresponds to a different training sample. Whenever we used the Stuart-
291 Landau model in our experiments (both when generating the dataset and when using it inside the
292 GOKU-nets), the time was rescaled by multiplying the right-hand side of Eqs. 1 by 20. Thus,
293 when integrating the equations with the used $dt = 0.05$, the input sequences of length 46 time steps
294 contain a few oscillations. The parameters a , ω and C were sampled from uniform distributions
295 within the following ranges

$$296 a \in [-0.2, 0.2]; \quad \omega \in [0.08\pi, 0.14\pi]; \quad C \in [0, 0.2]$$

297 while $G = 0.1$ and $\eta = 0.02$. On the other hand, the initial conditions for the six state variables were
298 sampled from uniform distributions within the ranges $[0.3, 0.4]$. For each set of parameters and initial
299 conditions, the system is integrated with the SOSRI solver, a Stability-optimized adaptive strong
300 order 1.5 and weak order 2.0 for diagonal/scalar Ito SDEs, from the DifferentialEquations.jl Julia
301 package (Rackauckas & Nie, 2017). The complete time span of the integration is 35 units of time
302 and the trajectories are saved every 0.05, resulting in 700 time points. The first 100 time steps are
303 trimmed, in order to remove possible initial transients. Afterwards, a random linear transformation
304 is independently applied to each of the 600 remaining time steps, in order to obtain 784 dimensions.
305 In other words, every state vector of length 6 from each sample is multiplied by the same 784×6
306 matrix, initialized randomly sampling from a uniform distribution in the range $[-1, 1]$. A training
307 dataset was created with 5000 samples, which serves as the source for the different training instances
308 using different sizes of training sets (see Figure 2a). A different test set with 900 samples was created
309 for the posterior evaluations of the model.

310 B.2 Empirical dataset generation

311 We used resting state fMRI data from 153 participants, obtained from the Track-On HD study (Klöp-
312 pel et al., 2015). The data underwent pre-processing, as described in Polosecki et al. (2020), and
313 a 20-component Canonical ICA (Varoquaux et al., 2010) was performed. Upon inspecting the re-
314 sulting 20 components, 9 were identified as artifacts and thus discarded, leaving 11 components for
315 further analysis in our experiments. Each subject contributed data from two visits, accumulating a
316 total of 306 data samples. Each sample comprised 160 time points, obtained at a temporal resolution
317 of 3 seconds.

318 For our investigation, we set aside approximately 20% of the data samples ($n=60$) for testing, while
319 ensuring balanced representation from sex, condition, and measurement site. The remaining data
320 samples ($n=246$) were allocated for training and validation. Specifically, the first 114 time points
321 from each of these samples were utilized for model training, with the remainder reserved for valida-
322 tion and early training termination. Finally, the training, validation, and test splits were all normal-
323 ized by the standard deviation of the training set.

324 B.3 Training settings

325 All the experiments underwent the same training procedure with identical hyperparameters, which
326 will be described here.

327 The input sequence length for all the models was 46 time steps, and the batch size was set at 64.
328 As described above, the full length of each sample in the training sets was 600 time steps for the
329 synthetic dataset and 114 for the fMRI dataset. The procedure for generating a batch of training data
330 is as follows: First, 64 samples that have not been used previously in the current training epoch are

331 randomly selected. Then, for each sample, a 46 time-step-long interval is randomly chosen within
332 the 600 or 114 time steps available in the full sample length.

333 The GOKU-net based models, contain the same Stuart-Landau differential equations as described
334 above, however, the allowed ranges of parameters differ from the ones used during the generation
335 of the synthetic dataset. In order to be closer to a real world use-case we allow for a wider range of
336 parameters than those actually used for generating the data, since in principle one would not know
337 the true range:

$$338 \quad a \in [-1, 1]; \quad \omega \in [0, 1]$$

339 while keeping, the other parameters the same except of the connectivity in the empirical fMRI
340 training, in which case it was allowed to be negative: $C \in [-0.2, 0.2]$. The differential equa-
341 tions definitions were optimized for higher computational performance with the help of Modeling-
342 Toolkit.jl (Ma et al., 2021). During training, they were solved with the SOSRI solver, a Stability-
343 optimized adaptive strong order 1.5 and weak order 2.0 for diagonal/scalar Ito SDEs, from the Dif-
344 ferentialEquations.jl Julia package (Rackauckas & Nie, 2017). The sensitivity algorithm used was
345 ForwardDiffSensitivity from the SciMLSensitivity.jl package (Rackauckas et al., 2020). The
346 models were defined and trained within the deep learning framework of the Flux.jl package (Innes
347 et al., 2018). The experiments were managed using DrWatson.jl package (Datseris et al., 2020).

348 The model was trained with Adam with a weight decay of 10^{-10} , and the learning rate was dynam-
349 ically determined by the following schedule. The learning rate begins with a linear growth (also
350 referred to as *learning rate warm-up*) from 10^{-7} , escalating up to 0.005251 across 20 epochs. After-
351 wards, it maintains that value until the validation loss stagnates (has not achieved a lower value for
352 50 epochs), at which point it starts a sinusoidal schedule with an exponentially decreasing amplitude.

353 For the multiple shooting training, all the presented experiments used a time window length of 10,
354 therefore partitioning 46-time-steps-long sequences into 5 windows with their endpoints overlap-
355 ing. The regularization coefficient in the loss function for the continuity constraint had a value of
356 2.

357 Since we found that models with variational versions of the GOKU-nets underperformed their
358 non-variational versions, all the results presented in this work were obtained using non-variational
359 GOKU-nets. This is, instead of sampling from normal distributions in the latent space as depicted
360 in Figure 1, we pass forward the mean values μ_{z_0} and μ_θ . Thus, the associated loss function does
361 not have the KL divergence term associated with the ELBO but retains the reconstruction loss given
362 by the mean squared error between the output of the model and the input, normalized by the mean
363 absolute value of the input. In addition, when multiple shooting training is employed, the extra
364 term regarding the continuity constraint is included in the loss function. This extra term consists of
365 the mean squared differences between the last point of a window and the initial from the next one,
366 divided by the number of junctions and multiplied by a regularization coefficient. Please, note that
367 this continuity regularization is performed in the state space of the differential equation and not in
368 the input space.

369 C Reconstruction plots

370 To provide a visual representation of the model’s performance, this section presents trajectories
371 from both the synthetic and empirical fMRI test sets, along with their corresponding reconstructions
372 by GOKU-UI and the original GOKU-nets (lacking attention mechanisms and trained with single
373 shooting). The x-axis represents time steps in all cases. To display representative cases, samples
374 were selected based on their mean reconstruction RMSE being closest to the median error across
375 all samples. For the synthetic data, 11 components were randomly selected for display in Figures
376 3 and 4, due to the impracticality of displaying all 784 components. Each figure displays results
377 from different instances of models, all trained with 4800 samples but each initialized with a unique
378 random seed. For the fMRI data, all 11 ICA components are displayed in Figures 5 and 6.

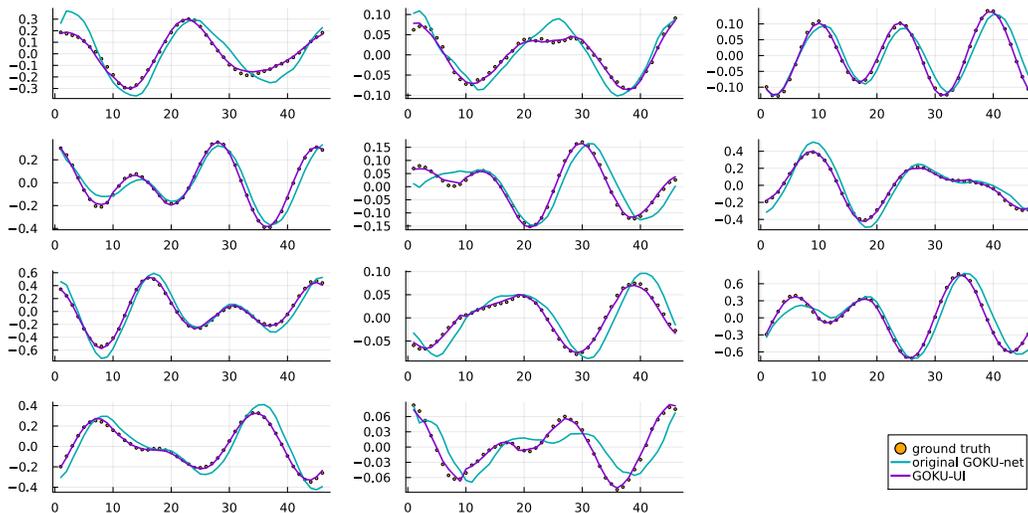


Figure 3: Representative example of a 46-time-step input sequence from the synthetic test set, accompanied by its reconstructions from both GOKU-UI and the original GOKU-nets (lacking attention mechanisms and trained with single shooting). The sample was selected so that its RMSE was the closest to the median error across all samples. 11 randomly selected components out of the 784 are displayed.

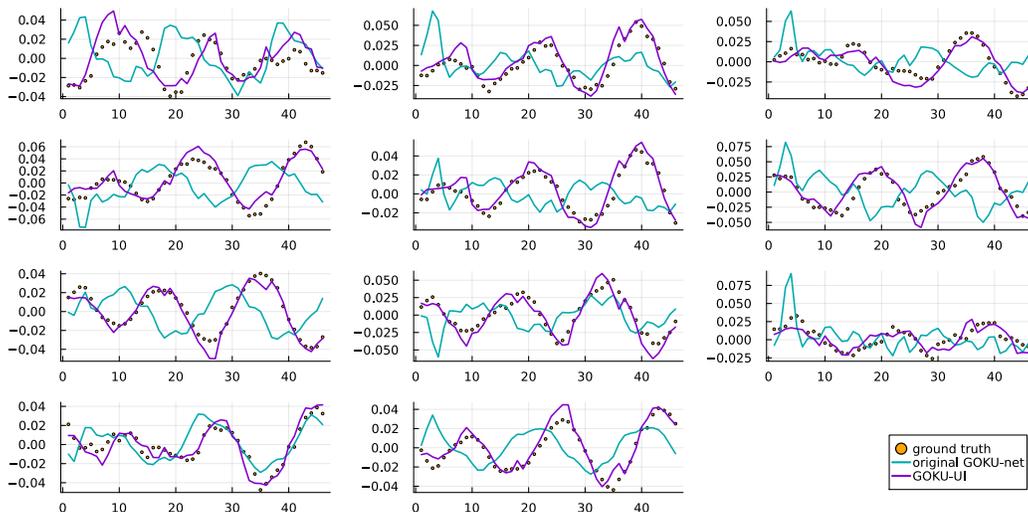


Figure 4: Representative example of a 46-time-step input sequence from the synthetic test set, accompanied by its reconstructions from both GOKU-UI and the original GOKU-nets (lacking attention mechanisms and trained with single shooting). The sample was selected so that its RMSE was the closest to the median error across all samples. 11 randomly selected components out of the 784 are displayed. This figure is similar to the previous one but presents results from different instances of the trained models, each initialized with a unique random seed.

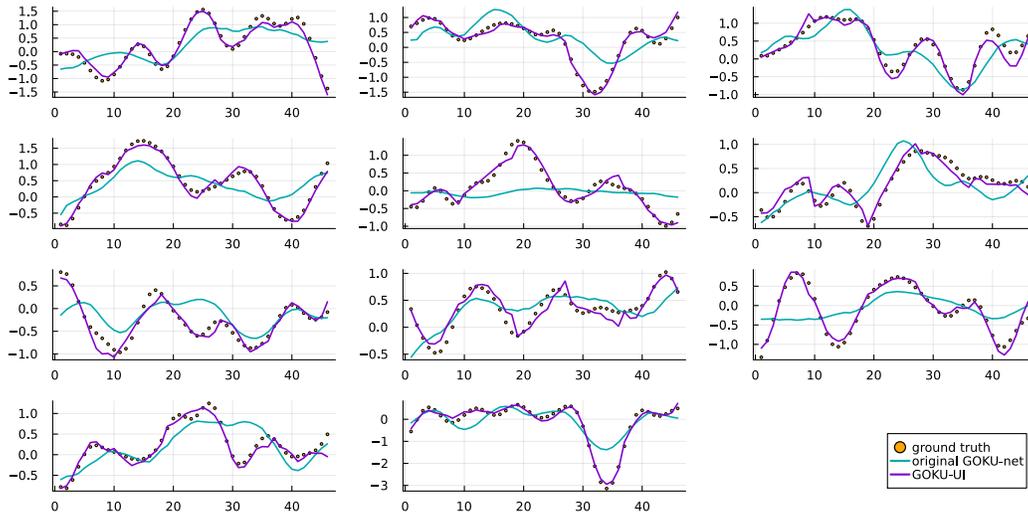


Figure 5: Representative example of a 46-time-steps input sequence for all considered ICA components from the empirical fMRI test set, accompanied by its reconstructions from both GOKU-UI and the original GOKU-nets (lacking attention mechanisms and trained with single shooting). The sample was selected so that its RMSE was closest to the median error across all samples. The x-axis represents time steps, each corresponding to 3 seconds.

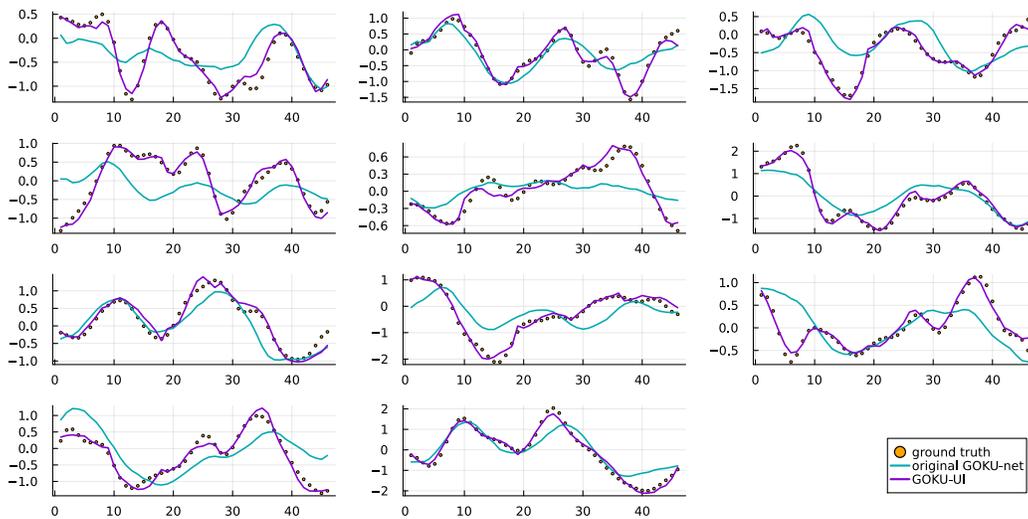


Figure 6: Representative example of a 46-time-steps input sequence for all considered ICA components from the empirical fMRI test set, accompanied by its reconstructions from both GOKU-UI and the original GOKU-nets (lacking attention mechanisms and trained with single shooting). The sample was selected so that its RMSE was closest to the median error across all samples. This figure is similar to the previous one but presents results from different instances of the trained models, each initialized with a unique random seed. The x-axis represents time steps, each corresponding to 3 seconds.