


# Replacing Neural Networks by Optimal Analytical Predictors for the Detection of Phase Transitions

Julian Arnold<sup>1,\*</sup> and Frank Schäfer<sup>1,2,†</sup>

<sup>1</sup>*Department of Physics, University of Basel, Klingelbergstrasse 82, 4056 Basel, Switzerland*

<sup>2</sup>*CSAIL, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA*

 (Received 21 February 2022; revised 2 July 2022; accepted 3 August 2022; published 28 September 2022)

Identifying phase transitions and classifying phases of matter is central to understanding the properties and behavior of a broad range of material systems. In recent years, machine-learning (ML) techniques have been successfully applied to perform such tasks in a data-driven manner. However, the success of this approach notwithstanding, we still lack a clear understanding of ML methods for detecting phase transitions, particularly of those that utilize neural networks (NNs). In this work, we derive analytical expressions for the optimal output of three widely used NN-based methods for detecting phase transitions. These optimal predictions correspond to the results obtained in the limit of high model capacity. Therefore, in practice, they can, for example, be recovered using sufficiently large, well-trained NNs. The inner workings of the considered methods are revealed through the explicit dependence of the optimal output on the input data. By evaluating the analytical expressions, we can identify phase transitions directly from experimentally accessible data without training NNs, which makes this procedure favorable in terms of computation time. Our theoretical results are supported by extensive numerical simulations covering, e.g., topological, quantum, and many-body localization phase transitions. We expect similar analyses to provide a deeper understanding of other classification tasks in condensed matter physics.

DOI: [10.1103/PhysRevX.12.031044](https://doi.org/10.1103/PhysRevX.12.031044)

Subject Areas: Computational Physics  
Condensed Matter Physics  
Quantum Physics

## I. INTRODUCTION

In recent years, machine learning (ML) has been used extensively to approach complex physics problems [1–3]. Among these applications, the task of classifying phases of matter and the identification of phase transitions is particularly exciting [4–7], as it could enable the autonomous discovery of novel phases of matter. Classical ML methods have successfully revealed the phase diagrams of a plethora of systems based on data from experimental measurements [8–12] and numerical simulations [4,5,13–39]. Many of the most powerful ML methods for detecting phase transitions utilize neural networks (NNs) at their core [4,5,14–26,28,29,31–34,36–39]. Prominent examples are supervised learning [4], the learning-by-confusion

scheme [5,22], and the prediction-based method [29,31,34], which are often applied in conjunction [10,23,31].

All three methods follow a similar work flow, which is illustrated in Fig. 1 (steps 1–3). They take as input samples that represent the state of a physical system at various values of a tuning parameter. The samples are processed by an NN whose parameters are tuned to minimize a specific loss function. By analyzing the NN predictions, one can compute a scalar quantity that highlights the critical value of the tuning parameter at which the system’s state changes most. As such, this quantity highlights phase boundaries and serves as an indicator for phase transitions. The decision whether the change corresponds to a crossover or a phase transition does, however, requires further analysis, such as finite-size scaling. The three methods differ in their choice of loss function, i.e., in the formulation of the underlying classification or regression task, and, thus, in the resulting indicator for phase transitions.

NNs are universal function approximators [42–45]. This fact makes supervised learning, the learning-by-confusion scheme, and the prediction-based method extremely powerful and has played a central role in the original conception of these methods [4,5,29]. Namely, the use of NNs for detecting phase transitions from data has been inspired by

\*julian.arnold@unibas.ch

†franksch@mit.edu

*Published by the American Physical Society under the terms of the Creative Commons Attribution 4.0 International license. Further distribution of this work must maintain attribution to the author(s) and the published article’s title, journal citation, and DOI.*

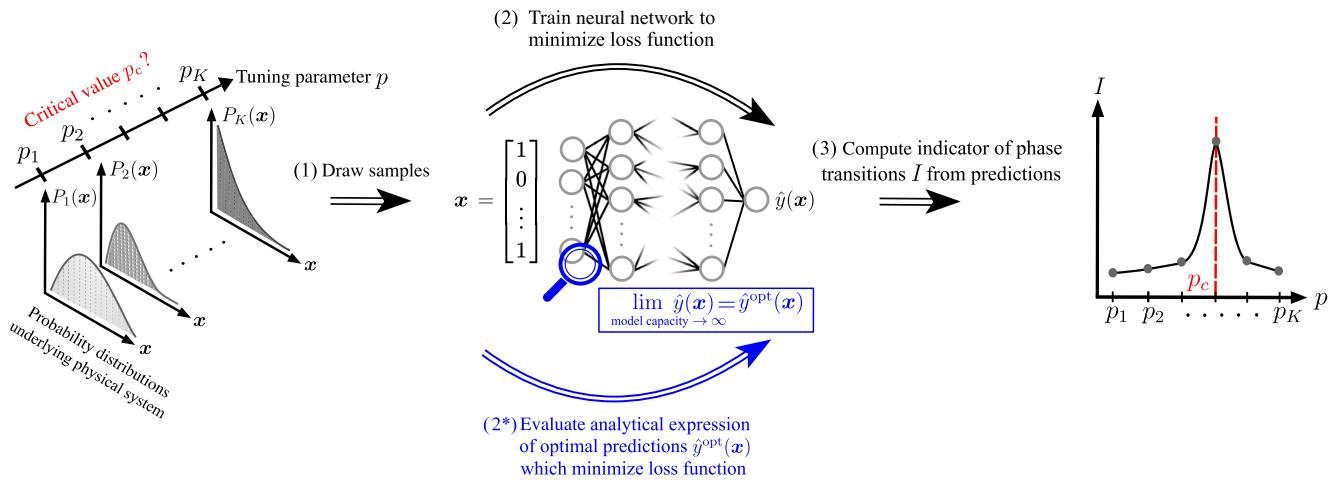


FIG. 1. Schematic representation of the setup and work flow of supervised learning, the learning-by-confusion scheme, and the prediction-based method for detecting phase transitions from data. The physical system under consideration is characterized by a tuning parameter  $p$ . The goal is to identify the critical value of the tuning parameter  $p_c$  at which the system transitions from one phase to another. In a first step (step 1), the state  $x$  of the physical system is (repeatedly) sampled at various values of the tuning parameter  $\{p_1, p_2, \dots, p_K\}$ , where  $\{P_1(x), P_2(x), \dots, P_K(x)\}$  are the corresponding probability distributions. Based on these samples, a neural network (NN) is trained to perform a particular classification or regression task; i.e., its tunable parameters are updated to minimize a particular loss function (step 2). The three ML methods for detecting phase transitions differ in their formulation of the underlying NN tasks. Having trained the NN, its predictions  $\hat{y}$  are used to compute the value of an indicator of phase transitions  $I$  at fixed values of the tuning parameter (step 3). Ideally, the indicator has a local maximum at  $p_c$ , where the largest change in the state of the system occurs. As a result, the ML methods then autonomously highlight phase boundaries along the chosen scanning range of the tuning parameter. Note that the indicators of phase transitions obtained with supervised learning, the learning-by-confusion scheme, and the prediction-based method differ. The contribution of our work is highlighted in blue: We derive analytical expressions for the optimal predictions  $\hat{y}^{\text{opt}}$  of the NNs used in these three methods. The optimal predictions minimize the corresponding loss function and can, thus be achieved by NNs whose capacity, i.e., ability to fit a wide variety of functions [40,41], is sufficiently high. The optimal predictions can solely be expressed in terms of the probability distributions underlying the physical system. Using the optimal predictions  $\hat{y}^{\text{opt}}$  in place of the NN predictions  $\hat{y}$ , we further obtain analytical expressions for the optimal indicators of phase transitions  $I^{\text{opt}}$  (step 2\*). Evaluating these analytical expressions provides an alternative path for computing indicators of phase transitions without *ever* training NNs; see Table I, where we compare the computation times of the two approaches.

the success of deep NNs (DNNs) in image recognition tasks [46]. The more expressive a ML model [40,47,48], such as an NN, the more resources are needed to train it, and the more difficult it is to interpret the underlying functional dependence of its prediction on the input [49,50]. Therefore, NNs typically act as black boxes that can correctly highlight phase transitions but whose internal workings remain opaque to the user. Since the proposal of supervised learning, the learning-by-confusion scheme, and the prediction-based method, there have been numerous attempts to understand their working principle, particularly through the extraction of order parameters. As an example, (kernel) support vector machines, which are easier to analyze than NNs due to their inherent linear nature, were used as predictive models [51–54]. Other approaches to improve interpretability rely on systematic input engineering, such that the objective function that the NN learns is approximately linearly [55], or on a systematic reduction of the NN expressivity [15]. Another set of works [56–59] analyze trained NNs using standard interpretability tools from ML, which rely on truncated Taylor expansions. Despite these efforts, we still understand little about the working principle

of ML methods for the detection of phase transitions based on NNs, when they fail or succeed, and how they differ [2]—in particular, when DNNs are used (i.e., in the limit of high model expressivity). These open questions reflect the general scarcity of rigorous theory in ML [35].

Here, we address these gaps in knowledge by pursuing a novel approach based on deriving analytical expressions for the optimal predictions of the NNs underlying supervised learning, learning by confusion, and the prediction-based method. The predictions are optimal in the sense that they minimize the target loss function; i.e., the corresponding model performs the desired task (as specified by the loss function) optimally. Based on the optimal predictions, we find analytical expressions for the optimal indicators of phase transitions of these three methods. The optimal indicators correspond to the output of the methods when using ideal high-capacity [40] predictive models, such as well-trained, highly expressive NNs. The inner workings of these methods are revealed through the dependence of the optimal indicators on the input data. Moreover, the analytical expressions make it possible to compute the optimal indicator directly from the input data without training NNs

(see step 2\* in Fig. 1), manifesting an alternative numerical routine to infer phase transitions. We demonstrate the procedure in a numerical study on a variety of models exhibiting, e.g., symmetry-breaking, topological, quantum, and many-body localization phase transitions.

This work is structured as follows: In Sec. II, we introduce the task of detecting phase transitions from data in an automated fashion, including supervised learning, the learning-by-confusion scheme, and the prediction-based method. Section III discusses the analytical expressions of their optimal indicators of phase transitions. A numerical study of the optimal predictions and indicators for the Ising model, Ising gauge theory,  $XY$  model,  $XXZ$  model, Kitaev model, and Bose-Hubbard model is presented in Sec. IV. Finally, the results are discussed in Sec. V, and conclusions are drawn in Sec. VI.

## II. AUTOMATED DETECTION OF PHASE TRANSITIONS FROM DATA

In this section, we formally introduce the task of automatically detecting phase transitions from data and how supervised learning (SL), learning by confusion (LBC), and the prediction-based method (PBM) approach this problem. We consider the following scenario: The physical system to be analyzed is characterized by a tuning parameter  $p$  sampled equidistantly with a grid spacing  $\Delta p$ . In the following, we denote the points at the boundary of the sampled region as  $p_1$  and  $p_K$  with  $K \in \mathbb{N}$  sampled points in total ( $K = (p_K - p_1)/\Delta p + 1$ ). At each sampled point  $p_k$  ( $1 \leq k \leq K$ ), we draw  $M \in \mathbb{N}$  samples from the system's state  $\{\mathcal{S}_{jk}\}_{j=1}^M$  which constitute our available data. We allow for these data to be preprocessed via a mapping to a representation space  $\mathcal{R}: \mathcal{S} \rightarrow \mathbf{x}$ . At the core of each of the three methods for detecting phase transitions under consideration lies a predictive model  $m: \mathbf{x} \rightarrow \hat{y}$ , such as an NN, which takes the preprocessed data  $\mathcal{X} = \{\mathbf{x}_{jk} | 1 \leq j \leq M, 1 \leq k \leq K\}$  as input. We denote the available data at sampled point  $p_k$  as  $\mathcal{X}_k = \{\mathbf{x}_{jk}\}_{j=1}^M$ . Note that  $\mathcal{X}$  may contain duplicates. Let  $\tilde{\mathcal{X}}$  be the set of unique inputs obtained from  $\mathcal{X}$  by removing all duplicates. We assume that the system is present either in a single phase  $A$  or two distinct phases  $A$  and  $B$  across the sampled range of the tuning parameter  $\{p_k\}_{k=1}^K$ . If a system exhibits multiple distinct phases, the parameter range can (in principle) be analyzed in a piecewise fashion (for more details on this case, see Appendixes A 1 and A 2). The task is then to compute a scalar indicator  $I(p)$ , which peaks at the phase boundary if two distinct phases are present, i.e., has a local maximum, and does not exhibit a peak otherwise. More specifically, if the system is in phase  $A$  from  $p_1$  to  $p_c$  and phase  $B$  from  $p_c$  to  $p_K$  with critical point  $p_c$  (not necessarily a sampled point), the indicator  $I(p_k)$  should exhibit a local maximum at the sampled point closest to the critical point  $\operatorname{argmin}_{p_k} |p_c - p_k|$ .

### A. Supervised learning

In SL, a predictive model  $m$  is trained on the data available in regions near the two boundaries of the chosen parameter range denoted by I and II. Regions I and II are comprised of the set of sampled points  $\{p_k | 1 \leq k \leq r_I\}$  and  $\{p_k | l_{II} \leq k \leq K\}$ , respectively. Here,  $r_I, l_{II} \in \mathbb{N}$  denote the rightmost and leftmost parameter point in region I and II, respectively. In SL, we assume that there exist two distinct phases  $A$  and  $B$ , with the regions I and II being located deep within these phases. Without loss of generality, we assign the label  $y = 1$  and  $y = 0$  to data obtained in region I and II, respectively. The predictive model is trained to minimize a cross-entropy (CE) loss

$$\mathcal{L}_{\text{SL}} = -\frac{1}{M_{\mathcal{T}}} \sum_{\mathbf{x} \in \mathcal{T}} \{y(\mathbf{x}) \ln [\hat{y}(\mathbf{x})] + [1 - y(\mathbf{x})] \ln [1 - \hat{y}(\mathbf{x})]\}, \quad (1)$$

where the sum runs over all  $M_{\mathcal{T}}$  data points in the training set  $\mathcal{T} \subseteq \mathcal{X}$ ,  $\mathcal{T} = \{\mathbf{x}_{jk} | 1 \leq j \leq M, k \in \{1, \dots, r_I\} \cup \{l_{II}, \dots, K\}\}$ . Let us denote the set containing all unique inputs present in  $\mathcal{T}$  without repetition as  $\tilde{\mathcal{T}}$ . The output of the predictive model  $\hat{y}(\mathbf{x}) \in [0, 1]$  corresponds to the probability of input  $\mathbf{x}$  having the label  $y = 1$ , whereas  $1 - \hat{y}(\mathbf{x})$  is the probability that the input  $\mathbf{x}$  carries the label  $y = 0$ .

After training the predictive model to minimize the loss function in Eq. (1), it is evaluated on all available data  $\mathcal{X}$ . Averaging over the predictions  $\hat{y}(\mathbf{x})$  for all data  $\mathcal{X}_k$  at a given point  $p_k$  ( $1 \leq k \leq K$ ) yields a prediction as a function of the tuning parameter:

$$\hat{y}_{\text{SL}}(p_k) = \frac{1}{M} \sum_{\mathbf{x} \in \mathcal{X}_k} \hat{y}(\mathbf{x}). \quad (2)$$

The indicator for phase transitions in SL,  $I_{\text{SL}}$ , is then given by the negative derivative of the prediction with respect to the tuning parameter:

$$I_{\text{SL}}(p_k) = -\left. \frac{\partial \hat{y}_{\text{SL}}(p)}{\partial p} \right|_{p_k}. \quad (3)$$

The estimated critical value of the tuning parameter in SL corresponds to the location of the global maximum in its indicator [Eq. (3)], which can easily be determined in an automated fashion without human supervision. If one chooses to label data obtained in region I with  $y = 0$  and region II with  $y = 1$  instead, the same indicator signal can be recovered via a sign change  $I_{\text{SL}}(p_k) \rightarrow -I_{\text{SL}}(p_k)$ . Note that it is also common to identify the estimated critical value of the tuning parameter in SL as  $\operatorname{argmin}_{p_k} |\hat{y}(p_k) - 0.5|$ ; see Appendix D 1 for a comparison motivating our choice.

Intuitively, if there is a transition from one phase to another (phase  $A$  to phase  $B$ ) when varying the tuning parameter  $p$ , the mean predictions  $\hat{y}_{\text{SL}}(p)$  should drop from

$\hat{y}_{\text{SL}}(p_1) = 1$  (deep within phase *A*) to  $\hat{y}_{\text{SL}}(p_K) = 0$  (deep within phase *B*) as  $p$  is increased. If the transition is sharp, the predictions should also change abruptly. Such a change results in a peak in the negative derivative of the predictions, i.e., in the indicator for phase transitions. In that case, the predictive model acts as an order parameter that approaches 1 (0) deep within phase *A* (*B*). In general, one expects the predictions—and, thus, the indicator—to vary most strongly at the critical point  $p_c$ . If there is only a single phase, one expects the predictions to be approximately constant, resulting in a flat indicator  $I_{\text{SL}}(p)$ . Our derivation of the optimal indicator  $I_{\text{SL}}^{\text{opt}}$  provides a rigorous basis for these heuristic arguments underlying the SL method.

### B. Learning by confusion

In LBC, predictive models are trained on all available data  $\mathcal{X}$ . The labels are obtained by performing a split of the sampled parameter range into two neighboring regions labeled I and II. Each input  $\mathbf{x}$  drawn in region I or II carries the label  $y = 1$  or  $y = 0$ , respectively. The values of the tuning parameters which realize each of the  $K + 1$  possible bipartitions are given as  $p_k^{\text{bp}} = p_1 - \Delta p/2 + (k - 1)\Delta p$ , where  $1 \leq k \leq K + 1$ . For a given bipartition point  $p_k^{\text{bp}}$ , regions I and II are then comprised of the sampled points  $\{p_j | p_j \leq p_k^{\text{bp}}, 1 \leq j \leq K\}$  and  $\{p_j | p_j > p_k^{\text{bp}}, 1 \leq j \leq K\}$ , respectively. Note that for bipartitions 1 ( $p_1^{\text{bp}} = p_1 - \Delta p/2$ ) and  $K + 1$  ( $p_{K+1}^{\text{bp}} = p_K + \Delta p/2$ ), region I or II encompasses the entire sampled parameter range, and all data are assigned the label 1 or 0, respectively.

To each bipartition, i.e., choice of data labeling, we associate a distinct predictive model  $m_k$  ( $1 \leq k \leq K + 1$ ) which is trained to minimize a CE loss:

$$\begin{aligned}
 \mathcal{L}_{\text{LBC}} = & -\frac{1}{M_{\mathcal{X}}} \sum_{\mathbf{x} \in \mathcal{X}} \{y(\mathbf{x}) \ln [\hat{y}(\mathbf{x})] \\
 & + [1 - y(\mathbf{x})] \ln [1 - \hat{y}(\mathbf{x})]\}, \quad (4)
 \end{aligned}$$

where the sum runs over all  $M_{\mathcal{X}} = KM$  data points. Again, the output of the predictive model  $\hat{y}(\mathbf{x}) \in [0, 1]$  corresponds to the probability of input  $\mathbf{x}$  having the label  $y = 1$ , whereas  $1 - \hat{y}(\mathbf{x})$  is the probability of the input  $\mathbf{x}$  carrying the label  $y = 0$ .

Once a predictive model has been trained to minimize the loss function in Eq. (4) for a given bipartition, it is evaluated on all available data points. In particular, we can compute the mean classification accuracy as a function of the bipartition parameter  $p_k^{\text{bp}}$  ( $1 \leq k \leq K + 1$ ) as

$$I_{\text{LBC}}(p_k^{\text{bp}}) = 1 - \frac{1}{M_{\mathcal{X}}} \sum_{\mathbf{x} \in \mathcal{X}} |\theta[\hat{y}(\mathbf{x}) - 0.5] - y(\mathbf{x})|, \quad (5)$$

where  $\theta$  denotes the Heaviside step function. The predictions  $\hat{y}(\mathbf{x})$  are obtained from the predictive model  $m_k$

associated with the bipartition point  $p_k^{\text{bp}}$ , and  $y(\mathbf{x})$  are the corresponding labels.

Clearly, the mean classification accuracy  $I_{\text{LBC}}$  exhibits trivial local maxima at the points  $p_1^{\text{bp}} = p_1 - \Delta p/2$  and  $p_{K+1}^{\text{bp}} = p_K + \Delta p/2$ , where the entire data are assigned the label 0 or 1, respectively. Therefore, a predictive model effortlessly reaches a perfect accuracy of 1, because it simply needs to predict a single label regardless of the input. However, given that the underlying data can be separated into two distinct classes of similar character (i.e., phases) through appropriate bipartitioning of the parameter range at  $p_c$ , one also expects the classification accuracy to have a local maximum at  $p_c$ . At such a point, the predictive model is “least confused” by the choice of data labeling. Hence, the mean classification accuracy serves as the indicator for phase transitions within LBC. The estimated critical value of the tuning parameter in LBC corresponds to the location of the largest local maximum (excluding the points  $p_1^{\text{bp}}$  and  $p_{K+1}^{\text{bp}}$  at the boundary) in its indicator [Eq. (5)].

### C. Prediction-based method

In PBM, a predictive model  $m$  is trained on all available data  $\mathcal{X}$  to infer the value of the tuning parameter  $p_k$  ( $1 \leq k \leq K$ ) at which an input  $\mathbf{x}$  was generated. While SL and LBC constitute supervised *classification* tasks, PBM corresponds to a supervised *regression* task, where the label is given by the tuning parameter itself  $y(\mathbf{x}) = p_k \forall \mathbf{x} \in \mathcal{X}_k$ .

We train the predictive model  $m$  to minimize a mean-square-error (MSE) loss function

$$\mathcal{L}_{\text{PBM}} = \frac{1}{M_{\mathcal{X}}} \sum_{\mathbf{x} \in \mathcal{X}} [\hat{y}(\mathbf{x}) - y(\mathbf{x})]^2. \quad (6)$$

After training, the predictive model is evaluated on all available data points  $\mathcal{X}$ . Averaging over the predictions  $\hat{y}(\mathbf{x})$  for all data  $\mathcal{X}_k$  at a given point  $p_k$  yields a mean prediction as a function of the tuning parameter:

$$\hat{y}_{\text{PBM}}(p_k) = \frac{1}{M} \sum_{\mathbf{x} \in \mathcal{X}_k} \hat{y}(\mathbf{x}). \quad (7)$$

We then compute the deviation of the prediction from the true underlying value of the tuning parameter  $\delta y_{\text{PBM}}(p_k) = \hat{y}_{\text{PBM}}(p_k) - p_k$ . The indicator for phase transitions of PBM,  $I_{\text{PBM}}$ , is then given by the derivative of this deviation with respect to the tuning parameter:

$$I_{\text{PBM}}(p_k) = \left. \frac{\partial \delta y_{\text{PBM}}(p)}{\partial p} \right|_{p_k} = \left. \frac{\partial \hat{y}_{\text{PBM}}(p)}{\partial p} \right|_{p_k} - 1. \quad (8)$$

The estimated critical value of the tuning parameter in PBM corresponds to the location of the global maximum in its indicator [Eq. (8)].

Intuitively, if there is only a single phase, in which inputs cannot be distinguished well by the predictive model, one expects the mean predictions to be approximately constant. This results in the deviations  $\delta y_{\text{PBM}}$  varying approximately linear with the tuning parameter. Hence, the indicator  $I_{\text{PBM}}$  is approximately constant. However, if there is a transition from one phase to another as the tuning parameter is varied, the predictions and the corresponding deviations also vary sharply. This results in a peak in the derivative of the deviations, i.e., the indicator for phase transitions  $I_{\text{PBM}}$ . In particular, one expects that the predictions are most susceptible at the phase boundary. Thus, its derivative should vary most strongly at the critical point  $p_c$ .

In many standard applications of NNs, it is typical to split the available data into multiple sets, in particular, to avoid overfitting [40]. For example, suppose we aim to construct an accurate on-the-fly classifier of individual samples into distinct phases of matter. In this case, it may be beneficial to split the available data into a training and validation set to avoid overfitting if only a limited amount of data is available. In the case of PBM and LBC, we do not explicitly split the dataset  $\mathcal{X}$  into a training set and test set (as well as a potential validation set). This can be done, e.g., to assess sampling convergence by comparing the predictions obtained on the training set and test set or to perform early stopping with NNs (see Appendix B 2 for concrete examples). Note, however, that the task we consider here is the detection of phase transitions given the data at hand. As such, the dataset  $\mathcal{X}$  does not *necessarily* need to be split. In particular, in the limit of a sufficient number of samples, all splits of a dataset coincide, assuming that all samples are drawn independently from the same probability distributions underlying the physical system (see Fig. 1). Therefore, the predictions and indicators obtained by training NNs using multiple distinct datasets coincide with the values obtained using the entire dataset for training and evaluation up to deviations arising from finite-sample statistics. That is, in the limit of a sufficient number of samples, the results obtained in the two scenarios coincide [40,60,61]. Moreover, given a fixed amount of data  $\mathcal{X}$ , better statistics are achieved by utilizing the entire data for training and evaluation.

### III. OPTIMAL INDICATORS OF PHASE TRANSITIONS

In this section, we discuss the optimal indicators of phase transitions  $I^{\text{opt}}$  for each of the phase-classification methods presented in Sec. II. The optimal indicators can be directly calculated given the predictions  $\hat{y}^{\text{opt}}(\mathbf{x})$  of an optimal model  $m^{\text{opt}}$  which minimizes the corresponding loss function. The detailed proofs can be found in Appendix A 1. In the limit of sufficient data, i.e., given accurate estimates of the probability distributions underlying the physical system  $\{P_k(\mathbf{x})\}_{k=1}^K$ , such a model is also *Bayes optimal* [40,62], meaning no other statistical model can outperform it (on

average) on the classification or regression task at hand. In this case, the optimal loss value it achieves coincides with the *Bayes error* [40,62], i.e., the irreducible error inherent to the problem.

*Supervised learning.*—In SL (see Sec. II A), the optimal predictions are given as

$$\hat{y}_{\text{SL}}^{\text{opt}}(\mathbf{x}) = \frac{P_{\text{I}}(\mathbf{x})}{P_{\text{I}}(\mathbf{x}) + P_{\text{II}}(\mathbf{x})} \quad \forall \mathbf{x} \in \bar{\mathcal{T}}, \quad (9)$$

where

$$P_{\text{I}}(\mathbf{x}) = \sum_{k=1}^{r_{\text{I}}} P_k(\mathbf{x}) \quad (10)$$

and

$$P_{\text{II}}(\mathbf{x}) = \sum_{k=l_{\text{II}}}^K P_k(\mathbf{x}) \quad (11)$$

are the (unnormalized) probabilities of drawing an input  $\mathbf{x}$  in region I and II, respectively. Hence, the optimal prediction for a particular input corresponds to the probability of drawing that input in region I compared to region II. Here,  $P_k(\mathbf{x})$  denotes the (normalized) probability to draw the input  $\mathbf{x}$  at the sampled point  $p_k$ . Given a dataset  $\mathcal{X}_k$ , this probability is estimated as  $P_k(\mathbf{x}) \approx M_k(\mathbf{x})/M$ , where  $M_k(\mathbf{x})$  is the number of times the input  $\mathbf{x}$  is present in the dataset  $\mathcal{X}_k$ . While having access to an analytical expression for the underlying probability distributions  $\{P_k(\mathbf{x})\}_{k=1}^K$  may ease computation and enable additional insights, it is not strictly required to compute the optimal predictions (see Sec. IV for application to physical systems). An expression for the optimal value of the loss in SL,  $\mathcal{L}_{\text{SL}}^{\text{opt}}$ , can be obtained by replacing  $\hat{y}(\mathbf{x})$  with  $\hat{y}_{\text{SL}}^{\text{opt}}(\mathbf{x})$  in Eq. (1), where, by definition,  $\mathcal{L}_{\text{SL}}^{\text{opt}} \leq \mathcal{L}_{\text{SL}}$ .

Assuming that all inputs within the entire dataset  $\mathcal{X}$  are already present in the training set  $\mathcal{T}$ , i.e.,  $\bar{\mathcal{T}} = \bar{\mathcal{X}}$ , the mean optimal prediction at a given point  $p_k$  ( $1 \leq k \leq K$ ) is

$$\hat{y}_{\text{SL}}^{\text{opt}}(p_k) = \sum_{\mathbf{x} \in \bar{\mathcal{X}}} P_k(\mathbf{x}) \hat{y}_{\text{opt}}(\mathbf{x}). \quad (12)$$

This corresponds to the probability of finding an input drawn at that point  $p_k$  in region I compared to region II. We find this assumption to be (approximately) satisfied for all physical systems analyzed in this work and can estimate the errors arising from a violation; see Sec. IV and Appendix A 2. The optimal indicator of phase transitions in SL is then given as

$$I_{\text{SL}}^{\text{opt}}(p_k) = - \left. \frac{\partial \hat{y}_{\text{SL}}^{\text{opt}}(p)}{\partial p} \right|_{p_k}. \quad (13)$$

In general, the prediction in Eq. (12) changes most at a transition point and, thus, its derivative, the optimal indicator in Eq. (13), shows a peak.

*Learning by confusion.*—For a given bipartition of the parameter range into regions I and II, the optimal predictions of LBC (see Sec. II B) are given as

$$\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}) = \frac{P_{\text{I}}(\mathbf{x})}{P_{\text{I}}(\mathbf{x}) + P_{\text{II}}(\mathbf{x})} \quad \forall \mathbf{x} \in \bar{\mathcal{X}}, \quad (14)$$

which corresponds to the probability of drawing the input in region I compared to region II. This characteristic is inherent to the underlying classification task [compare Eqs. (9) and (14)]. The mean classification error associated with an input  $\mathbf{x}$  is given by  $\min\{\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}), 1 - \hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x})\}$ . This classification error arises from a “confusion” of the model: Different labels can be assigned to the same input due to an overlap of the underlying probability distributions. The mean classification error over the entire parameter range given a particular choice of bipartition, i.e., labeling of the data, then corresponds to

$$I_{\text{LBC}}^{\text{opt}} = 1 - \frac{1}{K} \sum_{k=1}^K \sum_{\mathbf{x} \in \bar{\mathcal{X}}} P_k(\mathbf{x}) \min\{\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}), 1 - \hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x})\}. \quad (15)$$

This forms the optimal indicator for phase transitions in LBC. An expression for the optimal value of the loss in LBC,  $\mathcal{L}_{\text{LBC}}^{\text{opt}}$ , can be obtained by replacing  $\hat{y}(\mathbf{x})$  with  $\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x})$  in Eq. (4). The critical point  $p_c$  is highlighted by a dip in the mean classification error, i.e., by a peak in the mean classification accuracy [Eq. (15)]. It corresponds to the bipartition point for which the probability distributions underlying the two regions have the least overlap (on average), resulting in the highest classification accuracy and the least confusion. While confusion can arise due to suboptimal predictions of models with restricted capacity (see Appendix B 2 for a concrete example), we find that confusion can persist even in the limit of high model capacity if it is inherent to the underlying data. Based on the analytical expressions, we thus gain an intuitive and rigorous understanding of the concept of confusion underlying LBC [5].

*Prediction-based method.*—The optimal predictions within PBM (see Sec. II C) are given as

$$\hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x}) = \frac{\sum_{k=1}^K P_k(\mathbf{x}) p_k}{\sum_{k=1}^K P_k(\mathbf{x})} \quad \forall \mathbf{x} \in \bar{\mathcal{X}}. \quad (16)$$

Here, the optimal prediction for a given input is obtained by a weighted sum over each sampled point in the parameter range, where the weight of a point corresponds to the probability of obtaining the input at that point compared to all other points along the parameter range. Therefore, the

prediction accuracy decreases if the same input can be drawn at multiple values of the tuning parameter, i.e., when the underlying probability distributions overlap. An expression for the optimal value of the loss in PBM,  $\mathcal{L}_{\text{PBM}}^{\text{opt}}$ , can be obtained by replacing  $\hat{y}(\mathbf{x})$  by  $\hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x})$  in Eq. (6). The mean prediction of an optimal model  $m^{\text{opt}}$  at a sampled point  $p_k$  is given by

$$\hat{y}_{\text{PBM}}^{\text{opt}}(p_k) = \sum_{\mathbf{x} \in \bar{\mathcal{X}}} P_k(\mathbf{x}) \hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x}). \quad (17)$$

Thus, the optimal indicator for phase transitions is

$$I_{\text{PBM}}^{\text{opt}}(p_k) = \left. \frac{\partial \delta y_{\text{PBM}}^{\text{opt}}(p)}{\partial p} \right|_{p_k}, \quad (18)$$

where  $\delta y_{\text{PBM}}^{\text{opt}}(p_k) = \hat{y}_{\text{PBM}}^{\text{opt}}(p_k) - p_k$ . Recall that, in PBM, phase transitions are detected by analyzing the dependence of the prediction error on the tuning parameter. The optimal indicator [Eq. (18)] highlights the value of the tuning parameter at which the mean predictions change most, i.e., where the overlap of the underlying probability distributions changes most. The optimal predictions and indicators of PBM have previously been derived in Ref. [34] but have neither been utilized in a numerical routine nor been used to explain previous studies.

The optimal predictions of SL, LBC, and PBM can *solely* be expressed in terms of the probability distributions  $\{P_k(\mathbf{x})\}_{k=1}^K$  governing the input data. Crucially, this means that the optimal predictions—and, thus, the optimal indicators of phase transitions—do not depend on the particular nature of an input or how similar it is to other inputs. Such notions of similarity form the basis of a large set of other phase-classification methods, e.g., based on principal component analysis [13], diffusion maps [27], or anomaly detection [32]. The analytical form of the optimal predictions indicates that SL, LBC, and PBM ultimately gauge changes in the probability distributions governing the data akin to probability metrics [63]. Note that the same optimal predictions and indicators are obtained for multiple choices of representations  $\mathcal{R}$  given that the same probability distributions can still describe the data in the representation space. Consequently, knowledge of the symmetries of the system can be utilized to calculate indicators of phase transitions more efficiently. We make use of this in Sec. IV.

### A. Demonstration on prototypical probability distributions

In this section, we compute the optimal indicators of SL, LBC, and PBM for a set of simple probability distributions governing the input data. The probability distributions governing the data in physical systems can be regarded as generalizations of the special cases discussed in this section. Thus, they serve as a reasonable basis for

understanding. We compare these results to the indicators obtained by numerical optimization of NNs. The details on the NN architecture and training, including the corresponding hyperparameters, can be found in Appendix B. This first demonstration shows how the analytical expressions can be used to calculate the optimal indicator directly from input data without NNs. Moreover, it confirms that the optimal predictive models can be recovered by training NNs with sufficient expressive power.

*Case 1.*—Let us first consider the case where the probability distribution governing the data is identical across the parameter range, i.e.,  $P_k(\mathbf{x}) = P(\mathbf{x}) \forall 1 \leq k \leq K$ . Clearly, in this case, all three methods should indicate the presence of a single phase. The optimal prediction in SL is

$$\hat{y}_{\text{SL}}^{\text{opt}}(p) = \frac{K_{\text{I}}}{K_{\text{I}} + K_{\text{II}}} = \text{const}, \quad (19)$$

corresponding to the relative size of region I compared to region II [see Fig. 2(b)]. Here,  $K_{\text{I}} = r_{\text{I}}$  and  $K_{\text{II}} = K - l_{\text{II}}$  correspond to the number of sampled parameter values in region I or II, respectively. Taking the derivative of Eq. (19) results in a flat indicator signal  $I_{\text{SL}}^{\text{opt}} = 0$ . In LBC, the optimal classification accuracy for a particular bipartition is given by  $I_{\text{LBC}}^{\text{opt}} = \max\{K_{\text{I}}/K, K_{\text{II}}/K\}$ . This results in a characteristic V shape [5], which has its minimum at the center of the parameter range under consideration; see Fig. 2(c). In PBM,

the optimal mean prediction is also placed at the center of mass  $\hat{y}_{\text{PBM}}^{\text{opt}}(p_k) = 1/K \sum_{k=1}^K p_k = \text{const}$ , which results in a constant indicator  $I_{\text{PBM}}^{\text{opt}} = -1$  [see Fig. 2(d)]. As such, all three methods yield optimal indicators that correctly signal the presence of a single phase, i.e., the absence of two distinct phases. For a concrete numerical demonstration, we consider the case of binary inputs  $\mathcal{X} = \{0, 1\}$  with equal probability  $P(0) = P(1) = 0.5$ . Figures 2(a)–2(e) show the results for all three methods using the analytical expressions as well as NNs. The analytical predictions and indicators can be approximated well using NNs as predictive models.

*Case 2.*—Next, we consider the case where the input data naturally separate into two distinct sets. That is, the underlying probability distributions result in a bipartition of the parameter range into two regions *A* and *B*, where each input can be drawn in only one of the two regions. In these regions, we choose the probability distributions to be identical:

$$P_k(\mathbf{x}) = \begin{cases} P_A(\mathbf{x}) & \forall k \leq c, \\ P_B(\mathbf{x}) & \forall k > c, \end{cases} \quad (20)$$

where  $1 \leq k, c \leq K$ . This is a prototypical example for the case where the physical system transitions from phase *A* to *B* when crossing a critical value of the tuning parameter  $p_c$ .

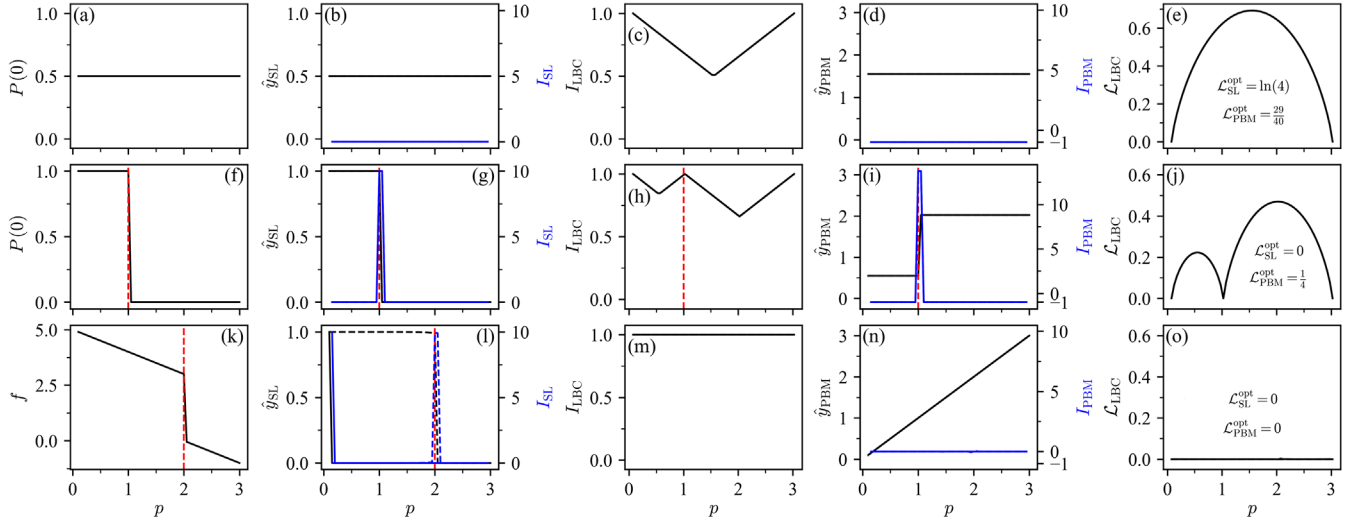


FIG. 2. Results for prototypical probability distributions in (a)–(e) case 1 with  $P_k(\mathbf{x}) = P(\mathbf{x}) \forall k$ , where  $P(0) = P(1) = 0.5$ , (f)–(j) case 2 given by Eq. (20) with  $P_A(0) = 1$ ,  $P_B(0) = 0$ , and  $p_c = 1$ , and (k)–(o) case 3 with Eqs. (26) and (27). The tuning parameter ranges from  $p_1 = 0.1$  to  $p_K = 3$  with  $\Delta p = 0.05$ . Critical values of the tuning parameter are highlighted with red dashed lines. For details on SL, LBC, and PBM using NNs, see Appendix B. (a), (f), (k) Illustration of the probability distributions underlying the data. (b), (g), (l) Mean prediction  $\hat{y}_{\text{SL}}(p)$  obtained using the analytical expression (black solid line) or an NN (black dashed line), as well as the corresponding indicator  $I_{\text{SL}}(p)$  (blue lines). Here, we choose  $r_{\text{I}} = 1$  and  $l_{\text{II}} = K$ . (c), (h), (m) The indicator of LBC,  $I_{\text{LBC}}$ , obtained using the analytical expression (black solid line) or an NN (black dashed line). (d), (i), (n) Mean prediction  $\hat{y}_{\text{PBM}}(p)$  of PBM obtained using the analytical expression (black solid line) or an NN (black dashed line), as well as the corresponding indicator  $I_{\text{PBM}}(p)$  (blue lines). (e), (j), (o) Value of the loss function in LBC,  $\mathcal{L}_{\text{LBC}}$ , for each bipartition point  $p^{\text{bp}}$  obtained using the analytical expression (black solid line) or evaluated after NN training (black dashed line). In addition, the optimal values of the loss function for SL and PBM obtained by evaluating the analytical expressions are reported. Note that, by definition,  $\mathcal{L}^{\text{opt}} \leq \mathcal{L}$  for all three methods.

Here,  $p_c$  corresponds to a sampled value of the tuning parameter, which may, in general, not be the case.

Using SL, the optimal strategy corresponds to

$$\hat{y}_{\text{SL}}^{\text{opt}}(p_k) = \begin{cases} 1 \forall k \leq c, \\ 0 \forall k > c. \end{cases} \quad (21)$$

This results in

$$I_{\text{SL}}^{\text{opt}}(p_k) = \begin{cases} 0 \forall k < c, \\ \frac{1}{2\Delta p} \forall k \in \{c, c+1\}, \\ 0 \forall k > c+1, \end{cases}$$

which diverges as  $\Delta p \rightarrow 0$  and exhibits a peak at the two points which constitute the boundary between regions  $A$  and  $B$  [see Fig. 2(g)]. Here, we approximate the derivative in Eq. (13) by a symmetric difference quotient:

$$I_{\text{SL}}^{\text{opt}}(p_k) \approx -\frac{\hat{y}_{\text{SL}}^{\text{opt}}(p_{k+1}) - \hat{y}_{\text{SL}}^{\text{opt}}(p_{k-1})}{2\Delta p}, \quad (22)$$

where  $2 \leq k \leq K-1$ .

In LBC, one can reach a perfect (error-free) classification when matching the natural bipartition present in the data. Let us denote the region between the bipartition point underlying the data,  $p_c$ , and the chosen bipartition point in the LBC scheme,  $p_k^{\text{bp}}$ , as III. The number of sampled parameter values within the smallest region between I, II, and III is  $K_k^m = \min\{K_I, K_{\text{II}}, K_{\text{III}}\}$ . Note that all input data drawn within one of these regions must be misclassified. Thus, the optimal strategy which yields the smallest classification error corresponds to misclassifying all input data drawn within the smallest region. The optimal classification accuracy is then given as

$$I_{\text{LBC}}^{\text{opt}}(p_k^{\text{bp}}) = 1 - \frac{K_k^m}{K}. \quad (23)$$

This results in a characteristic W shape of the indicator [5] [see Fig. 2(h)], where the middle peak occurs at the bipartition point  $p_k^{\text{bp}}$  closest to  $p_c$ .

In PBM, we have

$$\hat{y}_{\text{PBM}}^{\text{opt}}(p_k) = \begin{cases} \langle p \rangle_A = 1/c \sum_{j=1}^c p_j \forall k \leq c, \\ \langle p \rangle_B = 1/(K-c) \sum_{j=c+1}^K p_j \forall k > c, \end{cases} \quad (24)$$

where  $\langle p \rangle_{A/B}$  denotes the center of region  $A$  and  $B$ , respectively. This results in

$$I_{\text{PBM}}^{\text{opt}}(p_k) = \begin{cases} -1 \forall k < c, \\ \frac{\langle p \rangle_B - \langle p \rangle_A}{2\Delta p} \forall k \in \{c, c+1\}, \\ -1 \forall k > c+1, \end{cases} \quad (25)$$

where we approximate the derivative in Eq. (18) by a symmetric difference quotient [see Fig. 2(i)]. The expression in Eq. (25) diverges as  $\Delta p \rightarrow 0$  for  $k \in \{c, c+1\}$  and results in a peak at the two points which constitute the boundary between regions  $A$  and  $B$ . As such, the optimal indicators of all three methods correctly indicate the presence of two distinct sets of data, i.e., two distinct phases. The results obtained using the analytical expressions can be approximated well using NNs as predictive models. This is illustrated in Figs. 2(f)–2(j), where we consider the special case of binary inputs with  $P_A(0) = 1$ ,  $P_B(0) = 0$ , and  $p_c = 1$ .

*Case 3.*—Lastly, we consider the case where the probability distributions underlying the data do not overlap; i.e., the probability of drawing a given input at two distinct values of the tuning parameter vanishes. This situation can, for example, occur when dealing with large state spaces, which are prone to result in insufficient sampling statistics in practice. That is, even in scenarios where the ground-truth probability distributions underlying the data *do* overlap, the estimated probabilities  $P_k(\mathbf{x}) \approx M_k(\mathbf{x})/M$  based on the drawn dataset  $\mathcal{X}$  may not (see Appendix A 5 for a concrete physical example). Many image classification tasks encountered in traditional ML applications [64–70] *a priori* fall into this category. In particular, the probability distributions underlying the data are typically not known in these cases. Therefore, constructing optimal models, in particular, Bayes optimal models, largely remains conceptual in nature [62,71].

Here, an optimal predictive model is capable of distinguishing between samples obtained at distinct values of the tuning parameter with perfect accuracy. This results in  $I_{\text{LBC}}^{\text{opt}}(p_k^{\text{bp}}) = 1 \forall 1 \leq k \leq K+1$  for LBC [see Fig. 2(m)]. In the case of PBM, we have  $\hat{y}_{\text{PBM}}^{\text{opt}}(p_k) = p_k$  such that  $I_{\text{PBM}}^{\text{opt}}(p_k) = -1 \forall 1 \leq k \leq K$  [see Fig. 2(n)]. In both cases, the indicator signals the absence of two distinct sets of data, i.e., phases. The optimal predictions of SL for  $\mathbf{x} \in \bar{\mathcal{X}}$  are underdetermined: Only the predictions for inputs within the training data  $\mathbf{x} \in \bar{\mathcal{T}}$  are fixed after training, and the assumption that  $\bar{\mathcal{X}} = \bar{\mathcal{T}}$  is violated in this particular case [see Fig. 2(l)]. Note, however, that the predictions are, in principle, also unconstrained when using SL with NNs. For a simple numerical example, we consider the case where a single unique (scalar) input is drawn at each point along the parameter range:

$$P_k(x) = \begin{cases} 1 & \text{for } x = f(p_k), \\ 0 & \text{otherwise,} \end{cases} \quad (26)$$

with

$$f(p) = \begin{cases} 5-p \forall p \leq 2, \\ 2-p \forall p > 2, \end{cases} \quad (27)$$



where  $1 \leq k \leq K$ . The results are shown in Figs. 2(k)–2(o). In practice, NNs tend to predict similar outputs for similar inputs. The continuous nature of the NN results in SL highlighting the value of the tuning parameter  $p = 2$  where a discontinuity in the input data is present. We also observe this tendency for the NNs in LBC and PBM during training.

### B. Computational cost

We can use the analytical expressions to assess the computational cost associated with the evaluation of the mean optimal predictions and optimal indicators of SL, LBC, and PBM for a given set of input data (see Appendix A 3 for proofs). In our estimation, we neglect the overhead arising from the computation of the probability distributions  $\{P_k(\mathbf{x})\}_{k=1}^K$  which is identical for all three methods. In the case of SL, the computation of its optimal predictions (as a function of the tuning parameter) and indicator scales as  $O(M_{\bar{\chi}}K)$ . Here, we assume that the number of sampled values of the tuning parameter during training is small compared to the total number of sampled points  $K_I + K_{II} \ll K$ . For PBM and LBC, the computation scales as  $O(M_{\bar{\chi}}K^2)$  and  $O(M_{\bar{\chi}}K^3)$ , respectively. By saving the optimal predictions for each input  $\hat{y}_{\text{opt}}(\mathbf{x})$  instead of recomputing it, the computational cost can be reduced and scales as  $O(M_{\bar{\chi}}K)$ ,  $O(M_{\bar{\chi}}K)$ , and  $O(M_{\bar{\chi}}K^2)$ , in the case of SL, PBM, and LBC, respectively. Note the appearance of  $M_{\bar{\chi}}$ , which can result in an exponential scaling for quantum problems due to the exponential growth of the Hilbert space  $\mathcal{H}$  (and, thus, the state space  $M_{\bar{\chi}}$ ).

## IV. APPLICATION TO PHYSICAL SYSTEMS

In this section, we compute the optimal predictions and indicators of phase transitions of SL, LBC, and PBM, directly from data using the analytical expressions introduced in Sec. III for the Ising model, Ising gauge theory, XY model, XXZ model, Kitaev model, and Bose-Hubbard model. For the classical systems, namely, the Ising model, Ising gauge theory, and XY model, spin configurations are sampled from a thermal distribution at various temperatures  $T_k$  using the Metropolis-Hastings algorithm [72]. Here, the temperature serves as a tuning parameter. The probability that a system in equilibrium at inverse temperature  $\beta_k = 1/k_B T_k$  (where  $k_B$  is the Boltzmann constant) is found in a state with spin configuration  $\sigma$  is given by a Boltzmann distribution:

$$P_k(\sigma) = \frac{e^{-\beta_k H(\sigma)}}{Z_k}, \quad (28)$$

where  $Z_k = \sum_{\sigma} e^{-\beta_k H(\sigma)}$  is the partition function and  $H$  is the respective system Hamiltonian. In principle, one could use the raw spin configurations as input, i.e., estimate the underlying probability distributions as  $P_k(\sigma) = M_k(\sigma)/M$ . However, the probability of drawing a particular spin

configuration depends only on its energy [see Eq. (28)]. One can show that the optimal predictions and indicators remain identical when the energy is used as input instead of the raw configurations, i.e., when the probability distributions governing the data are given by

$$P_k(E) = \frac{g(E)e^{-\beta_k E}}{Z_k}, \quad (29)$$

where  $g(E)$  is the degeneracy factor (see Appendix A 4 for a proof). Using the energy as input instead of the raw configurations reduces both the input dimension and the size of the associated state space. This, in turn, reduces the cost of computing the optimal predictions and indicators. In general, one can take advantage of the symmetries of the system by adopting a symmetry-adapted representation.

In the quantum case, we are typically looking at a state associated with a Hamiltonian  $H(p)$  that depends on the tuning parameter  $p$ . This state could, for example, be the ground state or a state which has undergone unitary time evolution starting from a fixed initial state. Having chosen a complete orthonormal basis  $\{|j\rangle\}_{j=1}^d$  to study the system [ $d = \dim(\mathcal{H})$ ], the relevant quantum state at  $p_k$  can be written as  $|\Psi_k\rangle = \sum_{j=1}^d c_{jk} |j\rangle$ . Thus, the probability distribution  $P_k$  associated with a given value  $p_k$  of the tuning parameter is  $P_k(j) = |c_{jk}|^2$  with  $1 \leq j \leq d$  and  $1 \leq k \leq K$ . The value of  $P_k(j)$  corresponds to the probability of measuring the system in state  $|j\rangle$  given that the value of the tuning parameter is  $p_k$ . This corresponds to using the indices of the basis states  $|j\rangle$  ( $1 \leq j \leq d$ ) as inputs, which are governed by the probability distributions  $\{P_k(j)\}_{k=1}^K$ . For simplicity, we choose  $M_{\bar{\chi}} = d$ . In the case of spin systems, we use the  $S^z$  basis, whereas we choose the Fock basis for bosonic and fermionic systems. This choice of bases corresponds to experimentally accessible local measurements [73–79]. In this work, we obtain the ground states through exact diagonalization. Thus, we have direct access to the underlying probability distributions and do not rely on sampling. In Appendix A 5, we show that the optimal indicators can also be obtained from individual samples, i.e., measurement outcomes (similar to the classical case). As such, the procedure is *in principle* applicable to experimental scenarios.

In general, we can consider scenarios where a state  $|\Psi_i\rangle$  is drawn with probability  $a_k(i)$  at  $p_k$ . Then, the relevant quantum state is given by a classical probabilistic mixture  $\rho_k = \sum_i a_k(i) |\Psi_i\rangle \langle \Psi_i|$ ,  $i \in \mathbb{N}$ . The probability distribution associated with such a state is  $P_k(j) = \sum_i a_k(i) |c_{ij}|^2$ . This case is particularly relevant for the study of many-body localization phase transitions where disorder is naturally present (see Sec. IV F). Here, the tuning parameter  $p_k$  itself characterizes a distribution  $a_k$ . Further details on the data generation can be found in Appendix C.

Clearly, in the quantum case, there is an ambiguity in the choice of input or, equivalently, the choice of measurement basis. Changing the measurement basis may change the probability distributions underlying the data and, thus, the corresponding optimal predictors and indicators. In turn, the estimated critical value of the tuning parameter may change (which is difficult to assess *a priori* for a given system). In order to avoid an explicit choice of measurement basis, sampling over various classical projections can be performed. Classical representations of quantum states obtained via classical shadow tomography [35,80,81] are an example of this. Alternatively, measurements given by informationally complete positive operator-valued measures (IC-POVMs) can be used [82,83]. However, projective measurements in a single basis have been the most common choice, reflecting experimental constraints or prior knowledge of the system [10,11,31,37,84,85].

### A. Ising model

The two-dimensional square-lattice ferromagnetic Ising model is described by the following Hamiltonian:

$$H(\boldsymbol{\sigma}) = -J \sum_{\langle ij \rangle} \sigma_i \sigma_j, \quad (30)$$

where the sum runs over all nearest-neighboring sites (with periodic boundary conditions) and  $J$  is the interaction strength ( $J > 0$ ). At each lattice site  $k$ , there is a discrete spin variable  $\sigma_i \in \{+1, -1\}$ . This results in a state space of

size  $2^{L \times L}$  for a square lattice of linear size  $L$ . The system is completely characterized by its spin configuration  $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_{L \times L})$ . Two example spin configurations of the Ising model at different temperatures are shown in Fig. 3(a). The Ising model exhibits a symmetry-breaking phase transition at a critical temperature of [86]

$$T_c = \frac{2J}{k_B \ln(1 + \sqrt{2})}. \quad (31)$$

The system undergoes a transition between a paramagnetic (disordered) phase at high temperature and a ferromagnetic (ordered) phase at low temperature. Spontaneous magnetization occurs below the critical temperature  $T_c$ , where the interaction is sufficiently strong to cause neighboring spins to align spontaneously. This spontaneous symmetry breaking leads to a nonzero mean magnetization. Above  $T_c$ , thermal fluctuations dominate over spin alignment, resulting in a vanishing magnetization. Consequently, the phase transition can be characterized by the magnetization  $M(\boldsymbol{\sigma}) = \sum_{i=1}^{L^2} \sigma_i$ , where  $M/L^2$  serves as an order parameter that is zero within the paramagnetic phase and approaches one in the ferromagnetic phase; see Fig. 3(h). The phase transition can also be revealed by the heat capacity

$$C(T) = \frac{d\langle E \rangle_T}{dT} = \frac{\langle E^2 \rangle_T - \langle E \rangle_T^2}{k_B T^2}, \quad (32)$$

which diverges at  $T_c$  [see Fig. 3(g)].

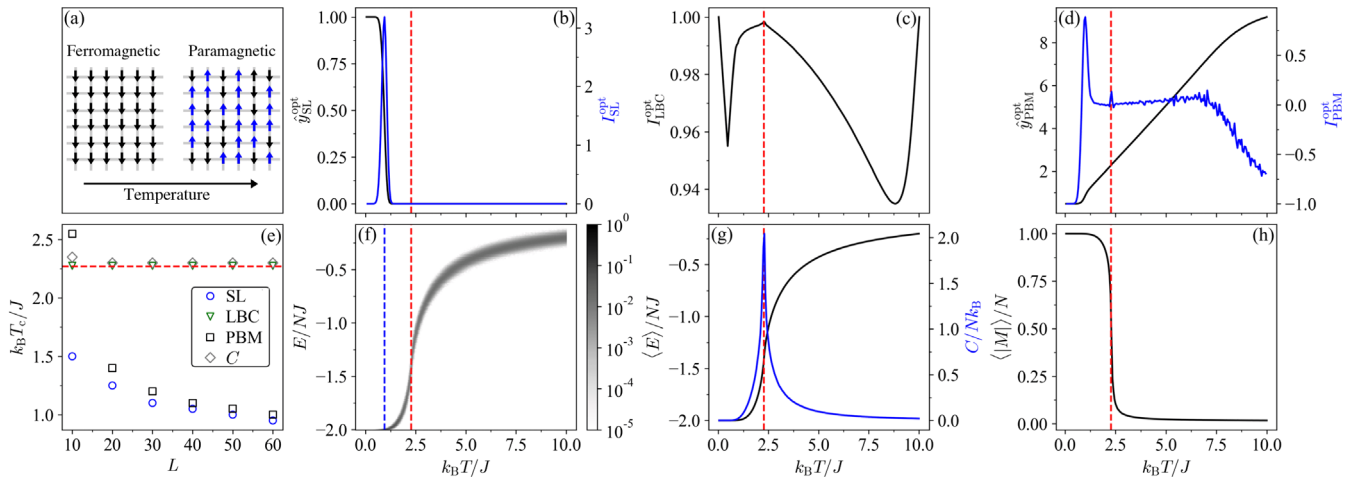


FIG. 3. Results for the Ising model ( $L = 60$ ) with the dimensionless temperature as a tuning parameter  $p = k_B T/J$ , where  $p_1 = 0.05$ ,  $p_K = 10$ , and  $\Delta p = 0.05$ . In SL, the data obtained at  $p_1$  and  $p_K$  constitute our training set, i.e.,  $r_1 = 1$  and  $l_{II} = K$ . The critical temperature [Eq. (31)] is highlighted by a red dashed line. (a) Illustration of the symmetry-breaking phase transition in the Ising model. (b) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL (black line) and the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$  (blue line). (c) Optimal indicator of LBC,  $I_{\text{LBC}}^{\text{opt}}$  (black line). (d) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM (black line) and the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$  (blue line). (e) Estimated critical temperatures based on  $I_{\text{SL}}^{\text{opt}}$  (SL),  $I_{\text{LBC}}^{\text{opt}}$  (LBC),  $I_{\text{PBM}}^{\text{opt}}$  (PBM), and heat capacity ( $C$ ) as a function of the lattice size  $L$ . The estimated critical temperature based on the heat capacity corresponds to the location of its maximum. (f) Probability distributions governing the input data (here, the energy) as a function of the tuning parameter, where the color scale denotes the probability. The blue dashed line highlights the predicted critical temperature of SL and PBM. (g) Average energy per site (black line) and associated heat capacity (blue line) as a function of the temperature, where  $N = L^2$ . (h) Average magnetization per site as a function of the temperature.

The results for the Ising model are shown in Fig. 3. Interestingly, SL fails to predict the correct critical temperature even for large lattices [see Figs. 3(b) and 3(e)]. In fact, we can further analyze the special case when the inputs are governed by Boltzmann distributions [Eq. (29)]: For training data obtained at  $T_1 = 0$  (region I) and  $T_K > 0$  (region II), the mean optimal prediction of SL at an intermediate temperature  $T_k$  is

$$\hat{y}_{\text{SL}}^{\text{opt}}(T_k) = \frac{P_k(E_{\text{gs}})}{1 + P_K(E_{\text{gs}})} \propto P_k(E_{\text{gs}}), \quad (33)$$

which approaches  $P_k(E_{\text{gs}})$  in the thermodynamic limit as  $T_K \rightarrow \infty$  (see Appendix A 4 for a proof). Here,  $E_{\text{gs}}$  denotes the ground-state energy. Therefore, in this case, the optimal indicator in SL peaks at the temperature at which the probability of drawing the ground state changes most [see the blue dashed line in Fig. 3(f)]. The location of the peak tends to zero as one approaches the thermodynamic limit; see Fig. 3(e).

The optimal indicator of PBM shows two distinct peaks. One coincides with the peak of the optimal indicator in SL, whereas the other coincides with the critical temperature of the Ising model [see Fig. 3(d)]. This observation suggests a deeper connection between SL and PBM. A similar indicator signal (with two distinct peaks) is observed in Ref. [29] with NNs after a sufficient number of training epochs. In principle, the finite-size scaling analysis allows one to identify the dominant peak as erroneous without prior knowledge of  $T_c$ , because it shifts toward  $T = 0$  as the lattice size is increased, whereas the small peak remains stable. In the same fashion, the output of SL can be identified to be erroneous. Note that the fluctuations present in the optimal indicator signal of PBM can be attributed to finite-sample statistics. A detailed study of the effect of finite-sample statistics on the optimal predictions and indicators can be found in Appendix A 5. Crucially, the analytical expression for the optimal indicator signal allows us to disentangle the stochasticity inherent to the NN training from other sources of noise, which was not rigorously possible in previous works.

In Ref. [4], SL with NNs is shown to predict the critical temperature of the Ising model for various lattice sizes correctly. In this case, small NNs with restricted expressive power in combination with  $\ell_2$  regularization are used. Similarly, using PBM in Ref. [29] a single, distinct peak at  $T_c$  is observed after a small number of training epochs with a second peak emerging after longer training. Training time, NN size, and explicit  $\ell_2$  regularization are all factors which influence the effective capacity of the resulting model and, thus, determine its ability to approximate the optimal predictive model [40,41], i.e., to realize the global minimum of the loss function corresponding to the optimal predictions and indicators. We recover the same behavior using NNs as in Refs. [4,29] by restricting the model

capacity, e.g., by choosing a small NN, stopping the training early, or using strong  $\ell_2$  regularization (see Appendix B 1 for details). As these restrictions are lifted, i.e., by choosing a larger NN, training for longer, or reducing the regularization strength, the NN-based predictions and indicators approach the corresponding optimal predictions and indicators displayed in Fig. 3. Thus, our analysis demonstrates that SL and PBM necessarily rely on models with restricted capacity and hyperparameter tuning to correctly predict the critical temperature of the Ising model.

Finally, the optimal indicator of LBC correctly highlights the critical temperature of the Ising model for various lattice sizes in accordance with Ref. [5]; see Figs. 3(c) and 3(e). Overall, the optimal indicators of all three methods show peaks at temperatures where the probability distribution underlying the data varies strongly. Recall the finding from Sec. III that all three methods gauge changes in the probability distributions underlying the data. Note that the results shown in Fig. 3 are stable against small perturbations of the chosen parameter range, including regions I and II in SL.

## B. Ising gauge theory

Wegner's Ising gauge theory (IGT) [88] is described by the following Hamiltonian:

$$H(\boldsymbol{\sigma}) = -J \sum_{\mathbb{P}} \prod_{i \in \mathbb{P}} \sigma_i, \quad (34)$$

where  $\mathbb{P}$  refers to plaquettes on the lattice; see Fig. 4(a). The IGT is a prototypical example of a classical system that exhibits a topological phase of matter [89]. It is a spin model ( $\sigma_i \in \{+1, -1\}$ ) defined on a square lattice of linear size  $L$  (with periodic boundary conditions) where the spins are placed on the lattice bonds [see Fig. 4(a)]. The IGT ground state is a degenerate manifold made up of all states which fulfill the condition that the product of spins on each plaquette is  $\prod_{i \in \mathbb{P}} \sigma_i = 1$  corresponding to a topological phase. These topological constraints can be violated at finite temperature, where the system leaves its ground state. Note that there is no phase transition at finite temperature: The critical temperature approaches zero in the thermodynamic limit. In finite-sized systems, however, the violations of local constraints are suppressed. Therefore, the system exhibits a crossover from the topological phase at low temperature to a phase with violated topological constraints at high temperature. The crossover temperature  $T_c$  is defined by the first appearance of a violated local constraint and scales as  $T_c \propto 1/\ln(2L^2)$  [87]. Figure 4(a), which shows typical spin configurations of the IGT, highlights that the phases of the IGT are hard to distinguish visually without prior knowledge of the local constraints or a dual representation [4,31]. Note that the heat capacity fails to identify the crossover; see Fig. 4(f). The topological

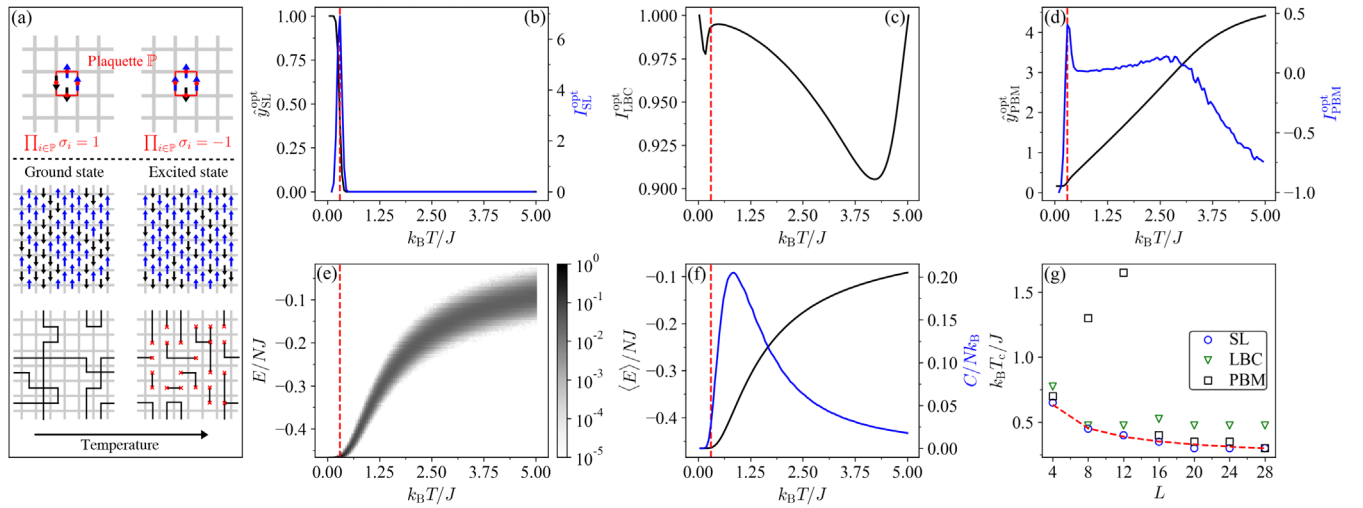


FIG. 4. Results for the IGT ( $L = 28$ ) with the dimensionless temperature as a tuning parameter  $p = k_B T/J$ , where  $p_I = 0.05$ ,  $p_K = 5$ , and  $\Delta p = 0.05$ . In SL, the data obtained at  $p_I$  and  $p_K$  constitute our training set, i.e.,  $r_I = 1$  and  $I_{\text{II}} = K$ . The crossover temperature is highlighted by a red dashed line and scales as  $k_B T_c/J \propto 1/\ln(2L^2)$  [87]. (a) Upper panels show examples of plaquettes  $\mathbb{P}$  where the topological constraint is met ( $\prod_{i \in \mathbb{P}} \sigma_i = 1$ ) and violated ( $\prod_{i \in \mathbb{P}} \sigma_i = -1$ ). Middle panels show examples of spin configurations within the topological ground-state phase (left) and phase with violated topological constraints at high temperature (right). Lower panels show the corresponding Wilson loops. (b) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL (black line) and the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$  (blue line). (c) Optimal indicator of LBC,  $I_{\text{LBC}}^{\text{opt}}$  (black line). (d) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM (black line) and the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$  (blue line). (e) Probability distributions governing the input data (here, the energy) as a function of the tuning parameter, where the color scale depicts the probability. (f) Average energy per site (black line) and associated heat capacity (blue line) as a function of the temperature, where  $N = 2L^2$ . Note that the heat capacity does not peak at the crossover temperature. (g) Estimated critical temperature based on  $I_{\text{SL}}^{\text{opt}}$  (SL),  $I_{\text{LBC}}^{\text{opt}}$  (LBC), and  $I_{\text{PBM}}^{\text{opt}}$  (PBM) as a function of the lattice size  $L$ .

character of the ground-state phase can be revealed through Wilson loops. These are formed by connecting edges with spins of the same orientation; see Fig. 4(a). In the ground-state phase, all such loops are closed. The violation of a plaquette constraint breaks a loop.

Recall that SL, LBC, and PBM are *a priori* sensitive to both phase transitions and crossovers. The results for the crossover in the IGT are shown in Fig. 4. The optimal indicator of SL [Fig. 4(b)] shows an appropriate scaling behavior. Moreover, the corresponding estimated critical temperature highlights the first appearance of violated local constraints; see Figs. 4(e) and 4(f). This can be confirmed explicitly, as SL can be shown to measure changes in the probability of drawing the ground state (cf. Sec. IV A). Observe that the underlying probability distribution undergoes a large change at the crossover temperature; see Fig. 4(e). In Refs. [4,31], SL and PBM are shown to correctly highlight the crossover temperature of the IGT using NNs. In fact, the optimal model underlying PBM for the IGT coincides with the physically motivated density-of-states-based model proposed in Ref. [31]; see Appendix D 2 for details. We find that the optimal indicator of PBM correctly marks the crossover temperature of the IGT except at small lattice sizes. As for the Ising model, the optimal indicator of PBM exhibits two peaks in this case. The peak located at the crossover temperature dominates for large lattice sizes. Note that for the IGT it is not

beneficial to reduce the model capacity when using PBM or SL, which leads to an erroneous peak closely matching the specific heat [see Fig. 4(f)], given that the corresponding optimal indicators correctly highlight the crossover temperature.

The optimal indicator of LBC correctly highlights the crossover temperature via its local maximum at small lattice sizes but shows slight deviations from the appropriate scaling behavior for large lattices. Reference [31] reports difficulties in identifying the crossover temperature using LBC due to a distorted W shape of its indicator. Choosing the same range for the tuning parameter, we can qualitatively reproduce their results using our analytical expression for the optimal indicator of LBC; see Appendix D 2. Using NNs, it is difficult to make concrete statements on whether a method succeeds or fails at identifying a given phase transition due to the inherent stochasticity arising during NN training and the choice of hyperparameters, such as the NN size. Our theoretical analysis allows for rigorous statements to be made about the optimal outcome when applying ML methods for detecting phase transitions to a given system (i.e., dataset). In this particular example, the analytical expressions allow us to determine that, when training highly expressive NNs for sufficiently long, the indicator signal of LBC is indeed ambiguous (as reported in Ref. [31]). Note that restricting the model capacity is not found to resolve this issue [31].

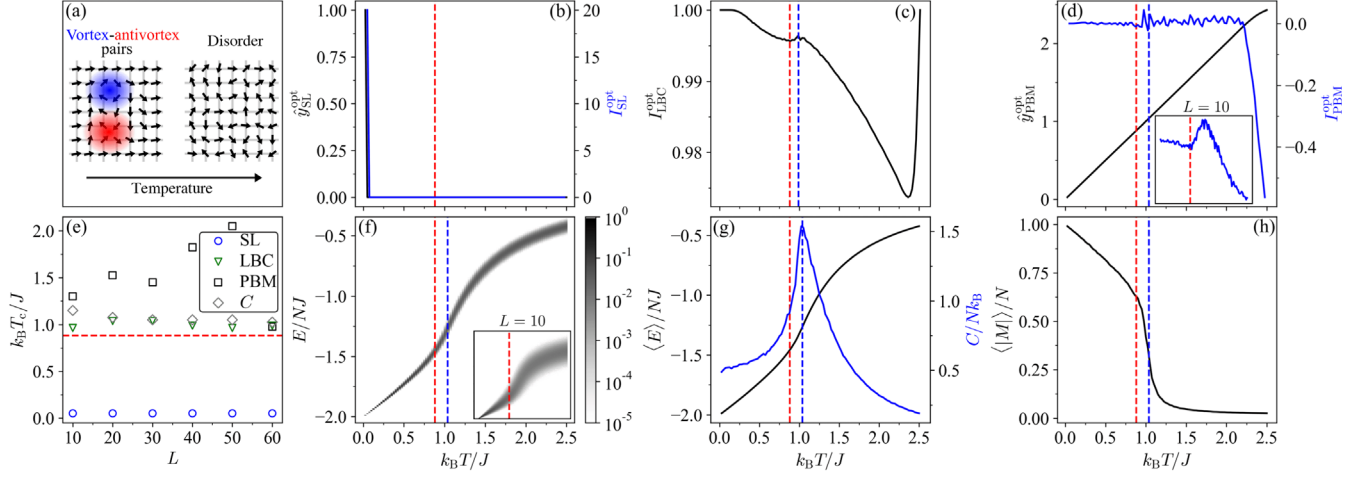


FIG. 5. Results for the XY model ( $L = 60$ ) with the dimensionless temperature as a tuning parameter  $p = k_B T / J$ , where  $p_I = 0.025$ ,  $p_K = 2.5$ , and  $\Delta p = 0.025$ . In SL, the data obtained at  $p_I$  and  $p_K$  constitute our training set, i.e.,  $r_I = 1$  and  $l_{II} = K$ . The BKT transition temperature  $k_B T_c / J \approx 0.8935$  [92] is highlighted by a red dashed line. The blue dashed line highlights the estimated critical temperature using LBC. (a) Illustration of the BKT phase transition in the XY model. (b) Mean optimal prediction  $\hat{y}_{SL}^{opt}$  in SL (black line) and the corresponding indicator  $I_{SL}^{opt}$  (blue line). (c) Optimal indicator of LBC,  $I_{LBC}^{opt}$  (black line). The blue dashed line highlights the predicted critical temperature of LBC. (d) Mean optimal prediction  $\hat{y}_{PBM}^{opt}$  in PBM (black line) and the corresponding indicator  $I_{PBM}^{opt}$  (blue line). The inset shows the optimal indicator signal of PBM for  $L = 10$ , which exhibits a peak near the location of the maximum in the heat capacity. (e) Estimated critical temperature based on  $I_{SL}^{opt}$  (SL),  $I_{LBC}^{opt}$  (LBC),  $I_{PBM}^{opt}$  (PBM), and heat capacity ( $C$ ) as a function of the lattice size  $L$ . The estimated critical temperature of the heat capacity corresponds to the location of its maximum. (f) Probability distributions governing the input data (here, the energy) as a function of the tuning parameter, where the color scale denotes the probability. The inset shows the probability distributions for  $L = 10$ . (g) Average energy per site (black line) and associated heat capacity (blue line) as a function of the temperature, where  $N = L^2$ . (h) Average magnetization per site as a function of the temperature.

### C. XY model

Next, we consider the two-dimensional classical XY model that exhibits a Berezinskii-Kosterlitz-Thouless (BKT) transition driven by the emergence of topological defects [90,91]. The model is described by the following Hamiltonian:

$$H = -J \sum_{\langle ij \rangle} \cos(\theta_i - \theta_j), \quad (35)$$

where  $\langle ij \rangle$  denotes the sum over nearest neighbors (with periodic boundary conditions) of a square lattice of linear size  $L$ . The angle  $\theta_i \in [0, 2\pi)$  corresponds to the orientation of the spin at site  $i$ . The formation of topological defects (i.e., vortices and antivortices) results in a quasi-long-range-ordered phase. The transition between the quasi-long-range-ordered phase at low temperature and a disordered phase at high temperature is a BKT transition, and the associated critical temperature is  $k_B T_c / J \approx 0.8935$  [92]. Below  $T_c$ , vortex-antivortex pairs form due to thermal fluctuations, but they remain bound to minimize their total free energy [see Fig. 5(a)]. At  $T_c$ , the entropic contribution to the free energy equals the binding energy of a pair which triggers vortex unbinding. These unbinding events drive the BKT phase transition. Note that the heat capacity has a peak at  $T > T_c$  which is associated with the entropy released when

most vortex pairs unbind [93,94]; see Fig. 5(g). Moreover, while the XY model has strictly zero magnetization for all  $T > 0$  in the thermodynamic limit, a nonzero value is found for systems of finite size [95]; see Fig. 5(h). Instead, the critical temperature can, for example, be estimated based on the helicity modulus [94,96] (see Appendix C).

The results for the XY model are shown in Fig. 5. Here, SL fails to predict the critical temperature correctly. This failure is linked to the fact that the optimal indicator of SL highlights changes in the probability to obtain the ground state (cf. Sec. IV A), which quickly vanishes with increasing temperature; see Fig. 5(f). In a similar spirit, in Ref. [23] it is found that “naive” SL (without engineering the features or NN architecture) fails to yield accurate estimates of the critical temperature. Here, we explicitly confirm that a classification based on detecting vortices does not correspond to the most optimal strategy. The peak in the optimal indicator of LBC matches the peak in the heat capacity at  $k_B T / J \approx 1$  [see Figs. 5(c) and 5(e)] and, thus, overestimates the critical temperature of the XY model. In Ref. [23], indicator signals of similar shape are obtained using LBC with NNs for the XY model. The rapid decrease in the optimal indicator of LBC for  $k_B T / J \gtrsim 1$  can be attributed to the increase in the overlap of the underlying probability distributions [Fig. 5(f)], which results in a higher classification error. Note that the overlap of the probability

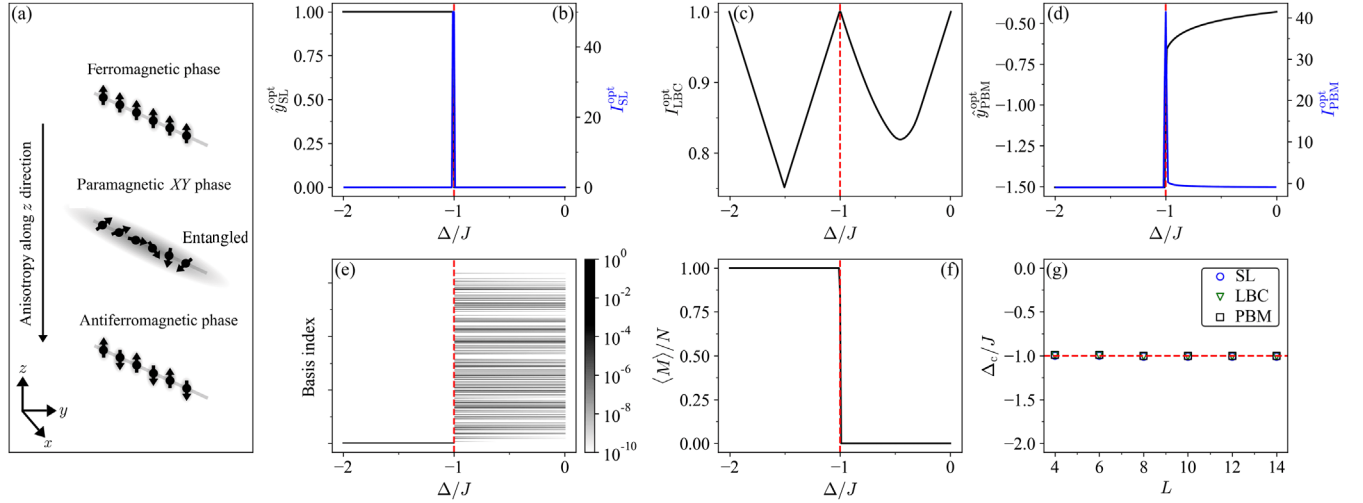


FIG. 6. Results for the XXZ chain ( $L = 14$ ) with the dimensionless anisotropy strength along the  $z$  direction as the tuning parameter  $p = \Delta/J$ , where  $p_1 = -2$ ,  $p_K = 0$ , and  $\Delta p = 0.01$ . In SL, the data obtained at  $p_1$  and  $p_K$  constitute our training set, i.e.,  $r_I = 1$  and  $l_{II} = K$ . The critical value of the tuning parameter  $\Delta/J = -1$  at which the phase transition between the ferromagnetic phase and paramagnetic XY phase occurs is highlighted by a red dashed line. (a) Illustration of the quantum phase transitions of the XXZ chain. (b) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL (black line) and the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$  (blue line). (c) Optimal indicator of LBC,  $I_{\text{LBC}}^{\text{opt}}$  (black line). (d) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM (black line) and the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$  (blue line). (e) Probability distributions governing the input data (indices of  $S^z$  basis states) as a function of the tuning parameter, where the color scale denotes the probability. The color scale is cut off at  $10^{-10}$  to improve visual clarity. (f) Average magnetization per site ( $N = L$ ). (g) Estimated critical value of the tuning parameter based on  $I_{\text{SL}}^{\text{opt}}$  (SL),  $I_{\text{LBC}}^{\text{opt}}$  (LBC), and  $I_{\text{PBM}}^{\text{opt}}$  (PBM) as a function of the chain length  $L$ .

distributions decreases with increasing lattice size; see Fig. 5(f). Hence, the indicator of PBM [Fig. 5(d)] shows a clear peak close to the location of the peak in the heat capacity for small lattice sizes. For systems of increasing size, the optimal predictions of PBM start to closely match the underlying tuning parameter, resulting in an increasingly linear behavior [see the black line in Fig. 5(d)]. This corresponds to an optimal indicator signal close to zero, where the variations in the predicted critical value of the tuning parameter [Fig. 5(e)] are due to small local fluctuations.

Overall, the behavior of the optimal indicators of all three methods closely resembles our previous example regarding perfectly distinguishable input data (see case 3 in Sec. III A). This can be traced back to the small overlap of the underlying probability distributions; see Fig. 5(f). The increase in the overlap with increasing temperature results in a decrease in the mean classification accuracy of LBC, i.e., its indicator [see Fig. 5(c)]. Evidently, in such a case, NNs with restricted expressive power and other phase-classification methods based on the similarity of input data [27] may provide more valuable insights. In particular, we find that the indicators peak close to the transition temperature, i.e., near the location of the peak in the heat capacity and drop in the magnetization, when restricting the model capacity, e.g., by stopping the NN training early (see Appendix B). Recall that this is also observed in the case of the Ising model (see Sec. IV A and Appendix B).

#### D. XXZ model

Having discussed classical models, we move on to the quantum case. First, we consider the spin-1/2 XXZ chain [97,98] with open boundary conditions whose Hamiltonian is given by

$$H = \sum_{i=1}^{L-1} J(S_{i+1}^x S_i^x + S_{i+1}^y S_i^y) + \Delta S_{i+1}^z S_i^z, \quad (36)$$

where  $J$  is the coupling strength along the  $x$  and  $y$  directions and  $\Delta$  is the coupling strength in the  $z$  direction. For  $\Delta/J < 1$ , the XXZ chain is in the ferromagnetic phase; see Fig. 6(a). The ground state is spanned by the two product states where all spins point in either the  $z$  or  $-z$  direction which have a magnetization of  $\langle M \rangle = 2\langle S_{\text{tot}}^z \rangle = \pm L$ . The ferromagnetic phase exhibits a broken symmetry: These states do not exhibit the discrete symmetry of spin reflection  $S_i^z \rightarrow -S_i^z$  under which the Hamiltonian is invariant. For  $\Delta/J > 1$ , the XXZ chain is in the antiferromagnetic phase with broken symmetry and two degenerate ground states. These are product states with vanishing magnetization. For  $-1 < \Delta/J < 1$ , the XXZ chain is in the paramagnetic XY phase characterized by uniaxial symmetry of the easy-plane type and vanishing magnetization.

Here, we restrict our analysis to the transition between the ferromagnetic phase and the paramagnetic XY phase. The ground states are obtained through exact

diagonalization. Figure 6 shows the results when the ground state with  $\langle S_{\text{tot}}^z \rangle = +L/2$  is selected in the ferromagnetic phase and  $S^z$  is chosen as a measurement basis. The quantum phase transition can be revealed by looking at the magnetization; see Fig. 6(f). The optimal indicators of all three methods correctly highlight the phase transition. Looking at the underlying probability distributions [see Fig. 6(e)], the problem closely resembles the prototypical case of a bipartitioned dataset (see case 2 in Sec. III A). Thus, the optimal predictions and indicators also qualitatively match the results obtained in this case. In particular, the optimal predictions of SL can be described by Eq. (33), where the ferromagnetic ground state takes the role of the ground-state energy (see Appendix A 4 for proof). We have verified that the optimal indicators also mark the phase transition when other states from the ground-state manifold are selected in the ferromagnetic phase and when measurements are performed in the  $S^x$  or  $S^y$  basis.

### E. Kitaev model

The Kitaev chain is a one-dimensional model based on  $L$  spinless fermions, which undergoes a quantum phase transition between a topologically trivial and nontrivial phase [99,100]. The Kitaev Hamiltonian is given by

$$H = \sum_{i=1}^{L-1} (\Delta c_{i+1} c_i - t c_{i+1}^\dagger c_i + \text{H.c.}) - \mu \sum_{i=1}^L n_i, \quad (37)$$

where we consider open boundary conditions,  $\mu$  is the chemical potential,  $t$  is the hopping amplitude, and  $\Delta$  is the induced superconducting gap. In the following, we set  $\Delta = -t$ . The ground state of this model features a quantum phase transition from a topologically trivial ( $|\mu/t| > 2$ ) to a nontrivial state ( $|\mu/t| < 2$ ); see Fig. 7(a). In the topological phase, Majorana zero modes [101] are present. Here, we restrict ourselves to  $\mu/t \leq 0$ . We compute the ground states through exact diagonalization. For results based on individual measurement outcomes (of projective measurements in the Fock basis), see Appendix A 5.

The topologically trivial and nontrivial phase can be distinguished through entanglement spectra and the corresponding entanglement entropy [102]. Consider the reduced density matrix  $\rho_A$  of a system in the pure state  $|\Psi\rangle$  obtained by subdividing the Hilbert space  $\mathcal{H}$  into two parts,  $A$  and  $B$ , and tracing out the degrees of freedom of  $B$ :

$$\rho_A = \text{Tr}_B |\Psi\rangle\langle\Psi|, \quad (38)$$

with  $\{\lambda_i\}$  the spectrum of  $\rho_A$  and  $\{-\ln(\lambda_i)\}$  the entanglement spectrum. Here, we consider the bipartition of the

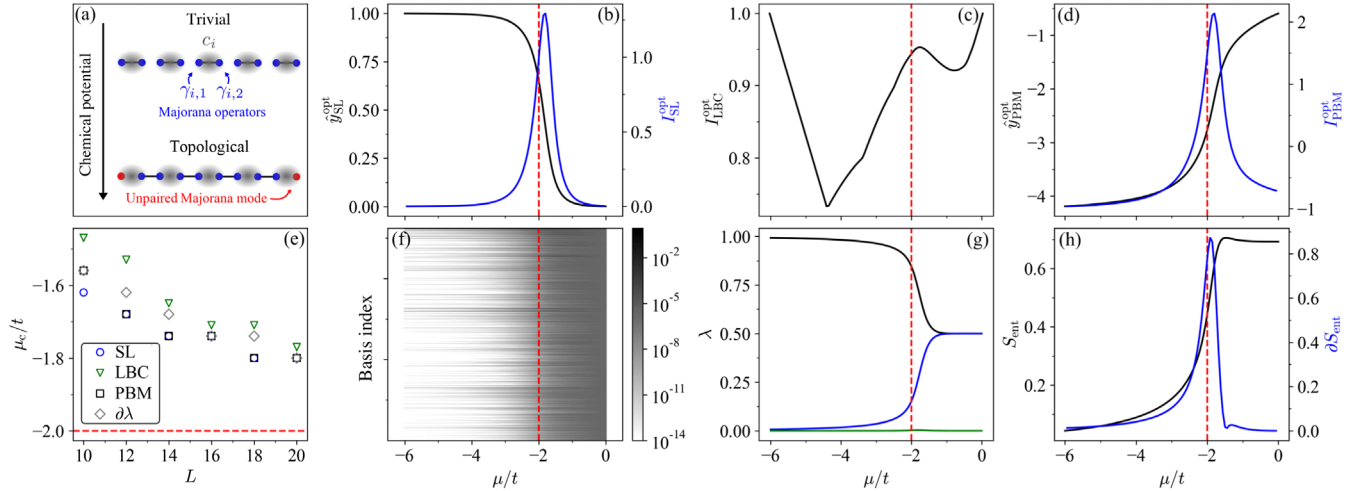


FIG. 7. Results for the Kitaev chain ( $L = 20$ ) with the dimensionless chemical potential as a tuning parameter  $p = \mu/t$ , where  $p_1 = -6$ ,  $p_K = 0$ , and  $\Delta p = 0.06$ . In SL, the data obtained at  $p_1$  and  $p_K$  constitute our training set, i.e.,  $r_1 = 1$  and  $l_{\text{II}} = K$ . The critical value  $\mu_c/t = -2$  is highlighted by a red dashed line. (a) Illustration of the phase transition in the Kitaev chain between a topological and trivial phase, where the Majorana operators  $\gamma_{i,1}$  and  $\gamma_{i,2}$  are defined by  $c_i = (\gamma_{i,1} + i\gamma_{i,2})/\sqrt{2}$ ,  $c_i^\dagger = (\gamma_{i,1} - i\gamma_{i,2})/\sqrt{2}$ . (b) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL (black line) and the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$  (blue line). (c) Optimal indicator of LBC,  $I_{\text{LBC}}^{\text{opt}}$  (black line). (d) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM (black line) and the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$  (blue line). (e) Estimated critical value of the tuning parameter based on  $I_{\text{SL}}^{\text{opt}}$  (SL),  $I_{\text{LBC}}^{\text{opt}}$  (LBC),  $I_{\text{PBM}}^{\text{opt}}$  (PBM), and the derivative of the largest eigenvalue of the reduced density matrix [see the black line in (g)] given by  $\partial\lambda/\partial p$  ( $\partial\lambda$ ), as a function of the chain length  $L$ . The estimated critical value of the tuning parameter denoted by  $\partial\lambda$  corresponds to the location of the minimum in  $\partial\lambda/\partial p$ . (f) Probability distributions governing the input data (indices of Fock basis states) as a function of the tuning parameter, where the color scale denotes the probability. The color scale is cut off at  $10^{-14}$  to improve visual clarity. (g) The three largest eigenvalues of  $\rho_A$  [Eq. (38)] as a function of the tuning parameter. (h) Entanglement entropy  $S_{\text{ent}}$  [Eq. (39)] (black line) and its derivative with respect to the tuning parameter  $\partial S_{\text{ent}}/\partial p$  (blue line).

chain into left and right halves with  $L_A = L_B = L/2$ . The entanglement entropy can then be computed as

$$S_{\text{ent}}(\rho_A) = -\sum_i \lambda_i \ln(\lambda_i). \quad (39)$$

The three largest eigenvalues of  $\rho_A$  are shown in Fig. 7(g), and the resulting entanglement entropy is shown in Fig. 7(h). Both the spectrum and entanglement entropy exhibit the largest change close to the critical value  $\mu_c/t = -2$ . The entanglement entropy approaches zero deep within the topologically trivial phase, signaling that the two halves of the ground state of the chain are not entangled. In the topological phase, the entanglement entropy approaches a value of  $\ln(2)$  characteristic of an entangled ground state.

Figure 7 shows the results of SL, LBC, and PBM. The location of the local maxima of the optimal indicators based on all three methods converges to the critical value of  $\mu_c/t = -2$  with increasing chain length. Considering the probability distributions governing the input data [see Fig. 7(f)], we observe that almost all basis states become occupied with non-negligible probability as the tuning parameter  $\mu/t$  is tuned across its critical value. Note that, in Ref. [5], the phase transition in the Kitaev model is successfully revealed using LBC with NNs where the entanglement spectrum of the

ground state serves as an input. The scaling behavior of the estimated critical value of the tuning parameter based on the optimal indicators of SL, LBC, and PBM is comparable to standard physical indicators, such as the eigenvalues of the reduced density matrix or the entanglement entropy [see Fig. 7(e)]. In the limit  $\mu/t \rightarrow -\infty$ , the ground state of the Kitaev chain corresponds to the Fock state with each site being occupied. Thus, in the limit  $\mu_1/t \rightarrow -\infty$ , the optimal predictions of SL follow Eq. (33), where the aforementioned Fock state takes the role of the ground-state energy (see Appendix A 4 for proof).

## F. Bose-Hubbard model

Finally, we consider the many-body localization (MBL) phase transition in the 1D Bose-Hubbard model (with open boundary conditions) following Refs. [10,75,76]. The system is described by the Hamiltonian

$$H = -J \sum_{i=1}^{L-1} (b_{i+1}^\dagger b_i + \text{H.c.}) + \sum_{i=1}^L \frac{U}{2} n_i(n_i - 1) + W h_i n_i, \quad (40)$$

where  $J$  is the hopping strength and  $U$  is the on-site interaction strength [see the top panel in Fig. 8(a)].

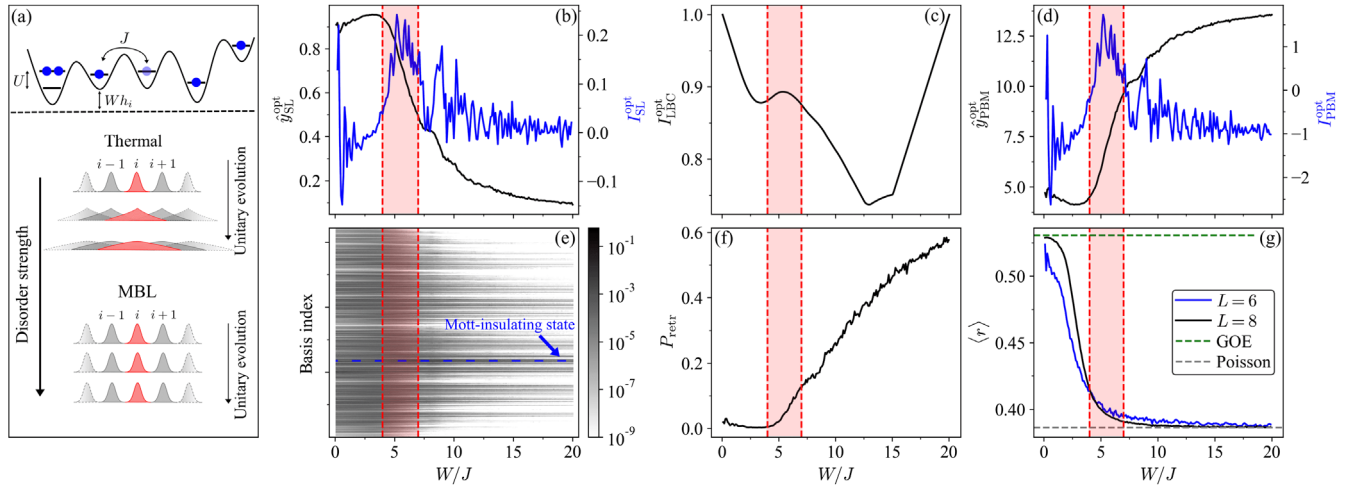


FIG. 8. Results for the MBL phase transition in the 1D Bose-Hubbard model ( $L = 8$ ) with the dimensionless disorder strength as a tuning parameter  $p = W/J$  ranging from  $p_1 = 0.1$  to  $p_K = 20$  in steps of  $\Delta p = 0.1$ . Here,  $1.1 \times 10^3$  different disorder realizations are considered. In SL, the data obtained at  $p_1$  and  $p_K$  constitute our training set, i.e.,  $r_1 = 1$  and  $l_{\text{II}} = K$ . The reference range for the critical value of the tuning parameter  $W_c/J \approx 4-7$  [10,76] at which the phase transition between the thermalizing and MBL phase occurs is highlighted in red. (a) Illustration of the 1D Bose-Hubbard model [Eq. (40)] (top) and the MBL phase transition (bottom), where the system is initialized in a Mott-insulating state. (b) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL (black line) and the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$  (blue line). (c) Optimal indicator of LBC  $I_{\text{LBC}}^{\text{opt}}$  (black line). (d) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM (black line) and the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$  (blue line). (e) Probability distributions governing the input data (indices of Fock basis states with  $N_b = 8$  particles) as a function of the tuning parameter, where the color scale denotes the probability. The color scale is cut off at  $10^{-9}$  to improve visual clarity. The blue dashed line highlights the initial Mott-insulating state. (f) Disorder-averaged retrieval probability  $P_{\text{retr}}$  as a function of the tuning parameter corresponding to the line cut marked in (e). (g) Average ratio of consecutive level spacings  $\langle r \rangle$  for a chain of length  $L = 6$  (blue line) and  $L = 8$  (black line) with reference values  $r_{\text{GOE}} = 0.5307$  (green dashed line) and  $r_{\text{Poisson}} = 2 \ln(2) - 1 \approx 0.3863$  (gray dashed line). We consider all eigenstates located in the middle one-third of the spectrum [76,103] restricted to subspace with  $N_b = L$  particles and additionally average over multiple disorder realizations ( $1 \times 10^4$  for  $L = 6$  and  $1.1 \times 10^3$  for  $L = 8$ ).



Here, we fix  $U/J = 2.9$ . The last term in Eq. (40) corresponds to a quasiperiodic potential  $h_i = \cos(2\pi\alpha i + \phi)$  mimicking on-site disorder with amplitude  $W$ , where we fix  $1/\alpha = 1.618$ . This system transitions to the MBL phase, where thermalization breaks down as the disorder strength is increased beyond a critical value  $W_c/J$ ; see the bottom panel in Fig. 8(a). We analyze the system in the long-time limit  $tJ = 100$  after unitary time evolution starting from a Mott-insulating state with one particle per site by solving the Schrödinger equation numerically. We average over different disorder realizations obtained by sampling the phase  $\phi \in [0, 2\pi)$  of the potential uniformly.

A popular way to differentiate between the thermalizing and MBL regimes relies on the study of spectral statistics using tools from random matrix theory [103–105]. In the thermal regime, the statistical distribution of level spacings is given by a Gaussian orthogonal ensemble (GOE), while a Poisson distribution is expected for localized states. The ratio of consecutive level spacings is

$$r_i = \frac{\min(\delta_i, \delta_{i+1})}{\max(\delta_i, \delta_{i+1})}, \quad (41)$$

with  $\delta_i = E_i - E_{i-1}$  at a given eigenenergy  $E_i$ . Averaging over the spectrum and multiple disorder realizations yields  $\langle r \rangle$ , which varies from  $r_{\text{GOE}} = 0.5307$  within the thermalizing phase to  $r_{\text{Poisson}} = 2 \ln(2) - 1 \approx 0.3863$  within the MBL phase; see Fig. 8(g).

The results are shown in Fig. 8. All three methods correctly identify the MBL phase boundary, where we take  $W_c/J \approx 4\text{--}7$  from Refs. [10,76] as a reference. This is in agreement with the spectral analysis: The crossover between the average ratio of consecutive level spacings for systems of size  $L = 6$  and  $L = 8$  is located at  $W_c/J \approx 4$ ; see Fig. 8(g). Moreover, the phase boundary marks the range of the tuning parameter in which the most significant change in the underlying probability distribution occurs [see Fig. 8(e)]. A line cut along the index corresponding to the initial Mott-insulating state is shown in Fig. 8(f). It corresponds to the disorder-averaged probability of retrieving the initial state after unitary time evolution. The MBL phase boundary is marked by the sudden increase in  $P_{\text{retr}}$  [75] which is correctly picked up by SL, LBC, and PBM.

Our results are also in agreement with Ref. [10], which examines the MBL phase transition within the same model using SL, PBM, and LBC with NNs on numerical and experimental data. As such, this example highlights the possibility of calculating optimal indicators directly from experimental data. Note that, in Ref. [10], the authors attempt to construct a simplified indicator for phase transitions when using LBC by subtracting the V-shaped indicator signal in the case of indistinguishable data (see case 1 in Sec. III A) as a baseline. However, we find that this procedure biases the peak of the optimal indicator signal of LBC toward the center of the parameter range

under consideration and is, thus, not a viable procedure; see Appendix D 3.

## V. DISCUSSION

In the previous section, we have demonstrated that the optimal indicators of SL, LBC, and PBM successfully detect phase transitions and crossovers in a variety of different classical and quantum systems based on numerical data. Recall that the optimal analytical predictors correspond to an optimal model that reaches the global minimum of the loss function. *A priori*, it is unclear if the optimal predictors can be recovered in practice when training NNs, because the employed NNs are of finite size and local optimization techniques are used. In Appendix B, we demonstrate that the optimal predictions and indicators of all six systems studied in Sec. IV can be recovered by training NNs. This reachability further underpins the practical relevance of our analysis for the case when using SL, LBC, and PBM with NNs.

In a traditional NN-based approach, one searches for the optimal model by iteratively updating the parameters of an NN in order to minimize a loss function (see step 2 in Fig. 1). In contrast, our numerical routine based on the derived analytical expressions allows for the optimal model to be constructed directly from data (see step 2\* in Fig. 1). As such, evaluating the analytical predictors also compares favorably to the NN-based approach in terms of computation time. For each of the three methods and across all six studied physical systems, we find that the time needed to train an NN of *minimal size* (one hidden layer with a single node) for a *single epoch* is of the same order of magnitude as the time needed to compute the optimal predictions, optimal indicator, and optimal loss (see Table I). Therefore, the computation time associated with constructing and evaluating an optimal model is *at worst* comparable with training and evaluating an NN-based model. In practice, however, the latter approach typically requires significantly more computation time, because larger NNs need to be used, the training takes many epochs, and hyperparameters need to be adjusted (see Appendix B for a detailed discussion). In particular, as the system size increases and the associated state space grows, converging to the global minimum of the loss function can become increasingly difficult. The convergence of the optimal model, on the other hand, is guaranteed *by construction*.

In Sec. IV, we have observed that the optimal indicator of a given method may fail to correctly highlight a phase transition. A failure can, for example, occur if only a limited amount of data is available and finite-sample statistics dominate. In this case, while the ground-truth probability distributions underlying the data show a significant overlap resulting in a peak in the indicator signal, the inferred probability distributions do not (see Appendix A 5 for a concrete example). However, even if the dataset is sufficiently large, i.e., the ground-truth

probability distributions are well approximated, the optimal model can fail (see classical systems in Sec. IV for examples). Both instances of failure can often be resolved by employing nonoptimal models. Such a model can be realized by an NN whose capacity, i.e., its ability to fit a wide variety of functions [40,41], is restricted. This can be achieved, e.g., by reducing the NN size, performing early stopping, or the explicit addition of  $\ell_2$  regularization (see Appendixes B 1 and B 2). In these instances, other phase-classification methods which are inherently based on the similarity of input data [13,27,32,34] are also expected to provide valuable insights. These methods stand in contrast to the optimal predictors of SL, LBC, and PBM, which are not explicitly based on learning order parameters, i.e., recognizing prevalent patterns or orderings. Instead, the optimal predictors gauge changes in the probability distributions governing the data. Contrary to popular opinion, the failure of optimal models, or, equivalently, high-capacity NNs, does not always correspond to overfitting in the traditional sense [40]: The gap between training and test loss vanishes in the limit of a sufficiently large dataset (which is available for the examples discussed in Sec. IV). Therefore, suboptimal models, such as NNs with insufficient capacity, are, in fact, *underfitting* the data. This signals a fundamental mismatch between the classification or regression task underlying a particular ML method, i.e., the corresponding loss function, and the goal of detecting phase transitions. In particular, it raises the intriguing question of whether one can adjust the learning task in SL, PBM, and LBC such that the corresponding optimal models also correctly highlight the phase transition in these problematic cases, e.g., through an appropriate modification of the underlying loss functions or by enforcing explicit constraints.

## VI. CONCLUSION AND OUTLOOK

The ML methods for detecting phase transitions from data given by SL, LBC, and PBM can be viewed under a unifying light: All three approaches have predictive models, such as NNs, at their heart which are trained to solve a given classification or regression task. Analyzing their predictions allows us to compute a scalar indicator that highlights phase boundaries. The power and success of these methods is largely attributed to the universal function approximation capabilities of their underlying NNs, which are often sacrificed in practice to regain interpretability [15,51–55,106]. Here, we take an alternative approach to cope with the interpretability-expressivity trade-off: By analyzing the class of predictive models that solve the classification and regression tasks underlying SL, LBC, and PBM optimally, we derive analytical expressions for the indicators of phase transitions of these three methods.

Our work establishes a solid theoretical foundation for SL, LBC, and PBM, based on which we are able to explain and understand the results of a variety of previous

studies [4,5,10,23,29,31]. We anticipate that similar analyses will be useful to gain an understanding of other methods for identifying phase transitions with NNs [21,28,32,36,38,107] and other classification tasks in condensed matter physics [85,108–113]. In these cases, the optimal models can also serve as benchmark solutions that enable future studies aimed at investigate the learning process of NNs and improving their design and update routines [11,85,114–117]. For example, in Refs. [4,16,20], it is shown that an NN trained to predict the phase transition in a given model using SL can successfully classify configurations generated from an entirely different Hamiltonian. An exciting prospect is to explore whether the success of this “transfer learning” can be rigorously explained based on our results.

The analytical expressions not only enable our understanding of the phase-classification methods under consideration, they also allow for the direct computation of their optimal predictions and indicators based on the input data *without* explicitly training NNs. We have demonstrated that this novel procedure can successfully reveal a broad range of different phase transitions in a numerical setting and is favorable in terms of computation time. Our results suggest a variety of avenues for further explorations. As a next step, one can consider whether tools from ML, especially for density estimation [118–121], can aid in the computation of the optimal indicators. In the quantum case, classical representations of quantum states obtained via classical shadow tomography [35,80] may help to evade the arising exponential complexity. We believe that optimal predictors will be a valuable tool to detect, interpret, and characterize phases of matter and their transitions from experimental data, particularly in the advent of digital quantum computers [122–126] and programmable quantum simulators [10,11,79,127–129].

The code for computing the optimal predictions and indicators of SL, LBC, and PBM utilized in this work is open source [130].

## ACKNOWLEDGMENTS

We thank Niels Lörch, Andreas Trabesinger, and Christoph Bruder for helpful suggestions on the manuscript. We thank Niels Lörch, Eliska Greplova, Eugene Demler, Florian Marquardt, and Christoph Bruder for stimulating discussions. We acknowledge financial support from the NCCR QSIT funded by the Swiss National Science Foundation (Grant No. 51NF40-185902), as well as from the Swiss National Science Foundation individual grant (Grant No. 200020\_200481). Computation time at sciCORE scientific computing core facility at the University of Basel is gratefully acknowledged. This material is based upon work supported by the National Science Foundation under Grants No. OAC-1835443, No. SII-2029670, No. ECCS-2029670, No. OAC-2103804, and No. PHY-2021825. We also gratefully acknowledge the U.S. Agency for International

Development through Penn State for Grant No. S002283-USAID. The information, data, or work presented herein was funded in part by the Advanced Research Projects Agency-Energy (ARPA-E), U.S. Department of Energy, under Grants No. DE-AR0001211 and No. DE-AR0001222. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof. This material was supported by The Research Council of Norway and Equinor ASA through Research Council project “308817—Digital wells for optimal production and drainage”. Research was sponsored by the United States Air Force Research Laboratory and the United States Air Force Artificial Intelligence Accelerator and was accomplished under Cooperative Agreement No. FA8750-19-2-1000. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the United States Air Force or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

## APPENDIX A: OPTIMAL PREDICTIONS AND INDICATORS

In this appendix, we provide detailed derivations of the optimal predictions and indicators of SL, LBC, and PBM. In particular, we discuss the assumptions underlying the derivation of the optimal predictions of SL and how the analytical predictors are evaluated in practice. This includes an analysis of the computational cost associated with constructing and evaluating the optimal models and the role of finite-sample statistics.

### 1. Derivation of optimal predictions and indicators

Here, we derive the form of the optimal predictions and indicators of phase transitions for SL, LBC, and PBM presented in Sec. II in the main text.

*Supervised learning.*—In SL, a predictive model  $m$  is trained to minimize the CE loss function given in Eq. (1). Now, consider a particular input contained within the training set  $\tilde{\mathbf{x}} \in \bar{\mathcal{T}}$ . We can determine the optimal model prediction  $\hat{y}_{\text{SL}}^{\text{opt}}(\tilde{\mathbf{x}})$  for this particular input by minimizing the loss function in Eq. (1) with respect to  $\hat{y}(\tilde{\mathbf{x}})$ , i.e., by solving the necessary condition

$$\frac{\partial \mathcal{L}_{\text{SL}}}{\partial \hat{y}(\tilde{\mathbf{x}})} = -\frac{1}{M_{\mathcal{T}}} \sum_{\tilde{\mathbf{x}} \in \bar{\mathcal{T}}} \left( \frac{y(\tilde{\mathbf{x}})}{\hat{y}(\tilde{\mathbf{x}})} - \frac{1 - y(\tilde{\mathbf{x}})}{1 - \hat{y}(\tilde{\mathbf{x}})} \right) = 0. \quad (\text{A1})$$

Using the explicit expressions for the labels ( $y = 1$  and  $y = 0$  for all inputs drawn in region I and II, respectively) in Eq. (A1), we have

$$\frac{\sum_{k=1}^{r_{\text{I}}} M_k(\tilde{\mathbf{x}})}{\sum_{k=1}^K M_k(\tilde{\mathbf{x}})} = \frac{M_{\text{I}}(\tilde{\mathbf{x}})}{M_{\text{II}}(\tilde{\mathbf{x}})} = \frac{\hat{y}(\tilde{\mathbf{x}})}{1 - \hat{y}(\tilde{\mathbf{x}})}. \quad (\text{A2})$$

Here,  $M_{\text{I/II}}(\tilde{\mathbf{x}})$  denotes the number of times the input  $\tilde{\mathbf{x}}$  is found in region I or II, respectively. In SL, the predictive model must, by definition, satisfy  $\hat{y}(\mathbf{x}) \in [0, 1] \forall \mathbf{x}$ . Thus, Eq. (A2) is satisfied given predictions of the form

$$\hat{y}_{\text{SL}}^{\text{opt}}(\tilde{\mathbf{x}}) = \frac{M_{\text{I}}(\tilde{\mathbf{x}})}{M_{\text{I}}(\tilde{\mathbf{x}}) + M_{\text{II}}(\tilde{\mathbf{x}})}. \quad (\text{A3})$$

The opposite choice of labeling ( $y = 0$  and  $y = 1$  for all inputs drawn in region I and II, respectively) is equally valid and results in

$$\hat{y}_{\text{SL}}^{\text{opt}}(\tilde{\mathbf{x}}) = \frac{M_{\text{II}}(\tilde{\mathbf{x}})}{M_{\text{I}}(\tilde{\mathbf{x}}) + M_{\text{II}}(\tilde{\mathbf{x}})}. \quad (\text{A4})$$

That is, the roles of  $\hat{y}_{\text{SL}}^{\text{opt}}(\tilde{\mathbf{x}})$  and  $1 - \hat{y}_{\text{SL}}^{\text{opt}}(\tilde{\mathbf{x}})$  are swapped. In this work, we stick to the former choice [Eq. (A3)]. The optimality of the predictions in Eq. (A3) can be confirmed by calculating the second derivative of the loss function:

$$\frac{\partial^2 \mathcal{L}_{\text{SL}}}{\partial \hat{y}(\tilde{\mathbf{x}})^2} = \frac{M_{\text{I}}}{M_{\mathcal{T}}} \frac{1}{\hat{y}(\tilde{\mathbf{x}})^2} + \frac{M_{\text{II}}}{M_{\mathcal{T}}} \frac{1}{[1 - \hat{y}(\tilde{\mathbf{x}})]^2} > 0. \quad (\text{A5})$$

The probability distribution governing the input data is denoted as  $P_k(\tilde{\mathbf{x}}) \approx M_k(\tilde{\mathbf{x}})/M$  ( $1 \leq k \leq K$ ). This allows for Eq. (A3) to be expressed as

$$\hat{y}_{\text{SL}}^{\text{opt}}(\tilde{\mathbf{x}}) = \frac{P_{\text{I}}(\tilde{\mathbf{x}})}{P_{\text{I}}(\tilde{\mathbf{x}}) + P_{\text{II}}(\tilde{\mathbf{x}})}, \quad (\text{A6})$$

where

$$P_{\text{I}}(\tilde{\mathbf{x}}) = \sum_{k=1}^{r_{\text{I}}} P_k(\tilde{\mathbf{x}}) \quad (\text{A7})$$

and

$$P_{\text{II}}(\tilde{\mathbf{x}}) = \sum_{k=r_{\text{I}}+1}^K P_k(\tilde{\mathbf{x}}) \quad (\text{A8})$$

are the (unnormalized) probabilities of drawing the input  $\tilde{\mathbf{x}}$  in region I and II, respectively. Repeating the above procedure for all inputs within the training set  $\bar{\mathcal{T}}$ , we obtain

$$\hat{y}_{\text{SL}}^{\text{opt}}(\mathbf{x}) = \frac{P_{\text{I}}(\mathbf{x})}{P_{\text{I}}(\mathbf{x}) + P_{\text{II}}(\mathbf{x})} \quad \forall \mathbf{x} \in \bar{\mathcal{T}}, \quad (\text{A9})$$

which matches Eq. (9) reported in the main text. Relaxations of the assumption in SL that there are only two distinct phases to be distinguished are discussed in Appendix A 2.

Note that the same optimal predictions are obtained when training on a MSE loss function:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{M_{\mathcal{T}}} \sum_{\mathbf{x} \in \mathcal{T}} [\hat{y}(\mathbf{x}) - y(\mathbf{x})]^2, \quad (\text{A10})$$

instead of a CE loss function. Again, consider a particular input  $\tilde{\mathbf{x}}$  contained within the training set  $\mathcal{T}$ . We can determine the optimal model prediction  $\hat{y}_{\text{SL}}^{\text{opt}}(\tilde{\mathbf{x}})$  for this input by minimizing the loss function in Eq. (A10) with respect to  $\hat{y}(\tilde{\mathbf{x}})$ , i.e., by solving

$$\frac{\partial \mathcal{L}_{\text{MSE}}}{\partial \hat{y}(\tilde{\mathbf{x}})} = \frac{2}{M_{\mathcal{T}}} \sum_{\tilde{\mathbf{x}} \in \mathcal{T}} [\hat{y}(\tilde{\mathbf{x}}) - y(\tilde{\mathbf{x}})] = 0. \quad (\text{A11})$$

Plugging the expression for the labels given by a one-hot encoding into Eq. (A11), we have

$$M_{\text{I}}(\tilde{\mathbf{x}})[1 - \hat{y}(\tilde{\mathbf{x}})] - M_{\text{II}}(\tilde{\mathbf{x}})\hat{y}(\tilde{\mathbf{x}}) = 0. \quad (\text{A12})$$

This coincides with the condition for the predictions given in Eq. (A2) obtained from a CE loss function. Their optimality can be confirmed via

$$\frac{\partial^2 \mathcal{L}_{\text{MSE}}}{\partial \hat{y}(\tilde{\mathbf{x}})^2} = \frac{2(M_{\text{I}} + M_{\text{II}})}{M_{\mathcal{T}}} > 0. \quad (\text{A13})$$

Therefore, in SL, the optimal predictions and indicators associated with optimal models trained on a CE or MSE loss function are identical.

*Learning by confusion.*—To reveal the phase transition by means of LBC, we perform several splits of the parameter range into two neighboring regions labeled I and II. For a fixed bipartition, we minimize a CE [Eq. (4)] or MSE loss function:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{M_{\mathcal{X}}} \sum_{\mathbf{x} \in \mathcal{X}} [\hat{y}(\mathbf{x}) - y(\mathbf{x})]^2. \quad (\text{A14})$$

Following the analysis of SL presented above, we obtain a similar expression for the optimal predictions:

$$\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}) = \frac{P_{\text{I}}(\mathbf{x})}{P_{\text{I}}(\mathbf{x}) + P_{\text{II}}(\mathbf{x})} \quad \forall \mathbf{x} \in \mathcal{X}, \quad (\text{A15})$$

with  $\mathcal{T} = \mathcal{X}$  in LBC. Thus, we recover Eq. (14) in the main text. Their optimality can be confirmed via

$$\frac{\partial^2 \mathcal{L}_{\text{LBC}}}{\partial \hat{y}(\tilde{\mathbf{x}})^2} = \frac{M_{\text{I}}}{M_{\mathcal{X}}} \frac{1}{\hat{y}(\tilde{\mathbf{x}})^2} + \frac{M_{\text{II}}}{M_{\mathcal{X}}} \frac{1}{[1 - \hat{y}(\tilde{\mathbf{x}})]^2} > 0 \quad (\text{A16})$$

or

$$\frac{\partial^2 \mathcal{L}_{\text{MSE}}}{\partial \hat{y}(\tilde{\mathbf{x}})^2} = \frac{2(M_{\text{I}} + M_{\text{II}})}{M_{\mathcal{X}}} > 0, \quad (\text{A17})$$

in the case of a CE or MSE loss, respectively. The value of the indicator in LBC for a given bipartition corresponds to the mean classification accuracy [Eq. (5)], where the continuous predictions  $\hat{y}(\mathbf{x}) \in [0, 1]$  are mapped to binary labels via  $\theta[\hat{y}(\mathbf{x}) - 0.5]$ . Using the optimal prediction in Eq. (A15), the mean classification error for a given input  $\mathbf{x}$  is  $\min\{\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}), 1 - \hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x})\}$ . Weighting the contribution of each input  $\mathbf{x}$  to the mean classification error by its probability  $P_k(\mathbf{x})$ , we arrive at Eq. (15) in the main text. Note that, in principle, the assumption in LBC that there are only two phases to be distinguished can be relaxed [5]. In this case, the optimal indicator may show multiple distinct peaks highlighting the different phase boundaries [131].

*Prediction-based method.*—In PBM, a predictive model  $m: \mathbf{x} \rightarrow \hat{y}(\mathbf{x})$  is trained to minimize the MSE loss function  $\mathcal{L}_{\text{PBM}}$  specified in Eq. (6). Consider a particular input  $\tilde{\mathbf{x}} \in \tilde{\mathcal{X}}$ . We can determine the optimal model prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}(\tilde{\mathbf{x}})$  for this input by minimizing the loss function in Eq. (6) with respect to  $\hat{y}(\tilde{\mathbf{x}})$ , i.e., by solving

$$\frac{\partial \mathcal{L}_{\text{PBM}}}{\partial \hat{y}(\tilde{\mathbf{x}})} = \frac{2}{KM} \sum_{k=1}^K M_k(\tilde{\mathbf{x}})[\hat{y}(\tilde{\mathbf{x}}) - p_k] = 0. \quad (\text{A18})$$

Solving Eq. (A18) yields

$$\hat{y}_{\text{PBM}}^{\text{opt}}(\tilde{\mathbf{x}}) = \frac{\sum_{k=1}^K P_k(\tilde{\mathbf{x}})p_k}{\sum_{k=1}^K P_k(\tilde{\mathbf{x}})}. \quad (\text{A19})$$

This prediction is indeed optimal, as

$$\frac{\partial^2 \mathcal{L}_{\text{PBM}}}{\partial \hat{y}(\tilde{\mathbf{x}})^2} = \frac{2}{K} \sum_{k=1}^K P_k(\tilde{\mathbf{x}}) > 0. \quad (\text{A20})$$

Repeating this procedure for all available inputs  $\mathbf{x} \in \tilde{\mathcal{X}}$  yields

$$\hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x}) = \frac{\sum_{k=1}^K P_k(\mathbf{x})p_k}{\sum_{k=1}^K P_k(\mathbf{x})} \quad \forall \mathbf{x} \in \tilde{\mathcal{X}}. \quad (\text{A21})$$

Thereby, we recover Eq. (16) in the main text. Note that this derivation can be generalized to higher-dimensional parameter spaces (which may host multiple distinct phases) in a straightforward manner (see Ref. [34]), resulting in

$$\hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x}) = \frac{\sum_k P_k(\mathbf{x})\mathbf{p}_k}{\sum_k P_k(\mathbf{x})}. \quad (\text{A22})$$

Here, the sum runs over all sampled points  $\mathbf{p}_k$  in parameter space. The optimal indicator is then given as a divergence:

$$I_{\text{PBM}}^{\text{opt}}(\mathbf{p}) = \nabla_{\mathbf{p}} \delta \mathbf{y}_{\text{PBM}}^{\text{opt}}(\mathbf{p}), \quad (\text{A23})$$

where  $\delta \mathbf{y}_{\text{PBM}}^{\text{opt}}(\mathbf{p}_k) = \sum_{\mathbf{x} \in \bar{\mathcal{X}}} P_k(\mathbf{x}) \hat{\mathbf{y}}_{\text{PBM}}^{\text{opt}}(\mathbf{x}) - \mathbf{p}_k$ .

## 2. Assumptions for supervised learning

Let us we review the assumption of  $\bar{\mathcal{X}} = \bar{\mathcal{T}}$  underlying the derivation for the optimal predictions and corresponding indicator of SL. In general, if  $\bar{\mathcal{X}} \neq \bar{\mathcal{T}}$ , the optimal predictions of SL can be expressed as

$$\hat{\mathbf{y}}_{\text{SL}}^{\text{opt}}(p_k) = \sum_{\mathbf{x} \in \bar{\mathcal{T}}} P_k(\mathbf{x}) \hat{\mathbf{y}}_{\text{SL}}^{\text{opt}}(\mathbf{x}) + \sum_{\mathbf{x} \notin \bar{\mathcal{T}}} P_k(\mathbf{x}) \hat{\mathbf{y}}_{\text{SL}}(\mathbf{x}). \quad (\text{A24})$$

The first contribution in Eq. (A24) comes from predictions for inputs contained in the training data, which are determined through minimization of the corresponding loss function [see Eq. (A1)]. The second contribution comes from predictions for inputs not contained in the training data, which are *a priori* restricted only to the unit interval  $\hat{\mathbf{y}}_{\text{SL}}(\mathbf{x}) \in [0, 1]$ . Therefore, this contribution to Eq. (A24) is bounded by the probability of drawing an input at  $p_k$  that is not present in the training data, which is given by  $\sum_{\mathbf{x} \notin \bar{\mathcal{T}}} P_k(\mathbf{x})$ . When using SL with NNs, the predictions for inputs not contained in the training data [second contribution in Eq. (A24)] are most susceptible to noise inherent to NN training and hyperparameter choices. As such, its physical relevance is questionable. It may be possible to obtain better bounds for this second contribution when using SL with NNs, e.g., based on the theory of neural tangent kernels [132].

Let us explicitly discuss the classical systems analyzed in this work, which are governed by Boltzmann distribution [Eqs. (28) and (29)]. Because the probability of drawing a particular configuration sample (or energy) at any nonzero temperature is nonzero, the assumption of  $\bar{\mathcal{X}} = \bar{\mathcal{T}}$  holds given a sufficient number of samples. When computing the optimal indicator of SL numerically, we work with a finite number of samples. Thus, it can happen that an input is encountered which is not part of the training data  $\mathbf{x} \notin \bar{\mathcal{T}}$ . In practice, we can verify *on the fly* whether this is the case. If so, we set  $y_{\text{SL}}(\mathbf{x}) = 0$  in Eq. (A24). Thereby, we effectively ignore the contribution to the predictions of SL from inputs not present in the training data. Note that, because these predictions correspond to inputs with low probability, they are also most susceptible to finite-sample statistics. This procedure is further justified by the fact that the optimal predictions  $\hat{\mathbf{y}}_{\text{SL}}^{\text{opt}}$  obtained in this manner track the ground-state probability with high accuracy [see Figs. 3(b), 4(b), and 5(b)]. That is, the optimal predictions closely match the expression in Eq. (33) valid in the case where deviations due to finite-sample statistics vanish.

In the quantum case, it is typically not straightforward to determine *a priori* whether the assumption of  $\bar{\mathcal{X}} = \bar{\mathcal{T}}$  is met for a given system and choice of basis. Here, when

calculating the optimal predictions and indicators numerically, we use the same procedure as described for the classical case. In our study, we find only cases where  $\mathbf{x} \notin \bar{\mathcal{T}}$  for the XXZ model. The error resulting from neglecting the second contribution in Eq. (A24) is marginal, as the probability of drawing such inputs across the parameter range is found to be small. Note that the optimal indicator of SL obtained in such a manner correctly reveals the quantum phase transition in the XXZ model (see Fig. 6). In fact, the optimal predictions calculated via this procedure correspond to the probability of measuring the ferromagnetic ground state (see Sec. IV D). For the above reasons, we expect that the optimal predictions of SL are capable of revealing phase transitions even if  $\bar{\mathcal{X}} \neq \bar{\mathcal{T}}$ .

A relevant scenario in which the assumption that  $\bar{\mathcal{X}} = \bar{\mathcal{T}}$  is violated occurs when the system transitions between multiple phases as the tuning parameter is varied. Then, inputs drawn in the phases present in the middle of the sampled range of the tuning parameter may not be present in the two boundary phases. By dropping the second contribution in Eq. (A24), we may still faithfully detect the transition between the first and second phase. However, all subsequent phase boundaries are then likely missed. In the future, it will be of interest to lift the assumption of  $\bar{\mathcal{X}} = \bar{\mathcal{T}}$  underlying the optimal predictions through appropriate interpolation schemes [31,35,132], which would allow for the generalization capabilities of SL to be explored.

## 3. Computational cost

Here, we derive the scaling of the computational cost with the number of unique inputs  $M_{\bar{\mathcal{X}}}$  and the number of sampled tuning parameter values  $N$  reported in Sec. III B in the main text. Note that we do not consider the overhead associated with computing the probability distributions  $\{P_k(\mathbf{x})\}_{k=1}^K \forall \mathbf{x} \in \bar{\mathcal{X}}$  from the data at hand (or any other constant overhead). The computation of the optimal predictions and indicators can be approached in two ways: Either the optimal predictions for a given input  $\hat{\mathbf{y}}^{\text{opt}}(\mathbf{x})$  are recomputed in each function call, or they are cached. We report the required number of floating-point operations in both instances, which can be counted based on the analytical expressions reported in Sec. III. This counting represents a rough, hardware-independent estimate of the required computational cost. In the following, we assume that the optimal indicators in SL and PBM are computed using a symmetric difference quotient; cf. Eq. (22).

*Supervised learning.*—The computation of  $\hat{\mathbf{y}}_{\text{SL}}^{\text{opt}}$  for all  $\mathbf{x} \in \bar{\mathcal{X}}$  requires  $M_{\bar{\mathcal{X}}} K_{\mathcal{T}}$  floating-point operations, where  $K_{\mathcal{T}} = K_{\text{I}} + K_{\text{II}}$  is the number of sampled values of the tuning parameter in the training regions I and II. Caching these values, the number of operations required to compute the mean optimal prediction  $\hat{\mathbf{y}}_{\text{SL}}^{\text{opt}}$  for all  $\{p_k\}_{k=1}^K$  is  $K(2M_{\bar{\mathcal{X}}} - 1) + M_{\bar{\mathcal{X}}} K_{\mathcal{T}}$ . Thus, computing the optimal indicator requires  $M_{\bar{\mathcal{X}}}(2K + K_{\mathcal{T}}) + K$  operations. Typically, in

SL we have  $K_{\mathcal{T}} \ll K$ . Under this assumption, the computation of the mean optimal predictions and the optimal indicators each require  $O(M_{\bar{x}}K)$  operations. If the values  $\hat{y}_{\text{SL}}^{\text{opt}}(\mathbf{x}) \forall \mathbf{x} \in \bar{\mathcal{X}}$  are not cached, computing the mean optimal prediction instead requires  $K[(2M_{\bar{x}} - 1) + M_{\bar{x}}K_{\mathcal{T}}]$  operations. Computing the optimal indicator then requires  $M_{\bar{x}}K(2 + K_{\mathcal{T}}) + K$  operations. For both quantities, this still corresponds to  $O(M_{\bar{x}}K)$  operations.

*Learning by confusion.*—The computation of  $\hat{y}_{\text{LBC}}^{\text{opt}}$  for all  $\mathbf{x} \in \bar{\mathcal{X}}$  requires  $M_{\bar{x}}K$  floating-point operations. Caching these values, the number of operations required to compute the optimal indicator is  $M_{\bar{x}}K^2(F_{\min} + 2)$ , where  $F_{\min}$  denotes the number of floating-point operations required to compute  $\min\{\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}), 1 - \hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x})\}$ . This corresponds to  $O(M_{\bar{x}}K^2)$  operations. Without caching, the optimal indicator requires  $M_{\bar{x}}K^3 + M_{\bar{x}}K^2(F_{\min} + 2) + K$  operations to compute, resulting in a scaling of  $O(M_{\bar{x}}K^3)$ .

*Prediction-based method.*—In PBM, the computation of  $\hat{y}_{\text{PBM}}^{\text{opt}}$  for all  $\mathbf{x} \in \bar{\mathcal{X}}$  requires  $M_{\bar{x}}(3K - 1)$  floating-point operations. Caching these values, the number of operations required to compute the mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  for all  $\{p_k\}_{k=1}^K$  is  $5M_{\bar{x}}K - K - M_{\bar{x}}$ . Computing the optimal indicator then requires  $M_{\bar{x}}(5K - 1) + K$  operations. The

computation of the mean optimal predictions and the optimal indicator each require  $O(M_{\bar{x}}K)$  operations. If the values  $\hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x}) \forall \mathbf{x} \in \bar{\mathcal{X}}$  are not cached, computing the mean optimal prediction instead requires  $3M_{\bar{x}}K^2 + K(M_{\bar{x}} - 1)$  operations. Computing the optimal indicator then requires  $3M_{\bar{x}}K^2 + KM_{\bar{x}} + K$  operations. For both quantities, this results in a scaling of  $O(M_{\bar{x}}K^2)$ .

*Numerical implementation.*—The measured computation times associated with calculating the optimal indicators of phase transitions of SL, LBC, and PBM for all six physical systems discussed in the main text (see Sec. IV) are reported in Table I. The corresponding code is open source [130]. Again, we do not consider the computational cost associated with generating samples and estimating the underlying probability distributions.

Overall, the computation times are remarkably low. For all systems, the optimal indicator of SL and PBM can be obtained in under a second and the optimal indicator of LBC in under a minute. We observe that the computation times of SL and PBM are comparable, with PBM being slightly slower than SL. In contrast, the computation times of LBC are 2 orders of magnitude larger. Note that these are the evaluation times corresponding to the largest system sizes under consideration. We find that the computation times

TABLE I. Measured computation times in seconds associated with constructing and evaluating optimal models,  $t^{\text{opt}}$ , or training an NN of minimal size (one hidden layer with a single node) for a single epoch,  $t^{\text{NN}}$ , for all three methods and six systems discussed in the main text (see Sec. IV). The linear system size  $L$  and the corresponding number of unique samples  $M_{\bar{x}}$  as well as the number of sampled values of the tuning parameter  $K$  for each system are also reported. The construction and evaluation of the optimal models yield the optimal predictions, optimal indicator, and optimal loss value. A training epoch is comprised of evaluating the NN at all  $M_{\bar{x}}$  unique samples, calculating the loss function, obtaining the gradient via backpropagation, and performing a single gradient step. For details on the NN architecture and training, see Appendix B. Note that, in LBC,  $t_{\text{LBC}}^{\text{NN}}$  corresponds to  $K + 1$  times the computation time of a training epoch for a single NN. All computation times are measured on a single CPU [Intel(R) Xeon(R) CPU E5-2630 v4 @ 2.20 GHz], and garbage collection times are subtracted from the total run-time. To gather statistics, for each method and system computations are run for 20 h. If  $10^5$  independent runs are completed in less than 20 h, the computations are stopped prematurely. The error corresponds to the observed standard deviation.

	Ising	IGT	XY	XXZ	Kitaev	Bose-Hubbard
$t_{\text{SL}}^{\text{opt}}$	0.0007 ± 0.0002	0.00007 ± 0.00002	0.00012 ± 0.00003	0.0049 ± 0.0009	0.17 ± 0.02	0.0044 ± 0.0009
$t_{\text{SL}}^{\text{NN}}$	0.00060 ± 0.00005	0.00030 ± 0.00002	0.00048 ± 0.00003	0.0060 ± 0.0009	0.14 ± 0.02	0.0023 ± 0.0003
$t_{\text{SL}}^{\text{NN}}/t_{\text{SL}}^{\text{opt}}$	0.9 ± 0.3	4.9 ± 1.3	4.0 ± 0.8	1.2 ± 0.3	0.9 ± 0.2	0.5 ± 0.1
$t_{\text{PBM}}^{\text{opt}}$	0.0016 ± 0.0004	0.00014 ± 0.00006	0.00021 ± 0.00008	0.019 ± 0.003	0.42 ± 0.05	0.009 ± 0.002
$t_{\text{PBM}}^{\text{NN}}$	0.0042 ± 0.0007	0.0005 ± 0.0001	0.00084 ± 0.00004	0.080 ± 0.006	1.2 ± 0.1	0.026 ± 0.004
$t_{\text{PBM}}^{\text{NN}}/t_{\text{PBM}}^{\text{opt}}$	2.7 ± 0.8	4.0 ± 2.1	4.0 ± 1.5	4.2 ± 0.8	2.8 ± 0.4	2.7 ± 0.6
$t_{\text{LBC}}^{\text{opt}}$	0.8 ± 0.1	0.042 ± 0.001	0.041 ± 0.004	3.7 ± 0.4	32.0 ± 1.7	1.4 ± 0.2
$t_{\text{LBC}}^{\text{NN}}$	1.11 ± 0.06	0.09 ± 0.01	0.12 ± 0.01	12.2 ± 1.2	93.9 ± 3.8	3.2 ± 0.4
$t_{\text{LBC}}^{\text{NN}}/t_{\text{LBC}}^{\text{opt}}$	1.3 ± 0.2	2.1 ± 0.2	2.8 ± 0.4	3.3 ± 0.5	3.0 ± 0.2	2.4 ± 0.5
$t_{\text{PBM}}^{\text{opt}}/t_{\text{SL}}^{\text{opt}}$	2.3 ± 0.9	2.0 ± 1.1	1.8 ± 0.7	3.8 ± 1.0	2.7 ± 0.5	2.1 ± 0.6
$t_{\text{LBC}}^{\text{opt}}/t_{\text{SL}}^{\text{opt}}$	1231 ± 374	629 ± 164	346 ± 77	751 ± 158	204 ± 33	308 ± 78
$L$	60	28	60	14	20	8
$M_{\bar{x}}$	1711	353	1000	16 384	524 288	6435
$K$	200	100	100	201	101	200

qualitatively agree with the complexity analysis described above (for the relevant case where caching is performed). An additional speedup can be gained through parallel execution. In particular, it is straightforward to compute optimal predictions (in the case of SL and PBM) and optimal indicators (in the case of LBC) at discrete values of the tuning parameter in parallel, e.g., via multithreading (which is implemented in Ref. [130]).

#### 4. Boltzmann-distributed inputs

Let us discuss the special case when the drawn inputs  $\mathbf{x}$ , such as spin configurations, follow a Boltzmann distribution

$$P_k(\mathbf{x}) = \frac{e^{-H(\mathbf{x})/k_B T_k}}{Z_k}. \quad (\text{A25})$$

The probability to draw a sample with energy  $E$  is, thus, given by

$$P_k(E) = \frac{g(E)e^{-E/k_B T_k}}{Z_k}, \quad (\text{A26})$$

where  $g(E)$  is the corresponding degeneracy factor

$$g(E) = \sum_{\mathbf{x} \in \mathcal{S}} \delta_{H(\mathbf{x}), E}. \quad (\text{A27})$$

Here,  $\mathcal{S}$  denotes the state space of the samples  $\mathbf{x}$ , i.e., the set of all unique samples without duplicates. Therefore, we have

$$P_k(\mathbf{x}) = P_k[H(\mathbf{x})]/g[H(\mathbf{x})]. \quad (\text{A28})$$

*Supervised learning.*—Plugging Eq. (A28) into Eq. (9), we immediately find that

$$\begin{aligned} \hat{y}_{\text{SL}}^{\text{opt}}(\mathbf{x}) &= \frac{P_{\text{I}}[H(\mathbf{x})]}{P_{\text{I}}[H(\mathbf{x})] + P_{\text{II}}[H(\mathbf{x})]} \\ &= \hat{y}_{\text{SL}}^{\text{opt}}[H(\mathbf{x})] \quad \forall \mathbf{x} \in \mathcal{S}, \end{aligned} \quad (\text{A29})$$

where we assume that  $\bar{\mathcal{T}} = \bar{\mathcal{X}} = \mathcal{S}$ . Using Eq. (12), we have

$$\begin{aligned} \hat{y}_{\text{SL}}^{\text{opt}}(p_k) &= \sum_{\mathbf{x} \in \mathcal{S}} P_k(\mathbf{x}) \hat{y}_{\text{opt}}(\mathbf{x}) \\ &= \sum_{\mathbf{x} \in \mathcal{S}} P_k[H(\mathbf{x})] \hat{y}_{\text{opt}}[H(\mathbf{x})]/g[H(\mathbf{x})] \\ &= \sum_{E \in \mathcal{S}_E} P_k(E) \hat{y}_{\text{opt}}(E), \end{aligned} \quad (\text{A30})$$

where  $\mathcal{S}_E$  is the set of unique energies corresponding to the state space  $\mathcal{S}$ . To obtain an expression for the optimal loss, we can rewrite Eq. (1) as

$$\begin{aligned} \mathcal{L}_{\text{SL}} &= -\frac{1}{r_{\text{I}} + (K - l_{\text{II}} + 1)} \sum_{k=1}^{r_{\text{I}}} \sum_{k=l_{\text{II}}}^K \sum_{\mathbf{x} \in \mathcal{S}} P_k(\mathbf{x}) \\ &\quad \times \{y(\mathbf{x}) \ln [\hat{y}(\mathbf{x})] + [1 - y(\mathbf{x})] \ln [1 - \hat{y}(\mathbf{x})]\}. \end{aligned} \quad (\text{A31})$$

Using Eq. (A29), we have

$$\begin{aligned} \mathcal{L}_{\text{SL}}^{\text{opt}} &= -\frac{1}{r_{\text{I}} + (K - l_{\text{II}} + 1)} \sum_{k=1}^{r_{\text{I}}} \sum_{k=l_{\text{II}}}^K \sum_{\mathbf{x} \in \mathcal{S}} P_k[H(\mathbf{x})] \\ &\quad \times \{y[H(\mathbf{x})] \ln [\hat{y}_{\text{SL}}^{\text{opt}}[H(\mathbf{x})]] \\ &\quad + [1 - y[H(\mathbf{x})]] \ln [1 - \hat{y}_{\text{SL}}^{\text{opt}}[H(\mathbf{x})]]\}, \end{aligned} \quad (\text{A32})$$

where we use the fact that  $y(\mathbf{x}) = y[H(\mathbf{x})]$ , i.e., the assigned labels remain identical. Equation (A32) can be simplified to

$$\begin{aligned} \mathcal{L}_{\text{SL}}^{\text{opt}} &= -\frac{1}{r_{\text{I}} + (K - l_{\text{II}} + 1)} \sum_{k=1}^{r_{\text{I}}} \sum_{k=l_{\text{II}}}^K \sum_{E \in \mathcal{S}_E} P_k(E) \\ &\quad \times \{y(E) \ln [\hat{y}_{\text{SL}}^{\text{opt}}(E)] + [1 - y(E)] \ln [1 - \hat{y}_{\text{SL}}^{\text{opt}}(E)]\}, \end{aligned} \quad (\text{A33})$$

using Eq. (A28).

*Learning by confusion.*—For a fixed bipartition in LBC, we can proceed in a similar manner. Plugging Eq. (A28) into Eq. (14) assuming  $\bar{\mathcal{X}} = \mathcal{S}$ , we have

$$\begin{aligned} \hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}) &= \frac{P_{\text{I}}[H(\mathbf{x})]}{P_{\text{I}}[H(\mathbf{x})] + P_{\text{II}}[H(\mathbf{x})]} \\ &= \hat{y}_{\text{LBC}}^{\text{opt}}[H(\mathbf{x})] \quad \forall \mathbf{x} \in \mathcal{S}. \end{aligned} \quad (\text{A34})$$

Using Eq. (15), this yields

$$\begin{aligned} I_{\text{LBC}}^{\text{opt}} &= 1 - \frac{1}{K} \sum_{k=1}^K \sum_{\mathbf{x} \in \mathcal{S}} P_k(\mathbf{x}) \min\{\hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x}), 1 - \hat{y}_{\text{LBC}}^{\text{opt}}(\mathbf{x})\} \\ &= 1 - \frac{1}{K} \sum_{k=1}^K \sum_{E \in \mathcal{S}_E} P_k(E) \min\{\hat{y}_{\text{LBC}}^{\text{opt}}(E), 1 - \hat{y}_{\text{LBC}}^{\text{opt}}(E)\}. \end{aligned} \quad (\text{A35})$$

To obtain an expression for the optimal loss, we follow the above procedure outlined for SL starting with Eq. (4) and eventually arrive at

$$\begin{aligned} \mathcal{L}_{\text{LBC}}^{\text{opt}} &= -\frac{1}{K} \sum_{k=1}^K \sum_{E \in \mathcal{S}_E} P_k(E) \\ &\quad \times \{y(E) \ln [\hat{y}(E)] + [1 - y(E)] \ln [1 - \hat{y}(E)]\}. \end{aligned} \quad (\text{A36})$$

*Prediction-based method.*—Plugging Eq. (A28) into Eq. (16) assuming  $\bar{\mathcal{X}} = \mathcal{S}$ , we find that

$$\begin{aligned} \hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x}) &= \frac{\sum_{k=1}^K P_k[H(\mathbf{x})] p_k}{\sum_{k=1}^K P_k[H(\mathbf{x})]} \\ &= \hat{y}_{\text{PBM}}^{\text{opt}}[H(\mathbf{x})] \quad \forall \mathbf{x} \in \mathcal{S}. \end{aligned} \quad (\text{A37})$$

Using Eq. (17), we have

$$\begin{aligned}\hat{y}_{\text{PBM}}^{\text{opt}}(p_k) &= \sum_{\mathbf{x} \in \mathcal{S}} P_k(\mathbf{x}) \hat{y}_{\text{PBM}}^{\text{opt}}(\mathbf{x}) \\ &= \sum_{E \in \mathcal{S}_E} P_k(E) \hat{y}_{\text{PBM}}^{\text{opt}}(E).\end{aligned}\quad (\text{A38})$$

To obtain an expression for the optimal loss, we rewrite Eq. (6) as

$$\mathcal{L}_{\text{PBM}} = \frac{1}{K} \sum_{k=1}^K \sum_{\mathbf{x} \in \mathcal{S}} P_k(\mathbf{x}) [\hat{y}(\mathbf{x}) - y(\mathbf{x})]^2. \quad (\text{A39})$$

Using Eq. (A37), we have

$$\mathcal{L}_{\text{PBM}}^{\text{opt}} = \frac{1}{K} \sum_{k=1}^K \sum_{\mathbf{x} \in \mathcal{S}} P_k(\mathbf{x}) \{\hat{y}_{\text{PBM}}^{\text{opt}}[H(\mathbf{x})] - y[H(\mathbf{x})]\}^2, \quad (\text{A40})$$

where  $y(\mathbf{x}) = y[H(\mathbf{x})]$ . With Eq. (A28), we finally get

$$\mathcal{L}_{\text{PBM}}^{\text{opt}} = \frac{1}{K} \sum_{k=1}^K \sum_{E \in \mathcal{S}_E} P_k(E) [\hat{y}_{\text{PBM}}^{\text{opt}}(E) - y(E)]^2. \quad (\text{A41})$$

This shows that the optimal predictions, indicators, and loss values of SL, LBC, and PBM remain identical when configuration samples which follow a Boltzmann distribution are used as input or when the corresponding energies are used as input instead. In practice, given a finite set of samples, the inferred probability distribution  $P_k(\mathbf{x}) \approx M_k(\mathbf{x})/M$  is only approximately Boltzmann, i.e.,  $\tilde{T}, \tilde{\mathcal{X}} \approx \mathcal{S}$ , and the two scenarios are equivalent only up to deviations due to finite-sample statistics. In particular, the inferred probability distribution  $P_k(\mathbf{x}) = M_k(\mathbf{x})/M$  based on raw configuration samples may not correspond to the inferred probability distribution  $P_k(E) = M_k(E)/M$  based on the corresponding energy, where the degeneracy factor for the conversion is inferred from the samples as

$$g(E) = \sum_{\mathbf{x} \in \tilde{\mathcal{X}}} \delta_{H(\mathbf{x}), E}. \quad (\text{A42})$$

However, using the energy as input instead of configuration samples yields a more accurate estimate of the ground-truth distribution. This is because the associated state space  $\mathcal{S}_E$  is significantly smaller compared to the entire configuration space  $\mathcal{S}$ , resulting in better statistics given a fixed number of samples. In the 2D Ising model, for example, the size of the configuration space is  $2^{L^2}$ , whereas there are  $L^2 - 1$  unique number of energies (for even  $L$ ). Therefore, the optimal predictions and indicators obtained using the energy as input converge significantly faster compared to the case where raw spin configurations are used. Note that the energy is readily available in numerical studies. However, in principle, one can obtain the same results without having access to the energy given that a sufficient number of raw configurations are

sampled. In the future, it will be of interest to employ more elaborate techniques for density estimation [118–121] in order to obtain a more accurate estimate of the underlying distribution given a reduced dataset size.

Finally, let us continue the analysis of the optimal predictions and indicators of SL in the case of Boltzmann-distributed inputs. We take region I to be composed of a single point  $T_1$ . Let  $T_1 \rightarrow 0$  such that

$$P_1(E) = \begin{cases} 1 & \text{if } E = E_{\text{gs}}, \\ 0 & \text{otherwise,} \end{cases} \quad (\text{A43})$$

where  $E_{\text{gs}}$  is the ground-state energy. Plugging into Eq. (A9) yields

$$\hat{y}_{\text{SL}}^{\text{opt}}(E) = \begin{cases} \frac{1}{1 + P_{\text{II}}(E_{\text{gs}})} & \text{if } E = E_{\text{gs}}, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{A44})$$

We calculate the mean prediction at a given temperature as

$$\hat{y}_{\text{SL}}^{\text{opt}}(T_k) = \sum_{E \in \mathcal{S}_E} P_k(E) \hat{y}_{\text{SL}}^{\text{opt}}(E). \quad (\text{A45})$$

Using Eq. (A44), this results in

$$\hat{y}_{\text{SL}}^{\text{opt}}(T_k) = \frac{P_k(E_{\text{gs}})}{1 + P_{\text{II}}(E_{\text{gs}})}. \quad (\text{A46})$$

Assuming region II is composed of a single point  $T_K$ , we have  $P_{\text{II}}(E_{\text{gs}}) = P_K(E_{\text{gs}})$  and recover Eq. (33) in the main text. For  $T_K \rightarrow \infty$ , we have  $P_K(E_{\text{gs}}) = g(E_{\text{gs}})/M_{\mathcal{S}}$ , where  $M_{\mathcal{S}}$  is the total number of unique system configurations. For the two-dimensional Ising model, for example,  $M_{\mathcal{S}} = 2^{L \times L}$ . Approaching the thermodynamic limit, this yields  $\hat{y}_{\text{SL}}^{\text{opt}}(T_k) \rightarrow P_k(E_{\text{gs}})$ .

Note that these results can be extended to non-Boltzmann distributions: Given that

$$P_1(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}^*, \\ 0 & \text{otherwise} \end{cases} \quad (\text{A47})$$

and following the same procedure as above, we have

$$\hat{y}_{\text{SL}}^{\text{opt}}(p_k) = \frac{P_k(\mathbf{x}^*)}{1 + P_{\text{II}}(\mathbf{x}^*)}. \quad (\text{A48})$$

In particular, Eq. (A48) can be used to qualitatively explain the optimal indicator signals of SL in the XXZ chain (Sec. IV D) and Kitaev chain (Sec. IV E). In this case,  $\mathbf{x}^*$  corresponds to a ground state which is one of the chosen basis states.

## 5. Finite-sample statistics

Finally, we investigate how the optimal predictions and indicators of SL, LBC, and PBM change as the number of



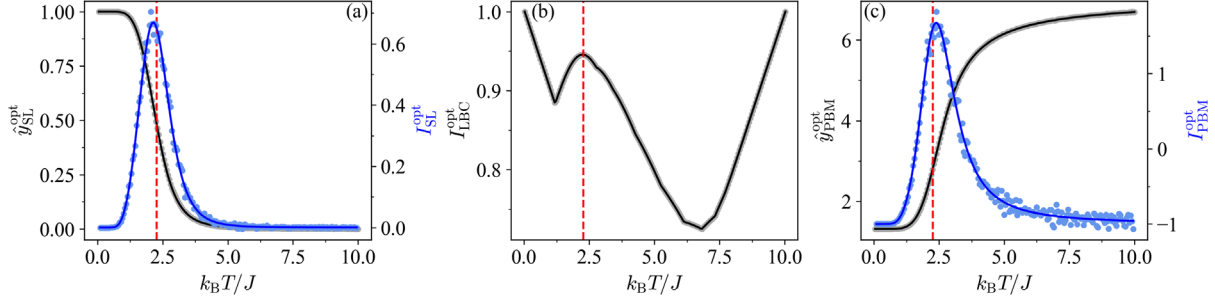


FIG. 9. Results for the Ising model ( $L = 4$ ) with the dimensionless temperature as a tuning parameter  $p = k_B T / J$ , where  $p_1 = 0.05$ ,  $p_K = 10$ , and  $\Delta p = 0.05$ . The critical temperature [Eq. (31)] is highlighted by a red dashed line. In SL, the data obtained at  $p_1$  and  $p_K$  constitute our training set, i.e.,  $r_1 = 1$  and  $I_{II} = K$ . The inputs are computed based on spin configurations obtained through exact enumeration (lines) or Monte Carlo sampling (points). (a) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL (black line) and the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$  (blue line). (b) Optimal indicator of LBC,  $I_{\text{LBC}}^{\text{opt}}$  (black line). (c) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM (black line) and the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$  (blue line).

data points  $M$  per sampled value of the tuning parameter is varied. Recall that the results for the classical systems displayed in the main text are obtained using the energy from Monte Carlo sampling as input, where  $M = 10^5$  spin configurations are drawn per temperature. For small lattice sizes, however, it is possible to enumerate all spin configurations explicitly. In Fig. 9, we compare the optimal predictions and indicators for the Ising model on a  $4 \times 4$  lattice when enumerating all  $2^{16} = 65536$  spin configurations explicitly or using Monte Carlo sampling with  $10^5$  number of configurations per sampled value of the tuning parameter. The results obtained based on the two distinct datasets are in good agreement, which is to be expected given that there are only 15 unique energies. The noise present in the indicator signals of SL and PBM when using Monte Carlo samples is absent when using exact enumeration. In the latter case, both indicators vary smoothly as a function of the temperature. As such, this noise can be attributed to finite-sample statistics.

In general, for both the classical and quantum systems, we observe that the overlap in the underlying probability

distributions leading to a peak in the indicator signals decreases as the number of samples  $M$  is decreased. However, meaningful results can already be obtained when only a fraction of the total state space is covered. In the case of the Ising model on a  $60 \times 60$  lattice, for example, we observe that the optimal predictions and indicators are already well converged for  $M = 10^2$ , i.e., matching the results obtained with  $M = 10^5$ . In particular, the key features in the indicators, i.e., the peak locations, can already be identified for  $M = 10$ . Compare this to the number of unique energies given by 3599.

Figure 10 shows the optimal predictions and indicators of SL, LBC, and PBM for the Kitaev chain of length  $L = 20$  given various values of  $M$ . Recall that the results for the quantum systems displayed in the main text (see Sec. IV) are obtained based on the “ground-truth” probability distributions from exact diagonalization. Here, we explicitly sample these probability distributions, i.e., perform projective measurements and infer the probability distribution based on the measurement results. In SL and PBM, accurate estimates for the critical value of the tuning

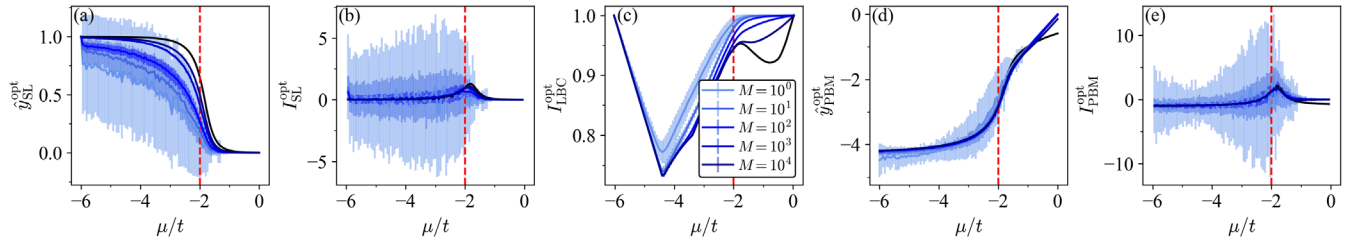


FIG. 10. Optimal predictions and indicators of SL, LBC, and PBM for the Kitaev chain ( $L = 20$ ) with a varying number of data points  $M$  per sampled value of the tuning parameter  $p = \mu / t$ , where  $p_1 = -6$ ,  $p_K = 0$ , and  $\Delta p = 0.06$ . In SL, the data obtained at  $p_1$  and  $p_K$  constitute our training set, i.e.,  $r_1 = 1$  and  $I_{II} = K$ . The critical value  $\mu_c / t = -2$  is highlighted by a red dashed line. The optimal predictions and indicators obtained based on the ground-truth probability distributions from exact diagonalization are shown in black. (a) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL and (b) the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$ . (c) Optimal indicator of LBC,  $I_{\text{LBC}}^{\text{opt}}$ . (d) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM and (e) the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$ . Here, we report results averaged over 100 independent datasets, where the error bars correspond to the standard deviation.

parameter can be obtained based on  $M = 10^3$  samples, whereas  $M = 10^4$  samples are required for a local maximum to emerge in LBC. This covers only a fraction of the total state space comprised of  $M_{\bar{\chi}} = 524288$  states. Notice that the indicator of LBC shows a plateau close to one in the topological phase for a small number of samples, which signifies the absence of “confusion” inherent to the data (see Fig. 10). Similarly, the optimal prediction of PBM is approximately linear in the topological phase for a small number of samples, corresponding to a model which can perfectly resolve the value of the tuning parameter associated with the input. This demonstrates the fact that, while the ground-truth probability distributions may have substantial overlap, estimated probabilities based on a drawn dataset may not.

The high level of uncertainty in the indicator of SL and PBM compared to LBC can be attributed to the symmetric difference quotient used to approximate the derivative. Moreover, in LBC, we associate a distinct optimal predictive model to each bipartition point, whereas the optimal indicator is extracted from a single optimal model in the case of SL and PBM. This leads to an additional suppression of fluctuations in the case of LBC. In the future, it will be of interest to enhance the quality of the optimal predictions and indicators based on finite data through improved derivative computations in the case of SL and PBM [133], as well as more elaborate techniques for density estimation [118–121].

## APPENDIX B: COMPUTATION USING NEURAL NETWORKS

In this appendix, we discuss the application of SL, LBC, and PBM to the six physical systems discussed in the main text (see Sec. IV) using NNs. First, we show that one can recover the optimal analytical predictions and indicators by training NNs. Next, we discuss the computational cost associated with training NNs compared to constructing and evaluating optimal models. Finally, we investigate the influence of NN size, early stopping, regularization, and finite-sample statistics on the results.

*Data preparation.*—For the classical systems (Ising model, IGT, and XY model), the energy  $H(\sigma)$  of the spin configurations  $\sigma$  sampled from Boltzmann distributions at various temperatures serves as an input. To counteract the effect of finite-sample statistics on the predictions in the case of SL due to inputs not contained in the training set  $\mathbf{x} \notin \bar{\mathcal{T}}$ , i.e.,  $\bar{\mathcal{X}} \neq \bar{\mathcal{T}}$ , we modify the corresponding probability distributions, such that  $P_K(\mathbf{x}) = 1/(M + M_{\notin \bar{\mathcal{T}}})$  as opposed to  $P_K(\mathbf{x}) = 0$ . Here,  $M_{\notin \bar{\mathcal{T}}}$  denotes the number of such inputs at  $p_K$ . That is, we add a single instance of each sample which does not appear at the boundary point  $p_K$  to the corresponding dataset  $\mathcal{X}_K$ . Alternatively, we could set these predictions to zero as discussed in Appendix A 2. While the NN-based indicator can change if no such

modifications are performed, this does not resolve the instances where the optimal indicator of SL fails to locate the phase transition (such as in the Ising model or XY model). For the quantum systems (XXZ chain, Kitaev chain, and Bose-Hubbard model), the index of the corresponding basis states serves as input. We use a physically motivated encoding, where the  $S^z$  eigenstate given by  $|\uparrow \downarrow \dots \uparrow\rangle$  and the Fock state  $|10\dots 1\rangle$  are encoded as a bit string  $\mathbf{x} = (10\dots 1)$ .

Before training the NNs, each input  $\mathbf{x} = \{x_i\}$  is standardized via the following affine transformation:

$$x'_i = \frac{x_i - \langle x_i \rangle}{\sigma_{x_i}}, \quad (\text{B1})$$

where  $\langle x_i \rangle$  and  $\sigma_{x_i}$  are the mean value and standard deviation of  $x_i$  across the training data, respectively. Standardization generally leads to a faster rate of convergence when applying gradient-based optimizers [134]. Note that this bijective mapping does not change the probability associated with each input, i.e.,  $P_k(\mathbf{x}) = P_k(\mathbf{x}') \forall 1 \leq k \leq K$ . Therefore, the optimal predictions and indicators remain unchanged.

*Neural network architecture.*—For simplicity, the NNs used in this work consist of a series of fully connected layers, where rectified linear units (ReLU),  $f(z) = \max(0, z)$ , are used as activation functions [40]. The NNs for SL and LBC have two output nodes, where a softmax activation function

$$f_i(z) = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad (\text{B2})$$

is used in the output layer to guarantee that  $\hat{y}(\mathbf{x}') \in [0, 1]$ . Here, the sum runs over all output nodes, and  $\hat{y}$  corresponds to the value of one of the output nodes after application of the softmax activation function. In PBM, no activation function is used for the output layer. The value of the single output node corresponds to  $\hat{y}(\mathbf{x}')$ , which is the estimated value of the tuning parameter at which the input  $\mathbf{x}'$  was drawn. For the prototypical probability distributions discussed in Sec. III A in the main text, we use a single hidden layer with 64 nodes. The number of hidden layers and nodes for all other models is reported in the corresponding figure captions.

*Training.*—The NNs are implemented using Flux in JULIA [135], where the weights and biases are optimized via gradient descent with Adam [136] to minimize the loss function over a series of training epochs. In SL and LBC, we train on a CE loss function [Eq. (1) and (4), respectively], whereas in PBM we train on a MSE loss function [Eq. (6)]. Gradients are calculated using back-propagation [40, 137, 138]. For the prototypical probability distributions discussed in Sec. III A, we train for 10 000 epochs with a learning rate of 0.001. The number of training epochs and learning rate for all other models are reported in the corresponding figure captions.

*Results.*—Figures 11 and 12 show the predictions and indicators of the three methods obtained using NNs (dashed lines) after long training for all six physical systems considered in the main text. Here, we choose the smallest system sizes for convenience. Overall, they are in excellent agreement with the corresponding optimal predictions and indicators (bold lines). As the system size is increased, it becomes increasingly difficult to approximate the corresponding optimal predictions and indicator with high accuracy, because the NN size has to be increased systematically; i.e., hyperparameters need to be adjusted more carefully. However, even for the largest system sizes considered in this work, qualitative agreement can still be achieved with moderate NN sizes; see Appendix B 1 for an explicit example.

*Computational cost.*—Finally, let us touch upon the computational cost of training NNs. Table I reports the measured computation times associated with training an NN with one hidden layer composed of a single node for

one epoch. A training epoch is comprised of evaluating the NN (or NNs in the case of LBC) at all  $M_{\bar{\chi}}$  unique samples (see Table I), calculating the loss function, obtaining the gradient via backpropagation, and performing a single gradient step. This represents a lower bound for the total computation time associated with obtaining NN-based predictions and indicators. In a typical application, however, larger NNs need to be used, the NNs need to be trained for multiple epochs, the NN parameters (or the corresponding predictions and indicator) need to be cached at regular intervals, hyperparameters need to be tuned, and, finally, the indicator needs to be computed based on the NN predictions. The computation time for a single epoch is also expected to increase if the data are processed in a batchwise fashion (albeit likely at the benefit of requiring fewer training epochs overall). We find that this lower bound on the training time is comparable with the evaluation time of the corresponding optimal predictions and indicators (and optimal loss) and the two times differ by less than an order

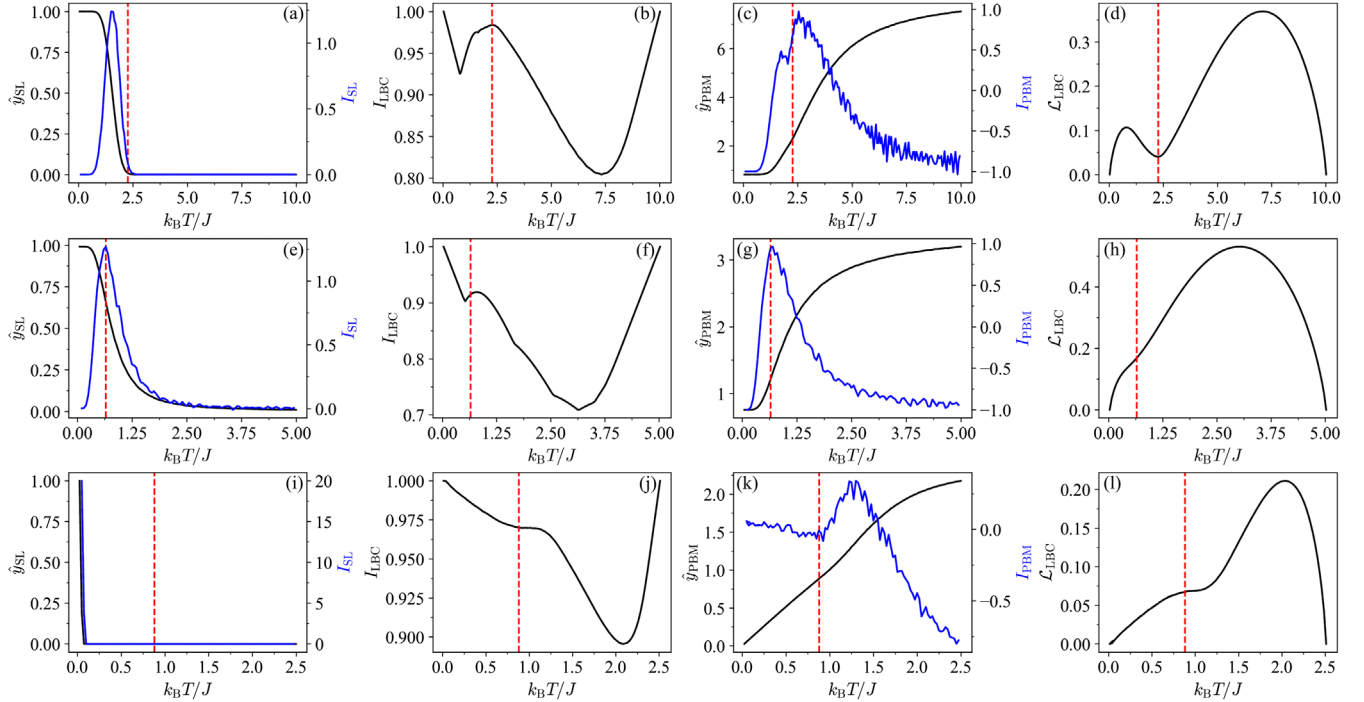


FIG. 11. (a)–(d) Results for the Ising model ( $L = 10$ ) using NNs. The NNs used in SL, LBC, and PBM are trained for 10 000, 1000, and 5000 epochs, respectively. The tuning parameter ranges from  $p_1 = 0.05$  to  $p_K = 10$  with  $\Delta p = 0.05$ . (e)–(h) Results for the IGT ( $L = 4$ ) using NNs. The NNs used in SL, LBC, and PBM are trained for 10 000, 1000, and 5000 epochs, respectively. The tuning parameter ranges from  $p_1 = 0.05$  to  $p_K = 5$  with  $\Delta p = 0.05$ . (i)–(l) Results for the XY model ( $L = 10$ ) using NNs. The NNs used in SL, LBC, and PBM are trained for 10 000, 1000, and 10 000 epochs, respectively. The tuning parameter ranges from  $p_1 = 0.025$  to  $p_K = 2.5$  with  $\Delta p = 0.025$ . The critical value of the tuning parameter  $p_c = k_B T_c/J$  is highlighted in red. (a),(e),(i) Mean prediction  $\hat{y}_{SL}(p)$  obtained using the analytical expression (black solid line) or an NN (black dashed line), as well as the corresponding indicator  $I_{SL}(p)$  (blue lines). Here, we choose  $r_I = 1$  and  $l_{II} = K$ . (b),(f),(j) The indicator of LBC,  $I_{LBC}(p)$ , obtained using the analytical expression (black solid line) or an NN (black dashed line). (c),(g),(k) Mean prediction  $\hat{y}_{PBM}(p)$  of PBM obtained using the analytical expression (black solid line) or an NN (black dashed line), as well as the corresponding indicator  $I_{PBM}(p)$  (blue lines). (d),(h),(l) Value of the loss function in LBC,  $\mathcal{L}_{LBC}$ , for each bipartition point  $p^{\text{bp}}$  obtained using the analytical expression (black solid line) or evaluated after NN training (black dashed line). In all three models, the NNs are comprised of three hidden layers with 64 nodes each, and the learning rate is set to 0.001.

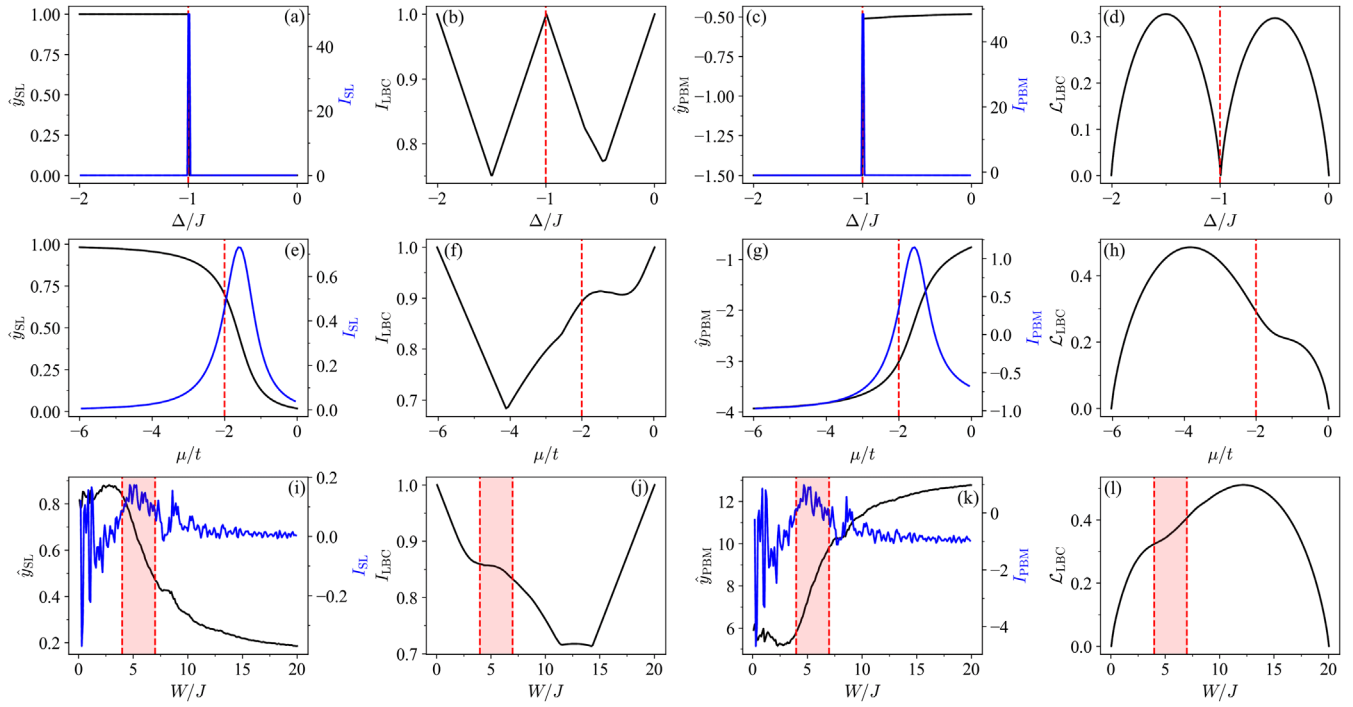


FIG. 12. (a)–(d) Results for the XXZ chain ( $L = 4$ ) using NNs. The NNs used in SL, LBC, and PBM are trained for 10 000, 1000, and 5000 epochs, respectively. The tuning parameter ranges from  $p_1 = -2$  to  $p_K = 0$  with  $\Delta p = 0.01$ . The critical value of the tuning parameter  $p_c = \Delta_c/J$  is highlighted in red. (e)–(h) Results for the Kitaev chain ( $L = 10$ ) using NNs. The NNs used in SL, LBC, and PBM are trained for 5000, 500, and 1000 epochs, respectively. The tuning parameter ranges from  $p_1 = -6$  to  $p_K = 0$  with  $\Delta p = 0.06$ . The critical value of the tuning parameter  $p_c = \mu_c/t$  is highlighted in red. (i)–(l) Results for the many-body localization phase transition in the Bose-Hubbard model ( $L = 6$ ) using NNs. The NNs used in SL, LBC, and PBM are trained for 10 000, 300, and 1000 epochs, respectively. The tuning parameter ranges from  $p_1 = 0.1$  to  $p_K = 20$  with  $\Delta p = 0.1$ . The critical value of the tuning parameter  $p_c = W_c/J$  is highlighted in red. (a),(e),(i) Mean prediction  $\hat{y}_{\text{SL}}(p)$  obtained using the analytical expression (black solid line) or an NN (black dashed line), as well as the corresponding indicator  $I_{\text{SL}}(p)$  (blue lines). Here, we choose  $r_1 = 1$  and  $l_{\text{II}} = K$ . (b),(f),(j) The indicator of LBC,  $I_{\text{LBC}}(p)$ , obtained using the analytical expression (black solid line) or an NN (black dashed line). (c),(g),(k) Mean prediction  $\hat{y}_{\text{PBM}}(p)$  of PBM obtained using the analytical expression (black solid line) or an NN (black dashed line), as well as the corresponding indicator  $I_{\text{PBM}}(p)$  (blue lines). (d),(h),(l) Value of the loss function in LBC,  $\mathcal{L}_{\text{LBC}}$ , for each bipartition point  $p^{\text{bp}}$  obtained using the analytical expression (black solid line) or evaluated after NN training (black dashed line). For the XXZ model, the NNs are comprised of three hidden layers with 64 nodes each. For the Kitaev chain and Bose-Hubbard model, we use two hidden layers with 128 nodes each, followed by three hidden layers with 64 nodes each. In all three cases, the learning rate is set to 0.001.

of magnitude across all six physical systems studied in the main text. This empirical finding can be explained as follows: To construct the optimal model, the probability of all inputs needs to be evaluated. Similarly, in each training epoch, the NN is evaluated at all inputs contained in the training dataset. The computation time associated with evaluating a small NN for a given input is comparable with evaluating the corresponding optimal model prediction, and the overhead associated with the gradient computation via backpropagation is of the same order of magnitude as the NN forward pass [139].

Suppose one is interested in the predictions and indicators of SL, PBM, and LBC, in the limit of a perfectly trained, highly expressive NNs. Evidently, based on the discussion above, the evaluation of the analytical expressions is generally more efficient in that case. The precise

timings depend on the particular implementation, as well as the choice of hyperparameters. However, even in the case where small NNs are trained for short times, the computation time associated with constructing and evaluating an optimal model is *at worst* comparable. Here, we neglect any overhead associated with constructing probability distributions based on drawn samples. In principle, when using NN one does not rely on the estimated probability distributions; i.e., one can directly work with the unprocessed dataset. Note, however, that in many scenarios (including this work) the overhead of estimated probability distributions from the dataset is negligible. When studying quantum systems using exact diagonalization, one has direct access to the underlying probability distributions. Similarly, when performing Monte Carlo studies, the energy statistics are readily available.

### 1. Controlling model capacity

Here, we investigate the effect of NN size, training time, and  $\ell_2$  regularization on the NN-based predictions and indicators and compare them with the corresponding optimal predictions and indicators. All three factors influence the capacity of the resulting model and, thus, determine its ability to approximate the optimal predictive model realizing the global minimum of the loss function corresponding to the optimal predictions and indicators [40,41]. As pointed out in the main text (see Sec. IV), there are instances where the optimal model does not correctly highlight the corresponding phase transition, whereas simpler models do.

As an example, let us consider the application of PBM to the Ising model. Figure 13 shows the results for a  $60 \times 60$  lattice obtained with NNs composed of a single hidden layer with a variable number of hidden nodes ranging from 2 to 2048. Figures 13(b) and 13(f) show the corresponding NN-based predictions and indicators after training for 10 000 epochs. For NNs with two and eight nodes, the indicator shows a clear peak at the critical value of the tuning parameter. As the number of nodes increases, the NN results start to resemble the optimal predictions and indicators (black) more closely. This reflects the fact that the expressivity of an NN increases as the number of nodes is increased. A similar behavior is also visible in Fig. 13(a), which shows the loss over time, where NNs with more than eight nodes achieve values close to the optimal loss (black), i.e., the global minimum.

Figures 13(c) and 13(g) show the predictions and indicators for the smallest NN (two hidden nodes) evaluated at various training epochs. Here, the indicator gradually converges toward its final form, which exhibits a peak at the critical value of the tuning parameter. Similarly, Figs. 13(d) and 13(h) shows the results for the largest NN (2048 hidden nodes). Here, early on during training the indicator is sharply peaked near the critical value of the tuning parameter. As the training progresses, the indicator signal starts to wash out and converge to the optimal indicator signal. The evolution of the global maximum of the indicator signal as a function of the training epoch for the various NN sizes is shown in Fig. 13(e). These results quantify how accurately the estimated critical value of the tuning parameter based on the optimal indicator (black) is reproduced for a given NN size and training time.

Figure 13(e) shows that even for the large NNs there seems to be an intermediate time period during training where the indicator peaks near the critical value of the tuning parameter correctly highlighting the phase transition. Looking at Fig. 13(a), during these intermediate time periods, the corresponding loss function starts to saturate and display a kink. This suggests a procedure for early stopping, where the training is stopped once a kink in the loss function is observed [40]. Early stopping based on the validation loss is discussed in the subsequent section (see Appendix B 2). During training, the model capacity increases as visible by the steady decrease in the

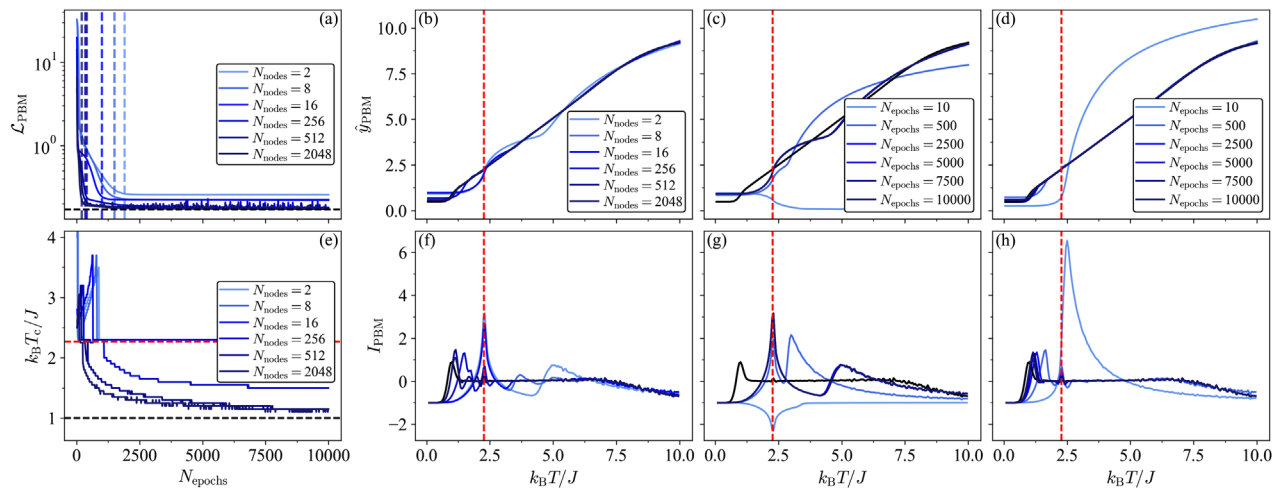


FIG. 13. Results for the Ising model ( $L = 60$ ) of PBM using NNs with a single hidden layer composed of different numbers of hidden nodes  $N_{\text{nodes}}$ . The learning rate is set to 0.01. The tuning parameter ranges from  $p_1 = 0.05$  to  $p_K = 10$  with  $\Delta p = 0.05$ . The critical value of the tuning parameter  $p_c = k_B T_c / J$  is highlighted in red. The optimal predictions, optimal indicator, optimal loss, and corresponding estimated critical value of the tuning parameter are highlighted in black. (a) Loss  $\mathcal{L}_{\text{PBM}}$  as a function of the number of training epochs  $N_{\text{epochs}}$ . The locations in kinks in the loss (as identified by eye) are marked by vertical dashed lines. (b),(f) Mean prediction  $\hat{y}_{\text{PBM}}(p)$  of PBM obtained using NNs after training for 10 000 epochs, as well as the corresponding indicator  $I_{\text{PBM}}(p)$ . (c),(g) Mean prediction  $\hat{y}_{\text{PBM}}(p)$  of PBM obtained using an NN with  $N_{\text{nodes}} = 2$  at various stages during training, as well as the corresponding indicator  $I_{\text{PBM}}(p)$ . (d),(h) Mean prediction  $\hat{y}_{\text{PBM}}(p)$  of PBM obtained using an NN with  $N_{\text{nodes}} = 2048$  at various stages during training, as well as the corresponding indicator  $I_{\text{PBM}}(p)$ . (e) Estimated critical value of the tuning parameter as a function of the number of training epochs.

corresponding loss [40,140,141]: Initially, the model cannot resolve anything; in the intermediate stages, it can resolve between the two phases leading to the sharp peak; and, eventually, it approaches the optimal predictive model (which, in this case, does not correctly highlight the phase transition). By stopping the training at the intermediate stage (i.e., selecting the corresponding NN parameters after the training is complete), a model of intermediate resolution can be obtained. Thus, early stopping acts as an implicit regularization [40,140,141]. In the case of PBM, stopping the training early yields an NN whose indicator peaks near the critical temperature of the Ising model. However, this is not always the case. In LBC, for example, the estimated critical temperature gradually improves during training, i.e., as the model capacity increases. Recall that the optimal indicator of LBC correctly highlights the phase transition. Qualitatively similar results can be obtained for the other methods and systems. In particular, in the Ising model and XY model, we find that the indicators of SL and PBM both show a clear peak near the critical transition temperature early on during training around the epochs marked by a kink in the loss function. The peak locations of the corresponding NN-based indicator signals coincide with the signals of physical indicators, such as the magnetization or heat capacity.

Lastly, we can also control the capacity of our model through explicit  $\ell_2$  regularization [40]

$$\mathcal{L} \rightarrow \mathcal{L} + \lambda_{\ell_2} \sum_i \theta_i^2, \quad (\text{B3})$$

where the sum runs over all tunable parameters  $\theta_i$  of the NN and  $\lambda_{\ell_2}$  is the regularization strength. Figure 14 shows the NN-based predictions and indicators of PBM for the Ising model after training with various regularization strengths. At large regularization strength, the resulting model cannot resolve any structure leading to a flat indicator signal. At an intermediate regularization strength, the resulting model can distinguish between the two phases,

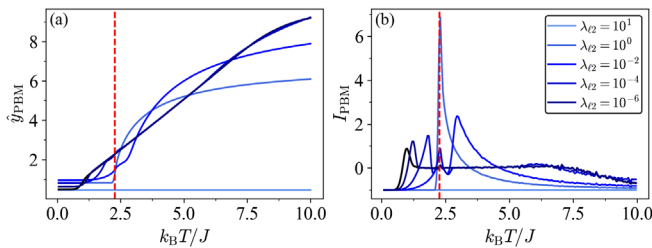


FIG. 14. (a) Mean prediction  $\hat{y}_{\text{PBM}}(p)$  and (b) the corresponding indicator  $I_{\text{PBM}}(p)$  of PBM for the Ising model ( $L = 60$ ) using NNs obtained after long training for various regularization strengths  $\lambda_{\ell_2}$  [cf. Eq. (B3)]. The tuning parameter ranges from  $p_1 = 0.05$  to  $p_K = 10$  with  $\Delta p = 0.05$ . The critical value of the tuning parameter  $p_c = k_B T_c / J$  is highlighted in red. The optimal predictions and indicator are highlighted in black. Each NN has a single hidden layer with 2048 nodes and is trained for 10 000 epochs with a learning rate of 0.01.

leading to a clear peak in the indicator signal at the critical temperature of the Ising model. As the regularization strength is decreased further, the resulting model becomes more complex and converges toward the optimal model that minimizes the loss function in the absence of regularization. Consequently, the predictions and indicators converge toward the optimal predictions and indicator. In the Ising model, we thus find that explicit regularization helps to construct a model of intermediate resolution whose indicator correctly highlights the critical temperature (similarly for SL). However, as mentioned above, models with restricted capacity may not always highlight the critical value of the tuning parameter correctly. In the IGT, for example, the indicator of regularized NNs tends to display an erroneous peak similar to the specific heat; see Fig. 4.

## 2. Finite-sample statistics: Splitting data into training, validation, and test sets

Here, we investigate NN-based predictions and indicators in the case where only a limited amount of data is available. In particular, we discuss the effect of splitting the data into a training, validation, and test set. Recall that, in the limit of sufficient data, the training, validation, and test set coincide, as they are all sampled independently from the same probability distribution underlying the physical system; see Sec. II. Therefore, in the limit of sufficient data, the training, validation, and test losses decrease in lock step during training. This is illustrated in Fig. 15, which shows the training, validation, and test loss of PBM for the Kitaev chain for different dataset sizes. For small datasets, the training, validation, and test sets can differ, resulting in differing training, validation, and test losses. In particular, one can observe a characteristic increase of the validation loss after a certain time period attributed to overfitting [40]. This allows one to perform early stopping such that the

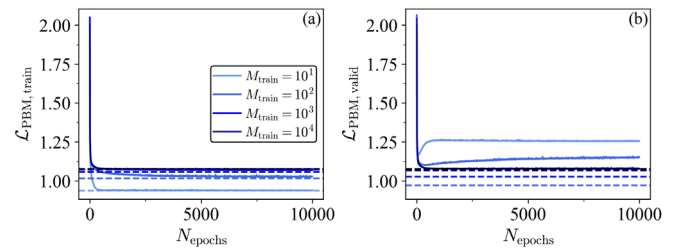


FIG. 15. (a) Training loss and (b) validation loss as a function of the number of training epochs of PBM for the Kitaev chain ( $L = 14$ ) using an NN composed of a single hidden layer with 128 nodes for various numbers of training samples  $M_{\text{train}}$  per parameter value, where  $M_{\text{valid}} = M_{\text{test}} = M_{\text{train}}/5$ . The corresponding optimal loss based on the training or validation dataset is highlighted by a colored dashed line. The optimal loss based on the ground-truth probability distributions is highlighted in black. The test loss shows the same behavior as the validation loss. Each NN is trained for 10 000 epochs with a learning rate of 0.01. The results are averaged over ten independent datasets.

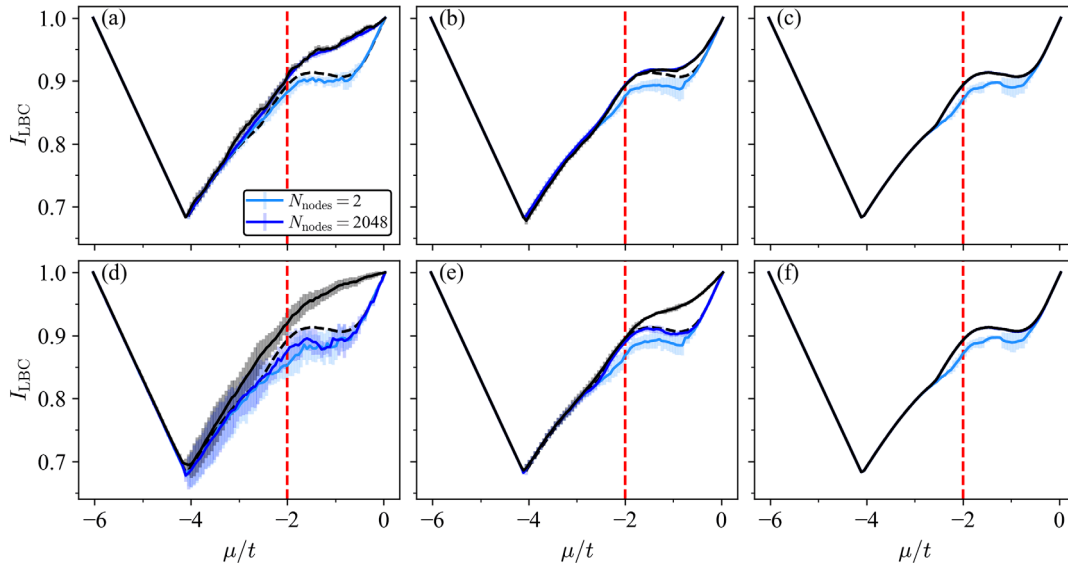


FIG. 16. Results of LBC for the Kitaev chain ( $L = 10$ ) using NNs composed of a single hidden layer with 2 or 2048 nodes for various numbers of training samples  $M_{\text{train}}$  per parameter value, where  $M_{\text{valid}} = M_{\text{test}} = M_{\text{train}}/5$ . The tuning parameter  $p = \mu/t$  ranges from  $p_1 = -6$  to  $p_K = 0$  with  $\Delta p = 0.06$ . The critical value  $\mu_c/t = -2$  is highlighted by a red dashed line. The optimal indicator obtained based on the corresponding dataset or the ground-truth probability distributions is highlighted by a black solid or dashed line, respectively. (a)–(c) Indicator  $I_{\text{LBC}}$  of LBC evaluated on the training set for (a)  $M_{\text{train}} = 10$ , (b)  $M_{\text{train}} = 10^2$ , and (c)  $M_{\text{train}} = 10^5$ , where the NN-based predictions are obtained after training. (d)–(f) Indicator  $I_{\text{LBC}}$  of LBC evaluated on the test set for (a)  $M_{\text{train}} = 10$ , (b)  $M_{\text{train}} = 10^2$ , and (c)  $M_{\text{train}} = 10^5$ , where early stopping is performed by minimizing the validation loss. Similar results are obtained when evaluating the NNs at the end of training instead. Each NN is trained for 10 000 epochs with a learning rate of 0.005. The results are averaged over ten independent datasets, and the error bars are given by the standard deviation.

minimum in the validation loss is realized [40]. Note that the location of the minima in the validation loss coincides with the kink in the corresponding training loss. The sharp local minimum in the validation loss fades as the dataset size is increased further, leaving only the corresponding kinks in the training loss as a signal for early stopping. The latter situation is discussed in Appendix B 1. Therefore, a splitting into training, validation, and test set may allow for a clearer signal to perform early stopping given a small dataset.

Another effect arising when a limited amount of data is available and finite-sample statistics play a role is best illustrated by investigating the Kitaev chain using LBC. Figure 16 shows the NN-based indicator signal of LBC obtained for training, test, and validation sets of various sizes. For small dataset sizes [see Figs. 16(a) and 16(d)] the optimal indicator (black solid line) shows no local maximum due to the negligible overlap in the inferred probability distribution. The NN-based indicator of a sufficiently large NN closely matches the optimal indicator on the training set after training [Fig. 16(a)], whereas a small NN is incapable of approximating the optimal indicator on the training set. However, interestingly, the indicator signal of the small NN qualitatively matches the optimal indicator signal based on the ground-truth probability distributions. In particular, it features a local maximum allowing for an estimate of the critical value of the tuning parameter to be obtained. This is another example illustrating how simple

models can lead to sharp indicator signals. While the inferred probability distribution has only a marginal overlap in the topological phase resulting in the absence of a local maximum in the optimal indicator signal (black), the data may be partially indistinguishable to a simple model. This illustrates how “confusion” can also arise due to models with restricted expressivity (see Sec. III). The same phenomenon can also be observed for the indicator signal of the large NN evaluated on the test set (or validation set); see Fig. 16(d). Here, the confusion arises because the predictions for the unseen data within the validation and test set are suboptimal. In the future, it will be of interest to investigate whether this effect can be mimicked through appropriate interpolation of the optimal predictions [31,35,132]. Figures 16(b) and 16(e) and Figs. 16(c) and 16(f) show how the discrepancy between the optimal indicator signal based on a finite dataset and the NN-based indicator vanishes for the large NN as the dataset size increases. This arises because, eventually, the training, validation, and test sets become indistinguishable. Note, however, that the discrepancy persists for the small NN.

### APPENDIX C: DATA GENERATION

In this appendix, we provide further details on the data-generation process for each of the physical systems analyzed in the main text (see Sec. IV). For the classical systems, given by the Ising model, IGT, and XY model,

we use the Metropolis-Hastings algorithm [72] to sample spin configurations from the thermal distribution at a given temperature  $T$ . The lattice is initialized in a state with all spins pointing up for the Ising model and a random spin configuration in the case of the IGT and XY model. The lattice is updated by drawing a random spin, which is flipped with probability  $\min(1, e^{-\Delta E/T})$ , where  $\Delta E$  is the energy difference resulting from the considered flip. In the XY model, instead of flipping a given spin, we add a perturbation  $\Delta\theta \in [-\pi, \pi]$ , which is drawn uniformly at random. To ensure that the systems are sufficiently thermalized, we sweep the complete lattice  $10^5$  times, where each lattice site is updated once per sweep. After the thermalization period, we collect  $10^5$  samples, which we find to be sufficient for achieving convergence (see Appendix A 5). In the Ising model and IGT, we increase the temperature gradually, whereas it is decreased in the XY model.

In the XY model, we can further validate the quality of the Monte Carlo samples by estimating the BKT transition point. One way to do this is to determine the temperature at which the helicity modulus  $\Upsilon(T)$  crosses  $2T/\pi$  [94,96]. The helicity modulus is also referred to as spin stiffness or spin rigidity and measures the response of the system to an in-plane twist of the spins. We find that the estimated BKT transition point based on our samples matches the literature value well; see Fig. 17. Note that, in the XY model, the angle of each spin can take on any value  $\theta \in [0, 2\pi]$ . This results in a continuum of states. Hence, we discretize the energy in practice, which serves as an input for the ML methods. This discretization eases computation and, more

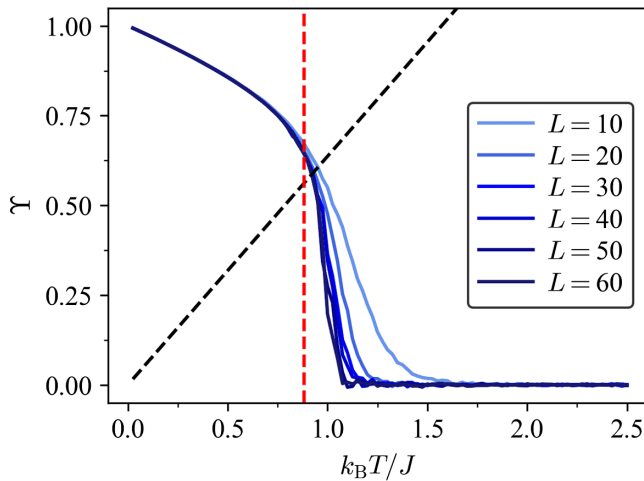


FIG. 17. Helicity modulus  $\Upsilon$  as a function of the tuning parameter  $p = k_B T/J$  for the two-dimensional XY model for various lattice sizes. The value of the BKT transition point from the literature,  $k_B T_c/J \approx 0.8935$  [92], is highlighted by a red dashed line. The estimated transition point based on our Monte Carlo samples at finite size corresponds to the point at which the helicity modulus crosses the line given by  $2k_B T/J\pi$  (black dashed line).

crucially, results in overlapping probability distributions given finite-sample statistics (see discussion on case 3 in Sec. III A). The discretization is performed through simple histogram binning using 1000 bins of equal size. The number of bins is increased systematically until a convergence of the optimal indicator signals is observed. In future works, histogram binning may be replaced by more elaborate techniques for density estimation [118–121].

Let us move on to the quantum case. To perform exact diagonalization and solve the Schrödinger equation, we use the QuSpin package [142,143] in PYTHON. Note that, when computing the ground state of the Kitaev chain through exact diagonalization, we restrict ourselves to the even-particle sector whose corresponding ground state has a lower energy within the topologically trivial phase. In the topological phase, the ground state is doubly degenerate, and the two states can be distinguished by their fermionic parity. This is because of the presence of the pairing term in the Kitaev chain Hamiltonian [Eq. (37)]. As a consequence,  $H$  does not conserve the total fermion number  $N_f = \sum_{i=1}^L n_i$ , i.e.,  $[H, N_f] \neq 0$ . However, the fermion number modulo 2 is conserved,  $[H, (-1)^{N_f}] = 0$  [144].

## APPENDIX D: COMPARISON TO OTHER WORKS

In this appendix, we provide additional material which facilitates comparison to other works.

### 1. Alternative approach toward supervised learning

Here, we review our approach to SL (see Sec. II A) and put it into context. In Ref. [4], the authors originally proposed to identify the estimated critical value of the tuning parameter in SL as  $\text{argmin}_{p_k} |\hat{y}(p_k) - 0.5|$ . In all systems analyzed in the main text (see Sec. IV), this yields similar results compared to our approach based on identifying the peak location of the mean prediction's derivative [Eq. (3)]. Note that the latter approach has, e.g., already been mentioned as an alternative in Ref. [20]. Looking at Fig. 8(b), we observe that these two procedures would yield slightly different estimated critical values for the MBL phase transition. This discrepancy is even more prominent for the Mott insulator to superfluid transition in the Bose-Hubbard model. Here, we investigate the two-dimensional Bose-Hubbard model whose Hamiltonian is given by

$$H = -J \sum_{\langle ij \rangle} (b_i^\dagger b_j + \text{H.c.}) + \sum_i \frac{U}{2} n_i (n_i - 1) - \mu n_i, \quad (\text{D1})$$

where  $J$  is the nearest-neighbor hopping strength,  $U$  is the on-site interaction strength, and  $\mu$  is the chemical potential. This model undergoes a quantum phase transition at zero temperature from a Mott insulating phase to a superfluid phase as the tuning parameter  $J/U$  is increased at a fixed



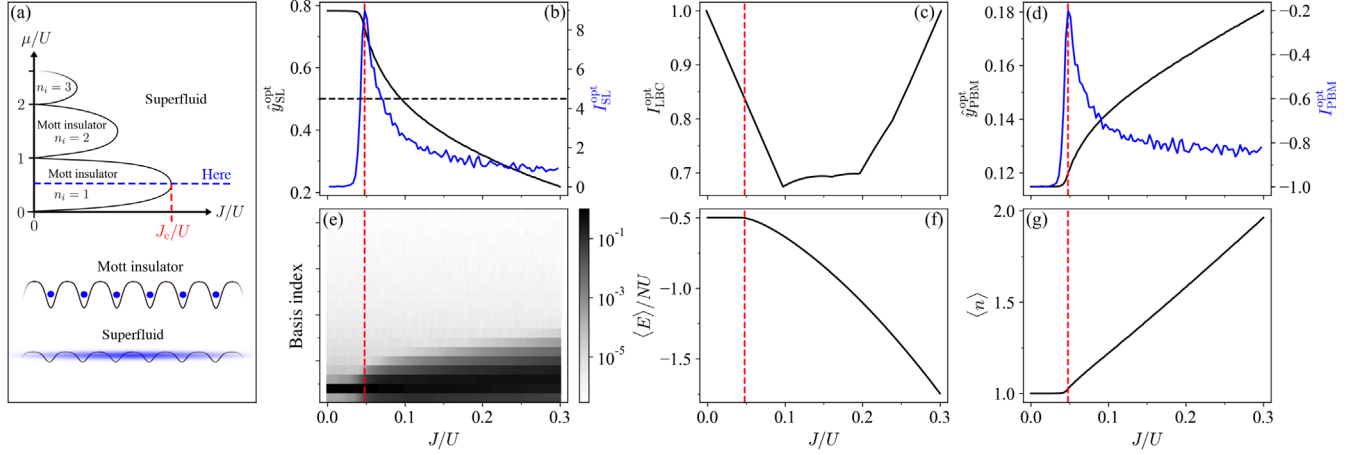


FIG. 18. Results for the Mott insulating to superfluid phase transition in the (two-dimensional) Bose-Hubbard model with the dimensionless coupling strength as a tuning parameter  $p = J/U$  ranging from  $p_1 = 0$  to  $p_K = 0.3$  in steps of  $\Delta p = 0.03$ , where  $\mu/U = 0.5$ . In SL, the data obtained at  $p_1$  and  $p_K$  constitute our training set, i.e.,  $r_I = 1$  and  $l_{II} = K$ . The reference value for the critical value of the tuning parameter  $J_c/U = 1/(5.8z)$  with  $z = 4$  [147] is highlighted by a red dashed line. (a) Illustration of the two-dimensional phase diagram of the Bose-Hubbard model containing three Mott lobes. Here, we analyze the quantum phase transition from a Mott insulating state to a superfluid state occurring at the tip of the first Mott lobe ( $\mu/U = 0.5$ ). A sketch of the two distinct phases is shown on the bottom. (b) Mean optimal prediction  $\hat{y}_{\text{SL}}^{\text{opt}}$  in SL (black solid line) and the corresponding indicator  $I_{\text{SL}}^{\text{opt}}$  (blue line). The value  $\hat{y}_{\text{SL}}^{\text{opt}} = 0.5$  is highlighted by a black dashed line. (c) Optimal indicator of LBC,  $I_{\text{LBC}}^{\text{opt}}$  (black line). (d) Mean optimal prediction  $\hat{y}_{\text{PBM}}^{\text{opt}}$  in PBM (black line) and the corresponding indicator  $I_{\text{PBM}}^{\text{opt}}$  (blue line). (e) Probability distributions governing the input data (indices of Fock basis states  $\{|n_i\rangle\}_{i=1}^{n_{\text{max}}}$ ) as a function of the tuning parameter, where the color scale denotes the probability. (f) Average energy per site ( $N = L$  sites in total) as a function of the tuning parameter. Notice the drop in the average energy as the system undergoes the quantum phase transition. (g) Average occupation number per site  $\langle n \rangle$  as a function of the tuning parameter.

chemical potential. This gives rise to the characteristic Mott lobes [145,146]; see Fig. 18(a).

We perform mean-field calculations based on a Gutzwiller ansatz in which the ground-state wave function is written as a product state

$$|\Psi_{\text{MF}}\rangle = \prod_i |\phi_i\rangle \quad (\text{D2})$$

with

$$|\phi_i\rangle = \sum_{n=0}^{n_{\text{max}}} f_n |n_i\rangle, \quad (\text{D3})$$

where  $|n_i\rangle$  denotes the Fock state with  $n$  bosons at site  $i$  [148]. We minimize the expectation value of the Hamiltonian with respect to the Gutzwiller coefficients  $\{|f_n|^2\}_{n=0}^{n_{\text{max}}}$  by means of simulated annealing [21,149] with a maximum number of bosons per site of  $n_{\text{max}} = 20$ . Here, the Gutzwiller coefficients  $\{|f_n|^2\}_{n=0}^{n_{\text{max}}}$  represent the relevant probability distributions governing the data. Note that the simulated annealing algorithm can get stuck in local energy minima. To counteract this noise, we average the Gutzwiller coefficients obtained from 500 independent simulated annealing runs.

At the tip of the first Mott lobe ( $\mu/U = 0.5$ ), the phase transition occurs at  $J_c/U = 1/(5.8z)$  [see Fig. 18(a)],

where  $z$  is the coordination number (here,  $z = 4$ ) [147]. The phase transition can be revealed by looking at the average boson number per site  $\langle n \rangle$ ; see Fig. 18(g). The Mott insulator is characterized by an integer density enforced by the Mott energy gap  $\propto U$ . As a result of the energy gap, the Mott insulator is incompressible. In contrast, the superfluid

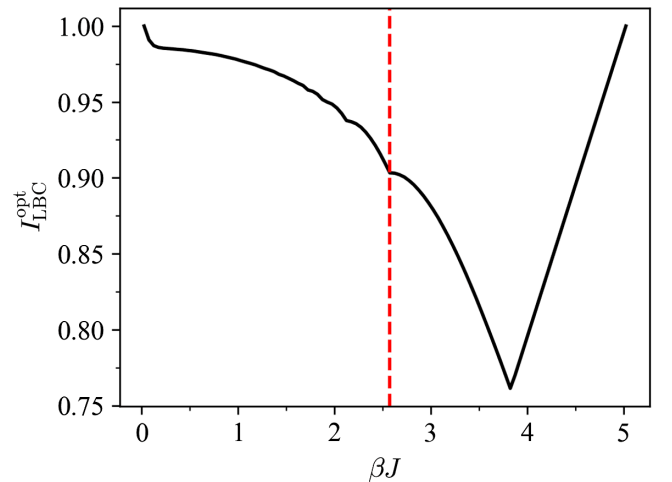


FIG. 19. Optimal indicator of LBC for the IGT ( $L = 12$ ) with dimensionless inverse temperature  $p = \beta J$  as a tuning parameter, where  $p_1 = 0.05$ ,  $p_K = 5$ , and  $\Delta p = 0.05$ . The critical value of the tuning parameter  $p_c = \beta_c J$  from Fig. 4 is highlighted in red.

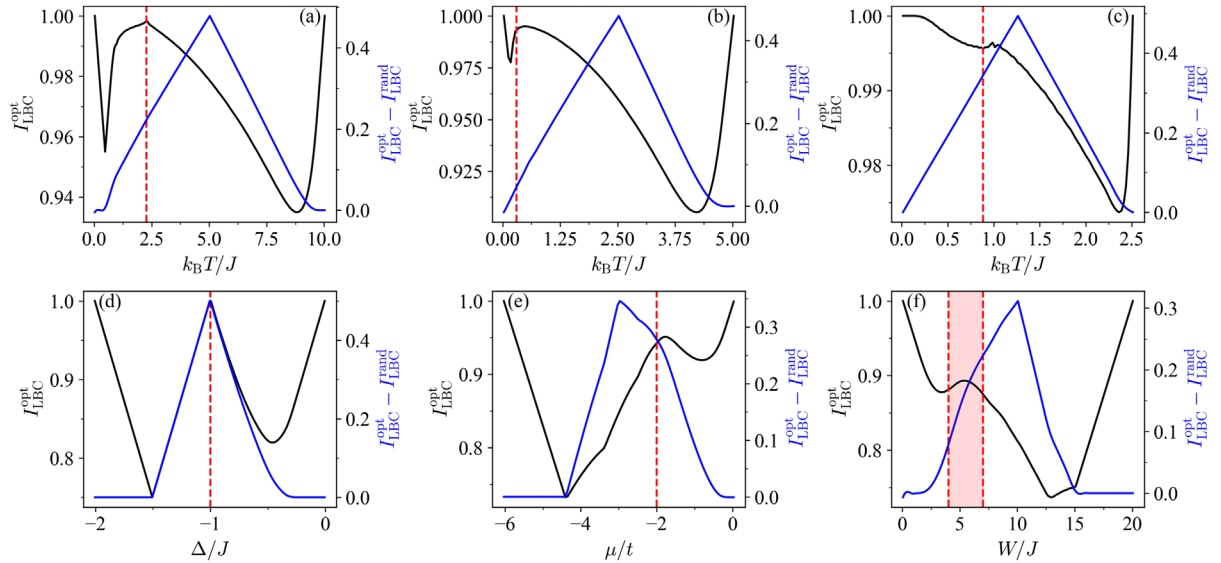


FIG. 20. Optimal indicator of LBC before (black line) and after (blue line) background subtraction for the (a) Ising model ( $L = 60$ ), (b) IGT ( $L = 28$ ), (c) XY model ( $L = 60$ ), (d) XXZ chain ( $L = 14$ ), (e) Kitaev chain ( $L = 20$ ), and (f) Bose-Hubbard model ( $L = 8$ ); see Sec. IV in the main text. The corresponding critical values of the tuning parameters are highlighted in red.

phase is compressible and is characterized by strong number fluctuations (even at low temperature).

Figure 18 shows the results of SL, LBC, and PBM. Here, both SL and PBM correctly identify the quantum phase transition, whereas LBC fails. Looking at Fig. 18(e), we see that a large change in the underlying probability distributions occurs at the quantum phase transition. In Ref. [22], LBC with NNs is shown to correctly highlight the Mott-insulating to superfluid transition in the Bose-Hubbard model. However, in this case, the Gutzwiller coefficients directly serve as input, whereas here the individual Fock basis states (i.e., their indices) constitute the input. Note that the phase transition would not be predicted with a high accuracy using SL if we estimate the predicted critical temperature as the value of the tuning parameter for which  $\hat{y}_{\text{SL}}^{\text{opt}} = 0.5$ ; see the black dashed line in Fig. 18(b). This motivates our approach to SL compared to the procedure originally proposed in Ref. [4]. However, both approaches for obtaining estimated critical values are directly applicable given optimal predictions.

## 2. Analysis of Ising gauge theory

Figure 19 shows the optimal indicator of LBC for the IGT with inverse temperature  $\beta$  as a tuning parameter. The signal qualitatively matches the indicator of LBC reported in Fig. C1 in Ref. [31] obtained with NNs, confirming that for high-capacity models the indicator signal of LBC is indeed ambiguous in this case.

In Ref. [31], the authors also investigate the IGT with PBM using NNs. They empirically find that the NN-based predictions agree well with a physical model based on the underlying density of states, which is proposed in an *ad hoc*

fashion guided by physical intuition. In our work, we explicitly confirm this physical intuition on what the NN learns by proving that the optimal prediction of PBM for a given configuration in the IGT corresponds to the most likely tuning parameter value based on the underlying Boltzmann distribution.

## 3. Background subtraction for learning by confusion

Figure 20 shows the optimal indicator in LBC for all physical systems considered in the main text, as well as a modified version where the V-shaped indicator signal characteristic of indistinguishable data is subtracted. Note that this V-shaped indicator signal is computed separately for each system, i.e., parameter range. For all systems, we find that the modified indicator peaks near the center of the parameter range under consideration, whereas the original indicator signal peaks near the phase transition (red dashed line). This bias arises because the subtracted signal is lowest near the center of the parameter range. As such, the bias can be easily missed if the transition point is indeed located in the center of the chosen parameter range; see Fig. 20(d).

- [1] V. Dunjko and H. J. Briegel, *Machine Learning & Artificial Intelligence in the Quantum Domain: A Review of Recent Progress*, *Rep. Prog. Phys.* **81**, 074001 (2018).
- [2] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, *Machine Learning and the Physical Sciences*, *Rev. Mod. Phys.* **91**, 045002 (2019).

- [3] A. Dawid, J. Arnold, B. Requena, A. Gresch, M. Płodzień, K. Donatella, K. Nicoli, P. Stornati, R. Koch, M. Büttner *et al.*, *Modern Applications of Machine Learning in Quantum Sciences*, [arXiv:2204.04198](#).
- [4] J. Carrasquilla and R. G. Melko, *Machine Learning Phases of Matter*, *Nat. Phys.* **13**, 431 (2017).
- [5] E. P. Van Nieuwenburg, Y.-H. Liu, and S. D. Huber, *Learning Phase Transitions by Confusion*, *Nat. Phys.* **13**, 435 (2017).
- [6] J. Carrasquilla, *Machine Learning for Quantum Matter*, *Adv. Phys. X* **5**, 1797528 (2020).
- [7] J. Carrasquilla and G. Torlai, *How to Use Neural Networks to Investigate Quantum Many-Body Physics*, *PRX Quantum* **2**, 040201 (2021).
- [8] B. S. Rem, N. Käming, M. Tarnowski, L. Asteria, N. Fläschner, C. Becker, K. Sengstock, and C. Weitenberg, *Identifying Quantum Phase Transitions Using Artificial Neural Networks on Experimental Data*, *Nat. Phys.* **15**, 917 (2019).
- [9] N. Käming, A. Dawid, K. Kottmann, M. Lewenstein, K. Sengstock, A. Dauphin, and C. Weitenberg, *Unsupervised Machine Learning of Topological Phase Transitions from Experimental Data*, *Mach. Learn.: Sci. Technol.* **2**, 035037 (2021).
- [10] A. Bohrdt, S. Kim, A. Lukin, M. Rispoli, R. Schittko, M. Knap, M. Greiner, and J. Léonard, *Analyzing Nonequilibrium Quantum States through Snapshots with Artificial Neural Networks*, *Phys. Rev. Lett.* **127**, 150504 (2021).
- [11] C. Miles, R. Samajdar, S. Ebadi, T. T. Wang, H. Pichler, S. Sachdev, M. D. Lukin, M. Greiner, K. Q. Weinberger, and E.-A. Kim, *Machine Learning Discovery of New Phases in Programmable Quantum Simulator Snapshots*, [arXiv:2112.10789](#).
- [12] Y. Yu, L.-W. Yu, W. Zhang, H. Zhang, X. Ouyang, Y. Liu, D.-L. Deng, and L.-M. Duan, *Experimental Unsupervised Learning of Non-Hermitian Knotted Phases with Solid-State Spins*, [arXiv:2112.13785](#).
- [13] L. Wang, *Discovering Phase Transitions with Unsupervised Learning*, *Phys. Rev. B* **94**, 195105 (2016).
- [14] S. J. Wetzel, *Unsupervised Learning of Phase Transitions: From Principal Component Analysis to Variational Autoencoders*, *Phys. Rev. E* **96**, 022140 (2017).
- [15] S. J. Wetzel and M. Scherzer, *Machine Learning of Explicit Order Parameters: From the Ising Model to SU(2) Lattice Gauge Theory*, *Phys. Rev. B* **96**, 184410 (2017).
- [16] K. Ch'ng, J. Carrasquilla, R. G. Melko, and E. Khatami, *Machine Learning Phases of Strongly Correlated Fermions*, *Phys. Rev. X* **7**, 031038 (2017).
- [17] T. Ohtsuki and T. Ohtsuki, *Deep Learning the Quantum Phase Transitions in Random Electron Systems: Applications to Three Dimensions*, *J. Phys. Soc. Jpn.* **86**, 044708 (2017).
- [18] F. Schindler, N. Regnault, and T. Neupert, *Probing Many-Body Localization with Neural Networks*, *Phys. Rev. B* **95**, 245134 (2017).
- [19] Y. Zhang and E.-A. Kim, *Quantum Loop Topography for Machine Learning*, *Phys. Rev. Lett.* **118**, 216401 (2017).
- [20] P. Broecker, J. Carrasquilla, R. G. Melko, and S. Trebst, *Machine Learning Quantum Phases of Matter beyond the Fermion Sign Problem*, *Sci. Rep.* **7**, 8823 (2017).
- [21] P. Huembeli, A. Dauphin, and P. Wittek, *Identifying Quantum Phase Transitions with Adversarial Neural Networks*, *Phys. Rev. B* **97**, 134109 (2018).
- [22] Y.-H. Liu and E. P. L. van Nieuwenburg, *Discriminative Cooperative Networks for Detecting Phase Transitions*, *Phys. Rev. Lett.* **120**, 176401 (2018).
- [23] M. J. S. Beach, A. Golubeva, and R. G. Melko, *Machine Learning Vortices at the Kosterlitz-Thouless Transition*, *Phys. Rev. B* **97**, 045207 (2018).
- [24] E. van Nieuwenburg, E. Bairey, and G. Refael, *Learning Phase Transitions from Dynamics*, *Phys. Rev. B* **98**, 060301(R) (2018).
- [25] P. Zhang, H. Shen, and H. Zhai, *Machine Learning Topological Invariants with Neural Networks*, *Phys. Rev. Lett.* **120**, 066401 (2018).
- [26] J. Venderley, V. Khemani, and E.-A. Kim, *Machine Learning Out-of-Equilibrium Phases of Matter*, *Phys. Rev. Lett.* **120**, 257204 (2018).
- [27] J. F. Rodriguez-Nieva and M. S. Scheurer, *Identifying Topological Order through Unsupervised Machine Learning*, *Nat. Phys.* **15**, 790 (2019).
- [28] P. Huembeli, A. Dauphin, P. Wittek, and C. Gogolin, *Automated Discovery of Characteristic Features of Phase Transitions in Many-Body Localization*, *Phys. Rev. B* **99**, 104106 (2019).
- [29] F. Schäfer and N. Lörch, *Vector Field Divergence of Predictive Model Output as Indication of Phase Transitions*, *Phys. Rev. E* **99**, 062107 (2019).
- [30] M. S. Scheurer and R.-J. Slager, *Unsupervised Machine Learning and Band Topology*, *Phys. Rev. Lett.* **124**, 226401 (2020).
- [31] E. Greplova, A. Valenti, G. Boschung, F. Schäfer, N. Lörch, and S. D. Huber, *Unsupervised Identification of Topological Phase Transitions Using Predictive Models*, *New J. Phys.* **22**, 045003 (2020).
- [32] K. Kottmann, P. Huembeli, M. Lewenstein, and A. Acín, *Unsupervised Phase Discovery with Deep Anomaly Detection*, *Phys. Rev. Lett.* **125**, 170603 (2020).
- [33] D. Zvyagintseva, H. Sigurdsson, V. Kozin, I. Iorsh, I. Shelykh, V. Ulyantsev, and O. Kyriienko, *Machine Learning of Phase Transitions in Nonlinear Polariton Lattices*, *Commun. Phys.* **5**, 8 (2022).
- [34] J. Arnold, F. Schäfer, M. Žonda, and A. U. J. Lode, *Interpretable and Unsupervised Phase Classification*, *Phys. Rev. Research* **3**, 033052 (2021).
- [35] H.-Y. Huang, R. Kueng, G. Torlai, V. V. Albert, and J. Preskill, *Provably Efficient Machine Learning for Quantum Many-Body Problems*, [arXiv:2106.12627](#).
- [36] W.-c. Guo and L. He, *Learning Phase Transitions from Regression Uncertainty*, [arXiv:2203.06455](#).
- [37] N. Maskara, M. Buchhold, M. Endres, and E. van Nieuwenburg, *Learning Algorithm Reflecting Universal Scaling Behavior near Phase Transitions*, *Phys. Rev. Research* **4**, L022032 (2022).
- [38] Z. Patel, E. Merali, and S. J. Wetzel, *Unsupervised Learning of Rydberg Atom Array Phase Diagram with Siamese Neural Networks*, [arXiv:2205.04051](#).
- [39] W. Zhang, H. Yang, and N. Wu, *Neural Network Topological Snake Models for Locating General Phase Diagrams*, [arXiv:2205.09699](#).

- [40] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning* (MIT Press, Cambridge, MA, 2016).
- [41] X. Hu, L. Chu, J. Pei, W. Liu, and J. Bian, *Model Complexity of Deep Learning: A Survey*, *Knowl. Inf. Syst.* **63**, 2585 (2021).
- [42] G. Cybenko, *Approximation by Superpositions of a Sigmoidal Function*, *Math. Control Signals Syst.* **2**, 303 (1989).
- [43] K. Hornik, *Approximation Capabilities of Multilayer Feedforward Networks*, *Neural Netw.* **4**, 251 (1991).
- [44] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang, *The Expressive Power of Neural Networks: A View from the Width*, in *Proceedings of Advances in Neural Information Processing Systems (NIPS 2017)*, Vol. 30, edited by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Curran Associates, Inc., Red Hook, NY, 2017).
- [45] D.-X. Zhou, *Universality of Deep Convolutional Neural Networks*, *Appl. Comput. Harmon. Anal.* **48**, 787 (2020).
- [46] A. Krizhevsky, I. Sutskever, and G. E. Hinton, *ImageNet Classification with Deep Convolutional Neural Networks*, in *Proceedings of Advances in Neural Information Processing Systems*, Vol. 25, edited by F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (Curran Associates, Inc., Red Hook, NY, 2012).
- [47] Y. Bengio and O. Delalleau, *On the Expressive Power of Deep Architectures*, in *Algorithmic Learning Theory*, edited by J. Kivinen, C. Szepesvári, E. Ukkonen, and T. Zeugmann (Springer, Berlin, 2011), pp. 18–36.
- [48] M. Raghu, B. Poole, J. Kleinberg, S. Ganguli, and J. Sohl-Dickstein, *On the Expressive Power of Deep Neural Networks*, in *Proceedings of the 34th International Conference on Machine Learning, PMLR*, Vol. 70, edited by D. Precup and Y. W. Teh (PMLR, Sydney, NSW, Australia, 2017), pp. 2847–2854.
- [49] P. Linardatos, V. Papastefanopoulos, and S. Kotsiantis, *Explainable AI: A Review of Machine Learning Interpretability Methods*, *Entropy* **23**, 18 (2021).
- [50] C. Molnar, *Interpretable Machine Learning*, 2nd ed., <https://christophm.github.io/interpretable-ml-book/>.
- [51] P. Ponte and R. G. Melko, *Kernel Methods for Interpretable Machine Learning of Order Parameters*, *Phys. Rev. B* **96**, 205146 (2017).
- [52] W. Zhang, L. Wang, and Z. Wang, *Interpretable Machine Learning Study of the Many-Body Localization Transition in Disordered Quantum Ising Spin Chains*, *Phys. Rev. B* **99**, 054208 (2019).
- [53] J. Greitemann, K. Liu, L. D. C. Jaubert, H. Yan, N. Shannon, and L. Pollet, *Identification of Emergent Constraints and Hidden Order in Frustrated Magnets Using Tensorial Kernel Methods of Machine Learning*, *Phys. Rev. B* **100**, 174408 (2019).
- [54] K. Liu, J. Greitemann, and L. Pollet, *Learning Multiple Order Parameters with Interpretable Machines*, *Phys. Rev. B* **99**, 104410 (2019).
- [55] Y. Zhang, P. Ginsparg, and E.-A. Kim, *Interpreting Machine Learning of Topological Quantum Phase Transitions*, *Phys. Rev. Research* **2**, 023283 (2020).
- [56] C. Casert, T. Viejira, J. Nys, and J. Ryckebusch, *Interpretable Machine Learning for Inferring the Phase Boundaries in a Nonequilibrium System*, *Phys. Rev. E* **99**, 023304 (2019).
- [57] A. Dawid, P. Huembeli, M. Tomza, M. Lewenstein, and A. Dauphin, *Phase Detection with Neural Networks: Interpreting the Black Box*, *New J. Phys.* **22**, 115001 (2020).
- [58] S. Blücher, L. Kades, J. M. Pawłowski, N. Strodthoff, and J. M. Urban, *Towards Novel Insights in Lattice Field Theory with Explainable Machine Learning*, *Phys. Rev. D* **101**, 094507 (2020).
- [59] A. Dawid, P. Huembeli, M. Tomza, M. Lewenstein, and A. Dauphin, *Hessian-Based Toolbox for Reliable and Interpretable Machine Learning in Physics*, *Mach. Learn.: Sci. Technol.* **3**, 015002 (2021).
- [60] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth, *Learnability and the Vapnik-Chervonenkis Dimension*, *J. Assoc. Comput. Mach.* **36**, 929 (1989).
- [61] V. N. Vapnik, *The Nature of Statistical Learning Theory*, 2nd ed. (Springer, New York, 1999).
- [62] L. Devroye, L. Györfi, and G. Lugosi, *The Bayes Error, in A Probabilistic Theory of Pattern Recognition* (Springer, New York, 1996), pp. 9–20.
- [63] L. Devroye, L. Györfi, and G. Lugosi, *Inequalities and Alternate Distance Measures, in A Probabilistic Theory of Pattern Recognition* (Springer, New York, 1996), pp. 21–37.
- [64] L. Fei-Fei, R. Fergus, and P. Perona, *Learning Generative Visual Models from Few Training Examples: An Incremental Bayesian Approach Tested on 101 Object Categories*, in *Proceedings of the 2004 Conference on Computer Vision and Pattern Recognition Workshop* (IEEE, New York, 2004), pp. 178–178.
- [65] Y. LeCun, F. J. Huang, and L. Bottou, *Learning Methods for Generic Object Recognition with Invariance to Pose and Lighting*, in *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (IEEE, New York, 2004), Vol. 2, pp. II–104.
- [66] G. Griffin, A. Holub, and P. Perona, *Caltech-256 object category dataset*, technical report 7694, 2007.
- [67] A. Krizhevsky, *Learning Multiple Layers of Features from Tiny Images*, Master's thesis, University of Toronto, Toronto, Ontario, 2009.
- [68] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, *ImageNet: A Large-Scale Hierarchical Image Database*, in *Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition* (IEEE, New York, 2009), pp. 248–255.
- [69] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, *Imagenet Large Scale Visual Recognition Challenge*, *Int. J. Comput. Vis.* **115**, 211 (2015).
- [70] F. A. Spanhol, L. S. Oliveira, C. Petitjean, and L. Heutte, *A Dataset for Breast Cancer Histopathological Image Classification*, *IEEE Trans. Biomed. Eng.* **63**, 1455 (2016).
- [71] G. James, D. Witten, T. Hastie, and R. Tibshirani, *Statistical Learning, in An Introduction to Statistical Learning: with Applications in R* (Springer, New York, 2013), pp. 15–57.
- [72] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, *Equation of State Calculations by*

- Fast Computing Machines*, *J. Chem. Phys.* **21**, 1087 (1953).
- [73] J. Simon, W. S. Bakr, R. Ma, M. E. Tai, P. M. Preiss, and M. Greiner, *Quantum Simulation of Antiferromagnetic Spin Chains in an Optical Lattice*, *Nature (London)* **472**, 307 (2011).
- [74] H. Bernien, S. Schwartz, A. Keesling, H. Levine, A. Omran, H. Pichler, S. Choi, A. S. Zibrov, M. Endres, M. Greiner *et al.*, *Probing Many-Body Dynamics on a 51-Atom Quantum Simulator*, *Nature (London)* **551**, 579 (2017).
- [75] A. Lukin, M. Rispoli, R. Schittko, M. E. Tai, A. M. Kaufman, S. Choi, V. Khemani, J. Léonard, and M. Greiner, *Probing Entanglement in a Many-Body-Localized System*, *Science* **364**, 256 (2019).
- [76] M. Rispoli, A. Lukin, R. Schittko, S. Kim, M. E. Tai, J. Léonard, and M. Greiner, *Quantum Critical Behaviour at the Many-Body Localization Transition*, *Nature (London)* **573**, 385 (2019).
- [77] P. N. Jepsen, J. Amato-Grill, I. Dimitrova, W. W. Ho, E. Demler, and W. Ketterle, *Spin Transport in a Tunable Heisenberg Model Realized with Ultracold Atoms*, *Nature (London)* **588**, 403 (2020).
- [78] P. N. Jepsen, W. W. Ho, J. Amato-Grill, I. Dimitrova, E. Demler, and W. Ketterle, *Transverse Spin Dynamics in the Anisotropic Heisenberg Model Realized with Ultracold Atoms*, *Phys. Rev. X* **11**, 041054 (2021).
- [79] S. Ebadi, T. T. Wang, H. Levine, A. Keesling, G. Semeghini, A. Omran, D. Bluvstein, R. Samajdar, H. Pichler, W. W. Ho *et al.*, *Quantum Phases of Matter on a 256-Atom Programmable Quantum Simulator*, *Nature (London)* **595**, 227 (2021).
- [80] H.-Y. Huang, R. Kueng, and J. Preskill, *Predicting Many Properties of a Quantum System from Very Few Measurements*, *Nat. Phys.* **16**, 1050 (2020).
- [81] H.-Y. Huang, M. Broughton, J. Cotler, S. Chen, J. Li, M. Mohseni, H. Neven, R. Babbush, R. Kueng, J. Preskill, and J. R. McClean, *Quantum Advantage in Learning from Experiments*, *Science* **376**, 1182 (2022).
- [82] M. A. Nielsen and I. L. Chuang, *Quantum Computation and Quantum Information: 10th Anniversary Edition* (Cambridge University Press, Cambridge, England, 2010).
- [83] J. Carrasquilla, G. Torlai, R. G. Melko, and L. Aolita, *Reconstructing Quantum States with Generative Models*, *Nat. Mach. Intell.* **1**, 155 (2019).
- [84] G. Torlai, B. Timar, E. P. L. van Nieuwenburg, H. Levine, A. Omran, A. Keesling, H. Bernien, M. Greiner, V. Vuletić, M. D. Lukin, R. G. Melko, and M. Endres, *Integrating Neural Networks with a Quantum Simulator for State Reconstruction*, *Phys. Rev. Lett.* **123**, 230504 (2019).
- [85] C. Miles, A. Bohrdt, R. Wu, C. Chiu, M. Xu, G. Ji, M. Greiner, K. Q. Weinberger, E. Demler, and E.-A. Kim, *Correlator Convolutional Neural Networks as an Interpretable Architecture for Image-like Quantum Matter Data*, *Nat. Commun.* **12**, 1 (2021).
- [86] L. Onsager, *Crystal Statistics. I. A Two-Dimensional Model with an Order-Disorder Transition*, *Phys. Rev.* **65**, 117 (1944).
- [87] C. Castelnovo and C. Chamon, *Entanglement and Topological Entropy of the Toric Code at Finite Temperature*, *Phys. Rev. B* **76**, 184442 (2007).
- [88] F. J. Wegner, *Duality in Generalized Ising Models and Phase Transitions without Local Order Parameters*, *J. Math. Phys. (N.Y.)* **12**, 2259 (1971).
- [89] J. B. Kogut, *An Introduction to Lattice Gauge Theory and Spin Systems*, *Rev. Mod. Phys.* **51**, 659 (1979).
- [90] J. M. Kosterlitz and D. J. Thouless, *Ordering, Metastability and Phase Transitions in Two-Dimensional Systems*, *J. Phys. C* **6**, 1181 (1973).
- [91] J. Kosterlitz, *The Critical Properties of the Two-Dimensional XY Model*, *J. Phys. C* **7**, 1046 (1974).
- [92] Y.-D. Hsieh, Y.-J. Kao, and A. W. Sandvik, *Finite-Size Scaling Method for the Berezinskii-Kosterlitz-Thouless Transition*, *J. Stat. Mech.* (2013) P09001.
- [93] P. M. Chaikin and T. C. Lubensky, *Principles of Condensed Matter Physics* (Cambridge University Press, Cambridge, England, 1995).
- [94] J. E. Van Himbergen and S. Chakravarty, *Helicity Modulus and Specific Heat of Classical XY Model in Two Dimensions*, *Phys. Rev. B* **23**, 359 (1981).
- [95] S. G. Chung, *Essential Finite-Size Effect in the Two-Dimensional XY Model*, *Phys. Rev. B* **60**, 11761 (1999).
- [96] P. Minnhagen and B. J. Kim, *Direct Evidence of the Discontinuous Character of the Kosterlitz-Thouless Jump*, *Phys. Rev. B* **67**, 172509 (2003).
- [97] U. Schollwöck, J. Richter, D. J. Farnell, and R. F. Bishop, *Quantum Magnetism* (Springer, Berlin, 2008), Vol. 645.
- [98] F. Franchini, *An Introduction to Integrable Techniques for One-Dimensional Quantum Systems* (Springer, Cham, 2017).
- [99] A. Y. Kitaev, *Unpaired Majorana Fermions in Quantum Wires*, *Phys. Usp.* **44**, 131 (2001).
- [100] J. Alicea, *New Directions in the Pursuit of Majorana Fermions in Solid State Systems*, *Rep. Prog. Phys.* **75**, 076501 (2012).
- [101] F. Wilczek, *Majorana Returns*, *Nat. Phys.* **5**, 614 (2009).
- [102] L. Amico, R. Fazio, A. Osterloh, and V. Vedral, *Entanglement in Many-Body Systems*, *Rev. Mod. Phys.* **80**, 517 (2008).
- [103] A. Pal and D. A. Huse, *Many-Body Localization Phase Transition*, *Phys. Rev. B* **82**, 174411 (2010).
- [104] V. Khemani, S. P. Lim, D. N. Sheng, and D. A. Huse, *Critical Properties of the Many-Body Localization Transition*, *Phys. Rev. X* **7**, 021013 (2017).
- [105] F. Alet and N. Laflorencie, *Many-Body Localization: An Introduction and Selected Topics*, *C. R. Phys.* **19**, 498 (2018).
- [106] J. Greitemann, K. Liu, and L. Pollet, *Probing Hidden Spin Order with Interpretable Machine Learning*, *Phys. Rev. B* **99**, 060404(R) (2019).
- [107] K. Kottmann, P. Corboz, M. Lewenstein, and A. Acín, *Unsupervised Mapping of Phase Diagrams of 2D Systems from Infinite Projected Entangled-Pair States via Deep Anomaly Detection*, *SciPost Phys.* **11**, 25 (2021).
- [108] A. Bohrdt, C. S. Chiu, G. Ji, M. Xu, D. Greif, M. Greiner, E. Demler, F. Grusdt, and M. Knap, *Classifying Snapshots of the Doped Hubbard Model with Machine Learning*, *Nat. Phys.* **15**, 921 (2019).
- [109] Y. Zhang, A. Mesaros, K. Fujita, S. Edkins, M. Hamidian, K. Ch'ng, H. Eisaki, S. Uchida, J. S. Davis, E. Khatami

- et al.*, *Machine Learning in Electronic-Quantum-Matter Imaging Experiments*, *Nature (London)* **570**, 484 (2019).
- [110] S. Pilati and P. Pieri, *Supervised Machine Learning of Ultracold Atoms with Speckle Disorder*, *Sci. Rep.* **9**, 5613 (2019).
- [111] S. Ghosh, M. Matty, R. Baumbach, E. D. Bauer, K. A. Modic, A. Shekhter, J. Mydosh, E.-A. Kim, and B. Ramshaw, *One-Component Order Parameter in URu<sub>2</sub>Si<sub>2</sub> Uncovered by Resonant Ultrasound Spectroscopy and Machine Learning*, *Sci. Adv.* **6**, eaaz4074 (2020).
- [112] T. Szoldra, P. Sierant, K. Kottmann, M. Lewenstein, and J. Zakrzewski, *Detecting Ergodic Bubbles at the Crossover to Many-Body Localization Using Neural Networks*, *Phys. Rev. B* **104**, L140202 (2021).
- [113] M. Gavreev, A. Mastiukova, E. Kiktenko, and A. Fedorov, *Learning Entanglement Breakdown as a Phase Transition by Confusion*, [arXiv:2202.00348](https://arxiv.org/abs/2202.00348).
- [114] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Babbush, and H. Neven, *Barren Plateaus in Quantum Neural Network Training Landscapes*, *Nat. Commun.* **9**, 4812 (2018).
- [115] T. Vieijra, C. Casert, J. Nys, W. De Neve, J. Haegeman, J. Ryckebusch, and F. Verstraete, *Restricted Boltzmann Machines for Quantum States with Non-Abelian or Anyonic Symmetries*, *Phys. Rev. Lett.* **124**, 097201 (2020).
- [116] M. Bukov, M. Schmitt, and M. Dupont, *Learning the Ground State of a Non-Stoquastic Quantum Hamiltonian in a Rugged Neural Network Landscape*, *SciPost Phys.* **10**, 147 (2021).
- [117] A. Valenti, E. Greplova, N. H. Lindner, and S. D. Huber, *Correlation-Enhanced Neural Networks as Interpretable Variational Quantum States*, *Phys. Rev. Research* **4**, L012010 (2022).
- [118] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)* (Springer, Berlin, 2006).
- [119] D. Wu, L. Wang, and P. Zhang, *Solving Statistical Mechanics Using Variational Autoregressive Networks*, *Phys. Rev. Lett.* **122**, 080602 (2019).
- [120] R. G. Melko, G. Carleo, and J. Carrasquilla, and J. I. Cirac, *Restricted Boltzmann Machines in Quantum Physics*, *Nat. Phys.* **15**, 887 (2019).
- [121] K. A. Nicoli, C. J. Anders, L. Funcke, T. Hartung, K. Jansen, P. Kessel, S. Nakajima, and P. Stornati, *Estimation of Thermodynamic Observables in Lattice Field Theories with Deep Generative Models*, *Phys. Rev. Lett.* **126**, 032001 (2021).
- [122] A. Smith, M. Kim, F. Pollmann, and J. Knolle, *Simulating Quantum Many-Body Dynamics on a Current Digital Quantum Computer*, *npj Quantum Inf.* **5**, 106 (2019).
- [123] F. Barratt, J. Dborin, M. Bal, V. Stojevic, F. Pollmann, and A. G. Green, *Parallel Quantum Simulation of Large Systems on Small NISQ Computers*, *npj Quantum Inf.* **7**, 79 (2021).
- [124] K. Satzinger, Y.-J. Liu, A. Smith, C. Knapp, M. Newman, C. Jones, Z. Chen, C. Quintana, X. Mi, A. Dunsworth *et al.*, *Realizing Topologically Ordered States on a Quantum Processor*, *Science* **374**, 1237 (2021).
- [125] J. Herrmann, S. M. Lima, A. Remm, P. Zapletal, N. A. McMahon, C. Scarato, F. Swiadek, C. K. Andersen, C. Hellings, S. Krinner *et al.*, *Realizing Quantum Convolutional Neural Networks on a Superconducting Quantum Processor to Recognize Quantum Phases*, [arXiv:2109.05909](https://arxiv.org/abs/2109.05909).
- [126] C. Noel, P. Niroula, D. Zhu, A. Risinger, L. Egan, D. Biswas, M. Cetina, A. V. Gorshkov, M. J. Gullans, D. A. Huse *et al.*, *Measurement-Induced Quantum Phases Realized in a Trapped-Ion Quantum Computer*, *Nat. Phys.* **18**, 760 (2022).
- [127] G. Semeghini, H. Levine, A. Keesling, S. Ebadi, T. T. Wang, D. Bluvstein, R. Verresen, H. Pichler, M. Kalinowski, R. Samajdar *et al.*, *Probing Topological Spin Liquids on a Programmable Quantum Simulator*, *Science* **374**, 1242 (2021).
- [128] P. Scholl, M. Schuler, H. J. Williams, A. A. Eberharter, D. Barredo, K.-N. Schymik, V. Lienhard, L.-P. Henry, T. C. Lang, T. Lahaye *et al.*, *Quantum Simulation of 2D Antiferromagnets with Hundreds of Rydberg Atoms*, *Nature (London)* **595**, 233 (2021).
- [129] E. Altman *et al.*, *Quantum Simulators: Architectures and Opportunities*, *PRX Quantum* **2**, 017003 (2021).
- [130] J. Arnold and F. Schäfer, *Replacing Neural Networks by Optimal Analytical Predictors for the Detection of Phase Transitions*, <https://github.com/arnoldjulian/Replacing-neural-networks-by-optimal-analytical-predictors-for-the-detection-of-phase-transitions>.
- [131] S. S. Lee and B. J. Kim, *Confusion Scheme in Machine Learning Detects Double Phase Transitions and Quasi-Long-Range Order*, *Phys. Rev. E* **99**, 043308 (2019).
- [132] A. Jacot, F. Gabriel, and C. Hongler, *Neural Tangent Kernel: Convergence and Generalization in Neural Networks*, in *Proceedings of Advances in Neural Information Processing Systems*, Vol. 31, edited by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Curran Associates, Inc., Red Hook, NY, 2018).
- [133] R. Chartrand, *Numerical Differentiation of Noisy, Nonsmooth Data*, *ISRN Appl. Math.* **2011**, 164564 (2011).
- [134] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, *Efficient Backprop*, in *Neural Networks: Tricks of the Trade* (Springer, New York, 2012), pp. 9–48.
- [135] M. Innes, *Flux: Elegant Machine Learning with Julia*, *J. Open Source Software* **3**, 602 (2018).
- [136] D. Kingma and J. Ba, *Adam: A Method for Stochastic Optimization*, [arXiv:1412.6980](https://arxiv.org/abs/1412.6980).
- [137] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, *Learning Representations by Back-propagating Errors*, *Nature (London)* **323**, 533 (1986).
- [138] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind, *Automatic Differentiation in Machine Learning: A Survey*, *J. Mach. Learn. Res.* **18**, 1 (2018).
- [139] É. Blayo, M. Bocquet, E. Cosme, and L. F. Cugliandolo, *Advanced Data Assimilation for Geosciences: Lecture Notes of the Les Houches School of Physics: Special Issue, June 2012* (Oxford University Press, New York, 2014).
- [140] C. M. Bishop, *Regularization and Complexity Control in Feed-Forward Networks*, in *Proceedings of the International Conference on Artificial Neural Networks ICANN'95* (unpublished), pp. 141–148, <https://publications.aston.ac.uk/id/eprint/524/>.

- [141] J. Sjöberg and L. Ljung, *Overtraining, Regularization and Searching for a Minimum, with Application to Neural Networks*, *Int. J. Control* **62**, 1391 (1995).
- [142] P. Weinberg and M. Bukov, *QuSpin: A PYTHON Package for Dynamics and Exact Diagonalisation of Quantum Many Body Systems Part I: Spin Chains*, *SciPost Phys.* **2**, 003 (2017).
- [143] P. Weinberg and M. Bukov, *QuSpin: A PYTHON Package for Dynamics and Exact Diagonalisation of Quantum Many Body Systems. Part II: Bosons, Fermions and Higher Spins*, *SciPost Phys.* **7**, 20 (2019).
- [144] H. Katsura, D. Schuricht, and M. Takahashi, *Exact Ground States and Topological Order in Interacting Kitaev/Majorana Chains*, *Phys. Rev. B* **92**, 115137 (2015).
- [145] M. P. A. Fisher, P. B. Weichman, G. Grinstein, and D. S. Fisher, *Boson Localization and the Superfluid-Insulator Transition*, *Phys. Rev. B* **40**, 546 (1989).
- [146] D. Jaksch, C. Bruder, J. I. Cirac, C. W. Gardiner, and P. Zoller, *Cold Bosonic Atoms in Optical Lattices*, *Phys. Rev. Lett.* **81**, 3108 (1998).
- [147] W. Zwerger, *Mott-Hubbard Transition of Cold Atoms in Optical Lattices*, *J. Opt. B* **5**, S9 (2003).
- [148] W. Krauth, M. Caffarel, and J.-P. Bouchaud, *Gutzwiller Wave Function for a Model of Strongly Interacting Bosons*, *Phys. Rev. B* **45**, 3137 (1992).
- [149] T. Comparin, `tcompa/bosehubbardgutzwiller v1.0.2`, <https://zenodo.org/record/1067968#.Yu1VQRzMKUk>.