

Interpretable Conservation Law Estimation

by Deriving the Symmetries of Dynamics from Trained Deep Neural Networks

本武 陽一

統計的機械学習研究センター 特任助教

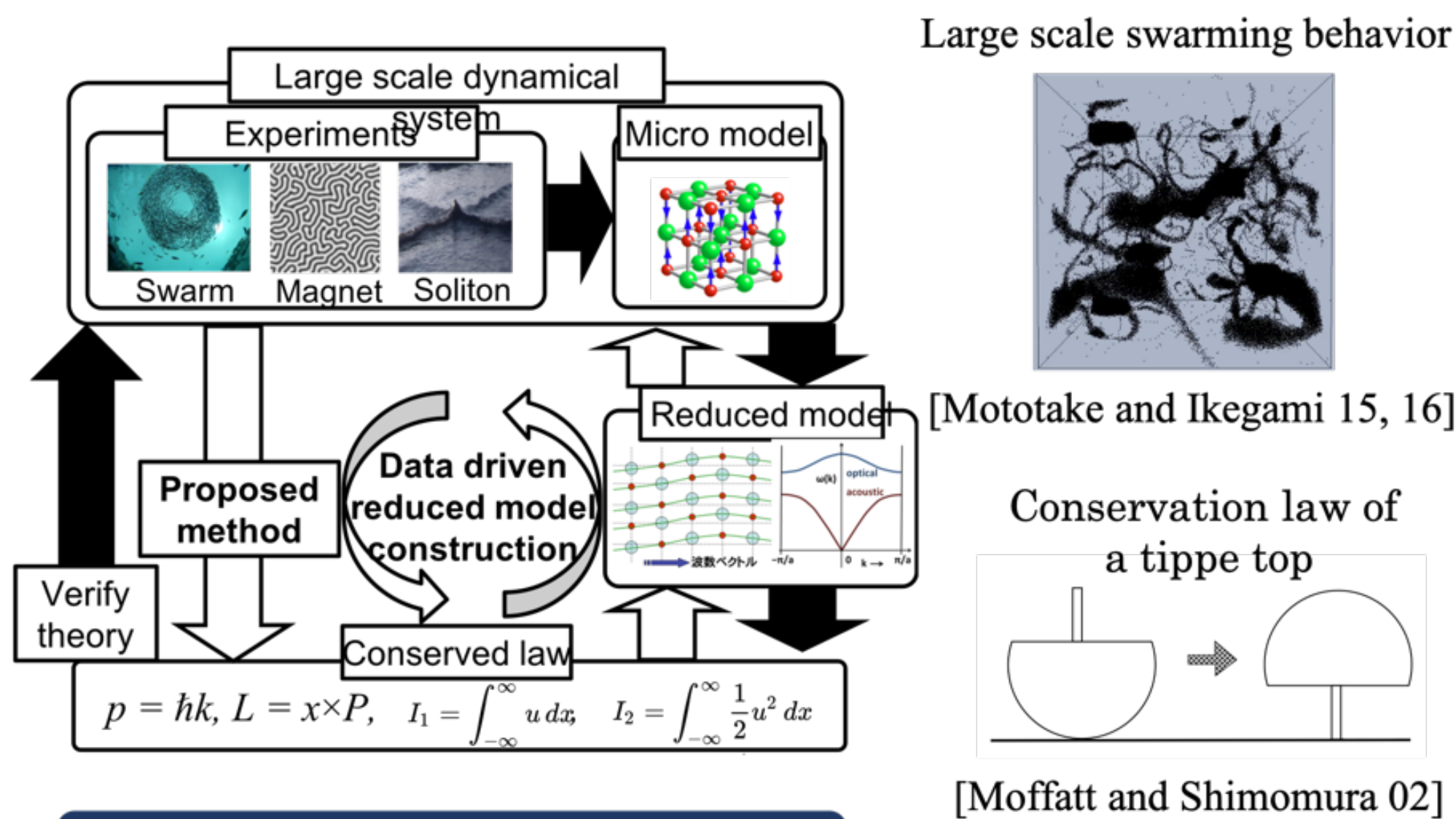
●Abstraction

We propose a framework to connect the deep neural networks (DNNs) to a study on conservation laws, the fundamental topic in physics. Our framework does not aspire to conduct physical data analysis using DNNs, but rather to find interpretable physical information from trained DNNs. Using Noether's theorem and an efficient sampling method, the inference of conservation laws was achieved through extracting symmetries from trained DNNs. It applies to a wide range of phenomena, including symmetries, because it only uses the general properties of DNNs and Hamiltonian mechanics. It will support physicists in constructing a minimal model of complex systems.

Paper: <https://arxiv.org/abs/2001.00111>Video: <https://www.youtube.com/user/MrKeaton2012>

●Introduction

Reduced model and conserved law

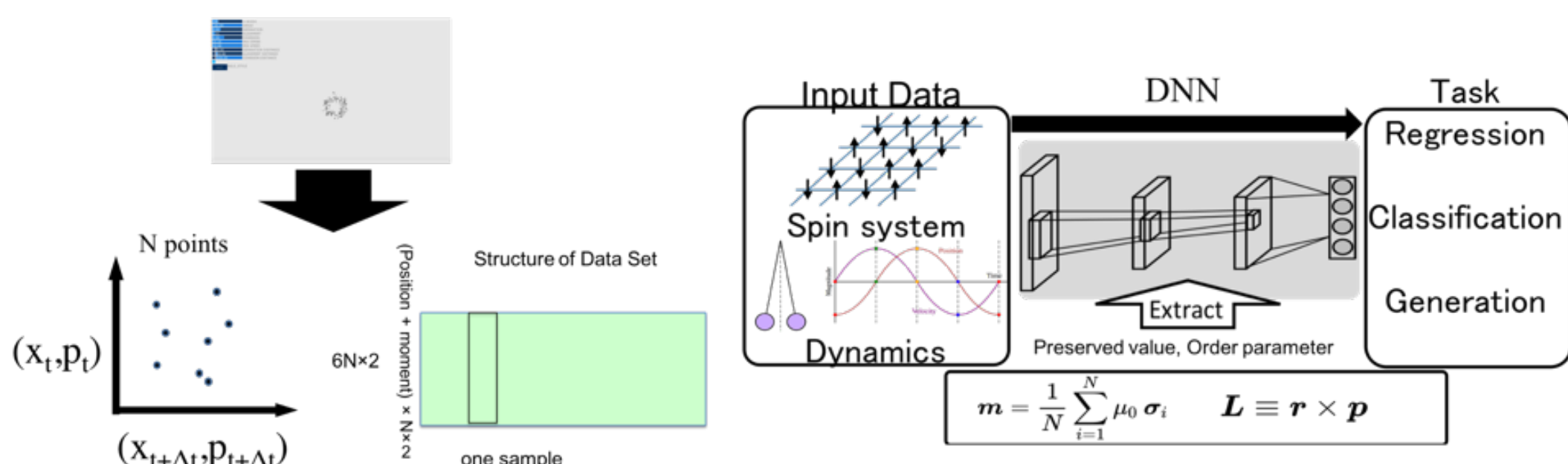


Noether's theorem

Noether's theorem connects the continuous symmetry of the Hamiltonian system and its conservation law. Considering the Hamiltonian systems in 2- d dimensional phase space (\mathbf{q}, \mathbf{p}) , let the system's Hamiltonian be $H(\mathbf{q}, \mathbf{p})$. Assuming that the Hamiltonian $H(\mathbf{q}, \mathbf{p})$ and the canonical equations (equations of motion), $\frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial \mathbf{q}} = -\dot{\mathbf{p}}$ and $\frac{\partial H(\mathbf{q}, \mathbf{p})}{\partial \mathbf{p}} = \dot{\mathbf{q}}$, are invariant for infinitesimal transformation, $(t', q'_i, p'_i) = (t + \delta t, q_i + \delta q_i, p_i + \delta p_i)$, where $i = 1 \sim d$. Then, based on Noether's theorem, the conserved value G satisfies the following equation:

$$(\delta q_j, \delta p_j) = \left(\frac{\partial G_\delta}{\partial p_j}, -\frac{\partial G_\delta}{\partial q_j} \right). \quad (1)$$

Data manifold of time series data



The candidate of invariant transformation of Hamiltonian system must be in transformations $(\mathbf{q}, \mathbf{p}) \mapsto (\mathbf{Q}, \mathbf{P}) := (\mathbf{Q}(\mathbf{q}, \mathbf{p}), \mathbf{P}(\mathbf{q}, \mathbf{p}))$ which satisfy below condition:

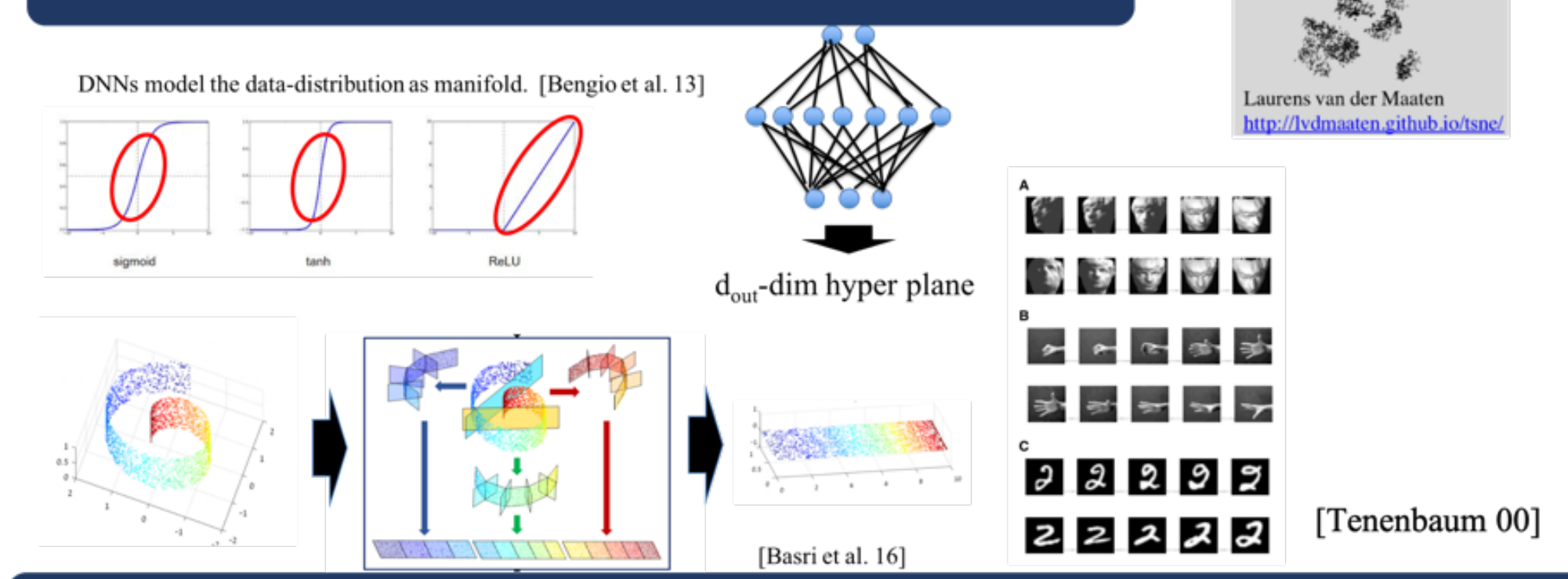
$$\left\{ \mathbf{q}_{t+\Delta t}, \mathbf{p}_{t+\Delta t}, \mathbf{q}_t, \mathbf{p}_t \mid H(\mathbf{q}_t, \mathbf{p}_t) = E, \mathbf{p}_{t+\Delta t} = \mathbf{p}_t - \frac{\partial H(\mathbf{q}_t, \mathbf{p}_t)}{\partial \mathbf{q}_t} \cdot \mathbf{q}_{t+\Delta t} = \mathbf{q}_t + \frac{\partial H(\mathbf{q}_t, \mathbf{p}_t)}{\partial \mathbf{p}_t} \cdot \mathbf{p}_{t+\Delta t} \right\} \quad (\mathbf{q}_t, \mathbf{p}_t)$$

$$= \left\{ \mathbf{Q}_{T+\Delta T}, \mathbf{P}_{T+\Delta T}, \mathbf{Q}_T, \mathbf{P}_T \mid H(\mathbf{Q}_T, \mathbf{P}_T) = E, \mathbf{P}_{T+\Delta T} = \mathbf{P}_T - \frac{\partial H(\mathbf{Q}_T, \mathbf{P}_T)}{\partial \mathbf{Q}_T} \cdot \mathbf{Q}_{T+\Delta T} = \mathbf{Q}_T + \frac{\partial H(\mathbf{Q}_T, \mathbf{P}_T)}{\partial \mathbf{P}_T} \cdot \mathbf{P}_{T+\Delta T} \right\} \quad (\mathbf{Q}_T, \mathbf{P}_T)$$

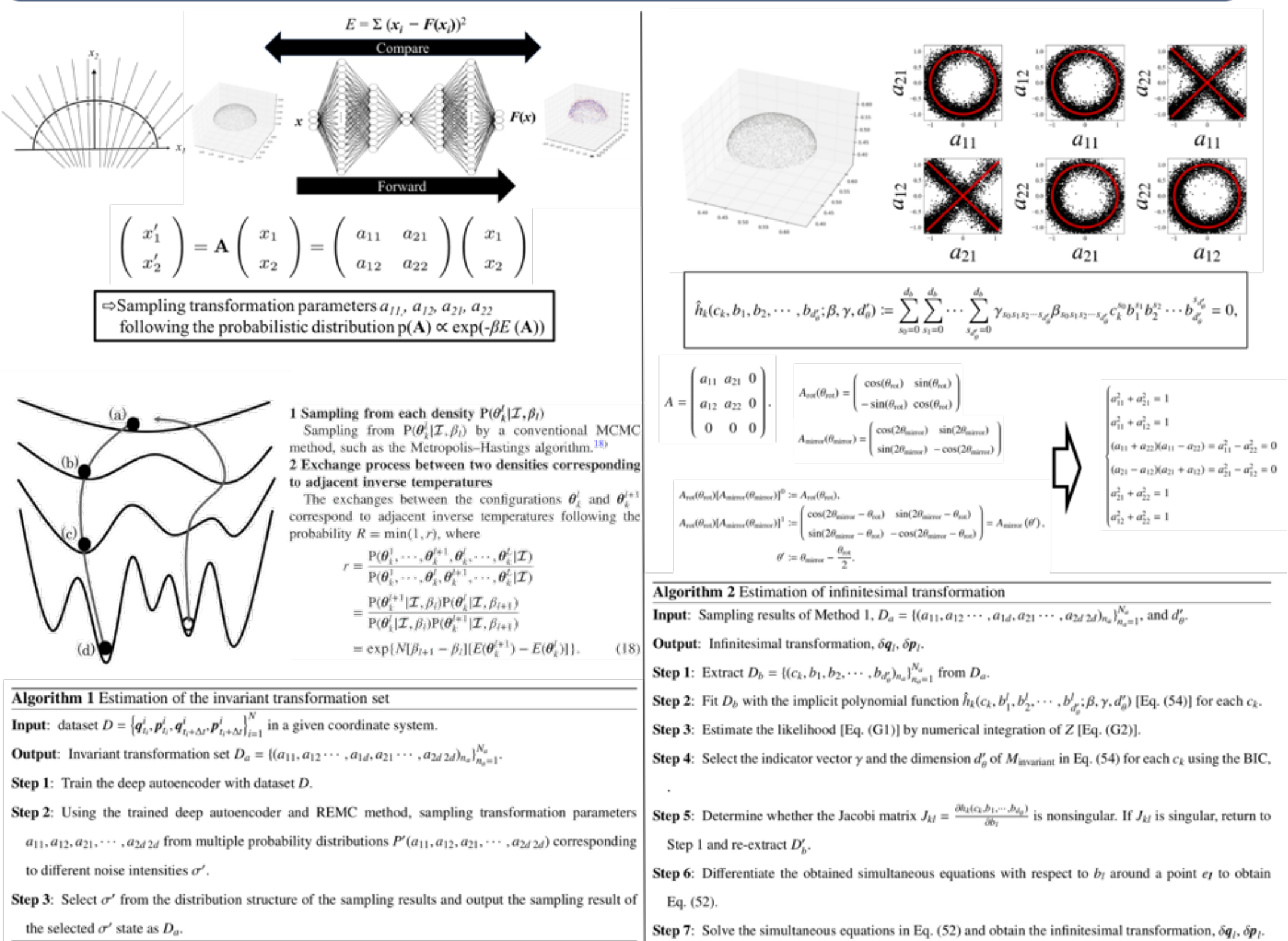
It's possible to estimate the conservation law directly from the time-series dataset!!

●Method

DNN and Manifold Hypothesis

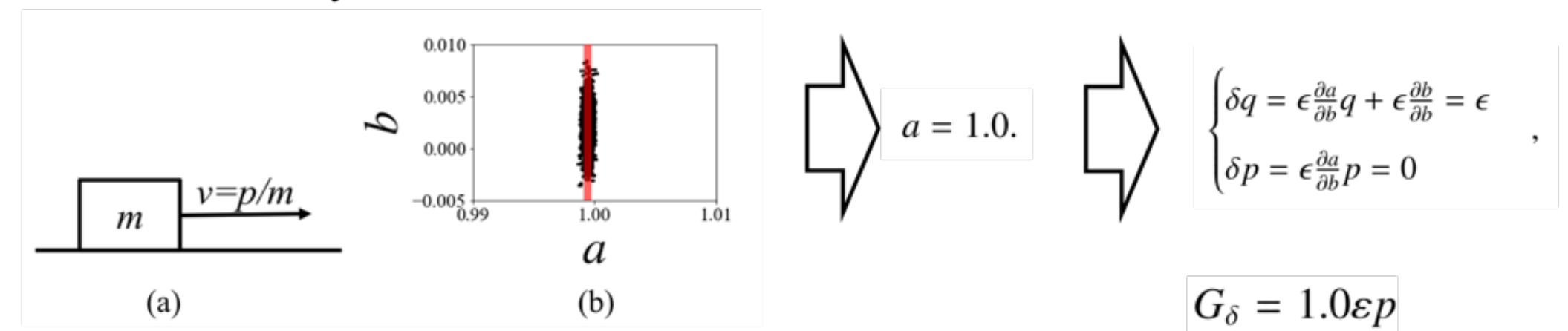


Method to estimate symmetry of manifold structure of time series data

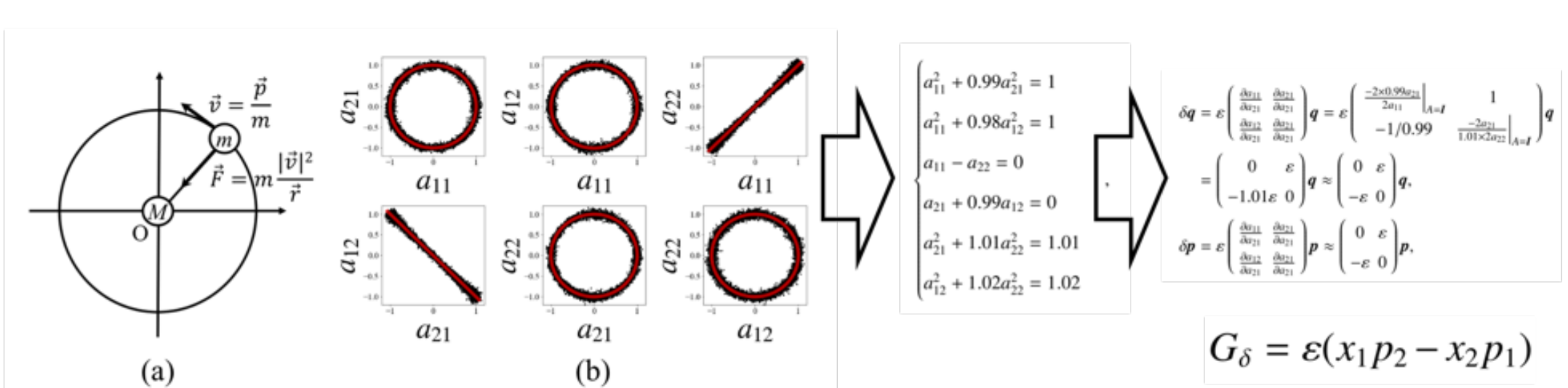


●Results

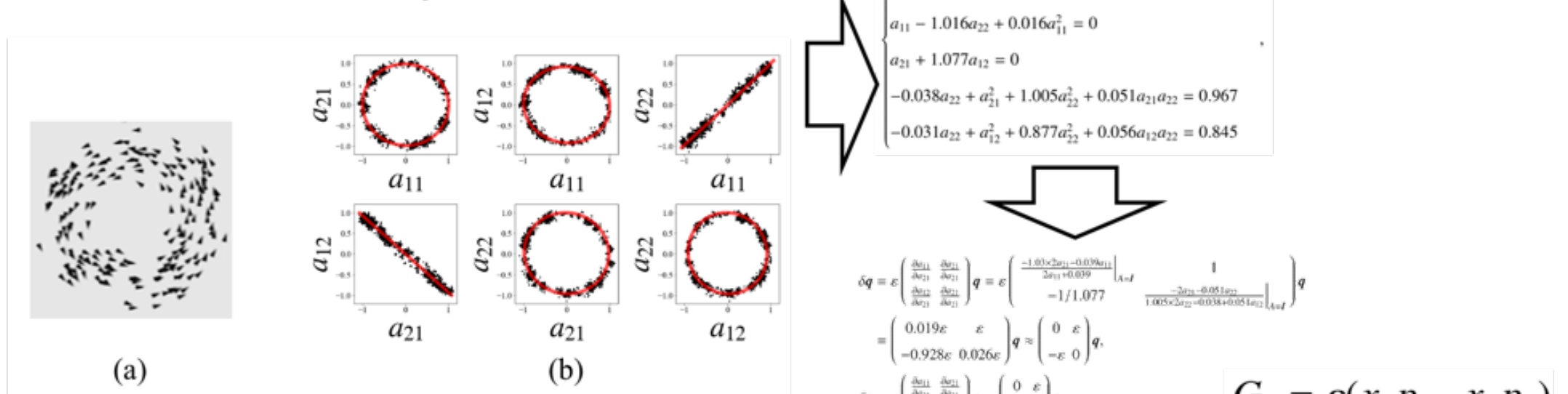
• Constant-velocity linear motion



• Two-dimensional central force system



• Collective motion system



●Summary and Discussion

Targeting the collective motion of living things, we show that our framework will be a powerful tool for physicists constructing a reduced model of complex systems. We believe that our work greatly contributes to the construction of a reduced model of complex systems such as a collective motion system