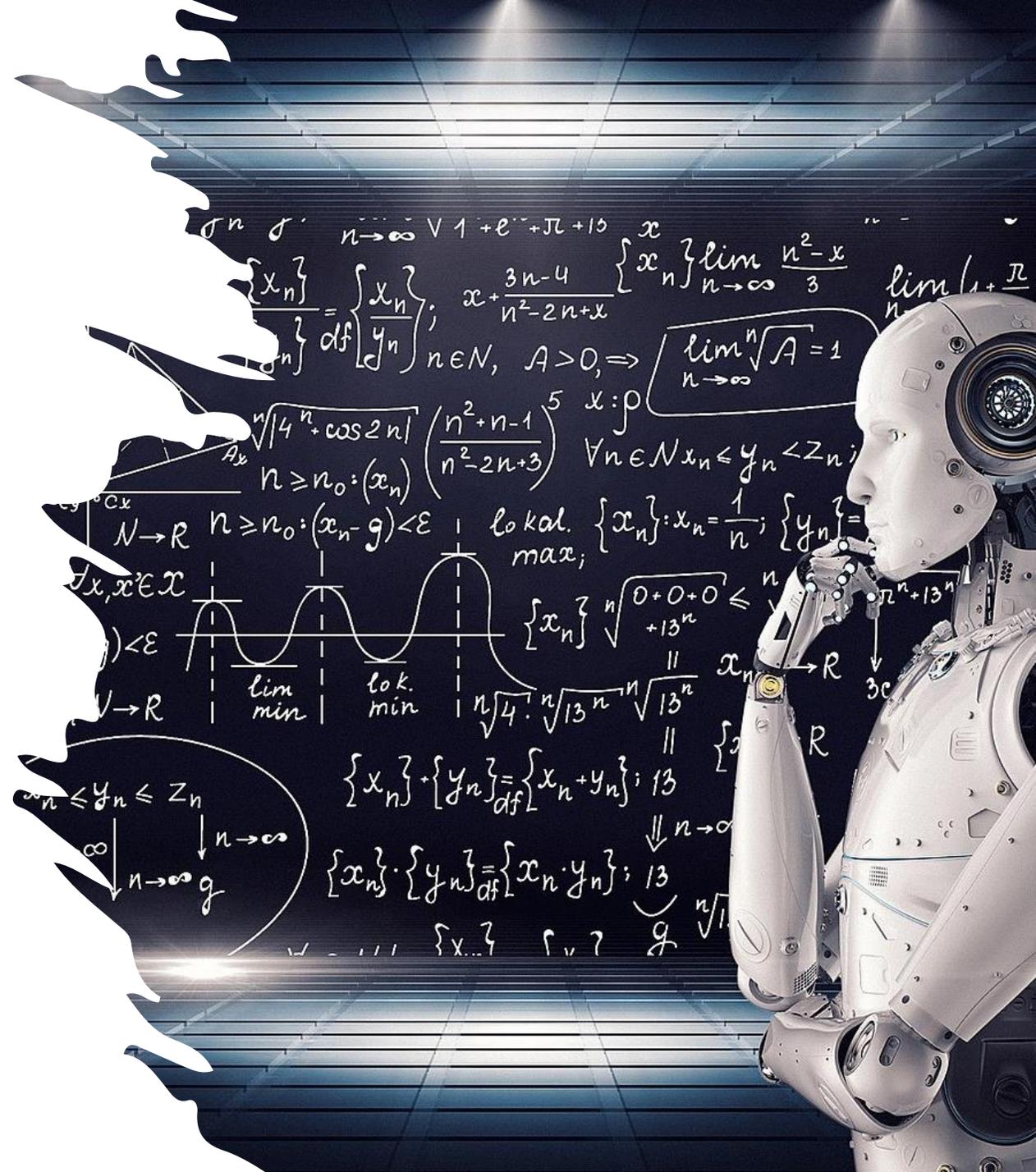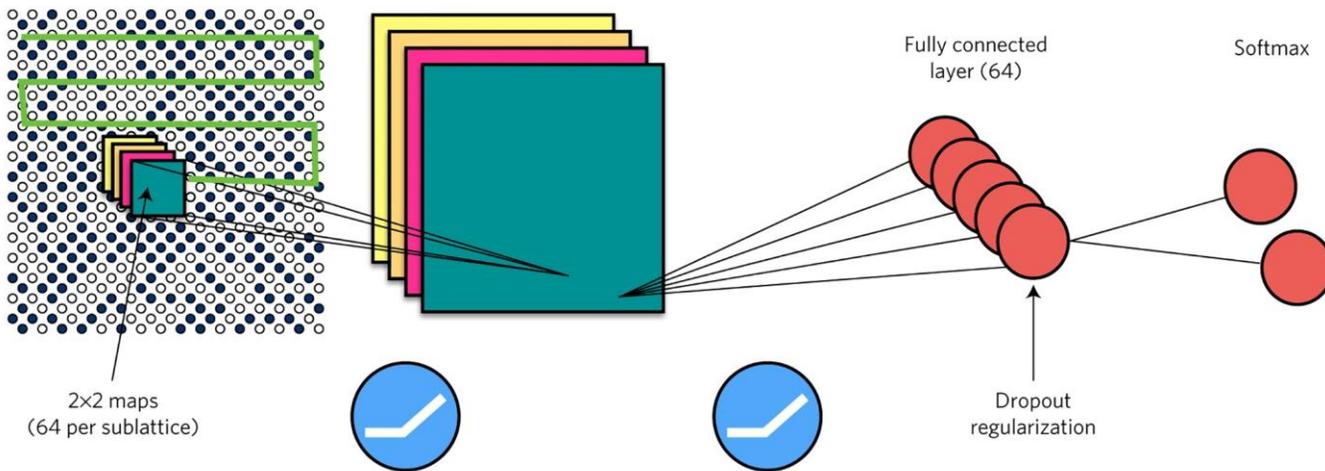# Towards interpretable and reliable machines learning physics
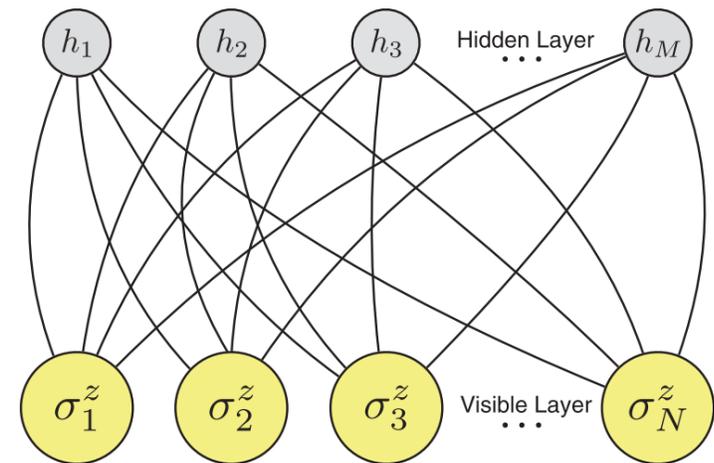
**Anna Dawid**

*Mach. Learn.: Sci. Technol. 3, 015002 (2022)*

Carrasquilla & Melko, Nat. Phys. **13**, 431-434 (2017)



Carleo & Troyer, Science **355**, 602-606, (2017)



1D TFIM

$\langle\sigma_i^z\sigma_j^z\rangle_{QMC}$

$\langle\sigma_i^z\sigma_j^z\rangle_{RBM}$

Torlai et al., Nat. Phys. **14**, 447–450 (2018)



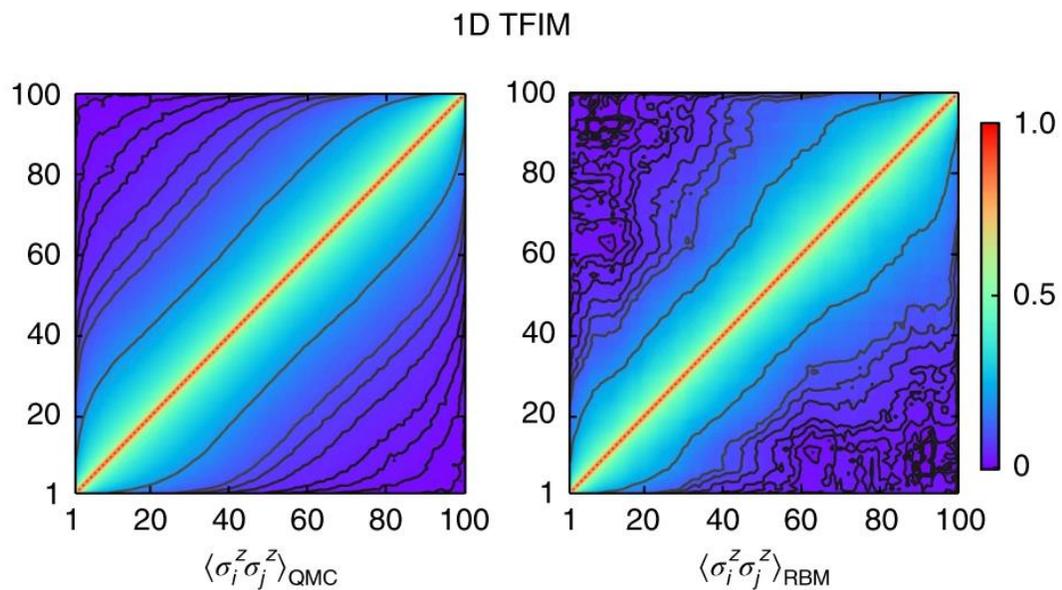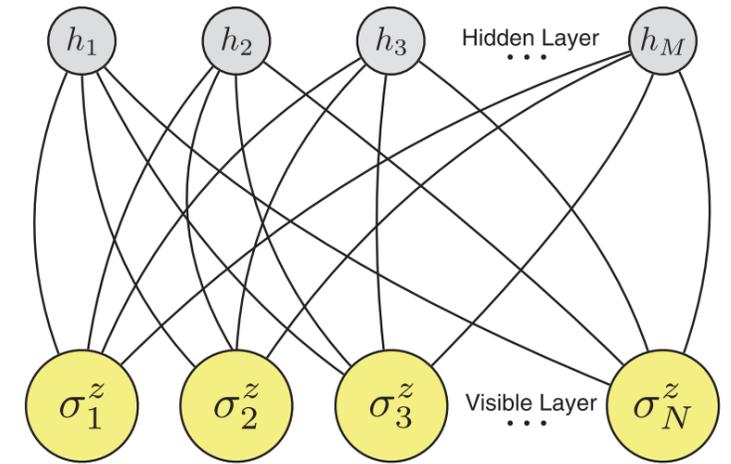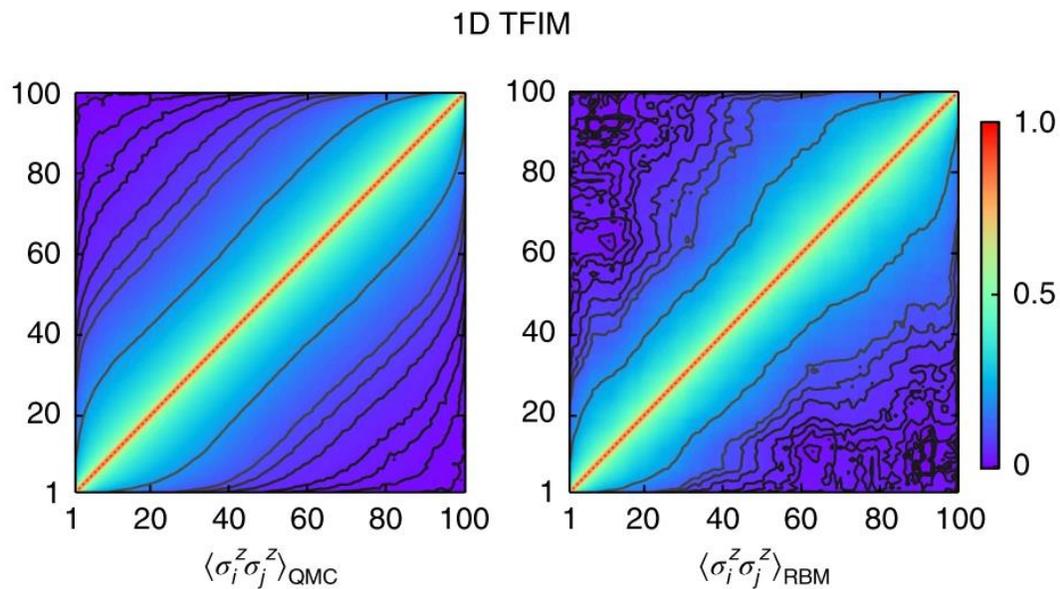Torlai et. al., Phys. Rev. Lett. **123**, 230504 (2019)

Carrasquilla & Melko, Nat. Phys. **13**, 431-434 (2017)

Carleo & Troyer, Science **355**, 602-606, (2017)

1D TFIM

$\langle \sigma_i^z \sigma_j^z \rangle_{QMC}$     $\langle \sigma_i^z \sigma_j^z \rangle_{RBM}$

Torlai et al., Nat. Phys. **14**, 447–450 (2018)

Torlai et. al., Phys. Rev. Lett. **123**, 230504 (2019)

# Machines in phase classification – open problems

→ **quantum many-body localization**

→ **topological phases of matter**

- disagreement of predicted critical exponents
- high sensitivity to hyperparameters describing the training process

- learning schemes trained on raw Monte Carlo configurations were found to be not effective

- pre-engineered features are often needed

→ mainly recovery of **known** results, but much cheaper

→ general problems with ML like…

*People worry that the computers will get too smart and take over the world, but the real problem is that they're too stupid and they've already taken over the world.*

Pedro Domingos "The Master Algorithm"

**!** even small invisible changes or a different background context can completely derail predictions

high error rates for faces from minority groups **!**

**!** the algorithm's hiring and insurance decisions are biased towards selecting men and white people

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

DATA

ANSWERS

Taken from: xkcd, A Webcomic of Romance, Sarcasm, Math, And Language, https://xkcd.com/1838/

# Some definitions

**_Interpretability_**
understanding what an ML model learns
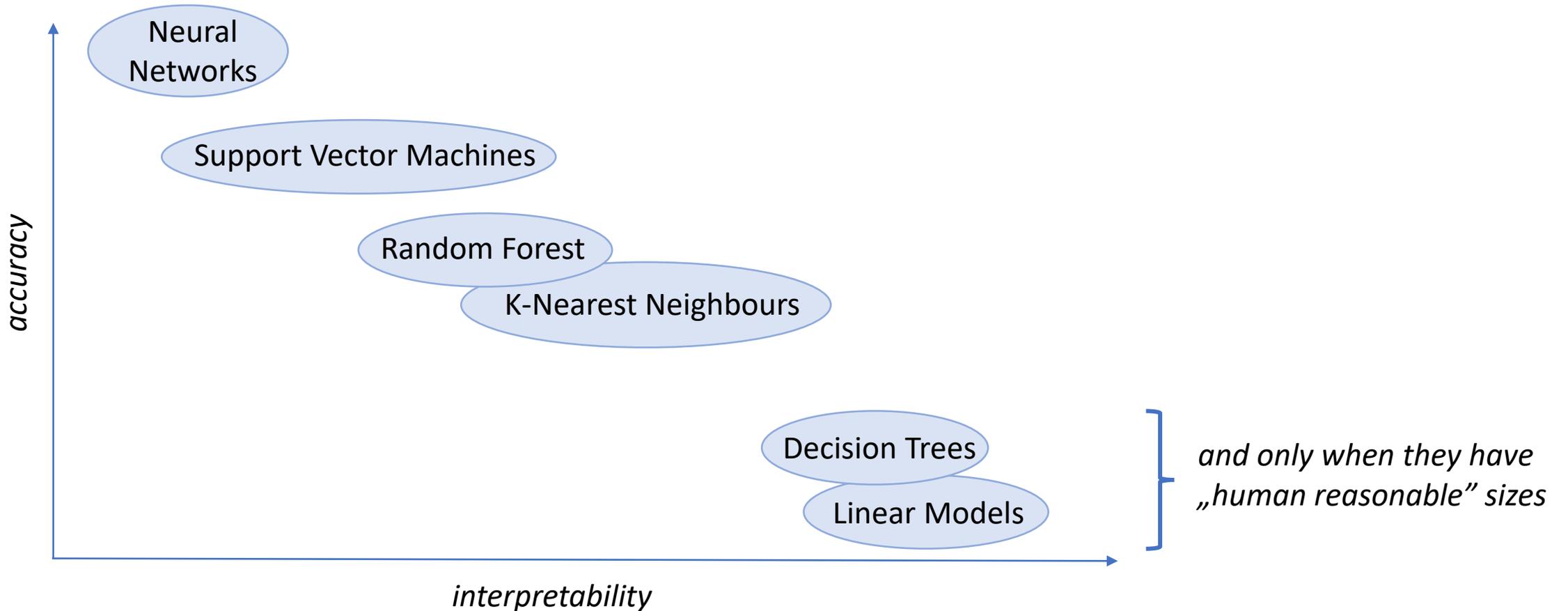and how it makes its predictions

**_Reliability_**
trusting our ML model predictions
(uncertainty)

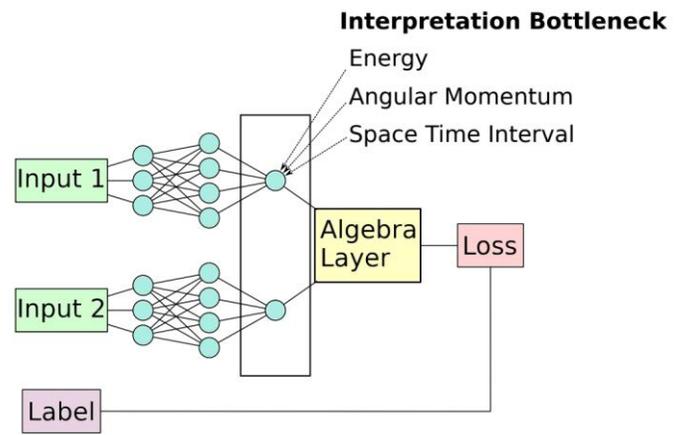_These two properties are closely
intertwined._
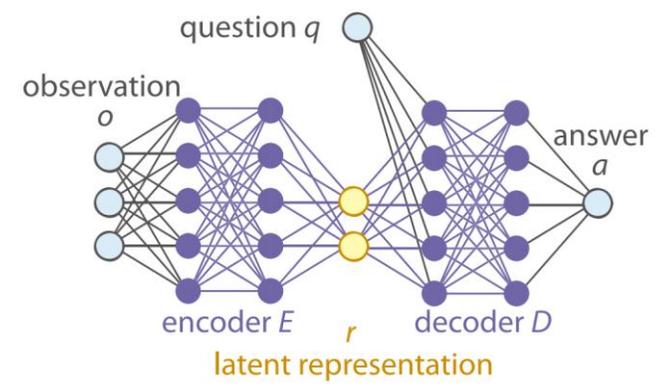
# Trade-off between complexity and interpretability
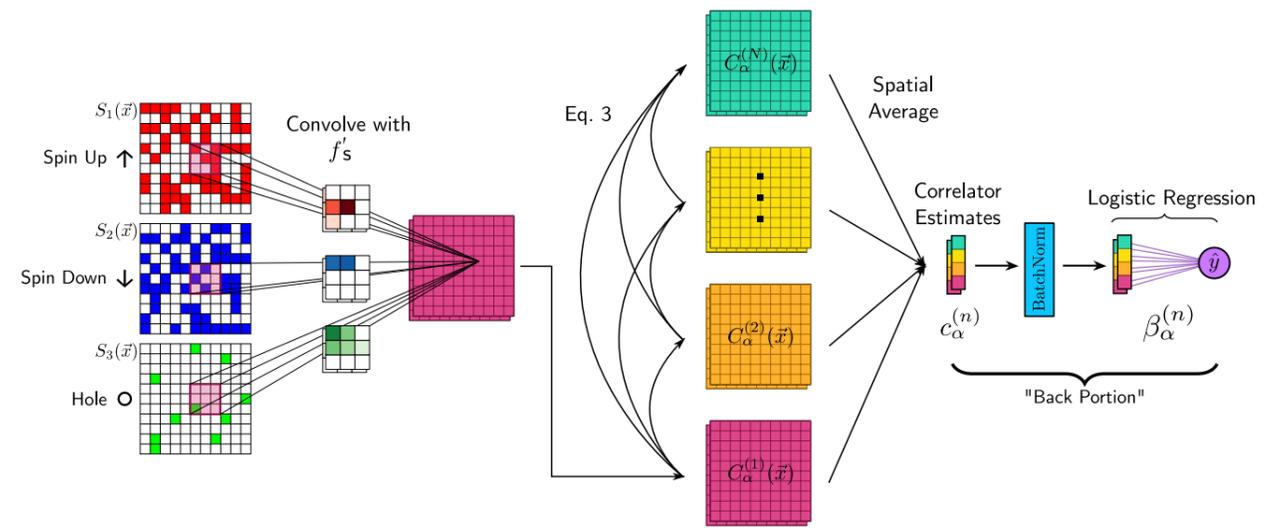
# Interpretation of ML in physics so far



**Interpretation Bottleneck**

Energy
Angular Momentum
Space Time Interval

Input 1
Input 2
Algebra Layer
Loss
Label

Phys. Rev. Research 2, 033499 (2020)

question $q$
observation $o$
encoder $E$
$r$
decoder $D$
answer $a$
latent representation

Phys. Rev. Lett. 124, 010508 (2020)

$S_1(\vec{x})$ Spin Up ↑
$S_2(\vec{x})$ Spin Down ↓
$S_3(\vec{x})$ Hole O
Convolve with $f'$s
Eq. 3
$C_\alpha^{(N)}(\vec{x})$
$C_\alpha^{(2)}(\vec{x})$
$C_\alpha^{(1)}(\vec{x})$
Spatial Average
Correlator Estimates
$c_\alpha^{(n)}$
BatchNorm
Logistic Regression
$\beta_\alpha^{(n)}$
$\hat{y}$
"Back Portion"

Nat. Commun. 12, 3905 (2021)

o Decision trees, kernel methods
o Bottleneck analysis

# Interpretation of ML in physics so far

o Decision trees, kernel methods
o Bottleneck analysis

**Interpretation Bottleneck**

Energy
Angular Momentum
Space Time Interval

Input 1

Input 2

Algebra Layer

Loss

Label

Phys. Rev. Re...

question $q$

observation $o$

answer $a$

decoder $D$

...entation

010508 (2020)

$S_1(\vec{x})$

Spin Up ↑

$S_2(\vec{x})$

Spin Down ↓

$S_3(\vec{x})$

Hole ○

$C_\alpha^{(2)}(\vec{x})$

$C_\alpha^{(1)}(\vec{x})$

$c_\alpha^{(n)}$

BatchNorm

Logistic Regression

$y$

$\beta_\alpha^{(n)}$

"Back Portion"

Nat. Commun. 12, 3905 (2021)

**Reliable** and **interpretable** ML which stays as smart as without these qualities, **independently** of model architecture and training
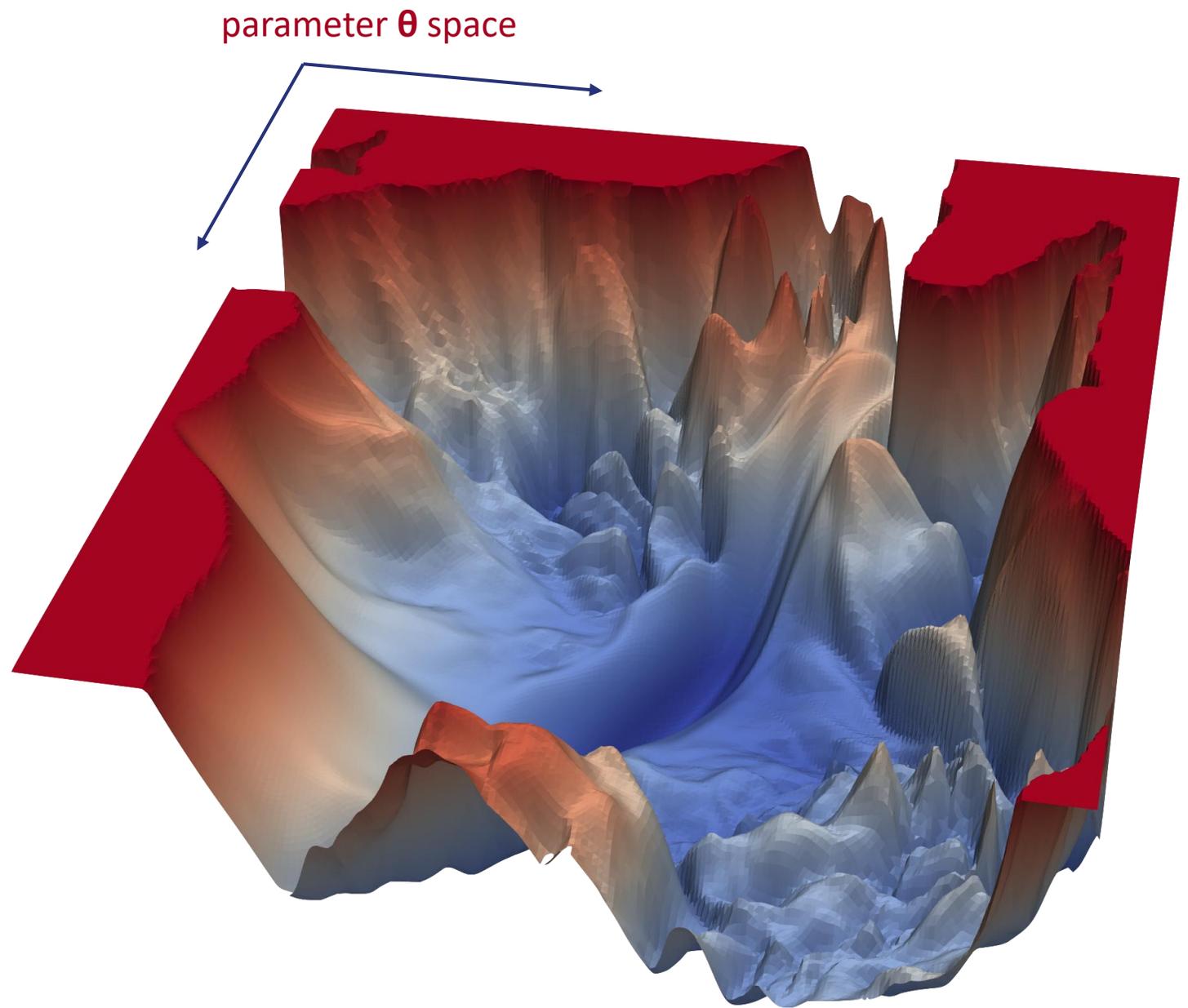
# Hessian-based toolbox

# „Minimum" of ML loss landscape

(# of classes $- 1$):     $\lambda_i > 0$
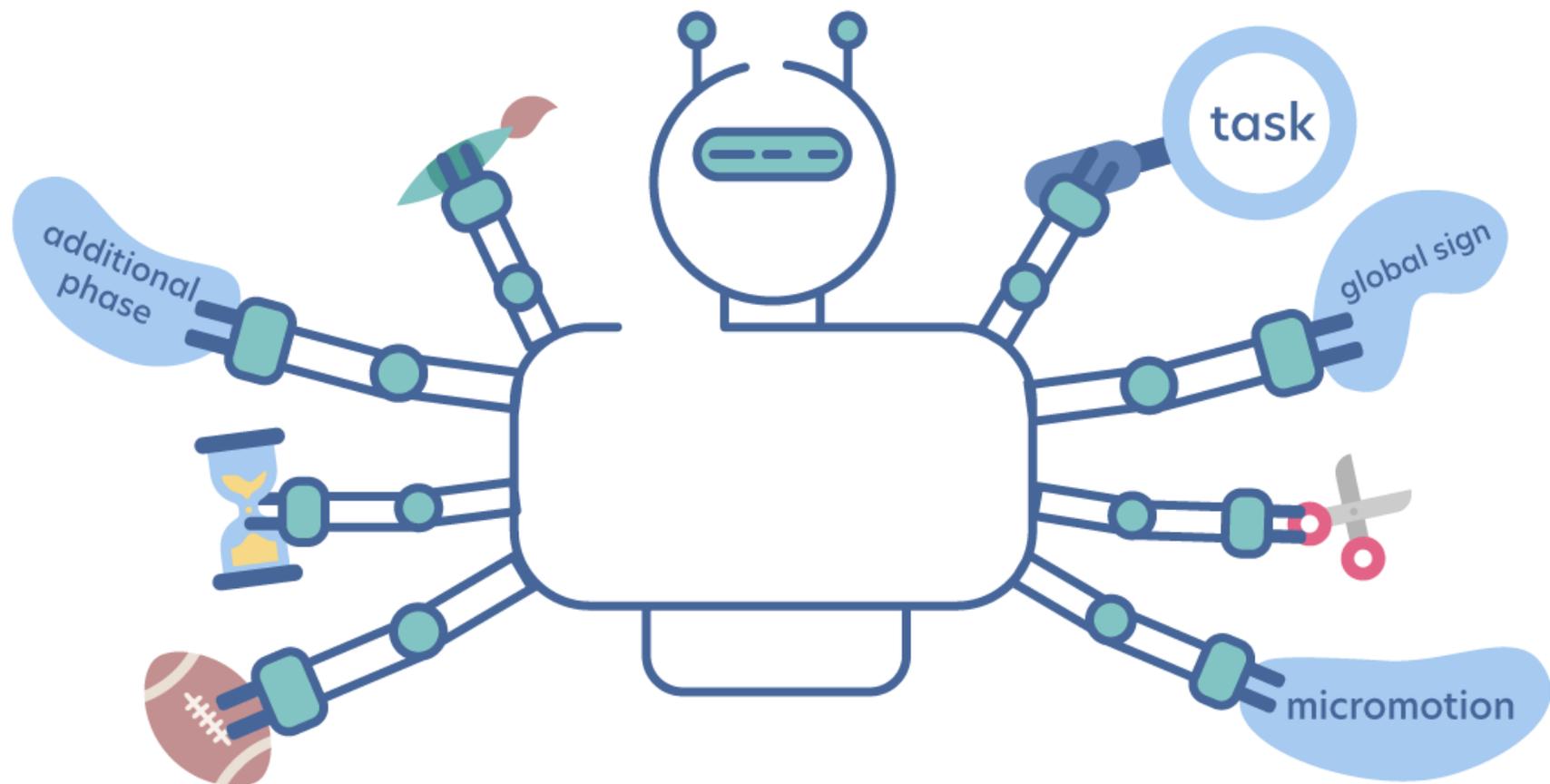
majority:     $\lambda_i \approx 0$

few:     $\lambda_i < 0$

$$H_{ij} = \frac{\partial^2 \mathcal{L}(D_{\text{train}}, \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j}\bigg|_{\boldsymbol{\theta} = \widetilde{\boldsymbol{\theta}}}$$
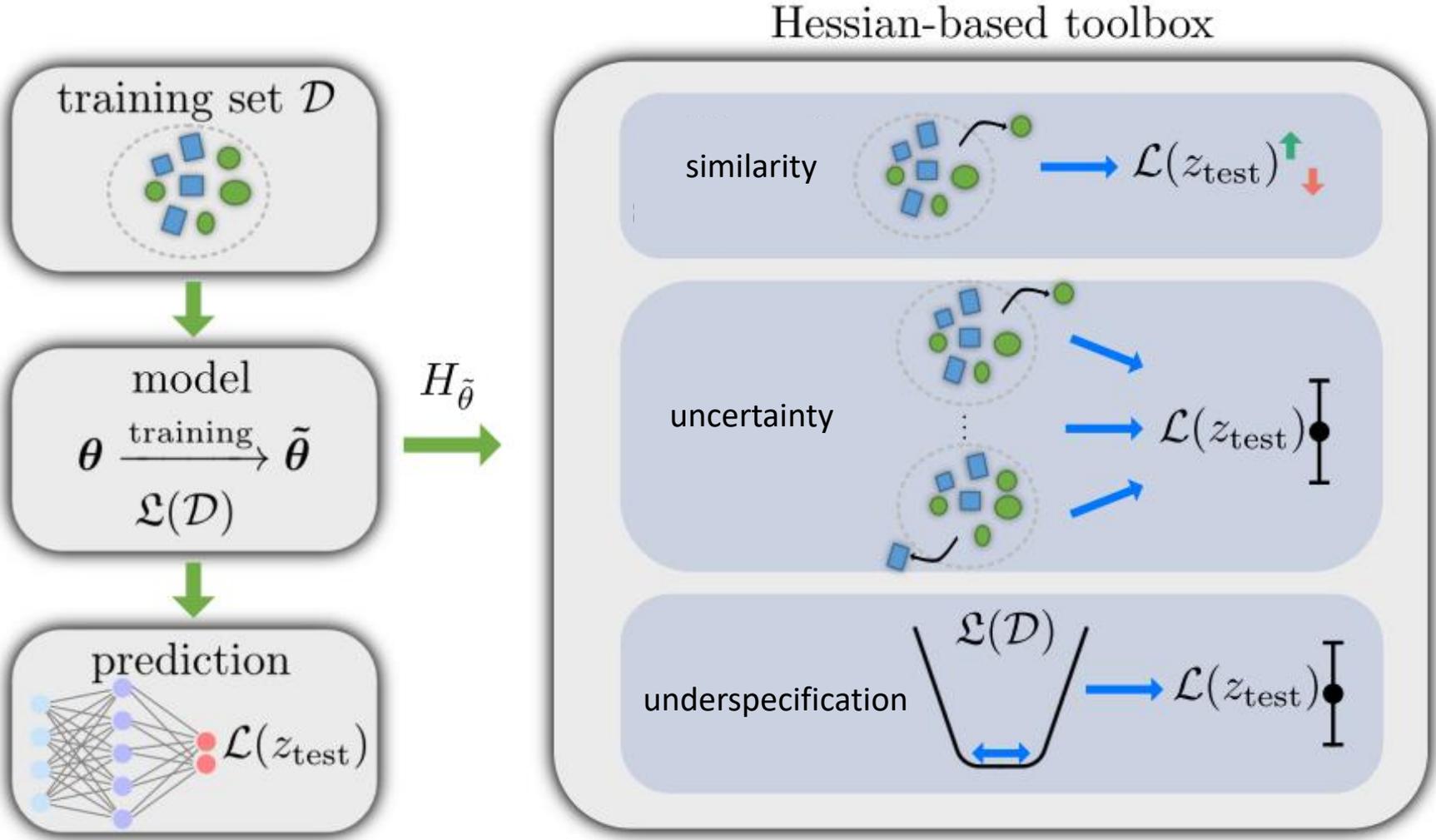


parameter $\boldsymbol{\theta}$ space

# Outline

1. Interpreting an ML model
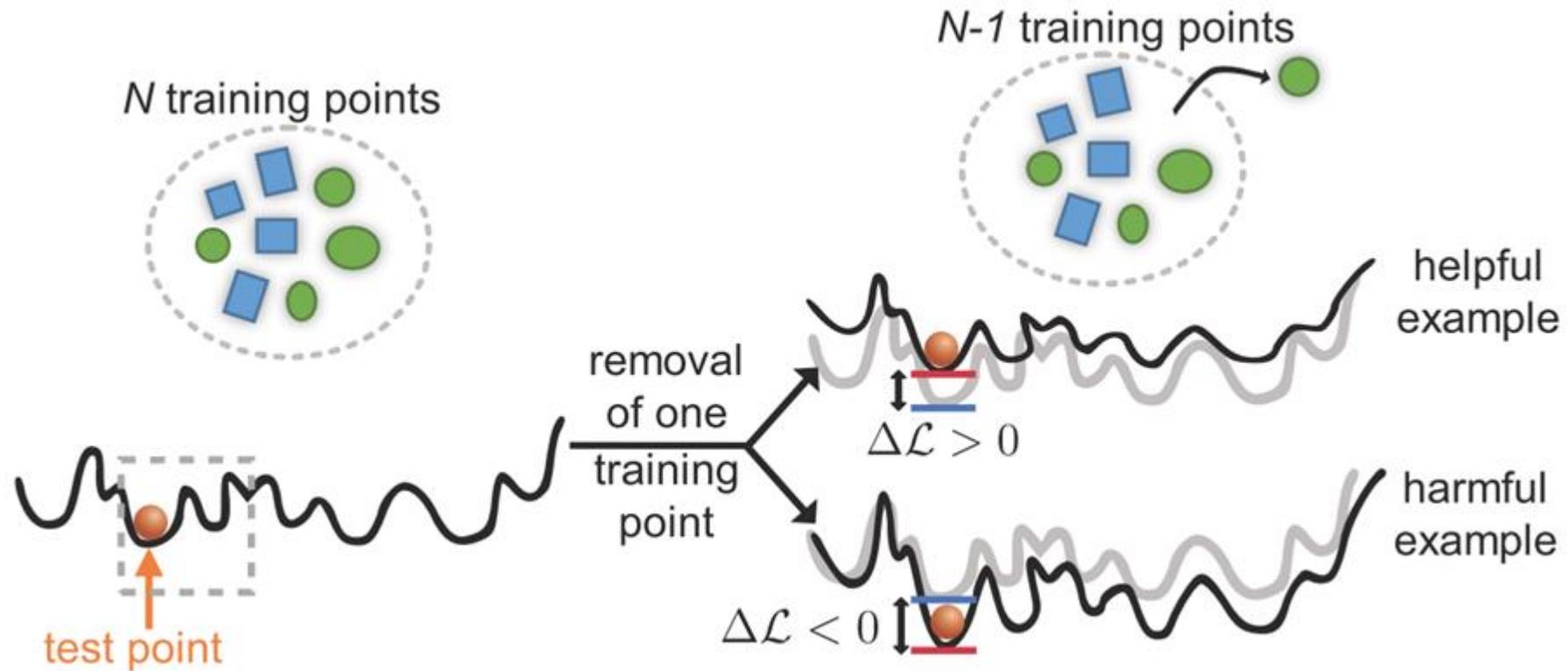2. Reliability methods

# Outline

1. Interpreting an ML model
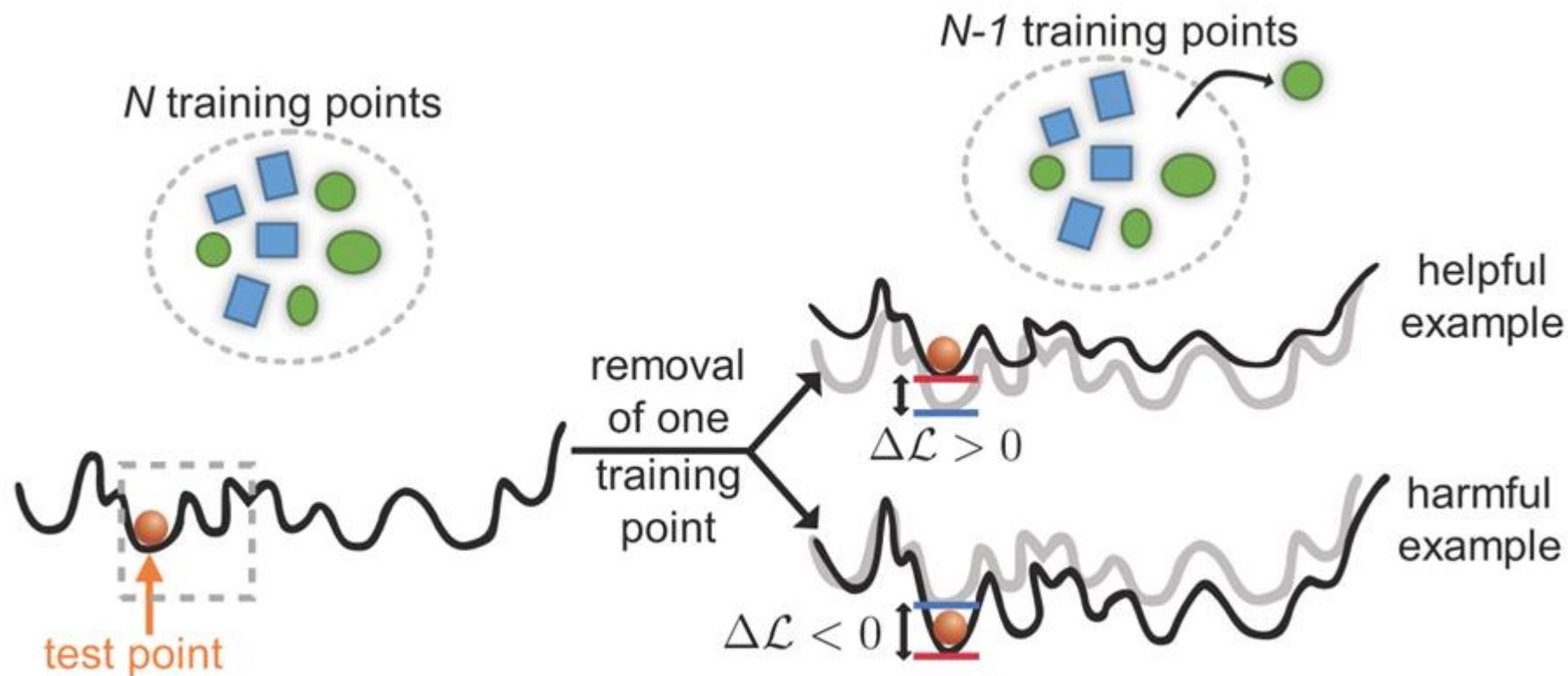2. Reliability methods

# Hessian-based toolbox



o **Influence functions**
Koh & Liang
arXiv:1703.04730

o Resampling
Uncertainty
Estimation (RUE)
Schulam & Saria
arXiv:1901.00403

o Local Ensembles
(LEs)
Madras, Atwood, D'Amour
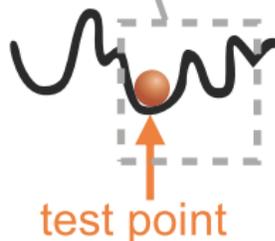arXiv:1910.09573
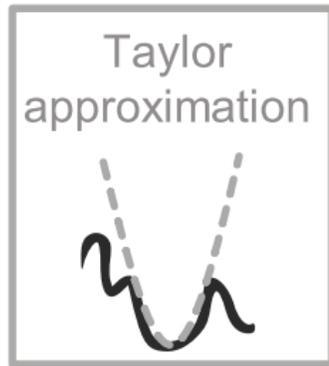
# Leave-one-out training



N training points

N-1 training points

removal of one training point

helpful example

$\Delta \mathcal{L} > 0$

harmful example

$\Delta \mathcal{L} < 0$

test point

# Leave-one-out training



prohibitively expensive!

# Influence functions



Taylor approximation

test point

**Analytical approximation for leave-one-out training**

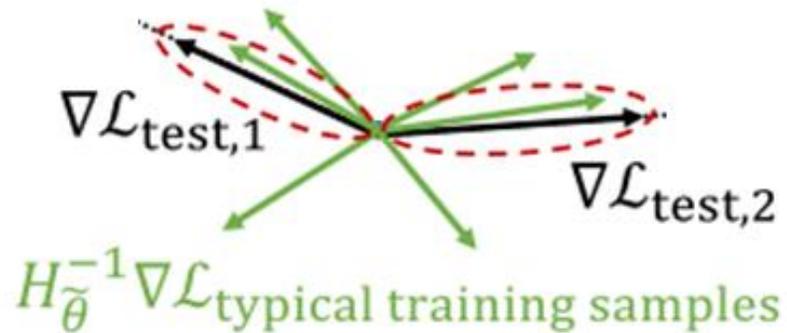$$\mathcal{I}(z_{\mathrm{r}}, z_{\mathrm{test}}) = \frac{1}{n} \nabla_\theta \mathcal{L}(z_{\mathrm{test}}, \hat{\theta})^T H_\theta^{-1}(\hat{\theta}) \nabla_\theta \mathcal{L}(z_{\mathrm{r}}, \hat{\theta})$$

approximated change in parameters due to removal of $z_{\mathrm{r}}$

Assumption: Hessian is positive-definite.
Generalization to non-convex models was done by Koh & Liang: arXiv:1703.04730, ICML 2017's best paper

# Geometrical interpretation



$$\mathcal{I}(z_{\mathrm{r}}, z_{\mathrm{test}}) = \frac{1}{n} \nabla_\theta \mathcal{L}(z_{\mathrm{test}}, \hat{\theta})^T H_\theta^{-1}(\hat{\theta}) \nabla_\theta \mathcal{L}(z_{\mathrm{r}}, \hat{\theta})$$

**notion of similarity
in the model
internal representation!**

it is a scalar product of gradient of a test point and the gradient of a training point, corrected by local curvature described by the Hessian

# Three messages

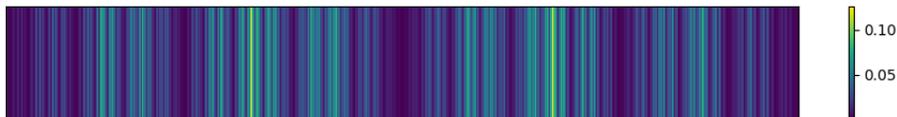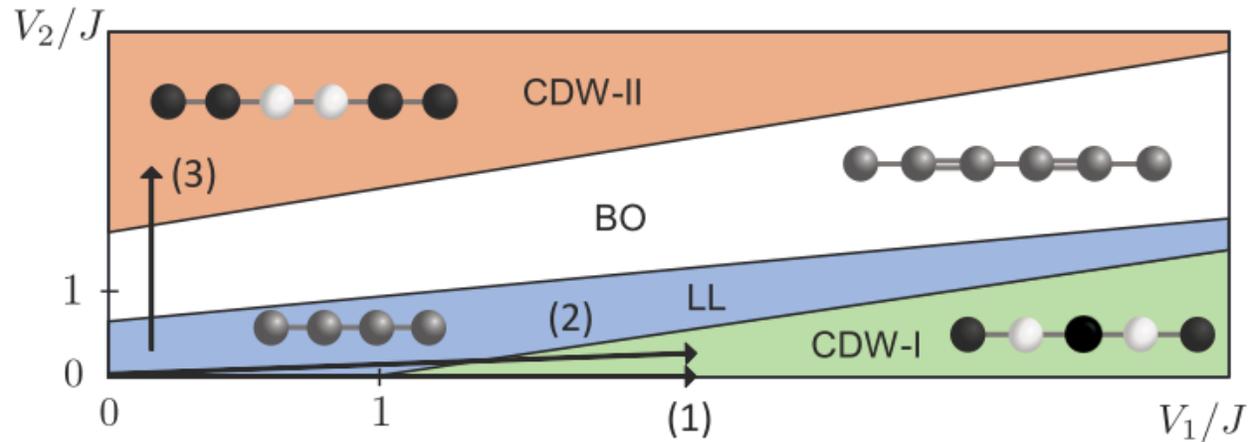- Detection additional phases
- Detecting influential data features
- Anomaly detection with influence functions
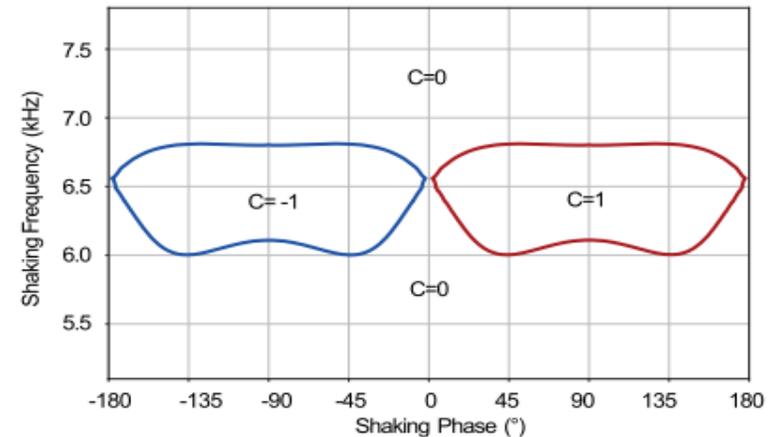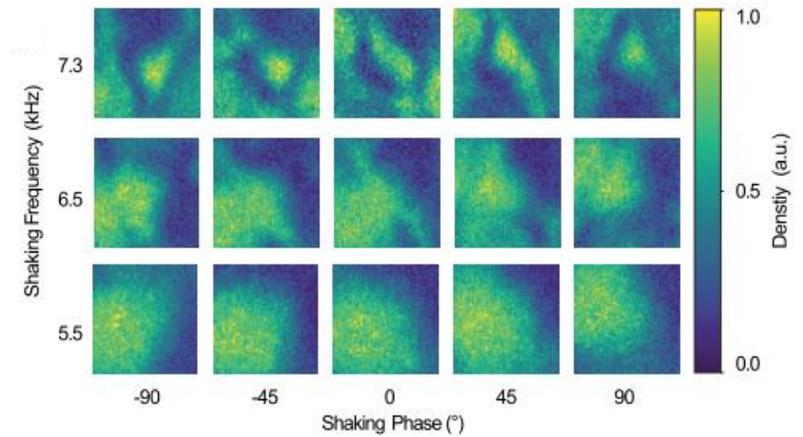
# Physical input data

## 1) simulated

spinless 1D Fermi-Hubbard model at half-filling



## 2) experimental

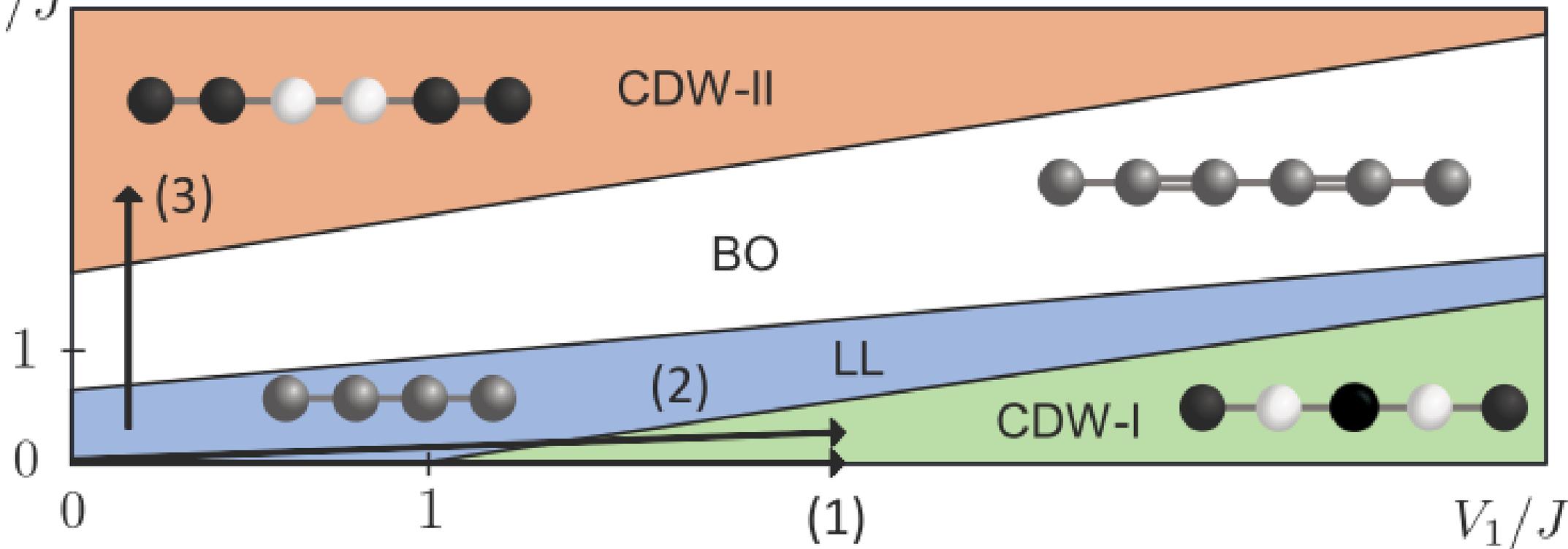topological Haldane model

**Three messages**

Detection additional phases

Detecting influential data features

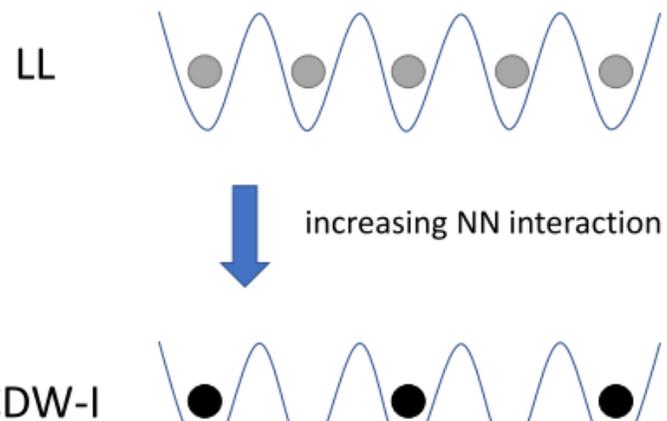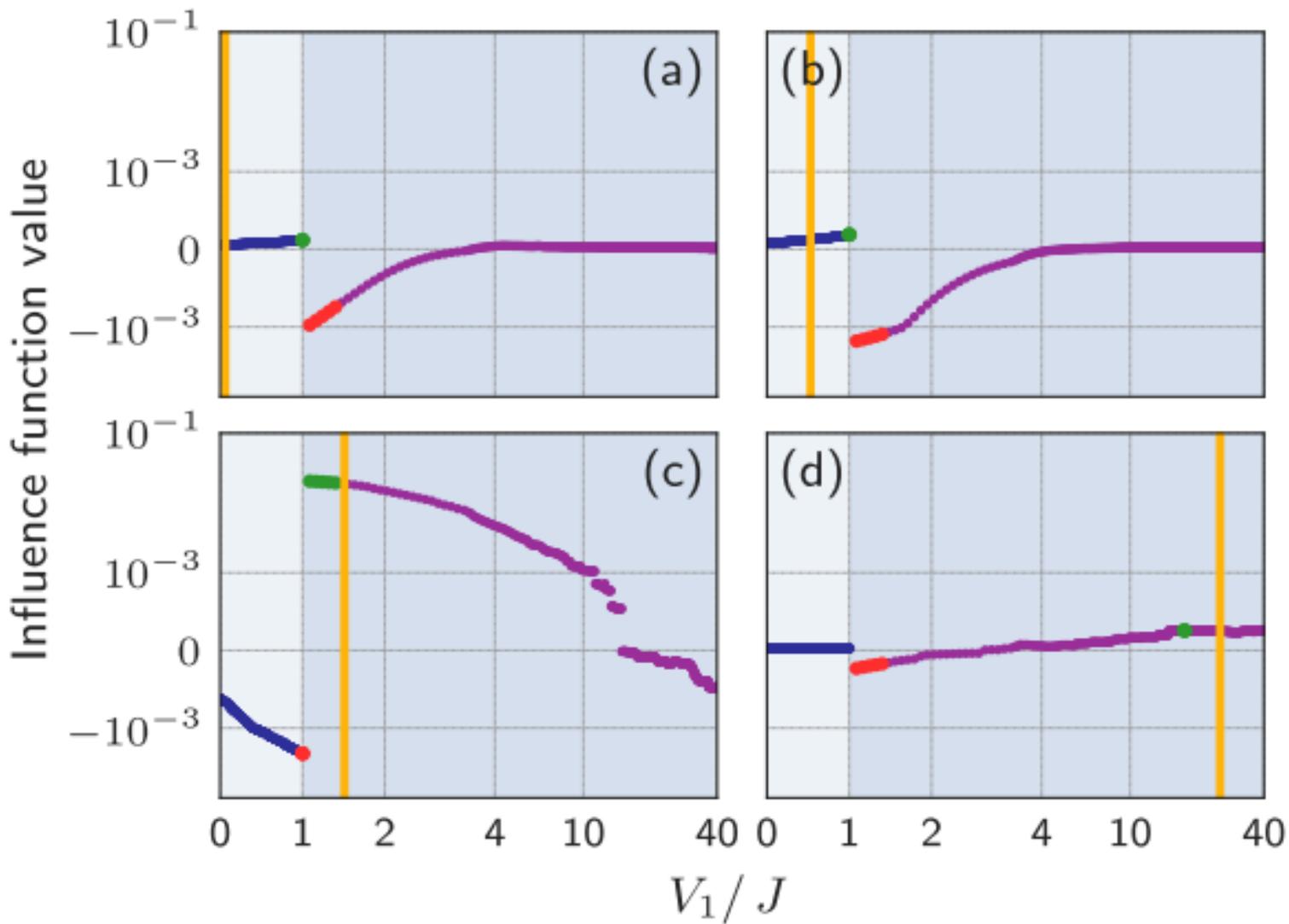Anomaly detection with influence functions

test point — training points from LL phase — training points from CDW-I phase

nearest neighbor interaction / hopping amplitude

A. Dawid et al, *New J. Phys.* **22** 115001 (2020)

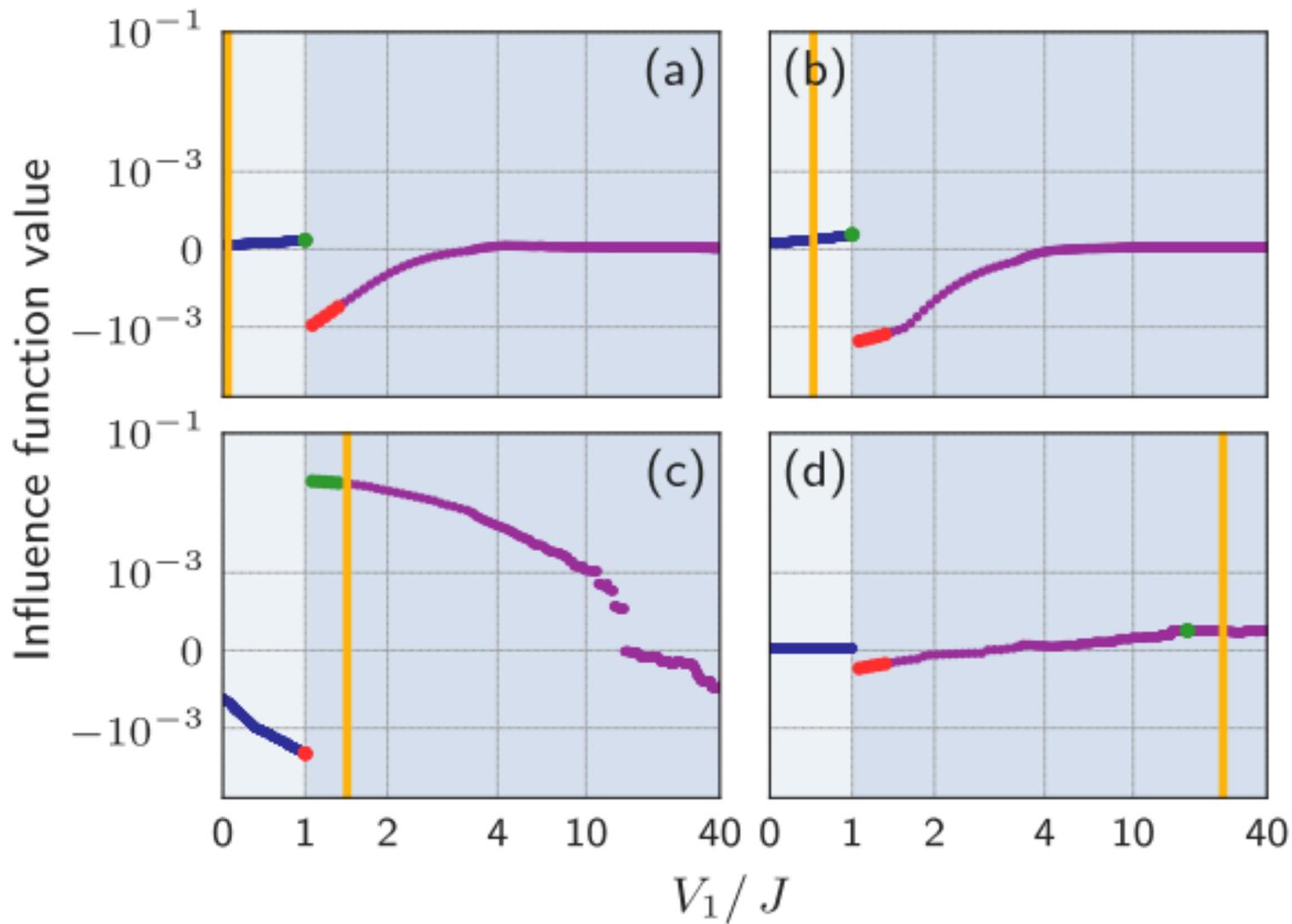test point ——— training points from LL phase ——— training points from CDW-I phase

nearest neighbor interaction / hopping amplitude

A. Dawid et al, *New J. Phys.* **22** 115001 (2020)

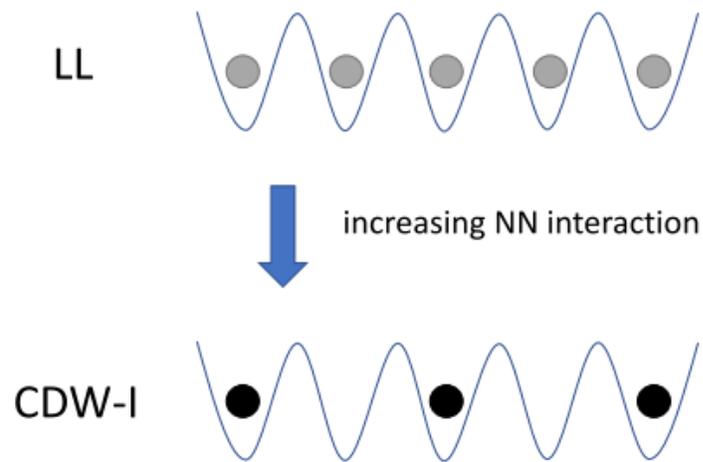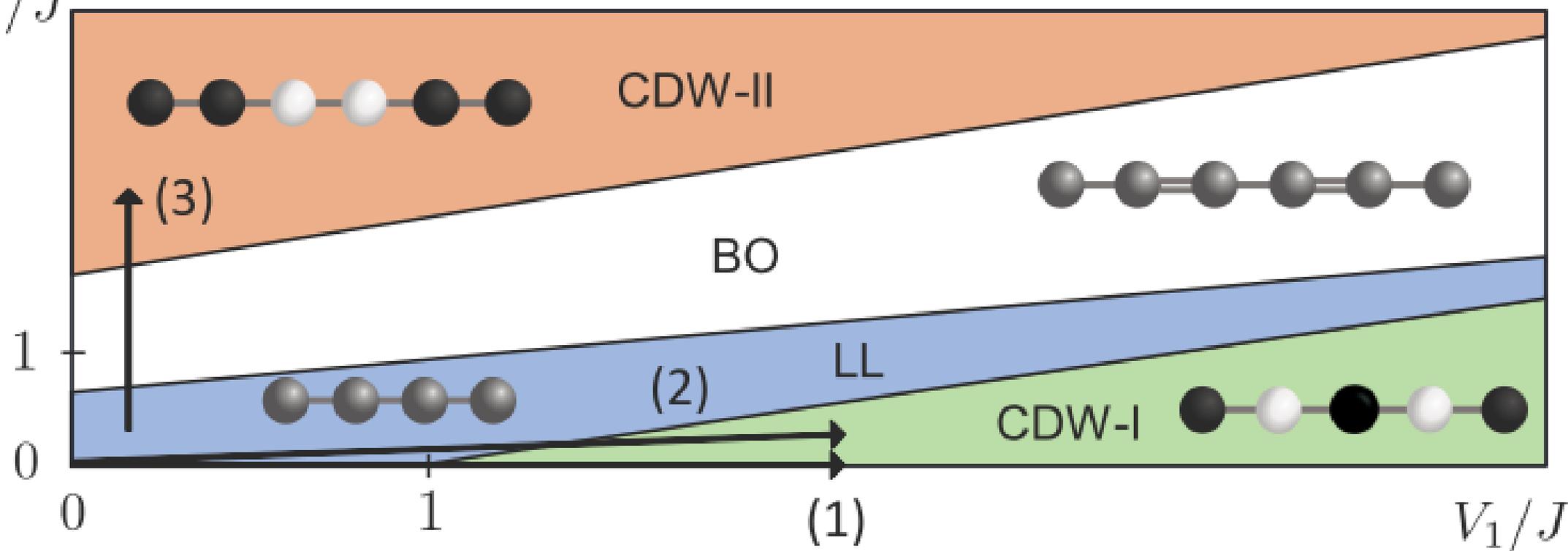test point     training points from LL phase     training points from BO and CDW-II phases

Influence function value

(a)     (b)

$V_2 / V_1$

nearest neighbor interaction / nearest neighbor interaction

A. Dawid et al, *New J. Phys.* **22** 115001 (2020)

It sees additional phase!

# Unsupervised machine learning of topological phase transition from experimental data



unsupervised approaches had troubles with distinghuishing between two topological phases...

N. Käming, A. Dawid et al., *MLST* **2** 035037 (2021)

Three messages

Detection additional phases

Detecting influential
data features

Anomaly detection
with influence functions

# Micromotion phase



b



The most influential points are localized around the same micromotion phase as test point

Same shaking frequency and shaking phase different micromotion phases



N. Käming, A. Dawid et al., *MLST* **2** 035037 (2021)

# Three messages

Detection additional phases

Detecting influential data features

Anomaly detection with influence functions

# Global sign

+ label: e.g., 0 for LL, 1 for CDW-I

We usually fix the global sign to +
Choice of global sign changes nothing in physics

global sign-imbalanced set
98% positive, 2% negative

global sign-balanced set
50% positive, 50% negative

test point    training points from LL phase    training points from CDW-I phase

Influence function value

$V_1 / J$

nearest neighbor interaction / hopping amplitude

Dawid et al., *MLST* **3**, 015002 (2022)

global sign-imbalanced set
98% positive, 2% negative

global sign-balanced set
50% positive, 50% negative

We can find outliers in the training data (according to the model's internal similarity measure!)

Influence function value

$V_1 / J$

nearest neighbor interaction / hopping amplitude

test point       training points from LL phase       training points from CDW-I phase

Dawid et al., *MLST* **3**, 015002 (2022)

test point ── training points from LL phase ── training points from CDW-I phase

**global sign-imbalanced set**
98% positive, 2% negative

**global sign-balanced set**
50% positive, 50% negative

(a)

(b)

Influence function value

$V_1 / J$

nearest neighbor interaction / hopping amplitude

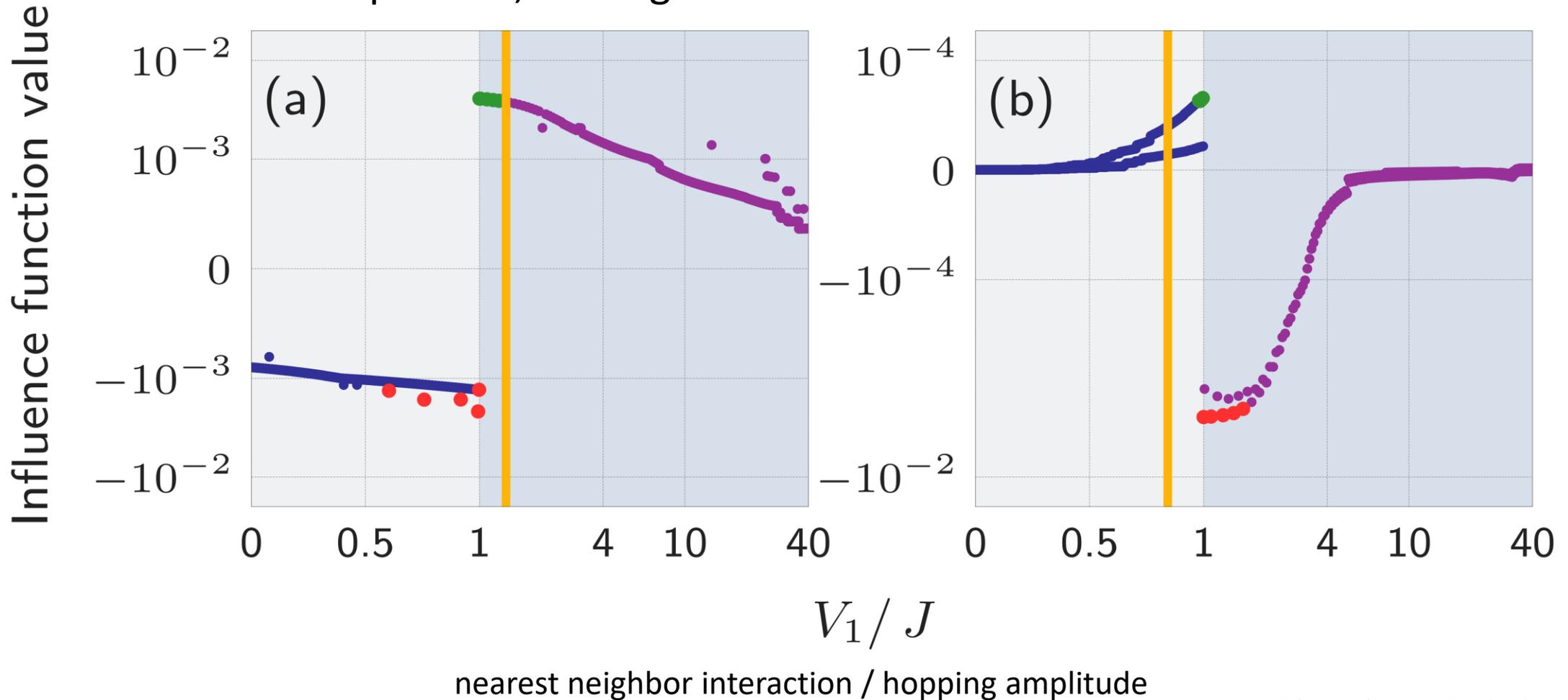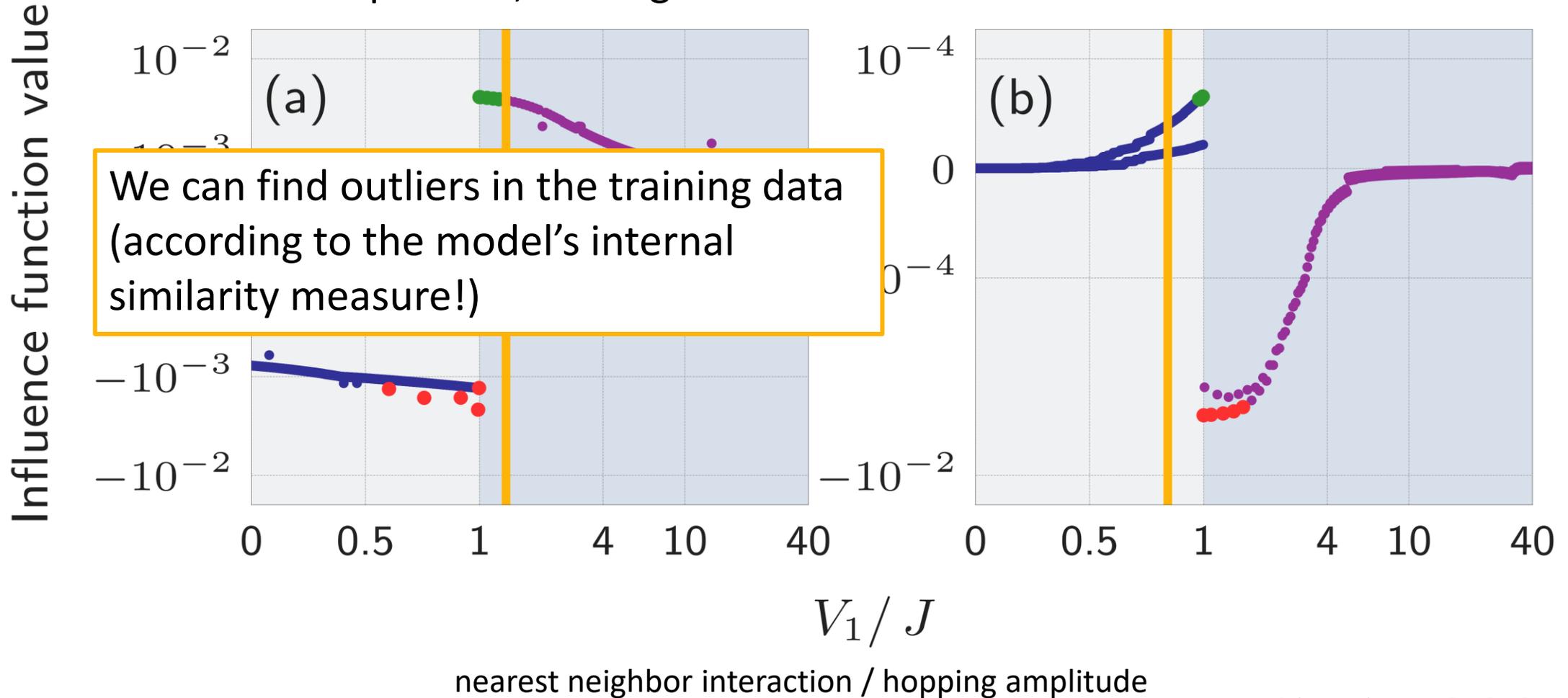Dawid et al., *MLST* **3**, 015002 (2022)

test point — training points from LL phase — training points from CDW-I phase

**global sign-imbalanced set**
98% positive, 2% negative

**global sign-balanced set**
50% positive, 50% negative

(a)

(b)

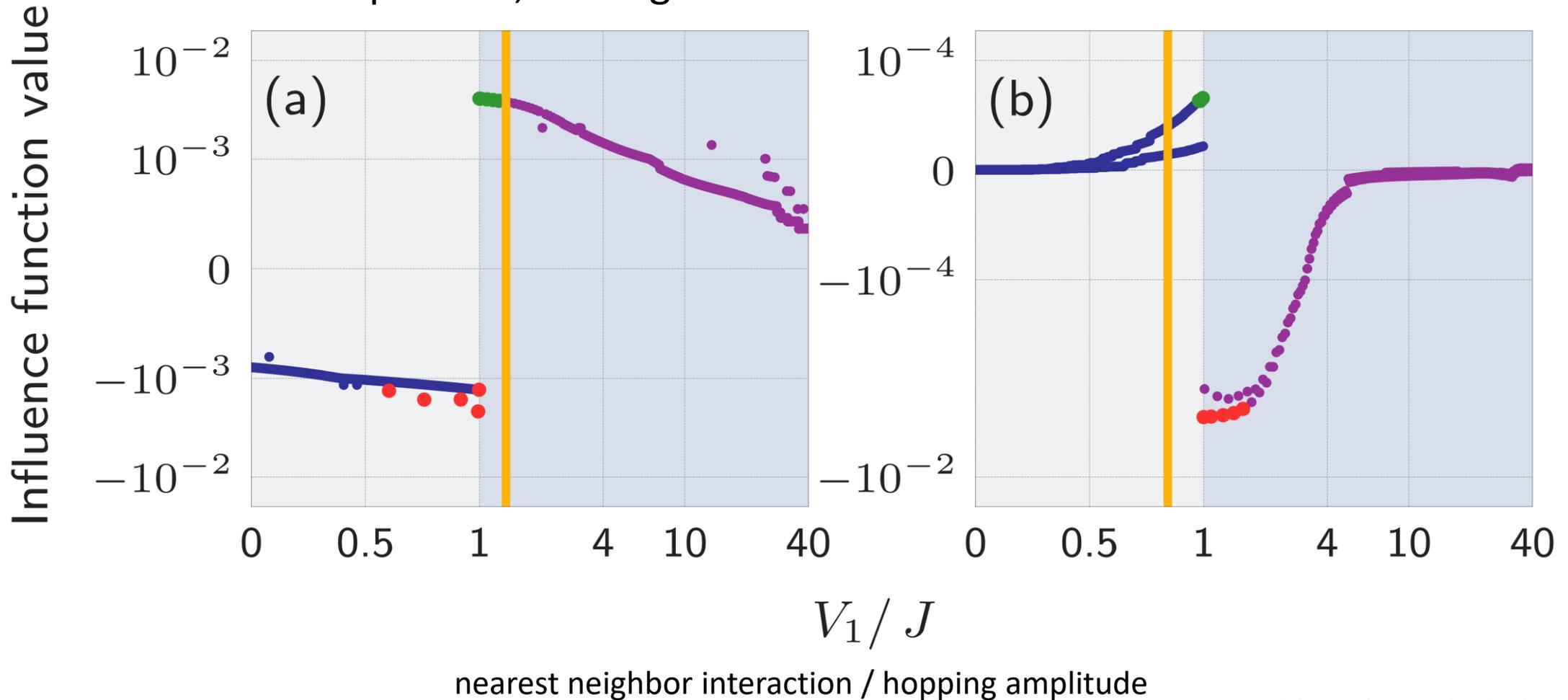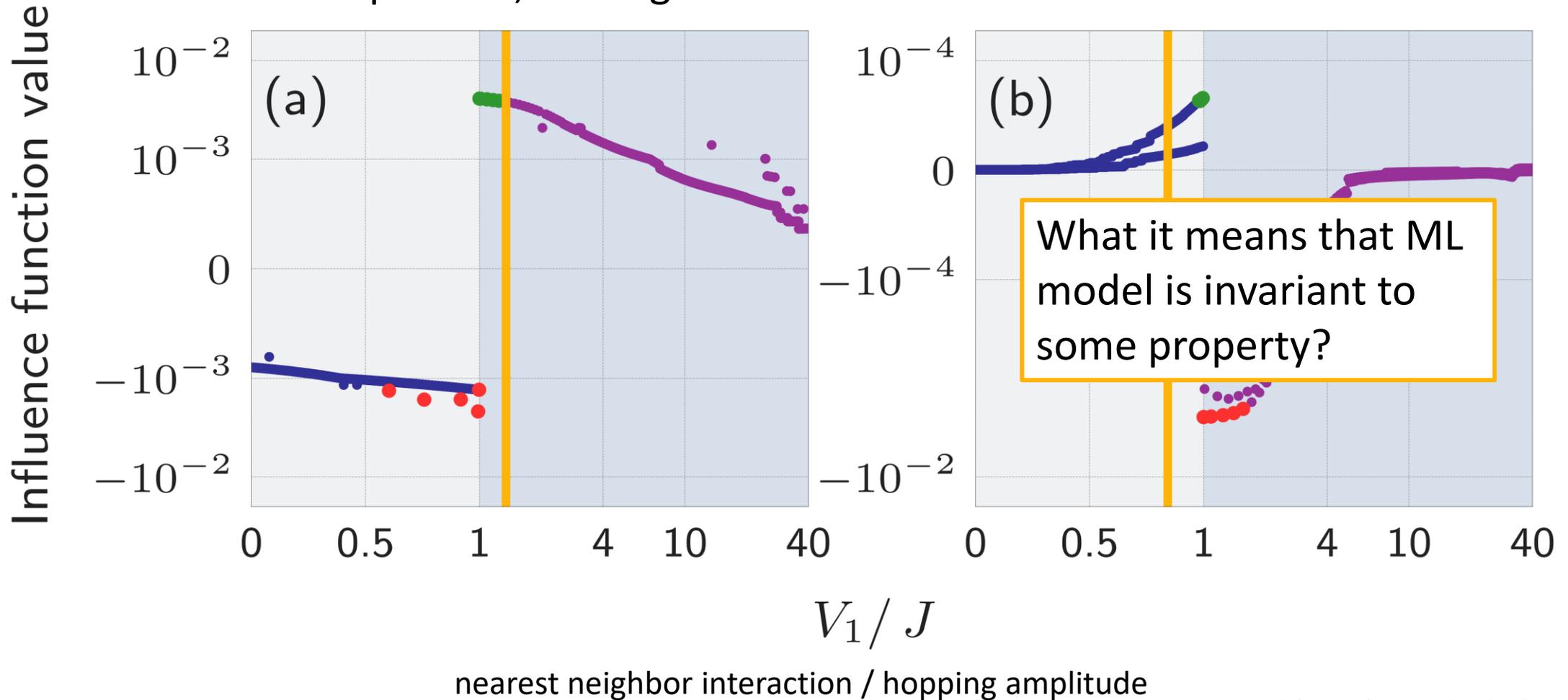What it means that ML model is invariant to some property?

Influence function value

$V_1 / J$

nearest neighbor interaction / hopping amplitude

Dawid et al., *MLST* **3**, 015002 (2022)

# Outline

1. Interpreting an ML model
2. Reliability methods

# Hessian-based toolbox



Hessian-based toolbox

training set $\mathcal{D}$

model
$\boldsymbol{\theta} \xrightarrow{\text{training}} \tilde{\boldsymbol{\theta}}$
$\mathfrak{L}(\mathcal{D})$

prediction
$\mathcal{L}(z_{\text{test}})$

$H_{\tilde{\theta}}$

similarity $\mathcal{L}(z_{\text{test}})$

uncertainty $\mathcal{L}(z_{\text{test}})$

underspecification $\mathfrak{L}(\mathcal{D})$ $\mathcal{L}(z_{\text{test}})$

o Influence functions
Koh & Liang
arXiv:1703.04730

o **Resampling Uncertainty Estimation (RUE)**
Schulam & Saria
arXiv:1901.00403

o Local Ensembles (LEs)
Madras, Atwood, D'Amour
arXiv:1910.09573

# Resampling Uncertainty Estimation (RUE)
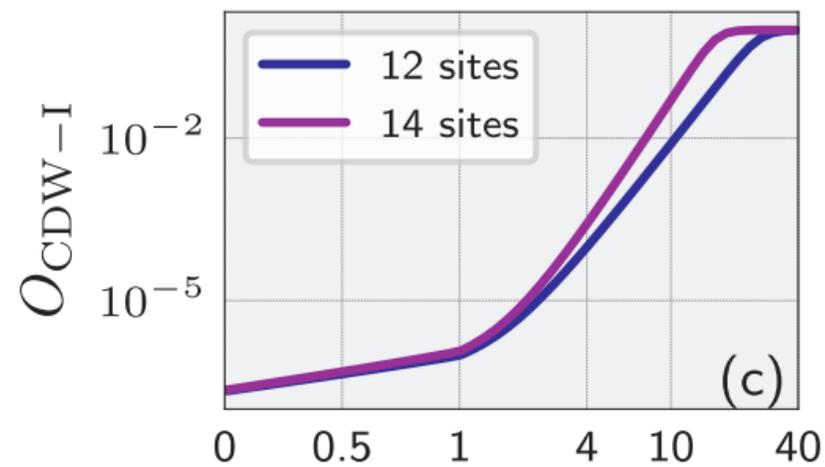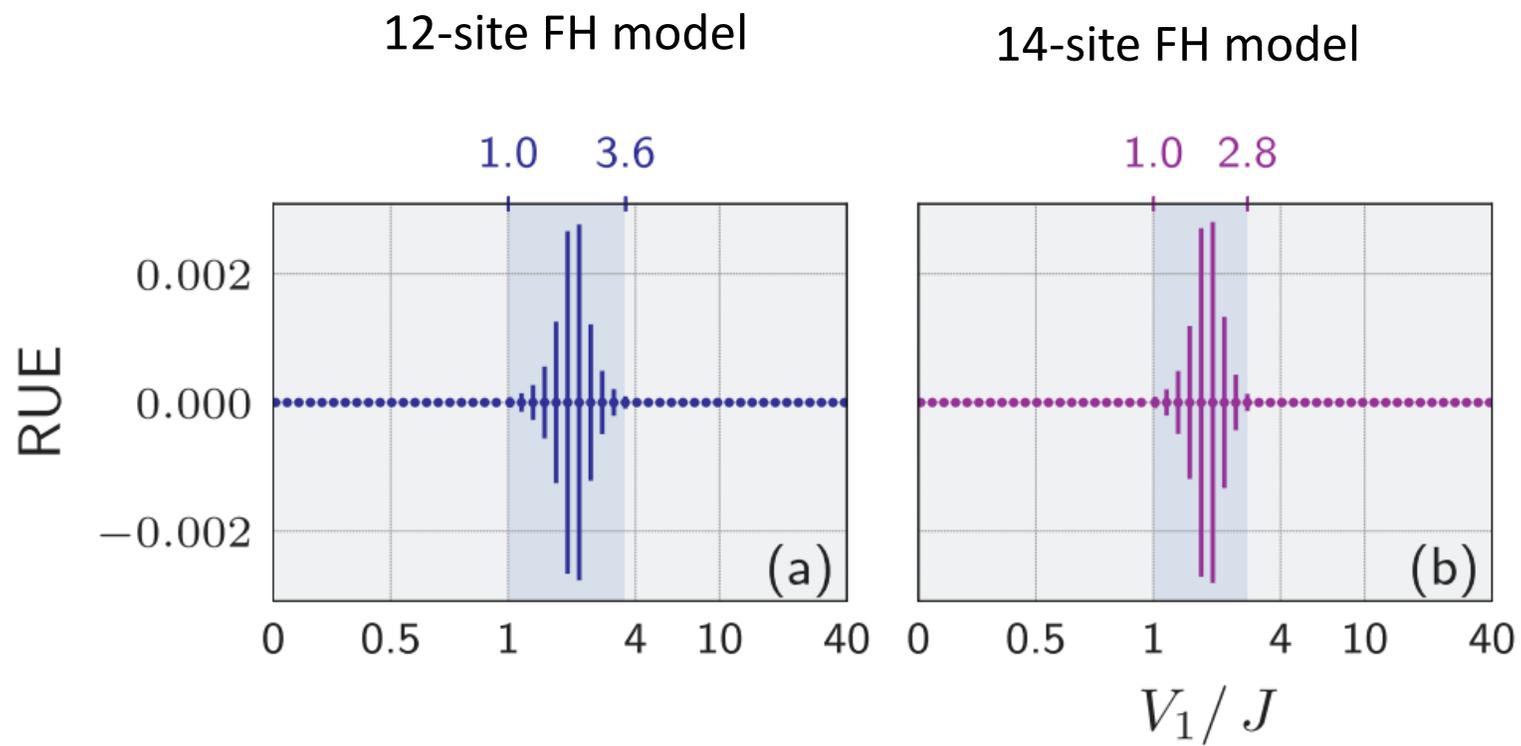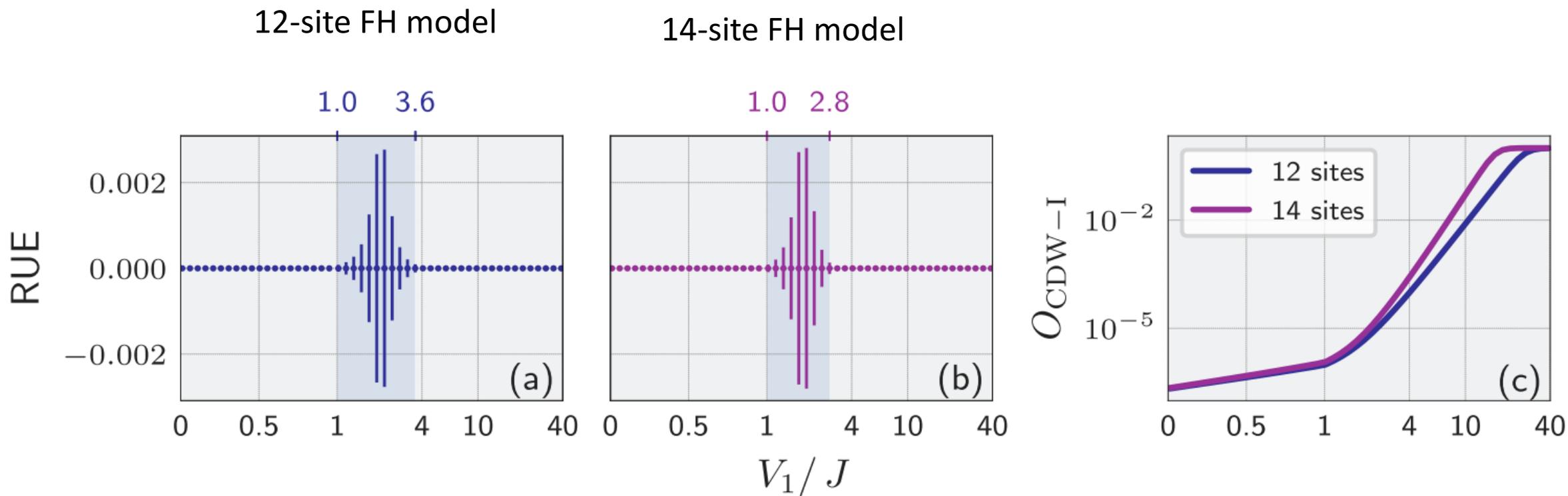
$b$ bootstrap samples

studied ML model

$D[1,1,1,1,...,1,1]$

training dataset, every
data point is taken once

$D_1[0,1,3,...,1,2]$

$D_2[1,4,0,...,0,1]$

$b$ predictions

variance

$D_b[2,0,3,...,1,0]$

P. Schulam and S. Saria, *Can you trust this prediction? Auditing pointwise reliability after learning*, AISTATS 2019 - 22nd Int. Conf. Artif. Intell. Stat. 89 (2020), arXiv:1901.00403v2.

12-site FH model     14-site FH model

Dawid et al., *MLST* **3**, 015002 (2022)

**12-site FH model**

**14-site FH model**

RUE indicates sharpness of the transition

Dawid et al., *MLST* **3**, 015002 (2022)

# Hessian-based toolbox



Hessian-based toolbox

training set $\mathcal{D}$

model
$\boldsymbol{\theta} \xrightarrow{\text{training}} \tilde{\boldsymbol{\theta}}$
$\mathfrak{L}(\mathcal{D})$

prediction
$\mathcal{L}(z_{\text{test}})$

$H_{\tilde{\theta}}$

similarity $\quad \mathcal{L}(z_{\text{test}})$

uncertainty $\quad \mathcal{L}(z_{\text{test}})$

underspecification $\quad \mathfrak{L}(\mathcal{D}) \quad \mathcal{L}(z_{\text{test}})$

o Influence functions
Koh & Liang
arXiv:1703.04730

o Resampling
Uncertainty
Estimation (RUE)
Schulam & Saria
arXiv:1901.00403

o **Local Ensembles
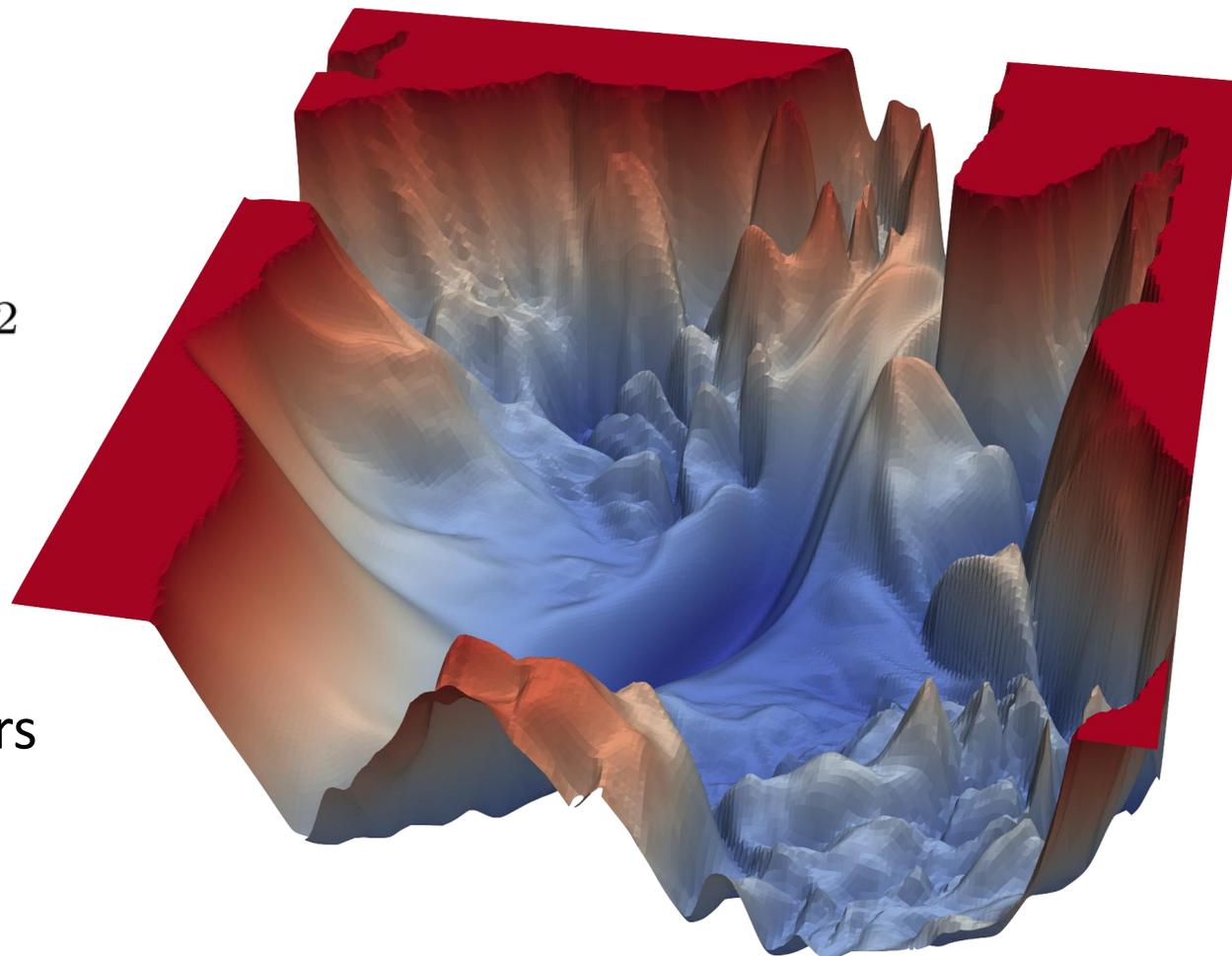(LEs)**
Madras, Atwood, D'Amour
arXiv:1910.09573

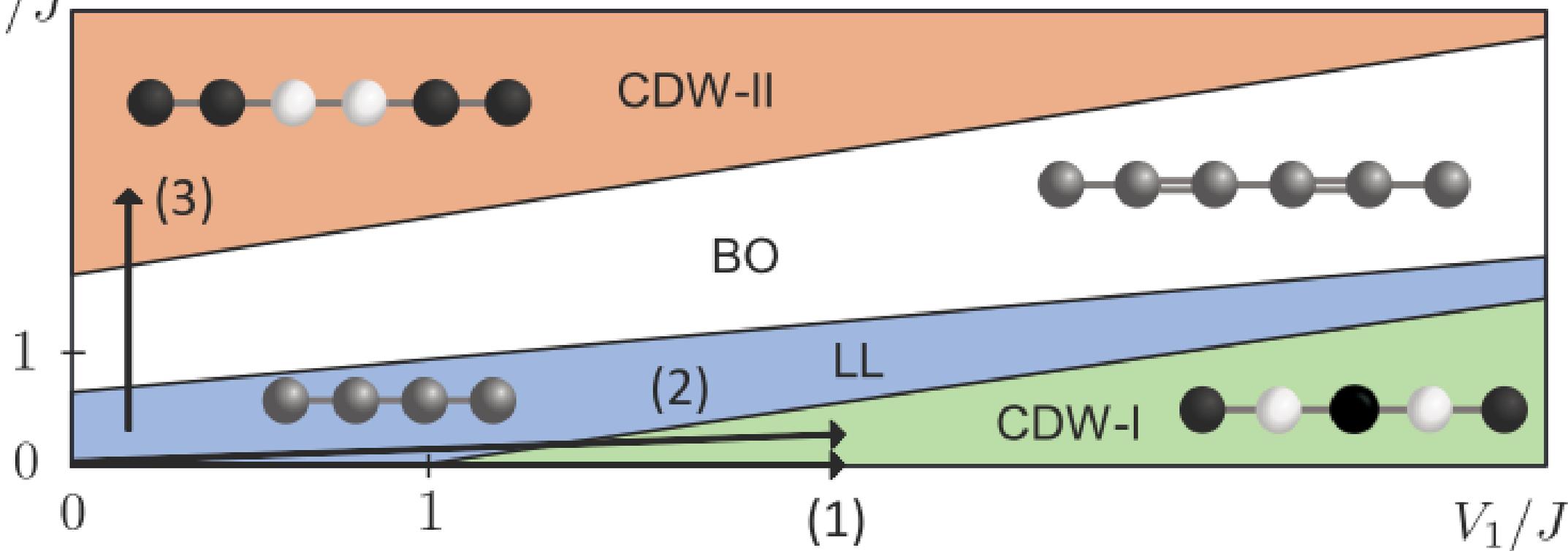# Local-Ensemble-based Extrapolation Score (LEES)

$$\mathcal{E}_m(x_{\text{test}}) = ||U_m^\top \nabla_\theta \mathcal{L}(z_{\text{test}}, \tilde{\theta})||_2$$

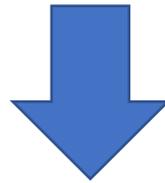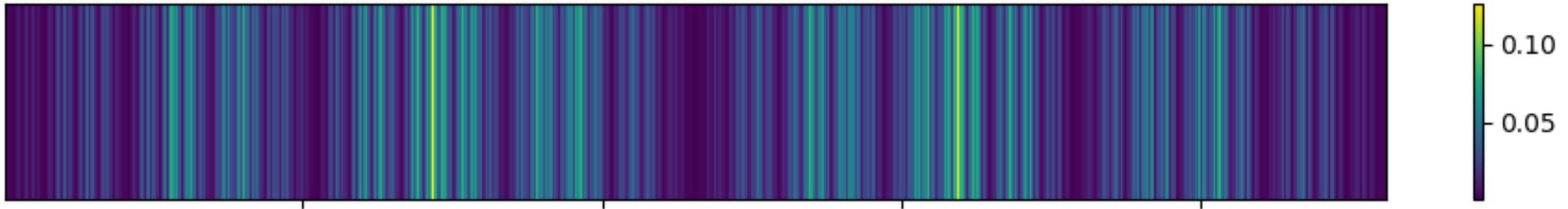matrix of (M – m) Hessian eigenvectors spanning a subspace of low curvature



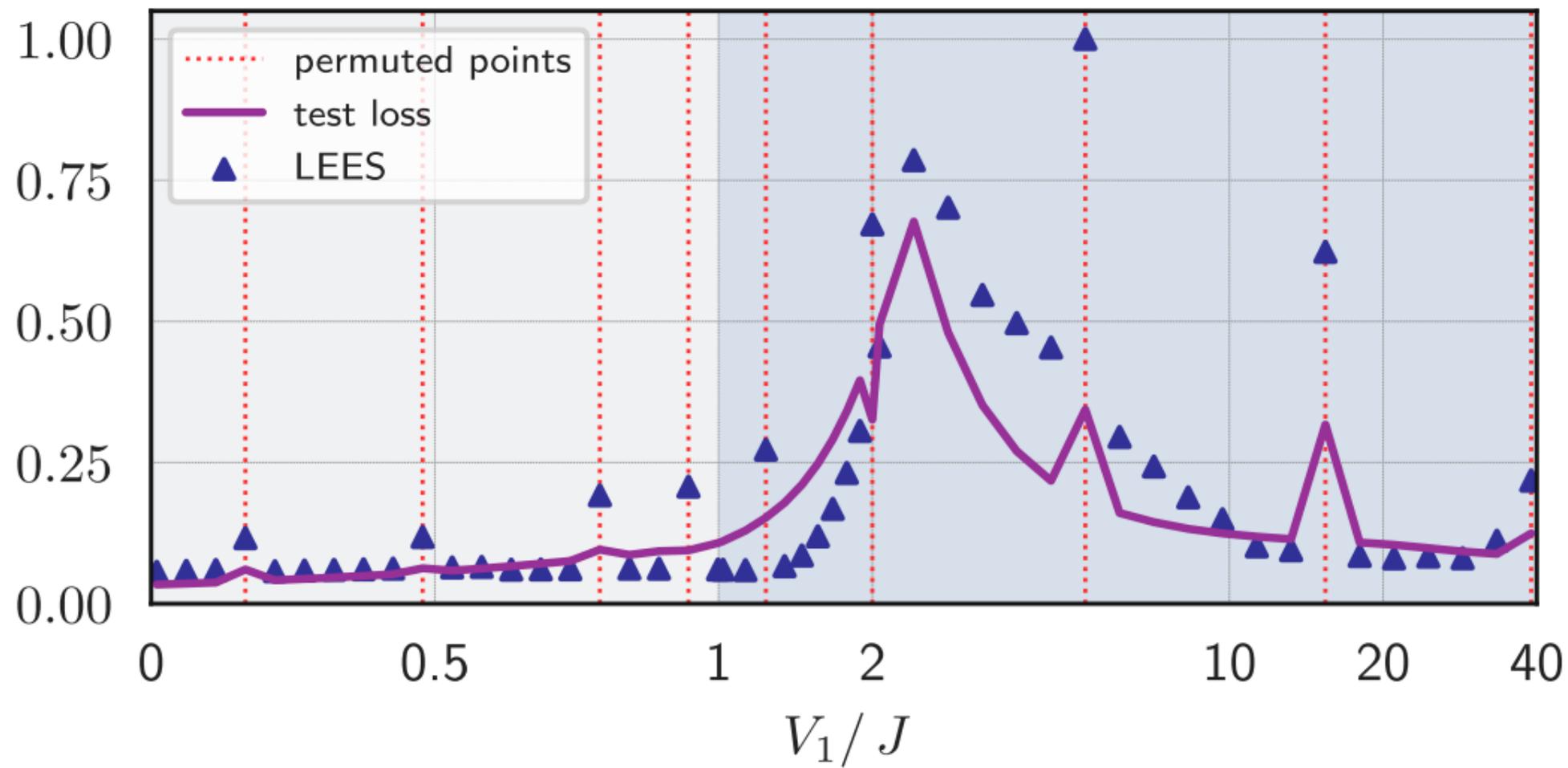D. Madras, J. Atwood, and A. D'Amour, *Detecting Extrapolation with Local Ensembles*. (2019) arXiv:1910.09573.

# Out-Of-Distribution (OOD) test points
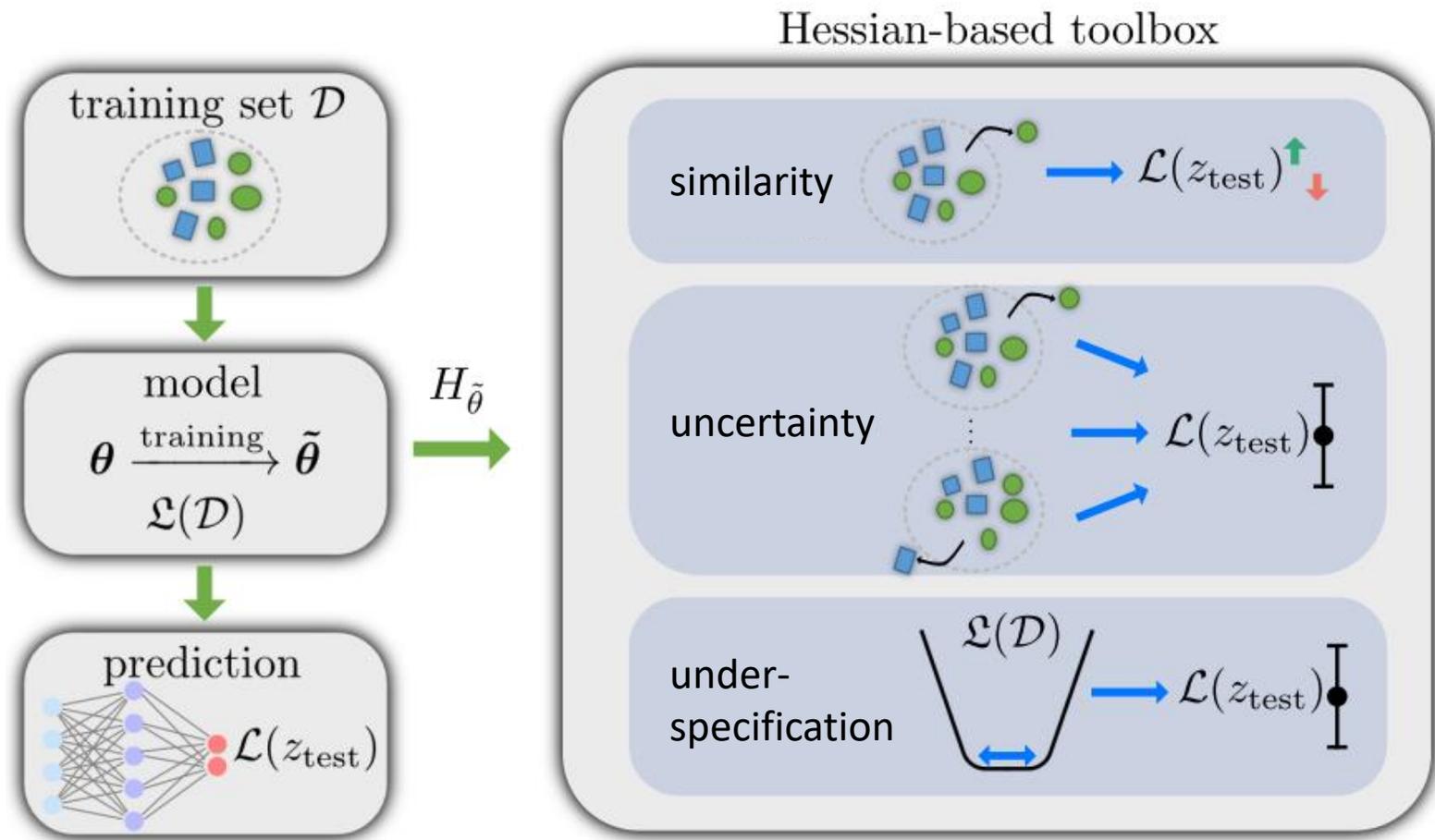


random permutation of eigenvector elements
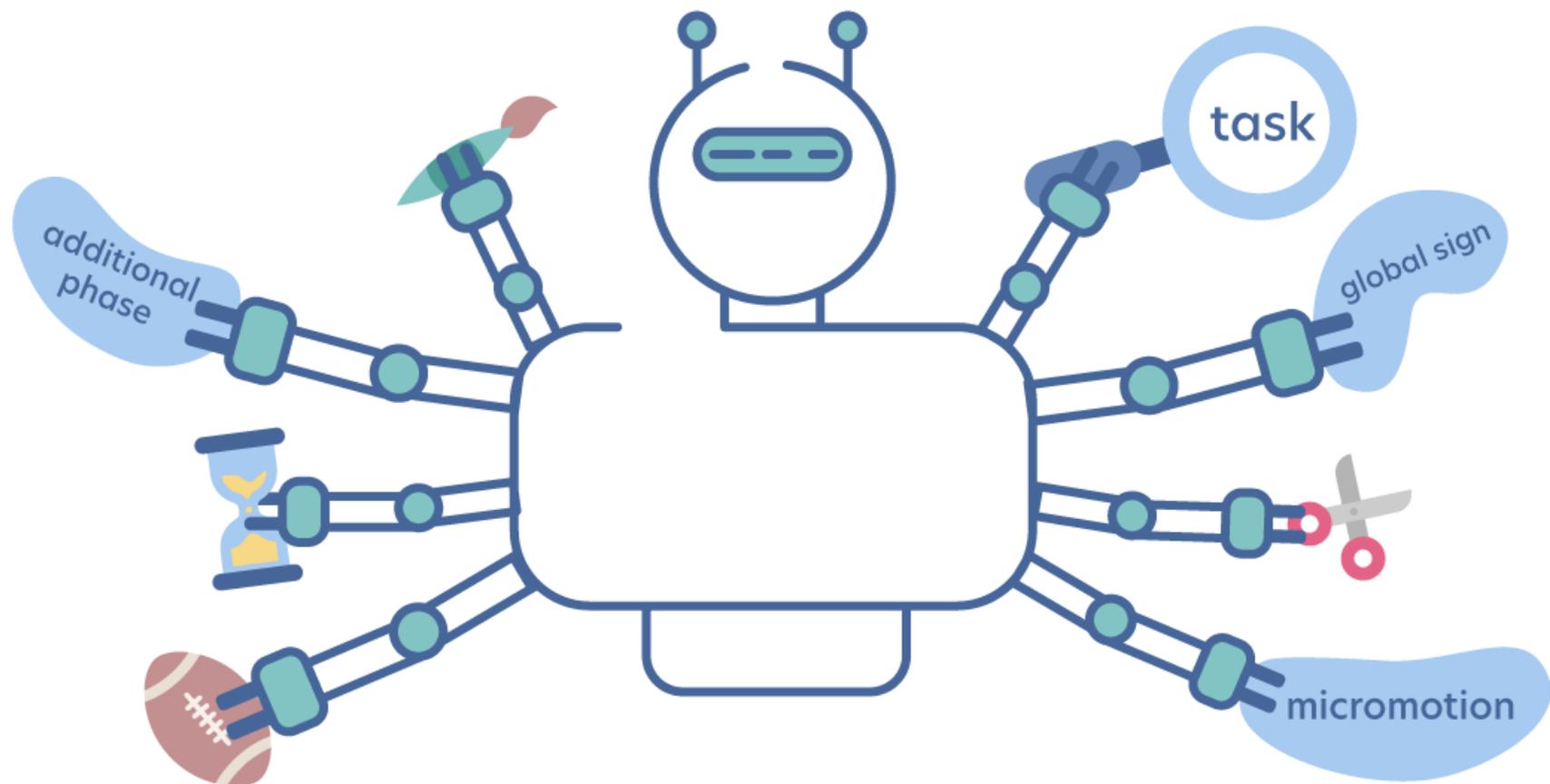
# Local Ensemble-based Extrapolation Score



Dawid et al., *MLST* **3**, 015002 (2022)

# Local Ensemble-based Extrapolation Score



LEES indicates perfectly OOD test points!

# Conclusions



training set $\mathcal{D}$

model
$\boldsymbol{\theta} \xrightarrow{\text{training}} \tilde{\boldsymbol{\theta}}$
$\mathfrak{L}(\mathcal{D})$

$H_{\tilde{\theta}}$

prediction
$\mathcal{L}(z_{\text{test}})$

## Hessian-based toolbox

similarity $\longrightarrow \mathcal{L}(z_{\text{test}})$

uncertainty $\longrightarrow \mathcal{L}(z_{\text{test}})$

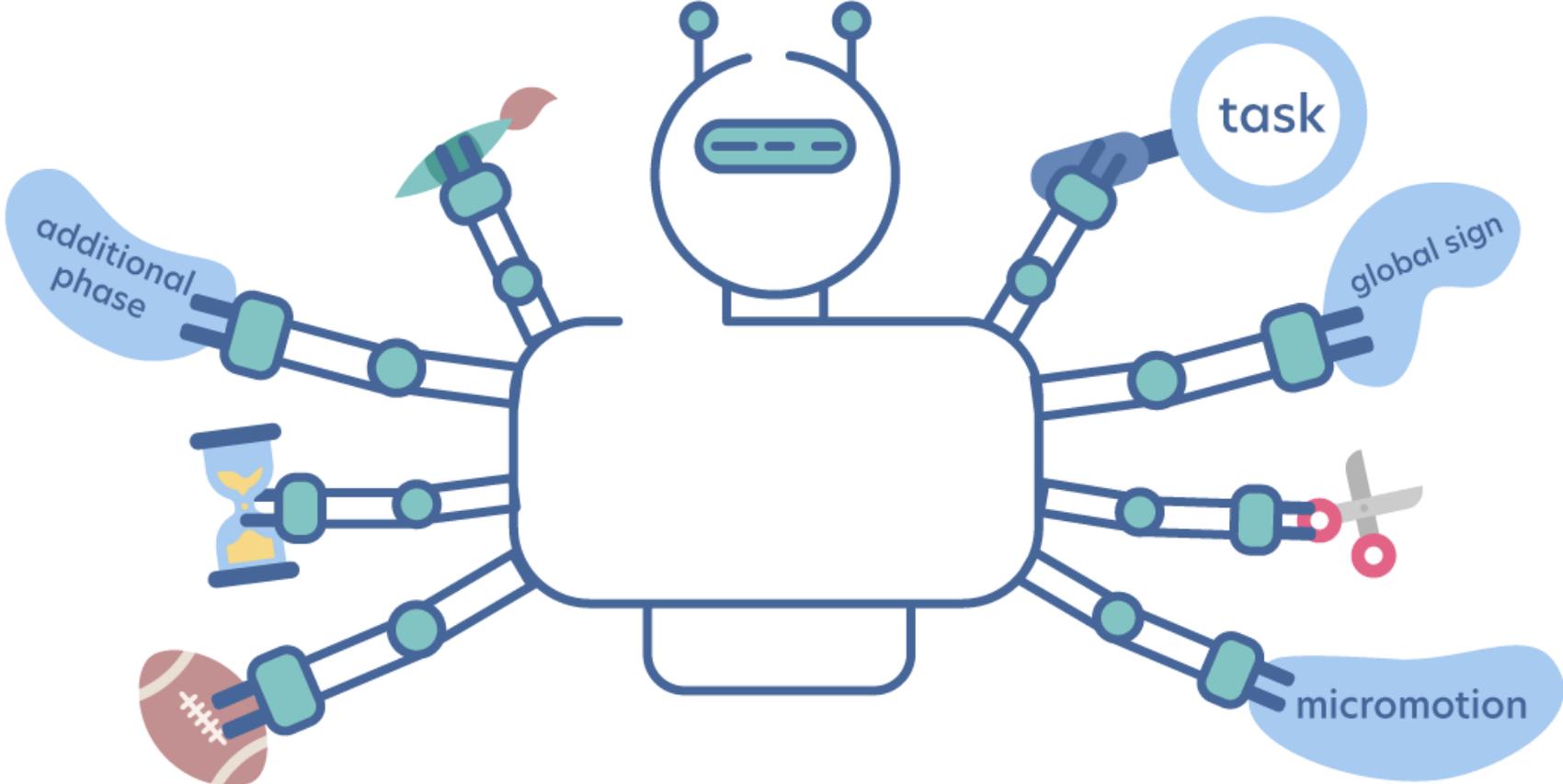under-specification $\mathfrak{L}(\mathcal{D}) \longrightarrow \mathcal{L}(z_{\text{test}})$

More likely to learn physical features
than spurious correlations?

UNIVERSITY OF WARSAW

ICFO
**The Institute of Photonic Sciences**

U·H Universität Hamburg

M. Tomza

P. Huembeli

K. Kottmann

N. Käming

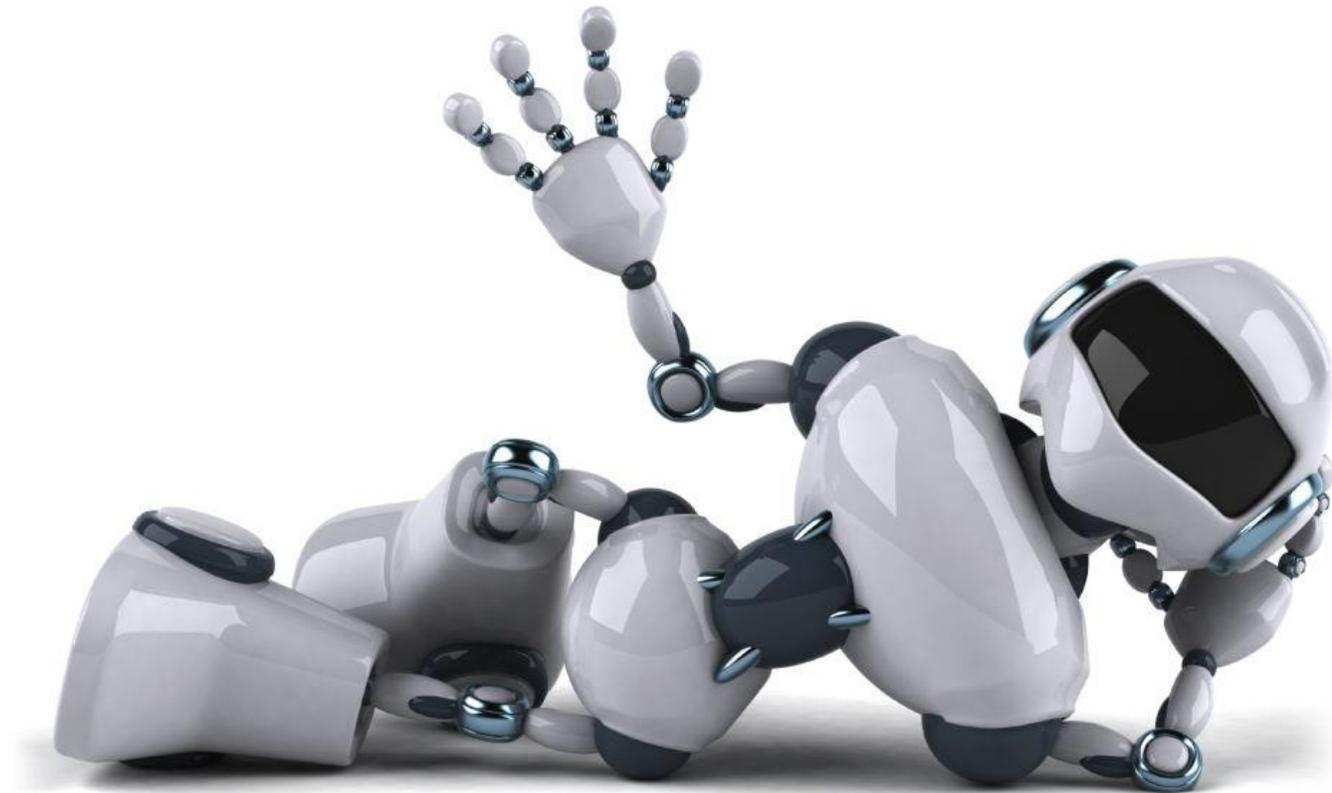A. Dauphin

M. Lewenstein

K. Sengstock

C. Weitenberg

# If you want to know more…

A. Dawid et al (2022)
*Mach. Learn.: Sci. Technol.* **3** 015002 (2022)

N. Käming, A. Dawid et al (2021)
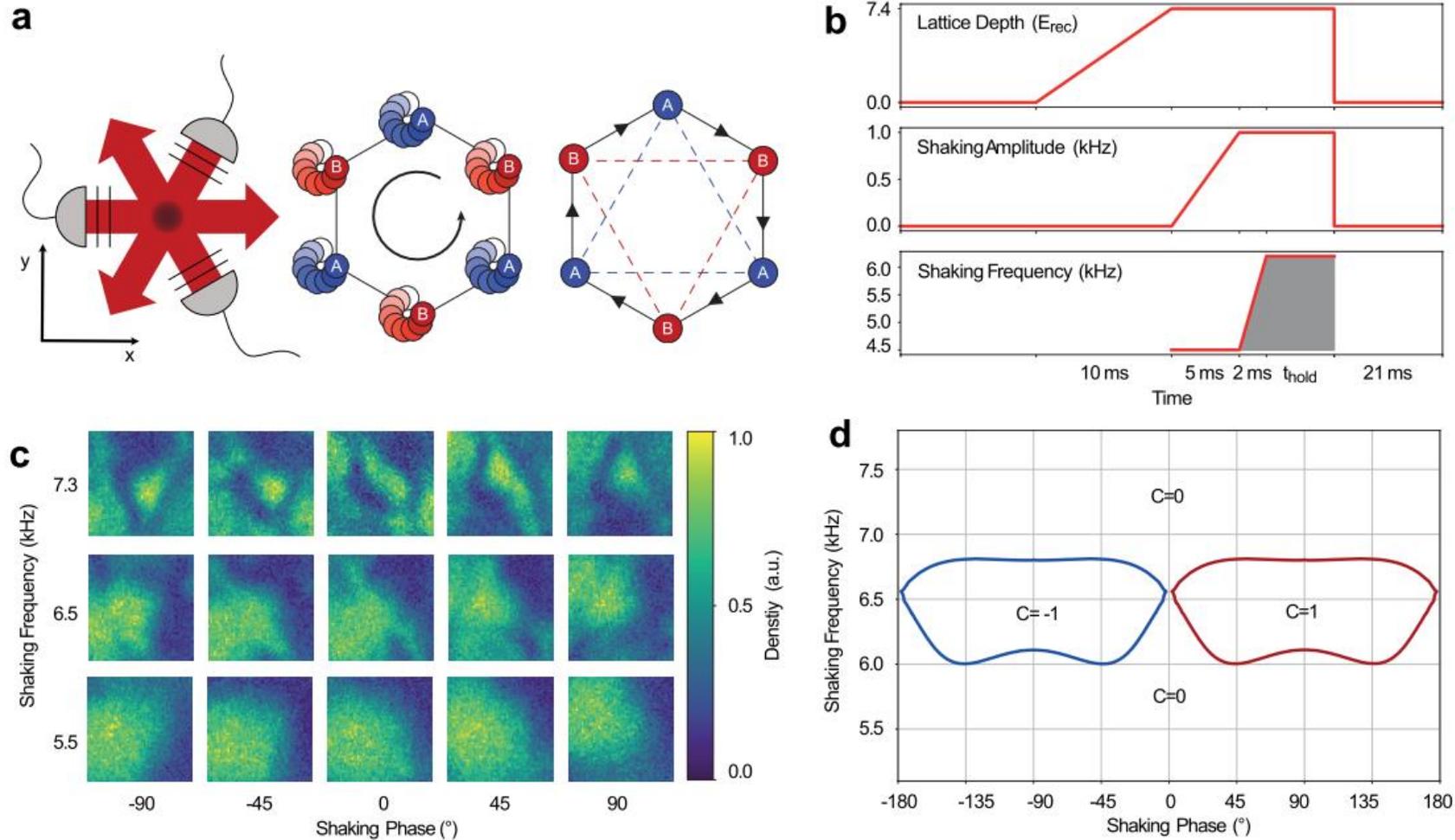*Mach. Learn.: Sci. Technol.* **2** 035037

A. Dawid et al (2020)
*New J. Phys.* **22** 115001

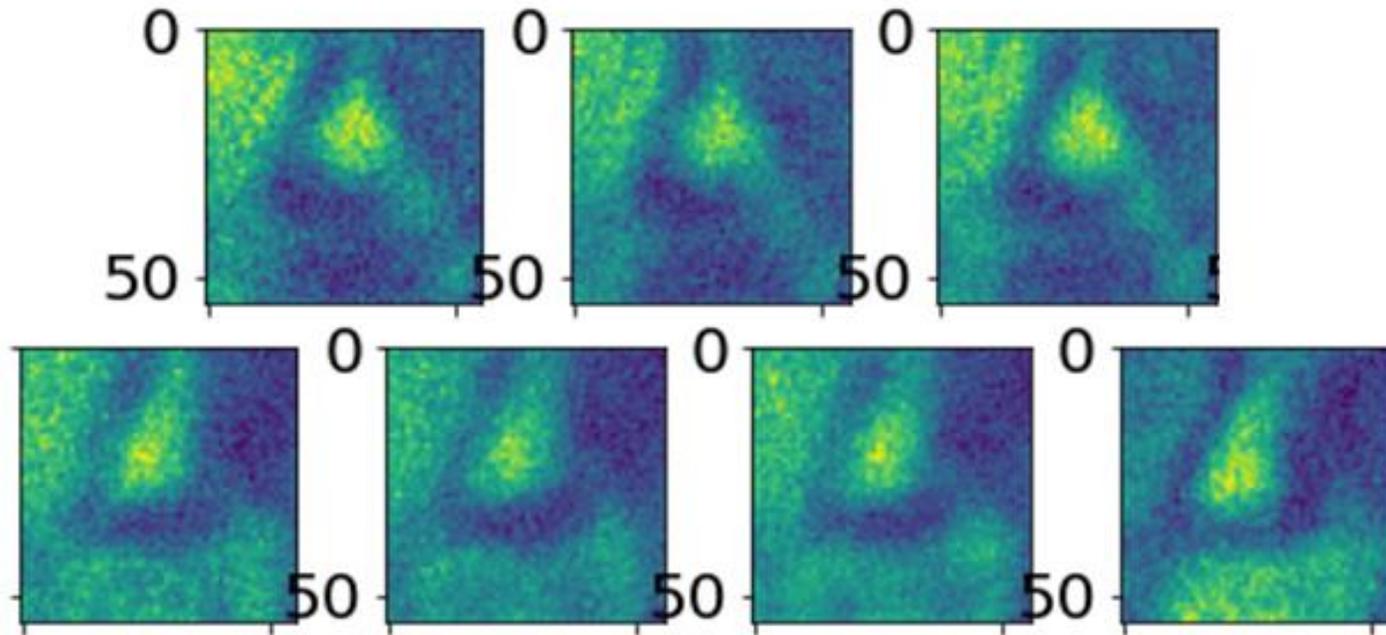/Shmoo137
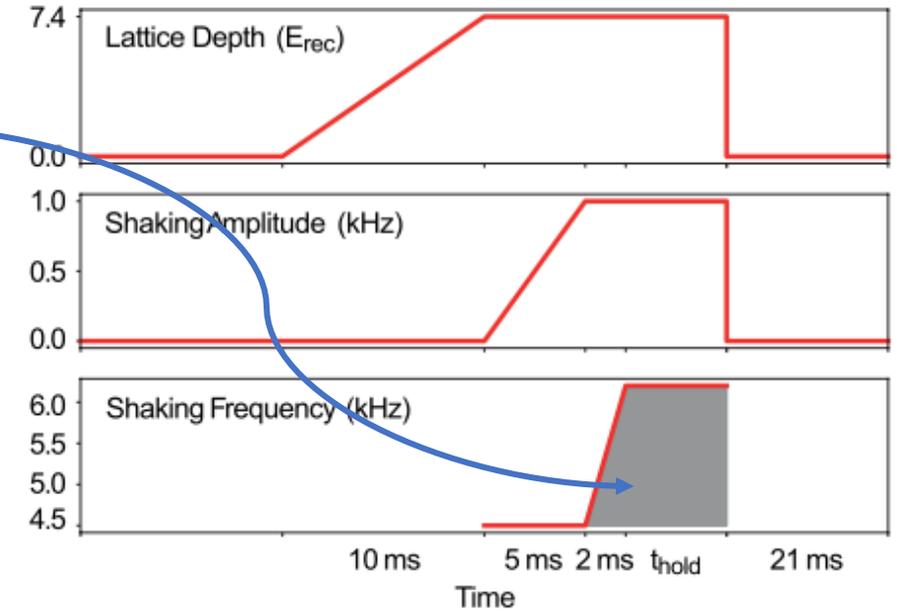
# Thank you for your attention!

# Topological Haldane model
## realized via Floquet-driving of ultracold fermions ($^{40}$K) in a honeycomb lattice

# Micromotion phase

shaking frequency = 7.4 Hz
shaking phase = 90°
different micromotion phases



N. Käming, A. Dawid et al., *MLST* **2** 035037 (2021)

# RUE vs LEES



OOD test points