

---

# DOMINANT MOTION IDENTIFICATION OF MULTI-PARTICLE SYSTEM USING DEEP LEARNING FROM VIDEO.

---

**Yayati Jadhav**  
Mechanical Engineering  
Carnegie Mellon University  
Pittsburgh, PA 15213  
yayatij@andrew.cmu.edu

**Amir Barati Farimani**  
Mechanical Engineering  
Carnegie Mellon University  
Pittsburgh, PA 15213  
barati@cmu.edu

April 27, 2021

## ABSTRACT

Identifying underlying governing equations and physical relevant information from high-dimensional observable data has always been a challenge in physical sciences. With the recent advances in sensing technology and available datasets, various machine learning techniques have made it possible to distill underlying mathematical models from sufficiently clean and usable datasets. However, most of these techniques rely on prior knowledge of the system and noise free data obtained by simulation of physical system or by direct measurements of the signals. Hence, the inference obtained by using these techniques is often unreliable to be used in the real world where observed data is noisy and requires feature engineering to extract relevant features. In this work, we provide a deep-learning framework that extracts relevant information from real-world videos of highly stochastic systems, with no prior knowledge and distills the underlying governing equation representing the system. We demonstrate this approach on videos of confined multi-agent/particle systems of ants, termites, fishes as well as a simulated confined multi-particle system with elastic collision interactions. Furthermore, we explore how these seemingly diverse systems have predictable underlying behavior. In this study, we have used computer vision and motion tracking to extract spatial trajectories of individual agents/particles in a system, and by using LSTM VAE we projected these features on a low-dimensional latent space from which the underlying differential equation representing the data was extracted using SINDy framework.

## 1 Introduction

Learning specific action cases from one systems and applying it to a different systems is one of the major research areas in machine learning. Nature provides a variety of well optimized action case scenarios, mimicking these natural systems and applying it to the other systems is a well researched topic in biological inspired computation. Distilling underlying dynamics of a system from data is a form of data-driven pattern recognition where similarities between the systems states over time are found. Recognition of this pattern not only helps in prediction of future states but it also helps to identify and quantify the deviation of a systems from ideal state due to disturbance. One of the major challenge in distilling underlying dynamics from data is the availability of specific features which define the states of the system. This problem is exacerbated in case of video data, where the data features are observable spatio-temporal data and required feature engineering and prior knowledge of the system to find a surrogate system model. Visually learning the dynamics of a system, say a group of ants foraging in space, in the form of a differential equation from videos of the systems not only allows for understanding of the system better but can also be used to transfer dynamics of this system to another system like swarm robots or optimization algorithms. Recent advancements in machine learning techniques has made it possible to approximate representative mathematical models directly from observable data [1–3] instead of building models from first principle physical laws.

Time-series data describing evolution of states of a system over time, is one of the most abundant form of data available. Underlining dynamics of a nonlinear dynamical system can be extracted from time series data using symbolic

regression [3, 10]. A more popular uses sparse regression, as demonstrated in SINDy framework, a predefined candidate library of all possible data features is created and sparse regression is used to select a parsimonious model that represents most prominent dynamics of the system [9, 41, 42]. However, since numerical approximation of derivatives from data is required for sparse regression, this approach is highly susceptible to noise. Moreover, this approach fails to handle spatio-temporal data where  $x$  and  $y$  co-ordinates from 1 feature [30] and becomes trickier with the increase in number of features/variables observed in the system, as it becomes difficult to objectively eliminate features and their interactions based on their contribution to the system dynamics. An extension of SINDy framework where auto encoder is used to find effective coordinate system representing the system in fewest possible feature [31], solves this dimensionality problem to some extent, however relatively clean and noise free data is still required for this approach. Besides sparse regression based methods, techniques such as Gaussian process is also used to find the governing equation from the data. In this process, a Gaussian prior is placed on an unknown coefficients which is inferred using maximum likelihood estimation [43, 47].

This problem of high dimensional noisy data is compounded in cases where there is limited sensor data and consequently less number of features, especially in cases such as quantifying animal movements, cellular movement, particle diffusion or in diverse fields of science requiring non-invasive forms of measurements. As in such cases it becomes difficult to find reliable sensor data and data is gathered from a series of images such as satellite images to track animal migration or microscopic images to track cellular movements. Although PDEs are analysed to understand and quantify various ecological systems, the difficulty to model or parameterize PDEs or form differential equations from data still remains. Various attempts have been made to form differential models from a given data using likelihood or probability density functions [38, 39], but these techniques have become increasingly complicated and time consuming with multiple overlapping or interacting population and large datasets [40]. Reaction-diffusion equations have been used to model animal behaviors and understand biological pattern formation [53, 54], however this approach still requires fitting the parameters to the PDEs.

Deep learning techniques have gained significant popularity in recent times partly due to their superior performance in various prediction and classification tasks [12]. Deep learning techniques use non-linear mapping to project input features and their interactions onto a output space thereby creating an approximated mathematical model for relating inputs to the outputs [4, 5]. Recent developments in deep learning to find data-drive solution of non-linear partial differential equations [11, 13] has opened an avenue for extracting the learned differential equation from the network. Physics informed neural networks (PINN) can parameterize a PDE of known structure, which limits their application for PDE discovery from data with no prior knowledge of the system. Convolutional neural networks (CNNs) which are known for automatic feature discovery have been used to identify the structure of unknown equation [44], however parsimony of the equation variables is not guaranteed.

One of the major drawback of neural networks is that it is a black box model. Though there have been various attempts at explaining and interpreting deep learning systems [7, 8] there has not been an agreed upon explanation [6] that links the weights of the neural networks and the function being approximated, since networks with different topology often provide same results. Furthermore, the obtained model is hidden, meaning it can only be viewed in terms of inputs and outputs which fails to provide any specific insights or physical interpretation of the function being approximated. This makes it extremely difficult to identify and extract specific function or equation being approximated by the network. These drawbacks of current techniques for model discovery demands for an all round approach which not only can handle noise, is easy to implement but also can provide interpretable models representing the system.

Videos or sequence of images offer rich visual information about the dynamics of a system. Understanding scene dynamics is one of the core problems shared between computer vision and physical science communities. Physical parameters such as mass and friction [55], rigid body motion [56] as well as Newtonian dynamics [57] can be estimated by observing interactions of objects in a video. However, these approaches rely on prior knowledge of event dynamics and hence can't be extended to infer physics of a stochastic system.

In this paper, we present a deep learning based approach to obtain a time dependent differential equation representing the dynamics of multi particle/agent stochastic systems directly from videos. We extract the spatio-temporal trajectory data of individual particle/agent from videos of multi particle systems. The obtained system states are then projected onto a low dimensional latent space representing the states, this latent vector is used to form a parsimonious model using sparse regression based SINDy framework. Furthermore, we explore the ability of the obtained model to identify anomalous behaviour in the system. Figure 1. provides the general workflow of the framework where the input is a video of multi-particle system from which trajectories of particles is extracted using object tracking algorithm. LSTM based variational autoencoder is used to create embeddings of system states and map the spatial features onto low dimensional latent vector. The extracted latent vector representing compressed states of the system is then distilled to find underlying system dynamics.

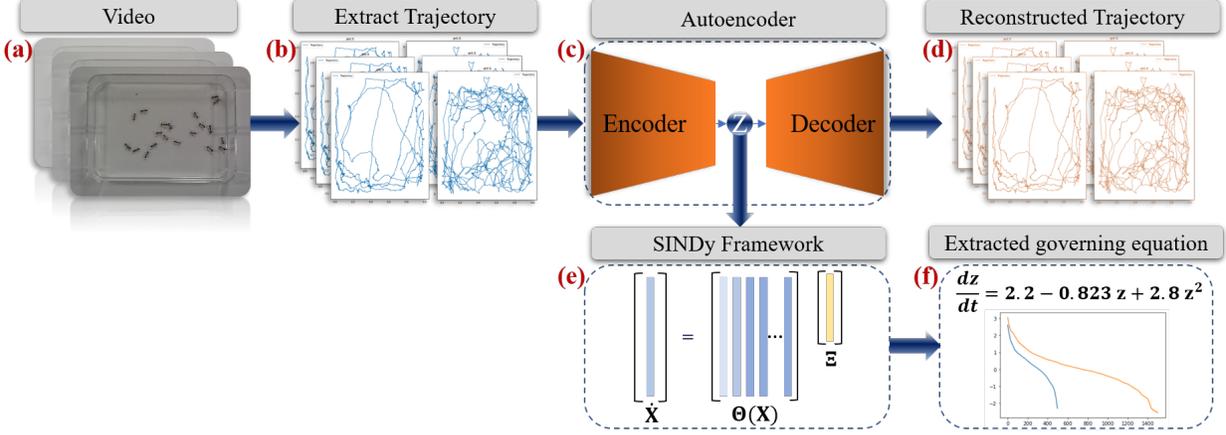


Figure 1: Schematic of framework. Frame in videos (a) are used to extract states (trajectory of particle) of the system at each time step (b). The extracted data is then fed into LSTM variational autoencoder (c). The network is trained till the reconstructed states or trajectories match the input states (d) and using sparse regression (e) the latent representation of states at each time steps is extracted is then used to model time differential equation of the system (f).

## 2 Methodology

### 2.1 Trajectory extraction:

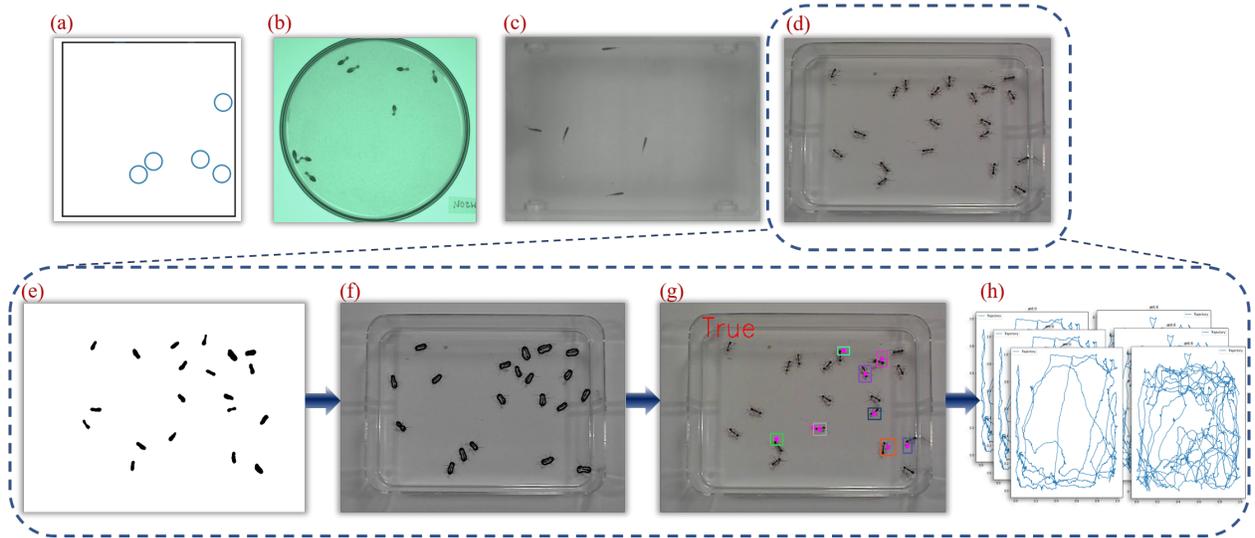


Figure 2: Examples of diverse multi agent system used in this study, (a) simulated multi particle system with elastic collision, (b) termites in confined space, (c) zebrafish behavior in box and (d) ants in a confined space. Trajectory extraction pipeline is demonstrated for ant system. Objects of interests are isolated based on color contrast and background removal (e). A contour is drawn around the object of interests (f) in each frame, the mid point of the drawn contour is then used to tag and track each object using CSRT algorithm [29].

In this study we explore our framework on diverse video datasets of 2 dimensional confined multi particle systems such as termite [27], zebrafish [27] and ants in a box [28]. We also analyze a system of simulated particles with elastic collision. Figure 2 (a),(b),(c),(d) shows various videos used in this study.

Tracking multiple objects or multi object tracking (MOT) in videos is an area of active research in the computer vision. Though there are various state of the art multi-object tracking systems available, since the priori of the all object appearance is known, all the particle are in-frame for the entire duration of the videos and there is no significant occlusion present, we have applied single object tracking model to solve MOT. In this study we use CSRT (Channel and

Spatial Reliability Tracking) [29] tracker available in OpenCV package to extract the trajectory of individual particle from the videos. CSRT tracker uses correlation filter with histogram of oriented gradients and colormnames as features to tag objects, and searches the regions around last known position of the object in previous frames. In order to initialize, the tracker requires initial position of the objects. CSRT algorithm doesnt involve motion tracking i.e. it doesnt estimate the position of the objects in successive frames, this makes the output trajectories noisy but it also help to remove any algorithmic bias in the trajectory data, thus helping the us learn the 'true' dynamics of the system.

Figure 2 shows the general workflow for automatic trajectory extraction from ants in box video. First the video is cropped to a specific region of interest and using Gaussian blur, noise and variations in lighting and contrast is smoothed. For object detection in the first frame, a color and contrast threshold based detection algorithm is used. Isolating the object of interest in the first frame is done by using background subtraction and color isolation, as shown in Figure 2e. Based on the this frame a contour is drawn around the objects of interests(2f) and it is tagged by drawing automated bounding boxes around the contours. The bounding box data is then used to initialize CSRT algorithm for tracking (2g).

One limitation of this approach is when the objects are interact with each other the CSRT algorithm sometimes fails to distinguish the two particle, in such cases the size of bounding box was modified to enable it to find distinctive features on each object or the objects are tagged manually.

## 2.2 Autoencoder: LSTM VAE

Recurrent neural networks (RNN) are powerful deep learning algorithms designed to analyze sequential data and are preferred over multilayered perceptrons (MLP) to map a sequence onto a sequence. The hidden states of RNNs are updated based on current as well as information from previous time steps, which make them appropriate to handle sequential time series data. LSTM [14] and GRU [15] are some of the most popular RNN architectures used to handle a variety of sequential data such as time series forecasting [18], forecasting chaotic systems [26], speech recognition [19, 20], translation [21, 22], human dynamics [25]. These architectures have shown to be promising compared to basic RNN structure which is difficult to train and has gradient vanishing problem [16, 17]. Moreover, LSTM networks have the ability to learn and reproduce long sequences and can handle spatio-temporal data like trajectories and have been extended to predict human trajectories in social settings [23, 24]. Auto encoders are typically unsupervised machine learning algorithms used to extract features, find compressed representation of the original data or to reconstruct and denoise data. LSTM autoencoders have be used in various fields [32–37] especially to extract spatial-temporal features from videos.

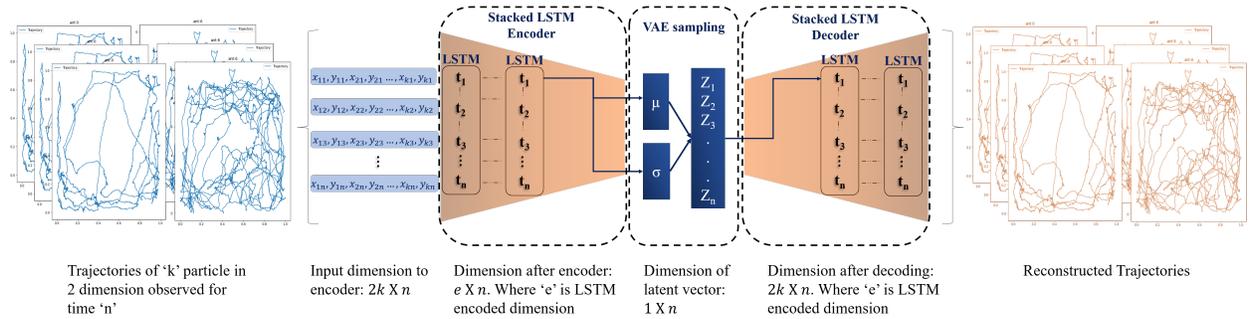


Figure 3: Structure of LSTM Variational autoencoder. The position of particles at time 't' or the system state at each time step is used as an input to the variational autoencoder(VAE). The LSTM encoder projects the system states at each time step onto a embedding of size  $e \times n$ . From this embedding a latent vector of size  $1 \times n$  is sampled. LSTM decoder then decodes this latent vector and reconstructs the input state.

The extracted trajectory data provides with an time ordered spatial states of all particles in the system. Given the stochastic nature of the data and noise from tracking algorithm, the LSTM Variational autoencoder(VAE) encodes the state of the system at a given time step and learns a probabilistic latent distribution to ensure that the 'dominant' latent vector is sampled from the encoded latent distribution.

LSTM VAE allows to map features at each time step to a sampled latent representation for the corresponding time step. Figure 3. shows the structure of LSTM-VAE. The trajectory data of each particle at time ' $t$ ' is in the given by  $x_i^t, y_i^t$  where  $i$  is the index of the particle. Hence ant system we have ' $k$ ' particles observed for ' $n$ ' time-steps will be a matrix of dimension  $(2k \times n)$ , where we have  $x_1, y_1, x_2, y_2, \dots, x_k, y_k$  features of  $k$  particles for  $n$  time-steps. These features are passed into the LSTM encoder, which maps these features to an encoded vector of predefined size. Here the feature extraction is done by LSTM encoder which finds features in spatial data at each time step. The encoded

vector from the encoder is further mapped to mean and standard deviation using a linear layer. During training a latent vector of size 1 is sampled from the encoder output. Hence, the  $2k$  input features are mapped to a latent vector of size 1 for each time-step. For decoding, the latent vector is passed through a linear layer to obtain initial state for decoder which projects the latent vector to the output size which is same as that of input. Since the system analysed in this study are stochastic and non-deterministic systems, more weight is given to the reconstruction loss which ensures that the reconstructed output is similar to the inputs i.e. the VAE is able to encode the trajectory of the particle to a single latent vector and that latent vector can reconstruct the original trajectories. To ensure the best reconstruction and subsequently the best latent representation of the features, MSE loss and SmoothL1 loss is used depending on which loss function gives best reconstruction.

### 2.3 SINDy Framework

The projected one dimensional latent space obtained from LSTM-VAE is considered to be representative of the spatial states of the system at corresponding time steps, leveraging this compressed representation we can describe the dynamics of the system in form of a time dependent differential equation. Sparse regression based approaches are the most commonly used to fit differential equation to measured data. In this work we have used SINDy framework which uses sparse regression to infer nonlinear dynamical systems from data.

It is assumed that the systems can be described by the equation of form:

$$\frac{dz}{dt} = f(z) \quad (1)$$

where,  $z$  is the obtained latent vector.

As per SINDy framework,  $z$  is known and time derivative of  $z$  is calculated from data. Since, SINDy framework becomes unstable with noise, the latent vector was smoothed using Savitzky–Golay filter. A candidate library of size  $n$  is constructed from the data as given by:

$$\Theta(z) = [1 \quad z \quad z^2 \quad z^3 \quad \dots \quad z^n] \quad (2)$$

A parsimonious model that includes most prominent and contributing terms is found by using sparse regression to solve:

$$\frac{dz}{dt} = \Theta(z) * \Xi \quad (3)$$

where  $\Xi$  is the library of coefficient corresponding to terms in  $\Theta(z)$  and is calculated by sparse regression.

In this study to find parsimonious model we have utilized, PySINDY [48] library, which is a sparse regression package.

## 3 Results

We demonstrate the workflow on diverse video sets of multi-particle, ranging from termites and ants in a box, to zebra fish. We also demonstrate the workflow on simulated 2d multi particle system with elastic collision particle for interaction.

### 3.1 Ants in a box

The video shows 20 ants in a container and all ants are in frame throughout the video hence we considered all 20 ants in this study. Machine learning algorithms perform poorly for unscaled features. In previous works data was simulated from known PDEs and hence did not demand transformation or scaling of the features. Since, the scaling transformation impacts the identified differential equation in this study, a simple coordinated transformation technique provided in scikit-learn module [49] is used which scales the features in range from 0 to 1. Figure?? shows the scaled trajectories of 10 selected ants extracted from the video.

Each observed state of the system corresponds to 1 frame of the video. Since, the video has a frame rate of 30 frames/sec, 30 observations were taken for each second in the video. Ant system has 20 particles (ants) and considering 500 time-steps (frames) a matrix of dimension of  $500 \times 40$  is obtained, where  $x_1, y_1, x_2, y_2, \dots, x_{20}, y_{20}$  are the spatial features of 20 particles. These features are passed into stacked LSTM encoder with 2 LSTM layers, which maps these 40 features to a vector of size 32 for each time-step from which a latent vector of size 1 is obtained. Hence for 500 time steps a sequential latent vector  $z$  of size  $500 \times 1$  is obtained.

Figure6 shows the reconstructed trajectories of 10 selected ants. It is observed that the neural network is able to approximate complex and stochastic trajectories of all ants. Based on these reconstructed trajectories we can say that

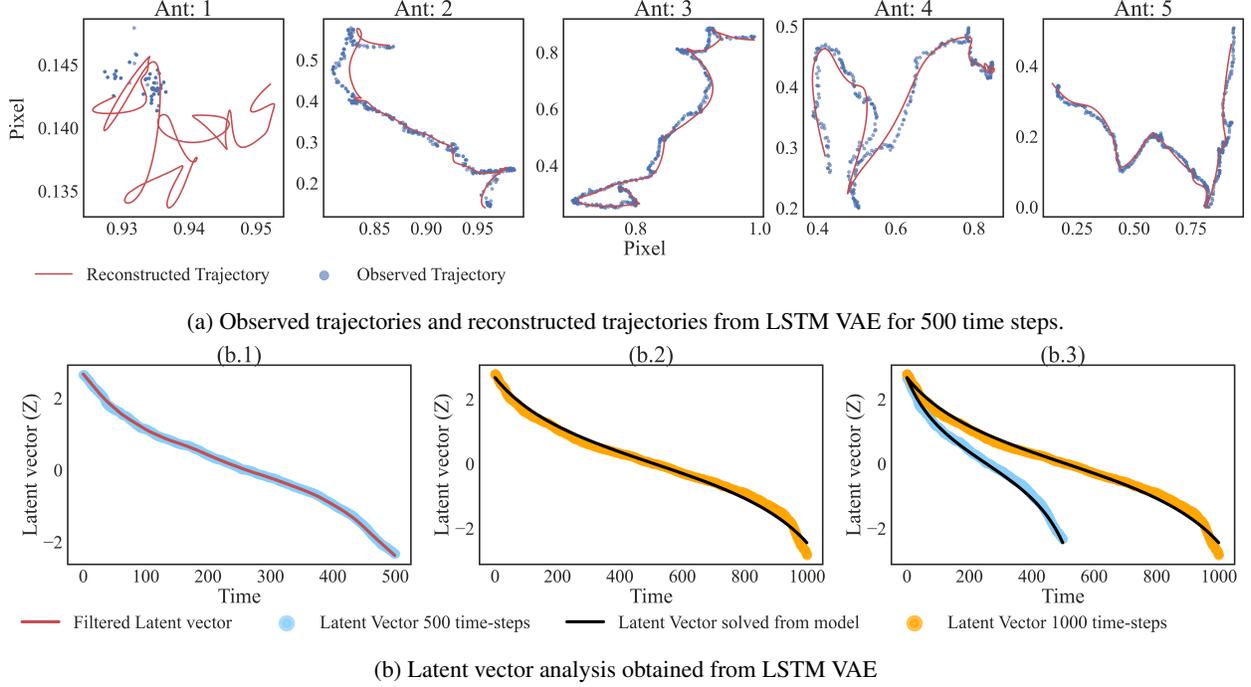


Figure 4: Latent vector obtained by training VAE for 500 time steps is filter using Savitzky–Golay (a). Differential equation is modeled from the filtered latent vector using SINDy framework, the obtained equation is then solved for 500 time steps (b) and 1000 time steps(c). The latent vector solved from model for 1000 time steps shows good correlation with the latent vector obtained by training the VAE for 1000 time steps.

latent vector obtained has the mappings of all the states of the system. Better approximations can be obtained by a deeper network or by data augmentation and filtering. Furthermore, since these trajectories correspond to a short real-world time, it can be said that network has better approximation once the dynamics of each particle in the system is established (see appendix).

To extract the underlining equation of the system represented by the latent vector, the latent vector 'z' is first denoised using Savitzky–Golay filter of window size 51 and polynomial order 1. Figure 7. shows the latent vector and smoother latent vector. Following SINDy framework, we obtain the following equation representing the dynamics of the system up till 500 time steps:

$$\frac{dz}{dt} = -1.232z^2 - 3.266 \quad (4)$$

In order to ascertain viability of the model obtained in Eq. 4 to make describe current system states and predict future states, the model is solved using ODEINT python package for 500 (current) time steps and 1000 (future) time steps. To compare the predicted latent vector obtained from our model, 'ground truth' latent vector is obtained by training the network for 1000 time steps. It can be observed from Figure 9. that the predicted latent vector obtained from solving the model defined by Eq. 4 shows good correlation with the ground truth latent vector.

### 3.2 Zebrafish

The zebrafish video set shows 5 fishes in a rectangular container and all the fishes are in frame for the entire duration of the video. Since the data obtained from videos had significant noise, the trajectories were filtered using Savitzky–Golay filter of window size 31 and polynomial order 2. Similar to ant dataset, the x and y spatial coordinates of 5 fishes form 10 features defining the state of the system at any given time. These features are passed through one layer of LSTM encoder, where these 10 features are mapped to a vector of size 64 and from this encoded vector 1 latent vector is sampled for corresponding time step.

Obtained latent vector for 500 time steps was filtered similar to technique used in ant data set and using SINDy following equation is obtained:

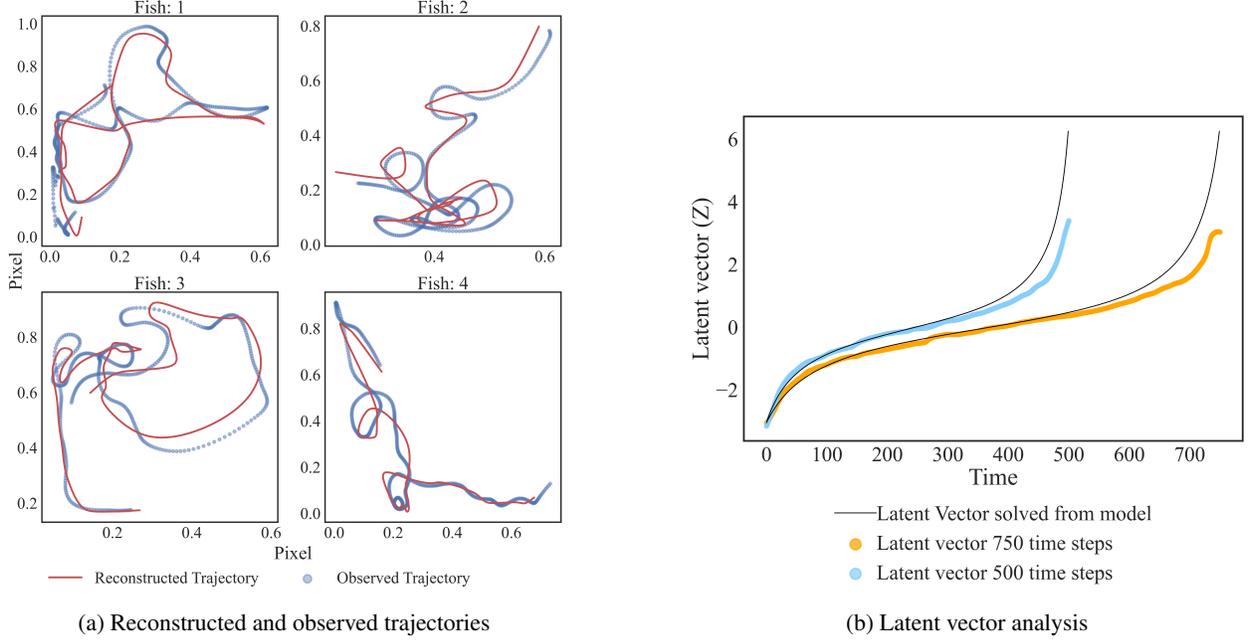


Figure 5: Dominant motion analysis on Zebrafish video dataset. Observed and reconstructed trajectories(a) show that the neural network is able to learn the dynamics of the system. The extrapolated latent vector model solved for 750 time steps shows good correlation with the latent vector latent vector obtained from training to 750 time steps (b).

$$\frac{dz}{dt} = 2.441 + 0.353z + 3.020z^2$$

The obtained model was validate by solving for 500 (current) time steps and for 750 (future) time steps. Comparing with the ground truth, as shown in Figure 5. it can be observed that the model can predict the latent state.

Similar analysis was done for termite dataset and simulated elastic collision, the results are tabulate in Figure6. (For details see SI:I and SI:I)

## 4 Discussion

### 4.1 Analysis of latent vector

Latent vectors across these diverse systems show surprising similarities as long as the reconstructed trajectories are approximate of the true trajectories. Figure below shows comparison of latent vector of all data set. Since the latent vector follows similar trend, the governing equation derived from SINDy framework also provides equations of a common structure. Based on this common structure of the extracted governing equation, it can be inferred that all these diverse multi-particle systems follow a similar governing principle.

To prove this and to find physical significance of latent vector, we need to look at long term trajectories of particles overtime. Figure below shows the 3D plot of trajectories of 2 randomly selected particles from all 4 dataset. It can be observed that all particles have a common trajectory behaviour, in which they start in the center of the frame and with time move towards to boundaries of the frame. Also no two particle (ants, termites and fishes in this case) can share same space at same time. The seemingly random behaviour when observed overtime is analogous to diffusion behaviour of fluid particles in a system in which the particles try to occupy the volume or the space available. In prior work animal movements are described by a class of PDEs called diffusion-taxis equations and methods methods have been developed to parametrize reaction-diffusion equations from animal location data [50–52].

Furthermore, this 'diffusive' behaviour of the particle becomes more prominent when we observer variation in probability density of these particle overtime. Figure 12 shows density distribution of ants calculated using Gaussian kernel estimation, at  $timestep = 0$  and  $timestep = 500$ . It can be observed that the probability of finding ants at the center of the frame at  $t = 0$  is more compared to at  $t = 500$ . Since the behaviour of ants is stochastic, comments can

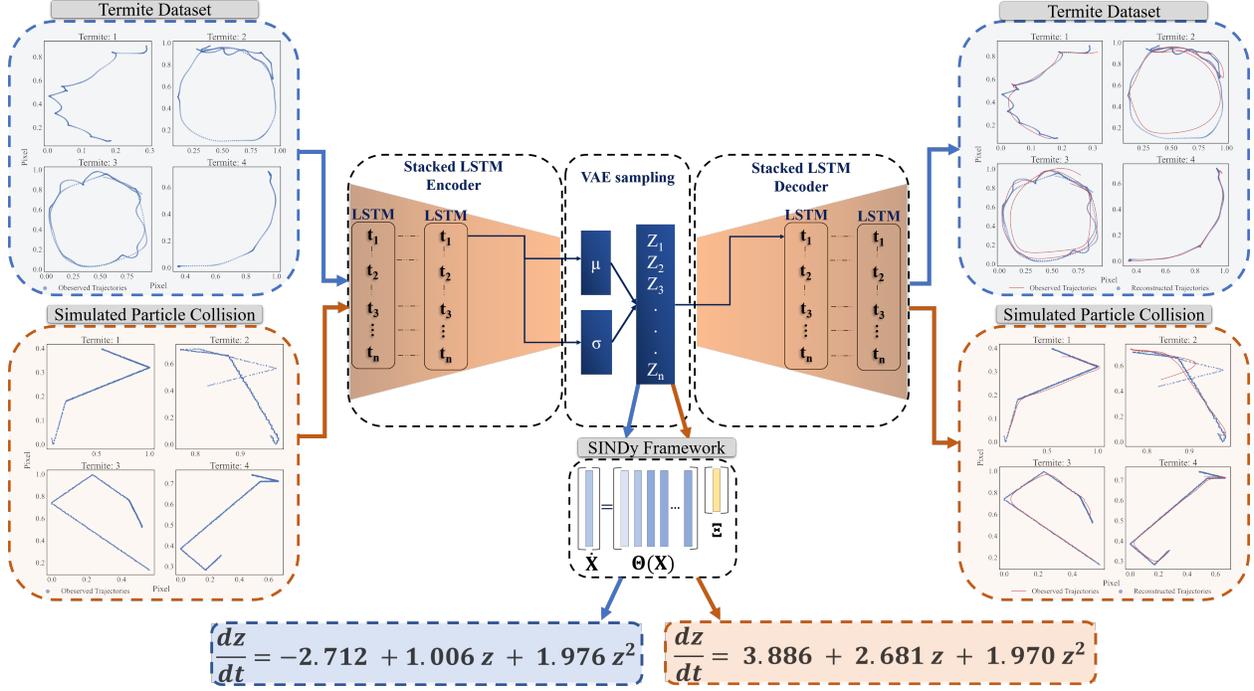


Figure 6: Dominant motion analysis for Termite and collision dataset.

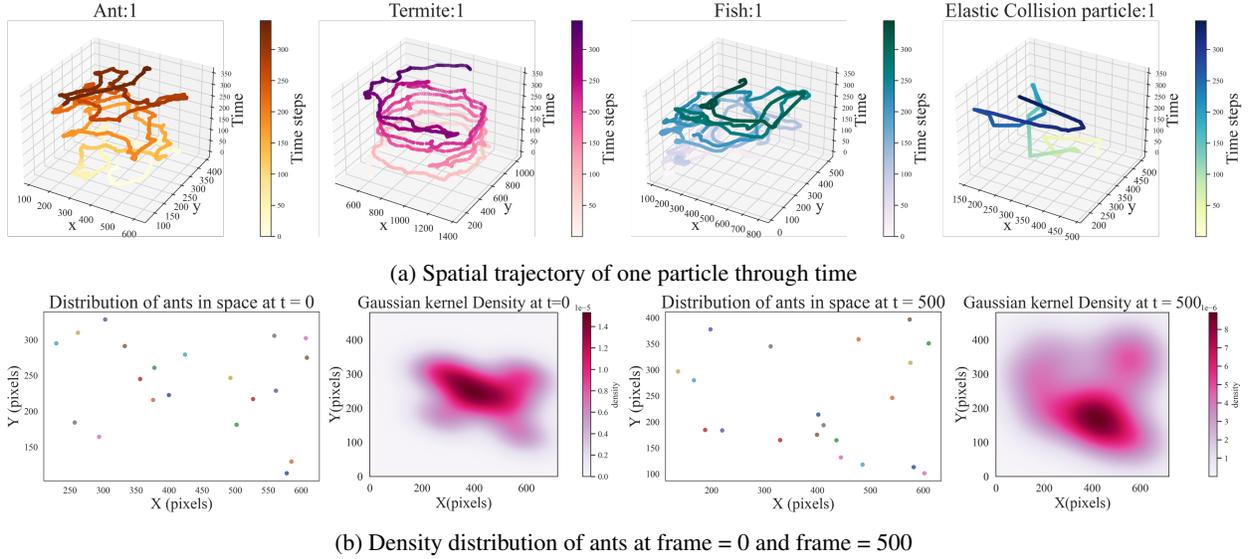


Figure 7: It is observed that the particle in each of the system move towards the boundaries of the system (a). The density plots show that the probability of finding an ant in space in reducing with time i.e. the ants are 'dispersing' and occupying unexplored space in the system.

only be made based on the observed timesteps, hence we can say that the stacked ants have diffused or occupied more space at  $t = 500$  compared to  $t = 0$ .

The neural network is approximating the states of the system at each time step by creating non-linear mappings of state to the latent vector. Based on the LSTM equation (See SI:I) the hidden state at time 't' ( $h_t$ ) can be written as a function of  $x_t$  and  $x_{t-1}$  i.e.

$$h_t = f(x_t, x_{t-1}) \tag{5}$$

Thus, it can be inferred that the latent vector learned by the neural network is a function of change in spatial states of the system, in other words the encoder is embedding the state of the system to a function of average displacement, which is analogous to mean square displacement function in diffusion. This makes our network/algorithm highly useful to study transport phenomena, molecular dynamics, etc. Furthermore, given the stochastic nature of the data used opens up the possibility to use this network to find governing dynamics of a Brownian motion system.

### 4.2 Potential applications of obtained governing equation

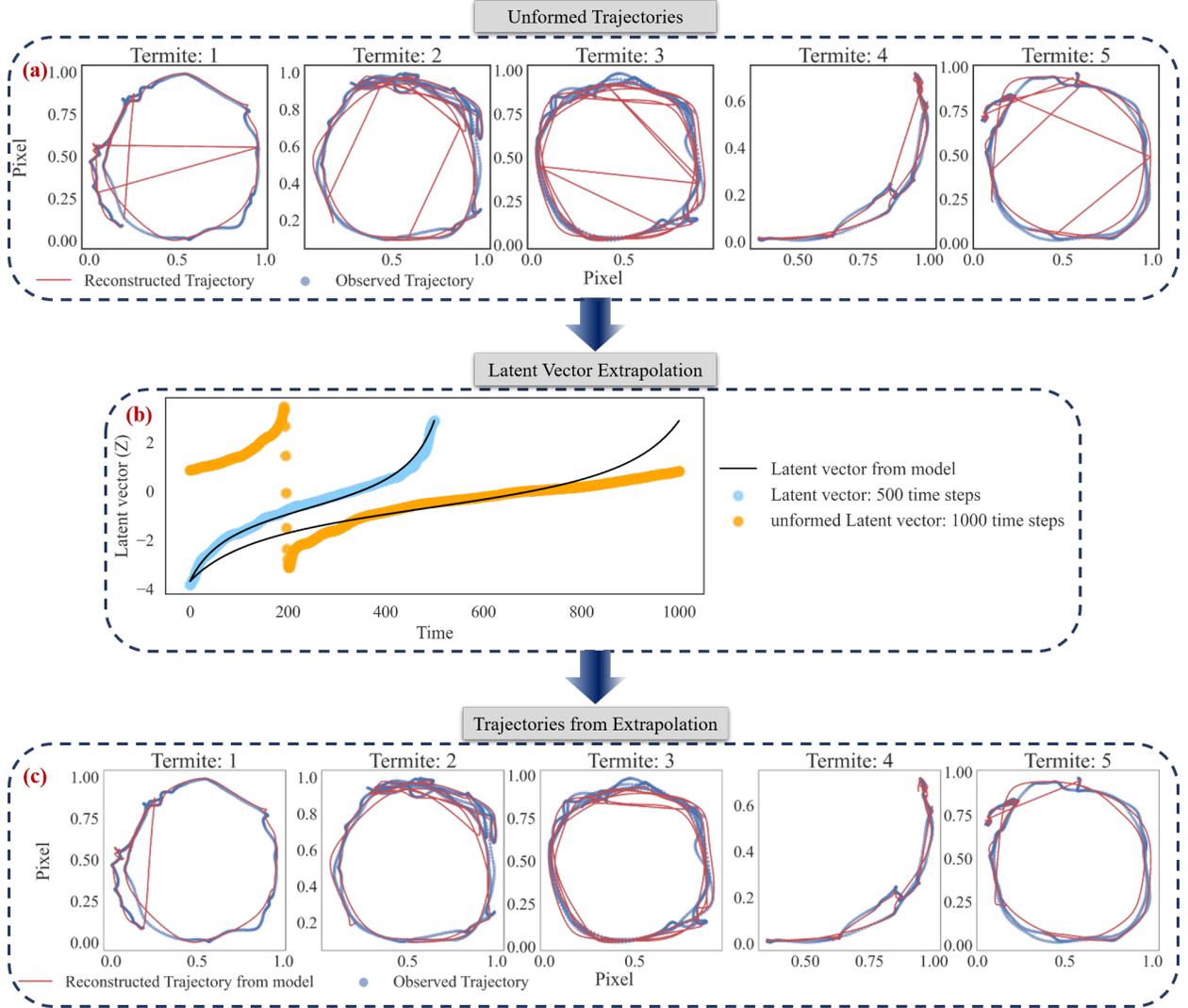


Figure 8: To demonstrate that the model can be used to reconstruct and extrapolate unformed trajectories first we initialize the LSTM VAE for 1000 time steps and stop it prematurely before the states are learned (a). The latent vector of unformed states show significant deviation from the latent vector of fully formed states (b). The latent vector model obtained from training the network 500 time step is solved for 1000 time steps. The extrapolated latent vector is then fed into the decoder to form fully reconstructed states for 1000 time steps(c).

One of the potential application of using the latent vector and consequently obtained governing equation is for anomaly detection. The nature of the latent vector doesn't change with time (see SI:I) unless there is a disturbance or sudden change to the system. Thus the obtained equation can be used to detect anomaly if there are significant changes in the system. Furthermore, it is also observed that the latent vector follows the trend defined by the model as long as

the reconstructed trajectory is similar to original trajectory. Thus, by changing the latent vector and by decoding the modified latent vector, different possible states of the system or simulate anomaly in the system can be simulated.

To prove both of these assertions, we trained the our network on termite dataset for 500 time steps and using SINDy framework extracted the differential equation representing state of the system from  $t = 0$  to  $t = 500$  from the latent vector. To simulate anomalous behaviour, we trained our network on same data set for 1000 time steps, this time stopping it prematurely so that the reconstructed trajectories are not full formed. Figure?? shows the latent vector for 500 time steps and unformed latent vector for 1000 time steps, along with the vector derived from the model. Looking at the reconstructed trajectories and original trajectories, it can be observed that there are significant mismatched states.

We solve the differential equation obtained from latent vector representing 500 time steps, for 1000 time steps. This new vector obtained by solving the equation for 1000 time steps is then passed into the decoder which was stopped prematurely. Based on the reconstructed trajectories, it can be observed in Figure?? that the mismatched states have been significantly reduced.

This proves that the vector obtained from the model can not only be used to represent normal system behaviour but can also be used to filter anomalous/disturbed conditions of the system and conversely be used to simulate anomaly. It should be noted that an encode-decoder model trained at 500 time steps cannot fully predict the future states when latent vector of  $t > 500$  is passed through the decoder, but a reconstructions from a unformed encode-decoder model can reformed by solving the equation derived from latent vector representing previous system states.

## References

- [1] Jordan, M.I.; Mitchell, T.M. Machine learning: Trends, perspectives, and prospects. *Science* 2015, 349, 255–260.
- [2] Marx, V. Biology: The big challenges of big data. *Nature* 2013, 498, 255–260.
- [3] M. Schmidt, H. Lipson, Distilling free-form natural laws from experimental data. *Science* 324, 81–85 (2009).
- [4] Rosenblatt, F. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychol. Rev.* 1958,65,86.
- [5] LeCun, Y., Bengio, Y. , Hinton, G. Deep learning. *Nature* 521, 436–444 (2015). <https://doi.org/10.1038/nature14539>
- [6] J. M. Benitez, J. L. Castro and I. Requena, "Are artificial neural networks black boxes?," in *IEEE Transactions on Neural Networks*, vol. 8, no. 5, pp. 1156-1164, Sept. 1997, doi: 10.1109/72.623216.
- [7] A. Adadi and M. Berrada, "Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI)," in *IEEE Access*, vol. 6, pp. 52138-52160, 2018, doi: 10.1109/ACCESS.2018.2870052.
- [8] Buhmester, V., Münch, D. and Arens, M., 2019. Analysis of explainers of black box deep neural networks for computer vision: A survey. *arXiv preprint arXiv:1911.12116*.
- [9] S. H. Rudy, S. L. Brunton, J. L. Proctor, and J. N. Kutz, Data-driven discovery of partial differential equations, *Science Advances* 26 Apr 2017: Vol. 3, no. 4, e1602614, DOI:10.1126/sciadv.1602614
- [10] J. Bongard H. Lipson, Automated reverse engineering of nonlinear dynamical systems. *Proceedings of the National Academy of Sciences* Jun 2007, 104 (24) 9943-9948; DOI: 10.1073/pnas.0609476104
- [11] Raissi, M., Perdikaris, P. and Karniadakis, G.E., 2017. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10561*.
- [12] Goodfellow, I., Bengio, Y., Courville, A. and Bengio, Y., 2016. *Deep learning* (Vol. 1, No. 2). Cambridge: MIT press.
- [13] Raissi, M., Perdikaris, P. and Karniadakis, G.E., 2018. Multistep neural networks for data-driven discovery of nonlinear dynamical systems. *arXiv preprint arXiv:1801.01236*.
- [14] Hochreiter, S. and Schmidhuber, J., 1997. Long short-term memory. *Neural computation*, 9(8), pp.1735-1780.
- [15] Chung, J., Gulcehre, C., Cho, K. and Bengio, Y., 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*.
- [16] Hochreiter, S., 1998. The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6(02), pp.107-116.
- [17] Dey, R. and Salemt, F.M., 2017, August. Gate-variants of gated recurrent unit (GRU) neural networks. In 2017 IEEE 60th international midwest symposium on circuits and systems (MWSCAS) (pp. 1597-1600). IEEE.
- [18] Elsworth, S. and Güttel, S., 2020. Time Series Forecasting Using LSTM Networks: A Symbolic Approach. *arXiv preprint arXiv:2003.05672*.

- [19] Graves, A., Jaitly, N. and Mohamed, A.R., 2013, December. Hybrid speech recognition with deep bidirectional LSTM. In 2013 IEEE workshop on automatic speech recognition and understanding (pp. 273-278). IEEE.
- [20] Graves, A., Mohamed, A.R. and Hinton, G., 2013, May. Speech recognition with deep recurrent neural networks. In 2013 IEEE international conference on acoustics, speech and signal processing (pp. 6645-6649). IEEE.
- [21] Cho, K., Van Merriënboer, B., Bahdanau, D. and Bengio, Y., 2014. On the properties of neural machine translation: Encoder-decoder approaches. arXiv preprint arXiv:1409.1259.
- [22] Luong, M.T., Pham, H. and Manning, C.D., 2015. Effective approaches to attention-based neural machine translation. arXiv preprint arXiv:1508.04025.
- [23] Alahi, A., Goel, K., Ramanathan, V., Robicquet, A., Fei-Fei, L. and Savarese, S., 2016. Social lstm: Human trajectory prediction in crowded spaces. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 961-971).
- [24] Gupta, A., Johnson, J., Fei-Fei, L., Savarese, S. and Alahi, A., 2018. Social gan: Socially acceptable trajectories with generative adversarial networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 2255-2264).
- [25] Fragkiadaki, K., Levine, S., Felsen, P. and Malik, J., 2015. Recurrent network models for human dynamics. In Proceedings of the IEEE International Conference on Computer Vision (pp. 4346-4354).
- [26] Vlachas, P.R., Byeon, W., Wan, Z.Y., Sapsis, T.P. and Koumoutsakos, P., 2018. Data-driven forecasting of high-dimensional chaotic systems with long short-term memory networks. Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences, 474(2213), p.20170844.
- [27] Sridhar, V.H., Roche, D.G. and Gingsins, S., 2019. Tracktor: Image-based automated tracking of animal movement and behaviour. Methods in Ecology and Evolution, 10(6), pp.815-820.
- [28] <http://www.eecs.qmul.ac.uk/~andrea/thdt.html>
- [29] Lukezic, A., Vojir, T., Čehovin Zajc, L., Matas, J. and Kristan, M., 2017. Discriminative correlation filter with channel and spatial reliability. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 6309-6318).
- [30] Rudy, S.H., Brunton, S.L., Proctor, J.L. and Kutz, J.N., 2017. Data-driven discovery of partial differential equations. Science Advances, 3(4), p.e1602614.
- [31] Champion, K., Lusch, B., Kutz, J.N. and Brunton, S.L., 2019. Data-driven discovery of coordinates and governing equations. Proceedings of the National Academy of Sciences, 116(45), pp.22445-22451.
- [32] Patraucean, V., Handa, A. and Cipolla, R., 2015. Spatio-temporal video autoencoder with differentiable memory. arXiv preprint arXiv:1511.06309.
- [33] Tu, J., Liu, H., Meng, F., Liu, M. and Ding, R., 2018, October. Spatial-temporal data augmentation based on LSTM autoencoder network for skeleton-based human action recognition. In 2018 25th IEEE International Conference on Image Processing (ICIP) (pp. 3478-3482). IEEE.
- [34] Marchi, E., Vesperini, F., Eyben, F., Squartini, S. and Schuller, B., 2015. A novel approach for automatic acoustic novelty detection using a denoising autoencoder with bidirectional LSTM neural networks. In Proceedings 40th IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP 2015 (pp. 5-pages).
- [35] Zhao, F., Feng, J., Zhao, J., Yang, W. and Yan, S., 2017. Robust lstm-autoencoders for face de-occlusion in the wild. IEEE Transactions on Image Processing, 27(2), pp.778-790.
- [36] Li, J., Luong, M.T. and Jurafsky, D., 2015. A hierarchical neural autoencoder for paragraphs and documents. arXiv preprint arXiv:1506.01057.
- [37] Srivastava, N., Mansimov, E. and Salakhudinov, R., 2015, June. Unsupervised learning of video representations using lstms. In International conference on machine learning (pp. 843-852).
- [38] Moorcroft, P.R., Lewis, M.A. and Crabtree, R.L., 2006. Mechanistic home range models capture spatial patterns and dynamics of coyote territories in Yellowstone. Proceedings of the Royal Society B: Biological Sciences, 273(1594), pp.1651-1659.
- [39] Moorcroft, P.R., Moorcroft, P. and Lewis, M.A., 2006. Mechanistic home range analysis. Princeton University Press.
- [40] Hays, G.C., Ferreira, L.C., Sequeira, A.M., Meekan, M.G., Duarte, C.M., Bailey, H., Bailleul, F., Bowen, W.D., Caley, M.J., Costa, D.P. and Eguíluz, V.M., 2016. Key questions in marine megafauna movement ecology. Trends in ecology & evolution, 31(6), pp.463-475.

- [41] S. L. Brunton, J. L. Proctor, J. N. Kutz, and W. Bialek, Discovering governing equations from data by sparse identification of nonlinear dynamical systems, *Proc. Natl. Acad. Sci. U. S. A.* 113, 3932 (2016)
- [42] H. Schaeffer, Learning partial differential equations via data discovery and sparse optimization, *Proc. R. Soc. A Math. Phys. Eng. Sci.* 473, (2017).
- [43] M. Raissi, P. Perdikaris, and G. E. Karniadakis, Machine learning of linear differential equations using Gaussian processes, *J. Computational Phys.* 348, 683 (2017).
- [44] Z. Long, Y. Lu, and B. Dong, PDE-Net 2.0: Learning PDEs from data with a numeric-symbolic hybrid deep network, *J. Comput. Phys.* 399, 108925 (2019).
- [45] K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991.
- [46] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems*, 2(4):303–314, Dec. 1989.
- [47] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Numerical Gaussian processes for time-dependent and nonlinear partial differential equations. *SIAM Journal on Scientific Computing*, 40(1):A172–A198, 2018.
- [48] de Silva et al., (2020). PySINDy: A Python package for the sparse identification of nonlinear dynamical systems from data. *Journal of Open Source Software*, 5(49), 2104, <https://doi.org/10.21105/joss.02104>
- [49] Scikit-learn: Machine Learning in Python, Pedregosa et al., *JMLR* 12, pp. 2825-2830, 2011
- [50] Hefley, T.J., Hooten, M.B., Russell, R.E., Walsh, D.P., Powell, J.A. (2017). When mechanism matters: Bayesian forecasting using models of ecological diffusion. *Ecol. Lett.*, 20, 640–650.
- [51] Garlick, M.J., Powell, J.A., Hooten, M.B. & McFarlane, L.R. (2011). Homogenization of large-scale movement models in ecology. *Bull. Math. Biol.*, 73, 2088–2108.
- [52] Thurfjell, H., Ciuti, S., Boyce, M. (2014). Applications of step-selection functions in ecology and conservation. *Mov. Ecol.*, 2, 4.
- [53] Kondo, S., Miura, T. (2010). Reaction-diffusion model as a framework for understanding biological pattern formation. *Science*, 329, 1616–1620
- [54] Li, X., Jiang, W., Shi, J. (2013). Hopf bifurcation and Turing instability in the reaction–diffusion holling–tanner predator–prey model. *IMA J. Appl. Math.*, 78, 287–306
- [55] Wu, J., Yildirim, I., Lim, J.J., Freeman, W.T., Tenenbaum, J.B.: Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. In: *NIPS* (2015)
- [56] Bhat, Kiran S., Seitz, Steven M., and Popovic, Jovan. Computing the physical parameters of rigidbody motion from video. In Heyden, Anders, Sparr, Gunnar, Nielsen, Mads, and Johansen, Peter (eds.), *ECCV* (1), volume 2350 of *Lecture Notes in Computer Science*, pp. 551–565. Springer, 2002. ISBN 3-540-43745-2.
- [57] Mottaghi, Roozbeh, Bagherinezhad, Hessam, Rastegari, Mohammad, and Farhadi, Ali. Newtonian image understanding: Unfolding the dynamics of objects in static images. *CoRR*, abs/1511.04048, 2015.

## Part I

# Supplementary information

Neural network architecture can be found here: <https://github.com/BaratiLab/LSTM-VAE-for-dominant-motion-extraction>

### SI 1: LSTM cell

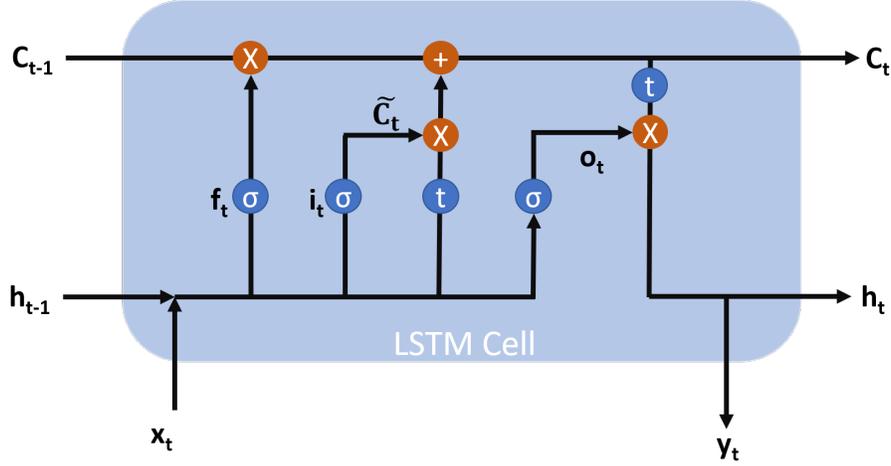


Figure 9: Spatial trajectory of one particle through time

The LSTM cell (Long short term memory) is a gated neural network, with 3 gates namely 1) input gate, 2) forget gate and 3) output gate. Figure9 shows a typical LSTM cell. The gates of LSTM cell have sigmoid activation. The output from the cell and the hidden state of the cell depends on hidden state and cell state at previous time step as well as on the input state at current time step.

LSTM gates are governed by the following set of equations:

$$\begin{aligned} i_t &= \sigma(w_i[h_{t-1}, x_t] + b_i) \\ f_t &= \sigma(w_f[h_{t-1}, x_t] + b_f) \\ o_t &= \sigma(w_o[h_{t-1}, x_t] + b_o) \end{aligned}$$

Where,  $i_t, f_t, o_t$  is the input, forget and output gates respectively.  $x_t$  is the input at time 't',  $h_{t-1}$  is hidden state at time 't-1' and  $w, b$  are the weights and bias of respective gates.

Based on the these gates, the outputs from the cell are:

$$\begin{aligned} C_t &= f_t * c_{t-1} + i_t * \tanh(w_c[h_{t-1}, x_t] + b_c) \\ h_t &= o_t * \tanh(C_t) \end{aligned}$$

Here,  $C_t$  is the cell state and  $h_t$  is the hidden state.

### SI 2: Simulated Elastic collision system

This system is a physical simulation of two-dimensional interaction of particles in a confined space. In this paper we have used 5 particles of same size which undergoes elastic collision when they interact with each other or with the walls. In case of the interaction the total kinetic energy remains unchanged and the momentum is conserved. Based on law of conservation of momentum:

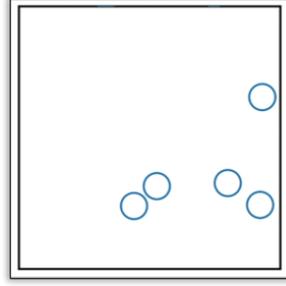


Figure 10: Simulated elastic collision video frame.

$$m_1 * u_1 + m_2 * u_2 = m_1 * v_1 + m_2 * v_2$$

where,  $m$  is the mass,  $u$  is initial velocity and  $v$  is final velocity after the collision. In this simulation the initial velocity at the start is randomly chosen. but since we have particle of same size therefore same mass the interaction of particle in this simulation is defined by:

$$u_1 + u_2 = v_1 + v_2$$

In order to track each particle we are using CSRT algorithm, bounding boxes were created manually around each particle. Form this video we get a total of 10 spatial features ( $x$ ,  $y$  coordinate for each particle) for each frame(times step). These 10 features are used as inputs to the LSTM VAE, where these features are compressed to a latent vector representation of size 1 for each corresponding time step. SINDy framework is then used to find the underlining equation represented by the latent vector.

### SI 3: Termite dataset

The termite dataset consist of 8 termites enclosed in a circular petri-dish as shown in Figure??b. Using CSRT tracker 64 spatial features ( $x,y$  coordinates) for each termite for each frame were captured. These 64 features were scaled in range 0-1 and then fed into LSTM VAE to obtain latent vector of size 1 for each time-step. The latent vector was smoothed using Savitzky-Golay filter of window size 51 and order 1. Using SINDy framework on the filtered latent vector the system dynamics was found.

### SI 4: Extracted Latent vector for multiple time steps

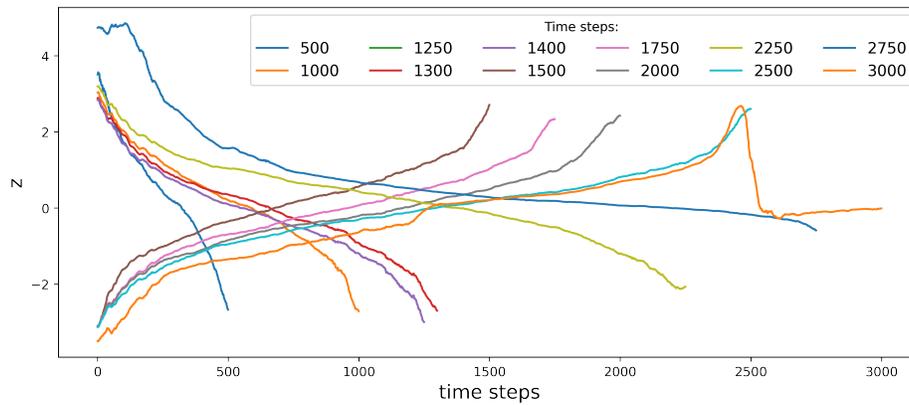


Figure 11: Latent vector for ants in box dataset extracted for multiple time steps.

In order to ascertain the nature of latent vectors remains the same for far in future time, we trained LSTM VAE and extracted latent vector for various time steps. For this we used the ground truth trajectory provided with the ant video to

compare and ensure that the nature of latent vector is same for the trajectory extracted using CSRT algorithm. The trajectory data has features for 10399 time steps. Figure 11 show the latent vector for upto 3000 time steps. It can be observed that the long term dynamics of the system follows similar structure as the short term dynamics. The ‘flip’ in the sign of latent vector observed after 1400 time step has no affect on the system dynamics except for the change in signs of the coefficients. It was also observed that the even without the changing the signs of coefficients when dynamical system was solved for future time steps and fed into the trained decoder, the reconstructed trajectories were similar to the ground truth trajectories.

Figure 12 shows the similarity in the latent vectors obtained by training the network on different datasets.

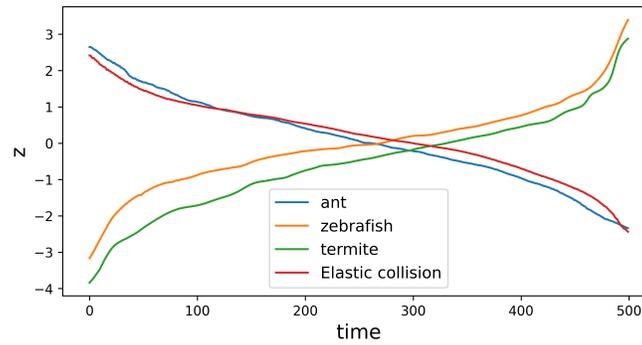


Figure 12: Latent vector of 500 time steps for ant, termite, elastic collision and fish dataset.