

Learning to see the hidden world: A perspective on causal representations

Francesco Locatello



What triggers grooming behaviors in ants as a collective hygiene policy?

Why causal models?

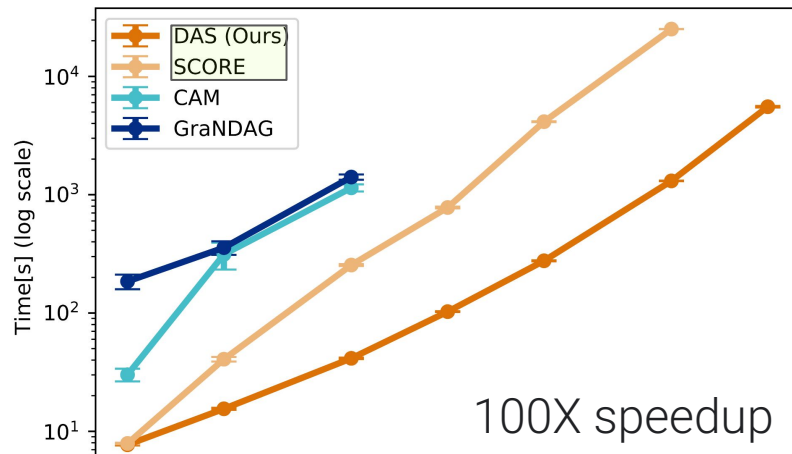
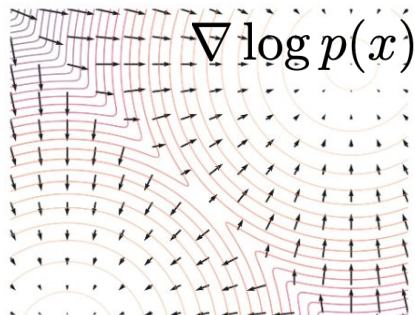
Observational data + **causal model** predict the result of **randomized studies** [1]

$$P(X_1, \dots, X_n) = \prod_i P(X_i | \mathbf{PA}_i)$$

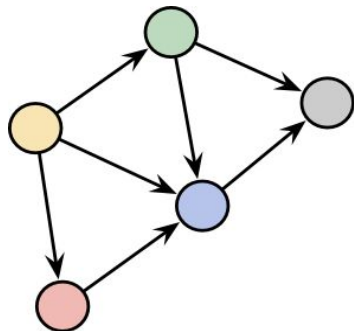
↓

$$P(X_1, \dots, X_n | do(X_j = x)) = \prod_{i \neq j} P(X_i | \mathbf{PA}_i) \delta(X_j = x)$$

Learning the graph



The discovery of the causal order converges linearly as $\exp(-\Theta(\frac{nC_m^2}{\log(m)}))$

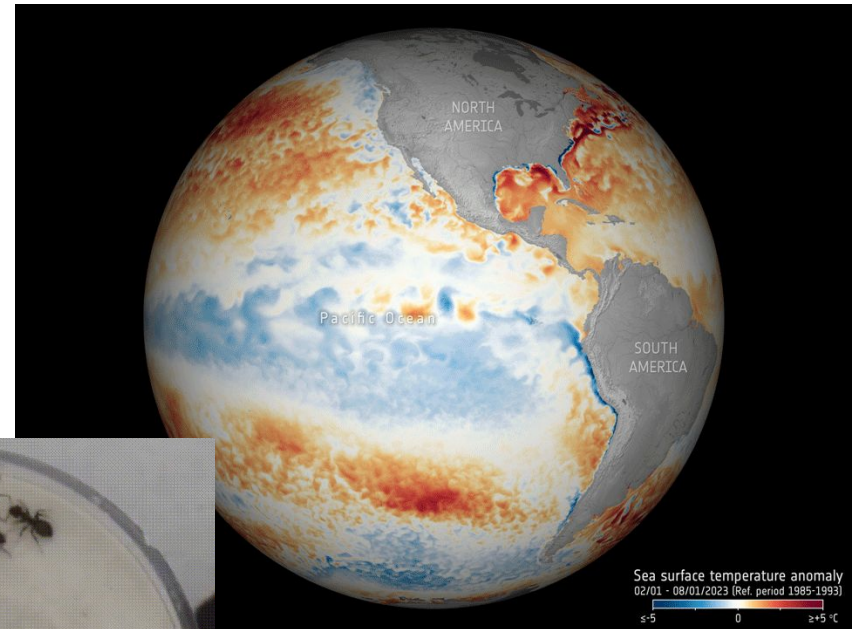
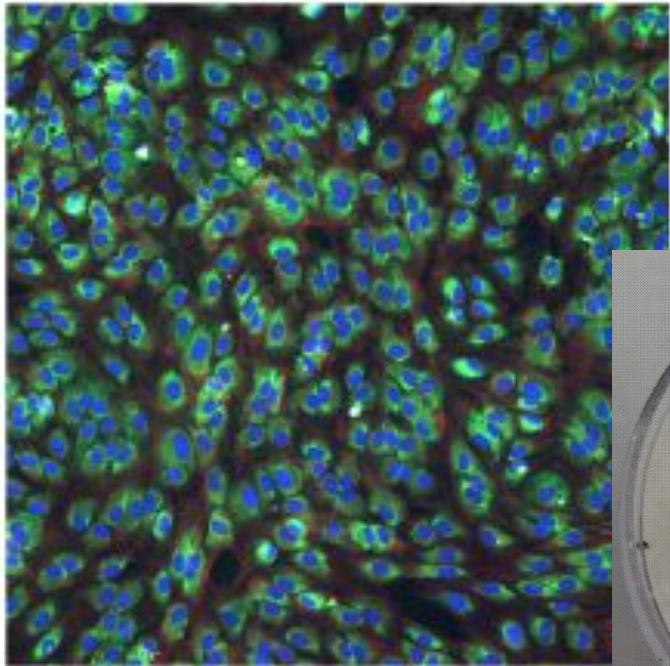


**Causal
Discovery**

infer

x_1	x_2	x_3	x_4	x_5
0.2	0.1	2.2	1.6	4.3
0.3	0.1	2.1	1.9	5.3
...

Data



- Prediction-powered causal inference
- Exploratory causal inference
- Beyond “standard” causal models

**WORLD OF
PERCEPTION**

**WORLD OF
KNOWLEDGE**

What do we want?

$$P(X_1, \dots, X_n) = \prod_i P(X_i | \mathbf{PA}_i)$$

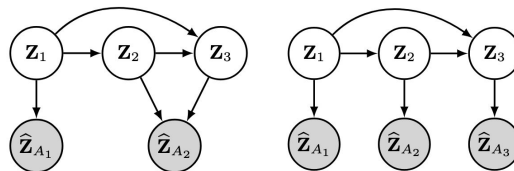
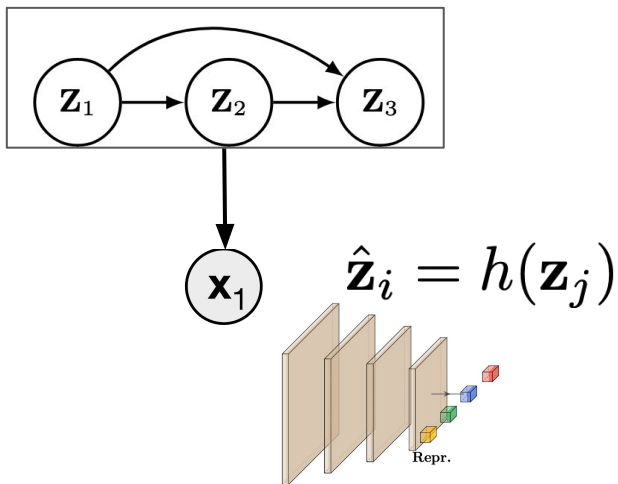
$$\begin{aligned} P(X_1, \dots, X_n | do(X_j = x)) &= \prod_{i \neq j} P(X_i | \mathbf{PA}_i) \delta(X_j = x) \\ &\downarrow \\ &= P(X_{V \setminus \{j, \mathbf{PA}_j\}} | X_j = x, X_{\mathbf{PA}_j}) P(X_{\mathbf{PA}_j}) \end{aligned}$$

Causal estimand that is statistically identified on causal variables is also identifiable from the representation

Idea: Representation should make it easier/possible to extract causal information with some downstream estimator

Measurement perspective of CRL and identifiability

When is a learned representation a valid proxy for a causal variable? “*measurement models*” [1]



Conditions for causal validity of downstream estimate [2]:

- Know $\hat{\mathbf{z}}_{A_j} \perp\!\!\!\perp \mathbf{z}_i \mid \mathbf{z}_{[N] \setminus \{i\}}$
- Estimate is invariant to h

[1] “Learning the structure of linear latent variable models” R. Silva, R. Scheines, C. Glymour, P. Spirtes, and D. M. Chickering. JMLR, 2006

[2] “The Third Pillar of Causal Analysis? A Measurement Perspective on Causal Representations”, Yao* and Huang*, Cadei, Zhang, L; NeurIPS 2025

[3] “Self-supervised Representation Learning Provably Isolates Content From Style”, von Kügelgen*, Sharma*, Gresele*, Brendel, Schölkopf+, Besserve+, L+ NeurIPS 2021

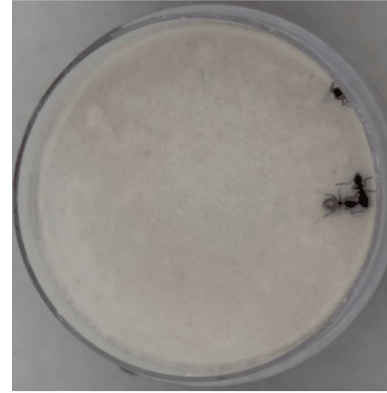
Exemplary pipeline in experimental ecology

Design and perform experiment

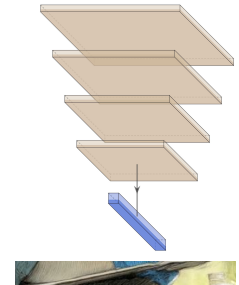
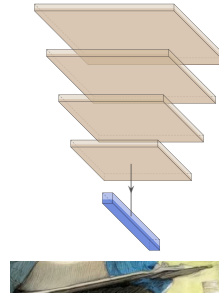


Collect data

Control Group
Batch: 2, Position: 4, Time: 2m30s



Treated Group
Batch: 1, Position: 6, Time: 2m45s



ISTAnt dataset



No grooming



Blue to Focal
grooming

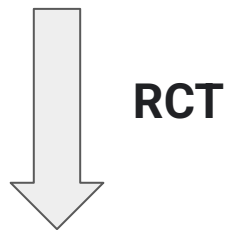


What? First benchmark to estimate causal effects from real world ecological videos collected in a randomized controlled trial.

Why is it unique?

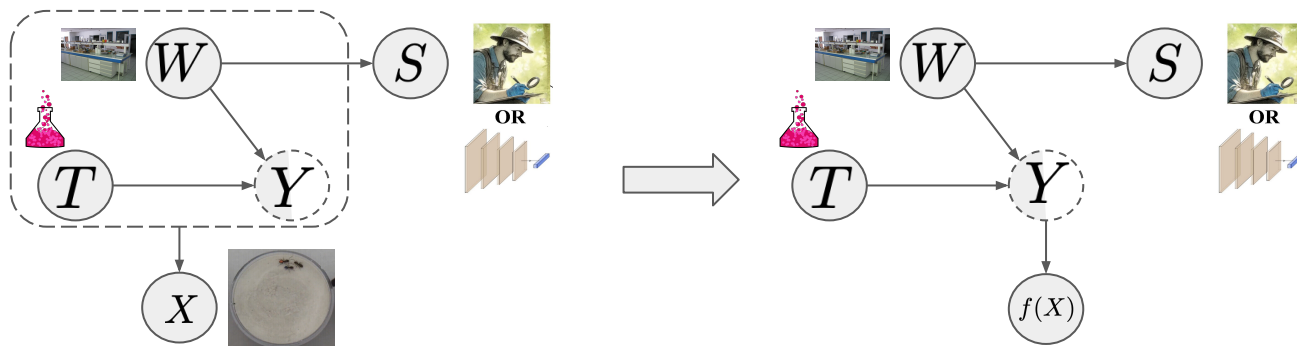
Other benchmarks stop at validating statistical accuracy. We estimate causal effects from real world ecological videos collected in a randomized controlled trial

$$ATE := \mathbb{E}[Y|do(T = 1)] - \mathbb{E}[Y|do(T = 0)].$$



$$AD := \mathbb{E}[Y|T = 1] - \mathbb{E}[Y|T = 0].$$

A measurement perspective on the problem

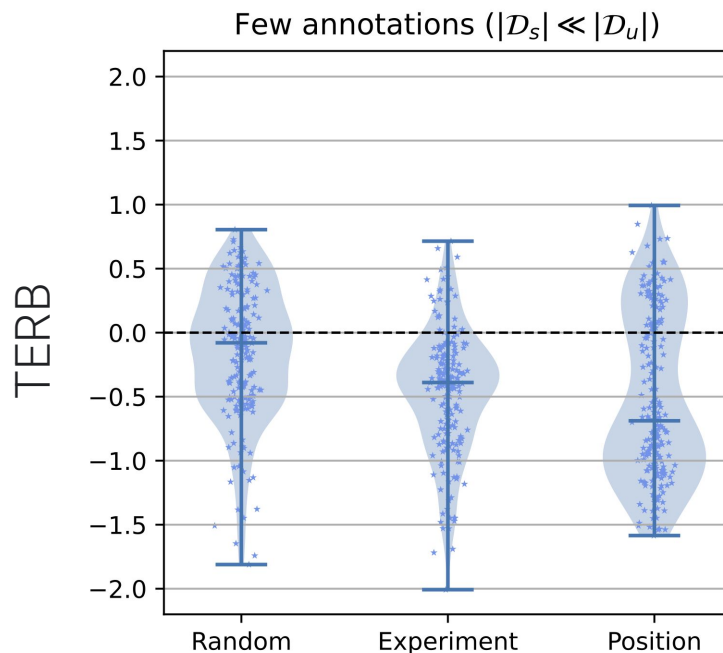


$$TEB := \left(\underbrace{\mathbb{E}_{\mathbf{X}|do(T=1)}[f(\mathbf{X})] - \mathbb{E}_{Y|do(T=1)}[Y]}_{\text{Interventional Bias under Treatment}} \right) - \left(\underbrace{\mathbb{E}_{\mathbf{X}|do(T=0)}[f(\mathbf{X})] - \mathbb{E}_{Y|do(T=0)}[Y]}_{\text{Interventional Bias under Control}} \right)$$

Errors come from:

- Selection bias: which samples are labelled?
- Pre-training data
- Discretization bias

Problem 1: What data to label?



Implication: Sampling choice matters, but random sampling is not always possible

Optimize for invariance in unified CRL

$$\iota(\mathbf{z}_A) = \iota(\tilde{\mathbf{z}}_A) \Leftrightarrow \mathbf{z}_A \sim_{\iota} \tilde{\mathbf{z}}_A.$$

Two vectors have the same projection onto the quotient induced by the equivalence relationship

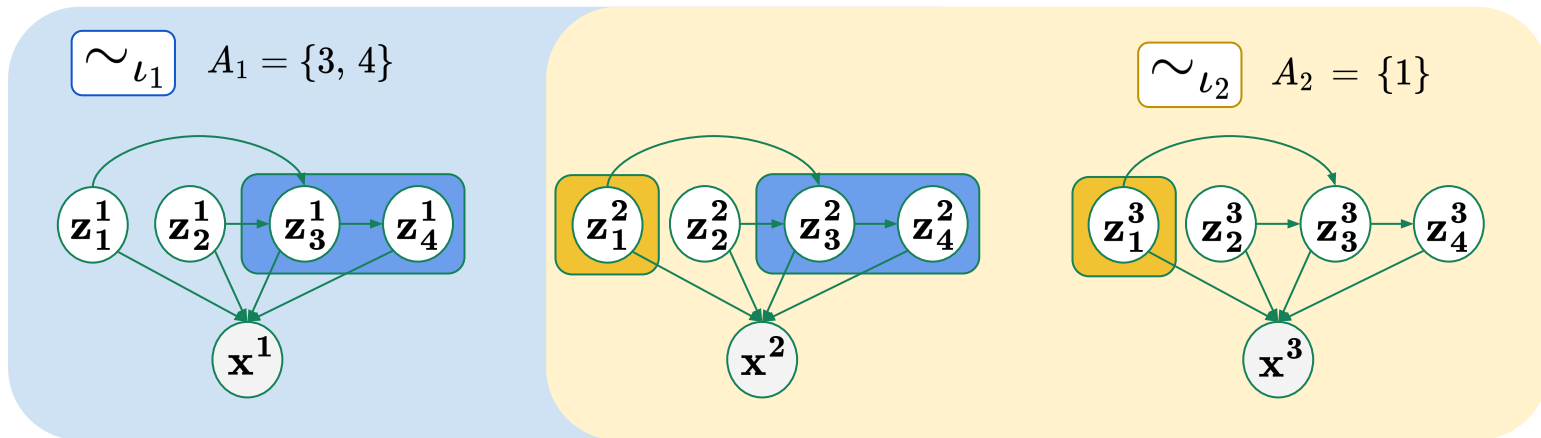
They belong to the same equivalence class

How? Assume access to multiple non-i.i.d. groups of sample. All samples in the same group are *equivalent* in some sense

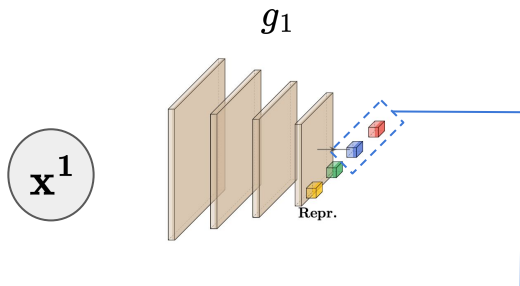
Why? Because you believe isolating these invariances is relevant to a task

Key idea for unified CRL: invariance principle

$$\iota(\mathbf{z}_A) = \iota(\tilde{\mathbf{z}}_A) \Leftrightarrow \mathbf{z}_A \sim_{\iota} \tilde{\mathbf{z}}_A.$$



Key idea for unified CRL: learning

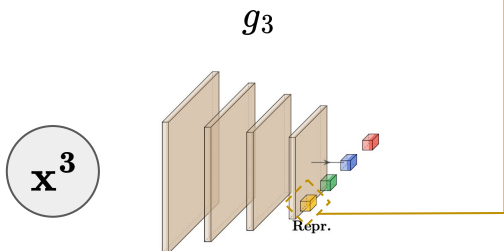


- Invariance Constraint:

$$g_1(\mathbf{x}^1) \sim_{\iota_1} g_2(\mathbf{x}^2)$$

$$g_2(\mathbf{x}^2) \sim_{\iota_2} g_3(\mathbf{x}^3)$$

Unify using a single “language” **31** different identification results from 28 paper!

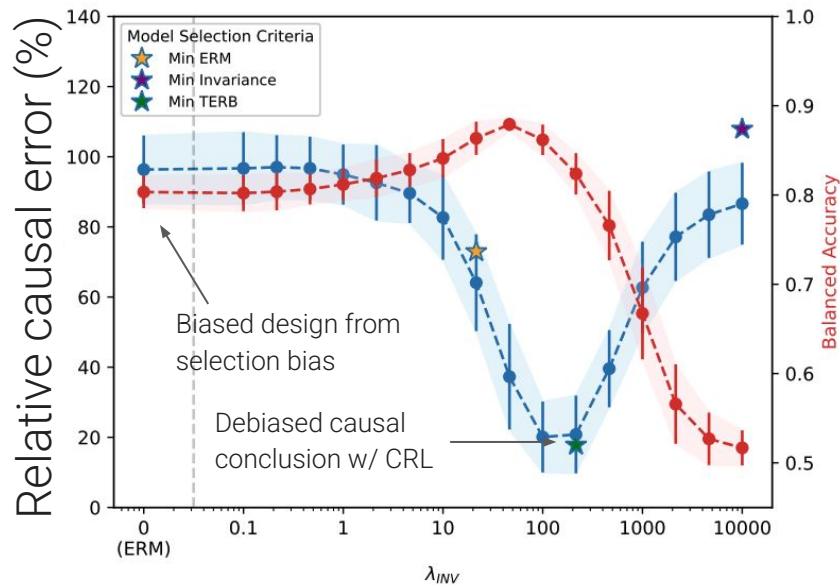


$$I(\mathbf{z}_{A_2}^k, g_k(\mathbf{x}^k)) = H(\mathbf{z}_{A_2}^k) \quad k = 2, 3$$

- **Result:** Smooth encoders satisfying the two constraints block-identify invariant components

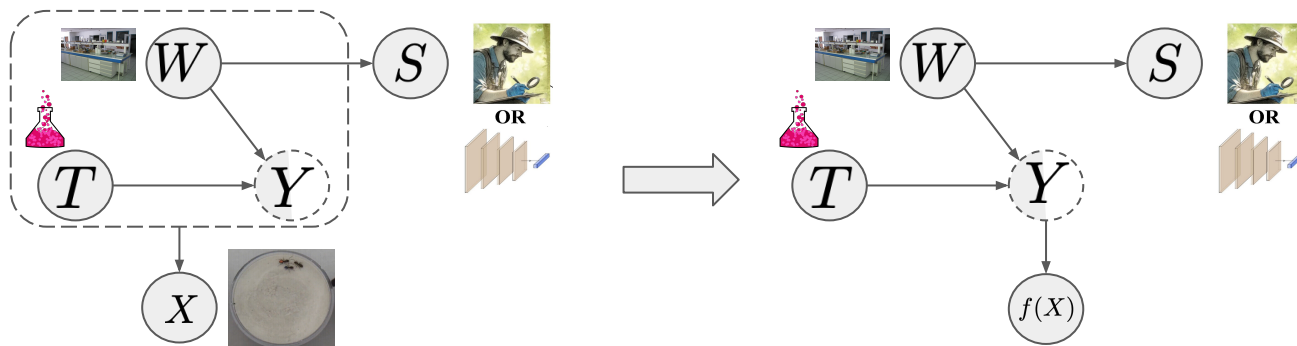
$$\hat{\mathbf{z}}_i = h(\mathbf{z}_j) \quad \mathbf{z}_{A_1}^{1,2} \quad \mathbf{z}_{A_2}^{2,3}$$

Debiasing with CRL and the invariance principle



Idea: Assume invariant representation across experiment settings

A measurement perspective on the problem

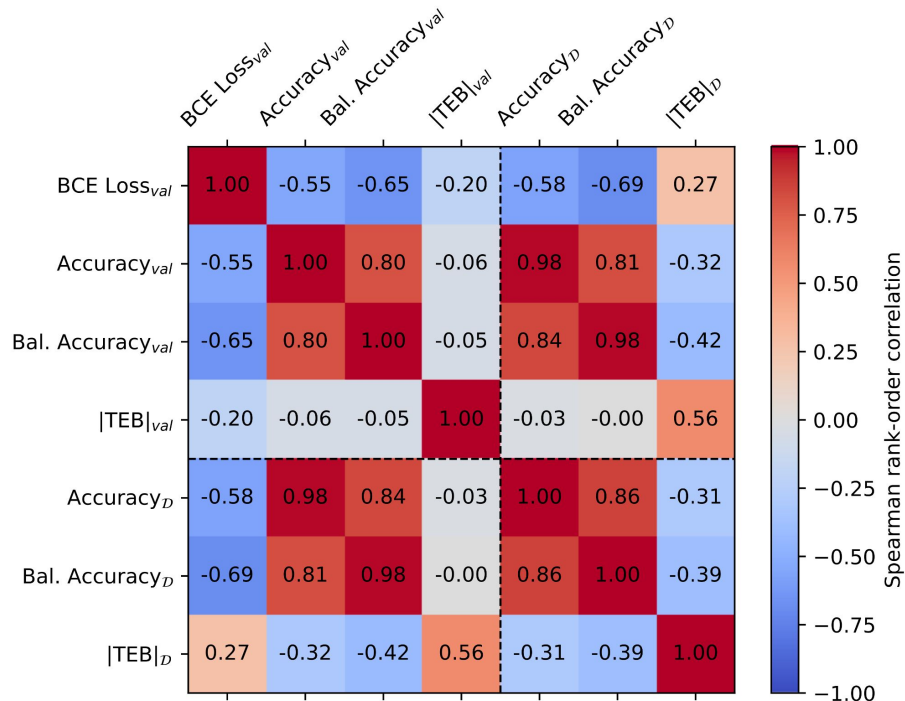
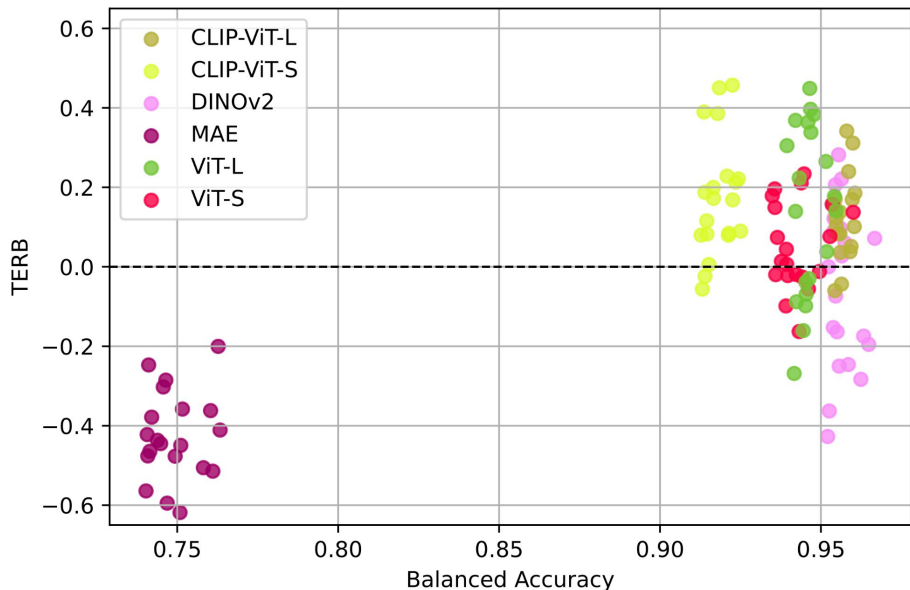


$$TEB := \left(\underbrace{\mathbb{E}_{\mathbf{X}|do(T=1)}[f(\mathbf{X})] - \mathbb{E}_{Y|do(T=1)}[Y]}_{\text{Interventional Bias under Treatment}} \right) - \left(\underbrace{\mathbb{E}_{\mathbf{X}|do(T=0)}[f(\mathbf{X})] - \mathbb{E}_{Y|do(T=0)}[Y]}_{\text{Interventional Bias under Control}} \right)$$

Errors come from:

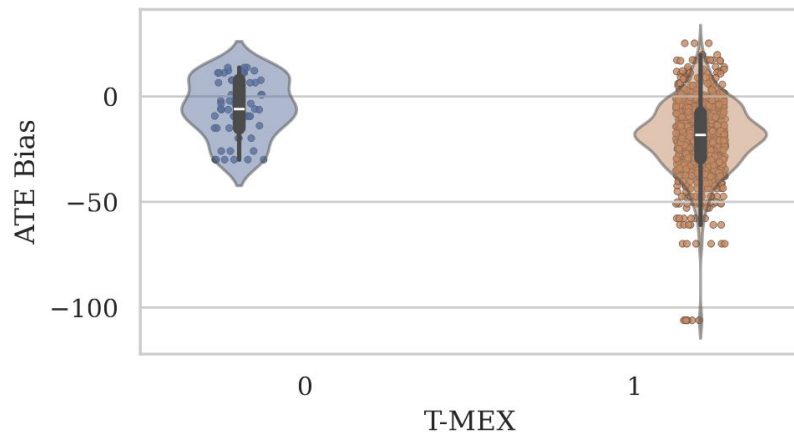
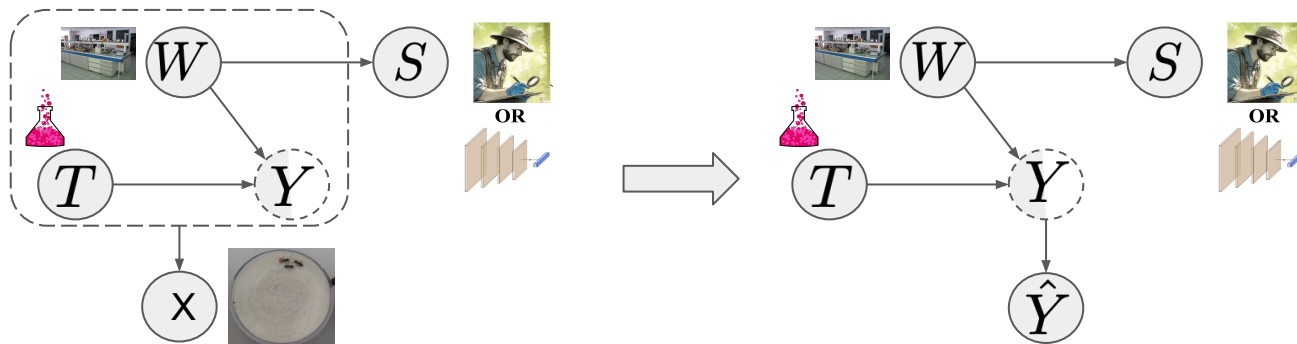
- Selection bias: which samples are labelled?
- Pre-training data
- Discretization bias

Problem 2: model choice



Implication: Models have different TERB and accuracy is not a good indicator of downstream causal performance. TEB on validation data works best

Selecting valid measurements with T-Mex



Conditions for causal validity of downstream estimator:

- $\hat{\mathbf{z}}_{A_j} \perp\!\!\!\perp \mathbf{z}_i \mid \mathbf{z}_{[N] \setminus \{i\}}$
- Estimator is invariant to h

Note: T-Mex doesn't need data to follow trial distribution (as long as the conditional indep. still holds)

Problem 3: Postprocessing

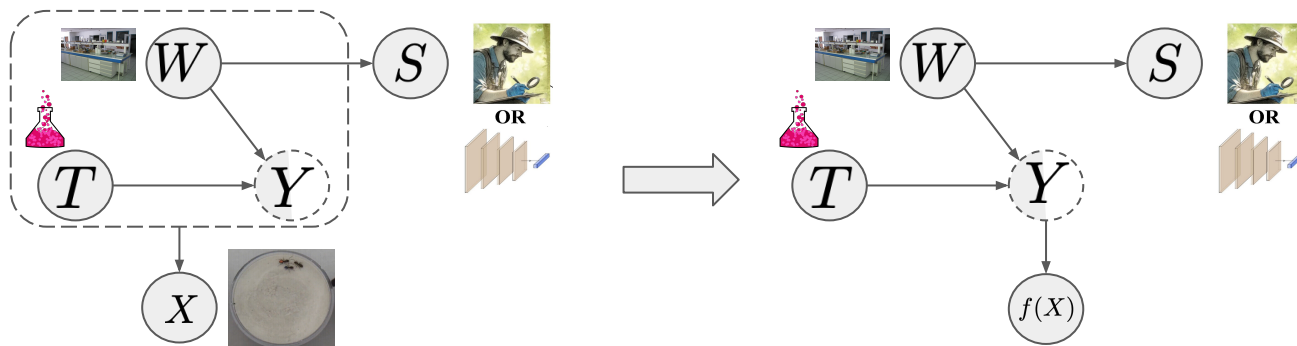
A common choice in ML is to threshold the predictions to make them binary.

Theorem [informal]: This choice can introduce bias.

$$\mathcal{H}_0 : \mathbb{E}[|\text{TEB}(f)|] = \mathbb{E}[|\text{TEB}(\mathbb{1}_{[0.5,1]}(f))|] \quad vs \quad \mathcal{H}_1 : \mathbb{E}[|\text{TEB}(f)|] < \mathbb{E}[|\text{TEB}(\mathbb{1}_{[0.5,1]}(f))|]$$

T-test rejects null hypothesis with $p \sim 10^{-25}$

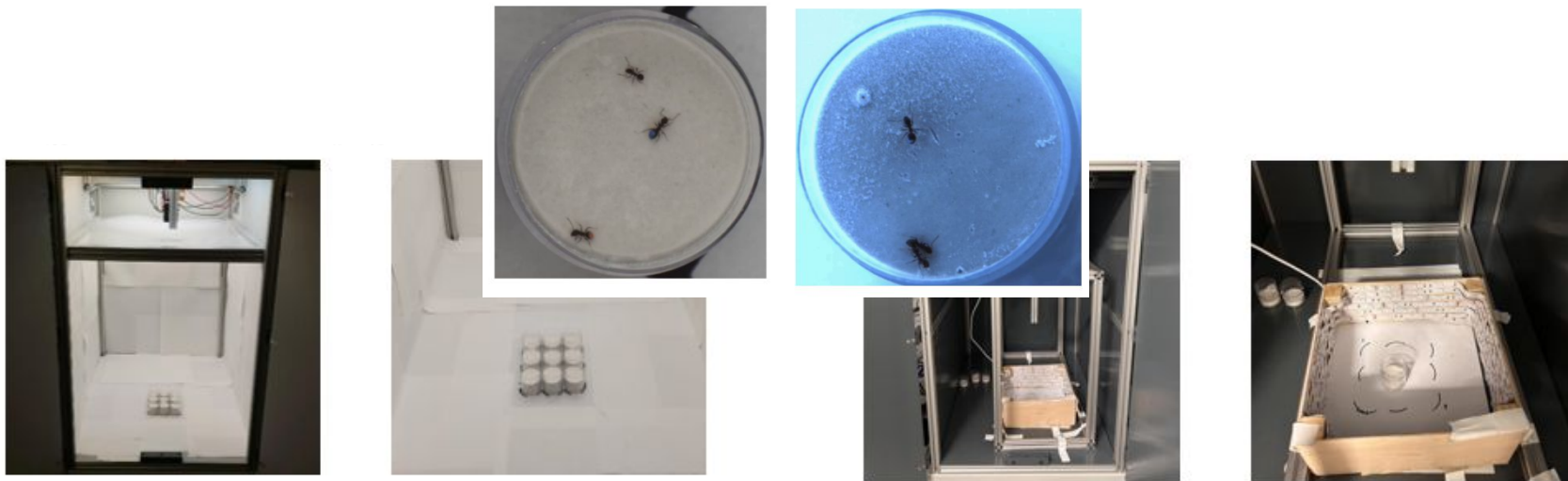
How do we optimize for validity



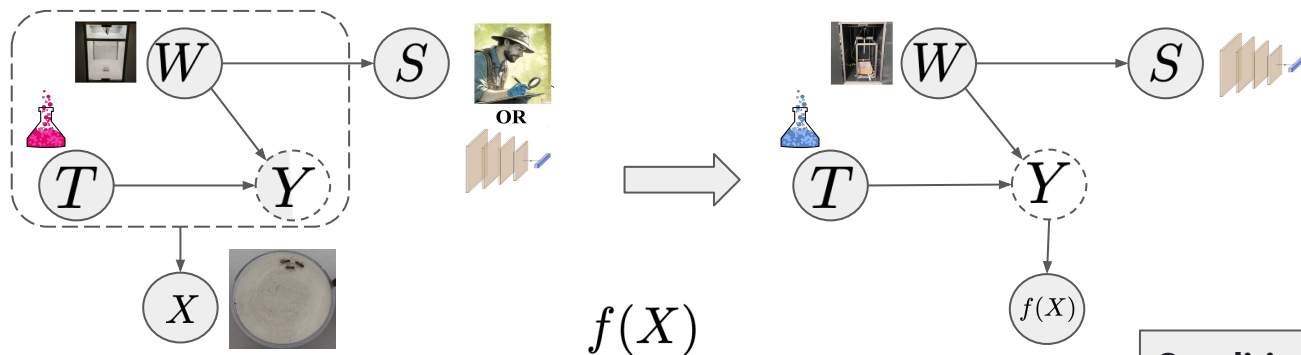
$$TEB := \underbrace{\left(\mathbb{E}_{\mathbf{X}|do(T=1)}[f(\mathbf{X})] - \mathbb{E}_{Y|do(T=1)}[Y] \right)}_{\text{Interventional Bias under Treatment}} - \underbrace{\left(\mathbb{E}_{\mathbf{X}|do(T=0)}[f(\mathbf{X})] - \mathbb{E}_{Y|do(T=0)}[Y] \right)}_{\text{Interventional Bias under Control}}$$

Thm: Conditional calibration $\mathbb{E}[Y - f(X)|W, T] = 0$ implies valid estimates with *correct confidence intervals* using AIPW. Conditional independence of the measurement model implies conditional calibration.

Zero-shot transfer across experiments



Validity across experiments



$$\min_{h, \phi} \mathbb{E}_{\mathbb{P}^{e_1}} [\mathcal{L}(Y, h \circ \phi(\mathbf{X}))]$$

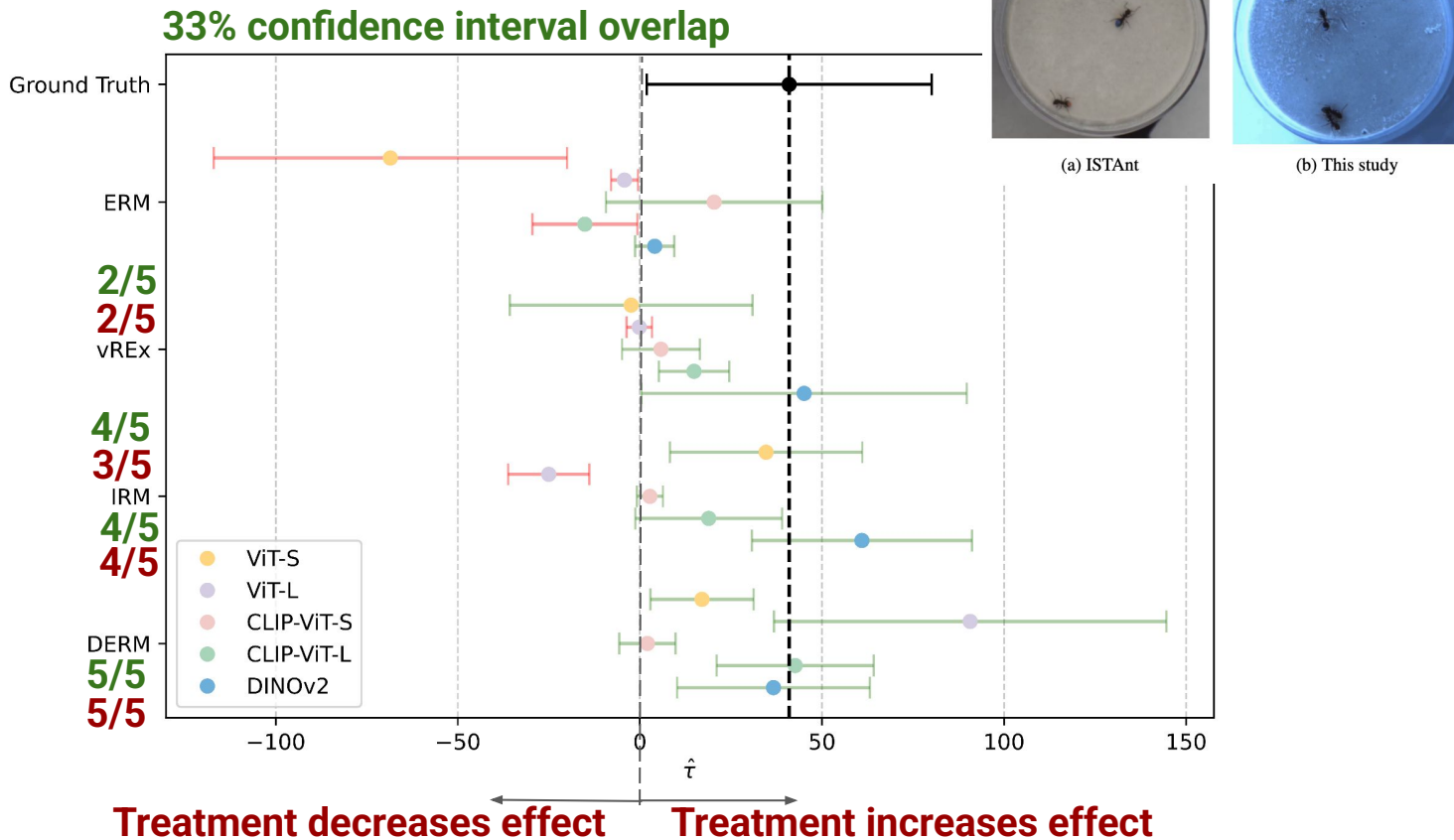
$$\text{s.t. } \phi(\mathbf{X}) \perp\!\!\!\perp \mathbf{Z} | Y = y \quad \forall y \in \mathcal{Y}$$

Conditions for causal validity of downstream estimate:

- Know $\hat{\mathbf{z}}_{A_j} \perp\!\!\!\perp \mathbf{z}_i | \mathbf{z}_{[N] \setminus \{i\}}$
- Estimate is invariant to h

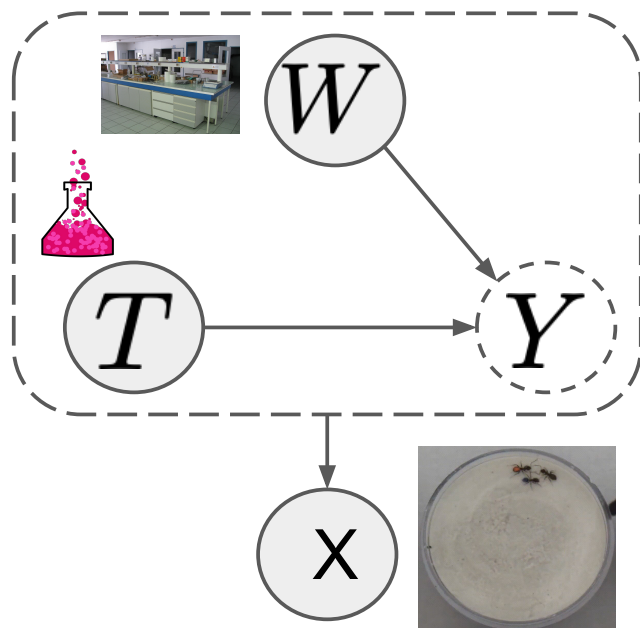
Thm: if a representation is valid on a training experiment and transfers to a target experiment while satisfying $\phi(\mathbf{X}) \perp\!\!\!\perp \mathbf{Z} | Y = y$, then it remains causally valid.

Results

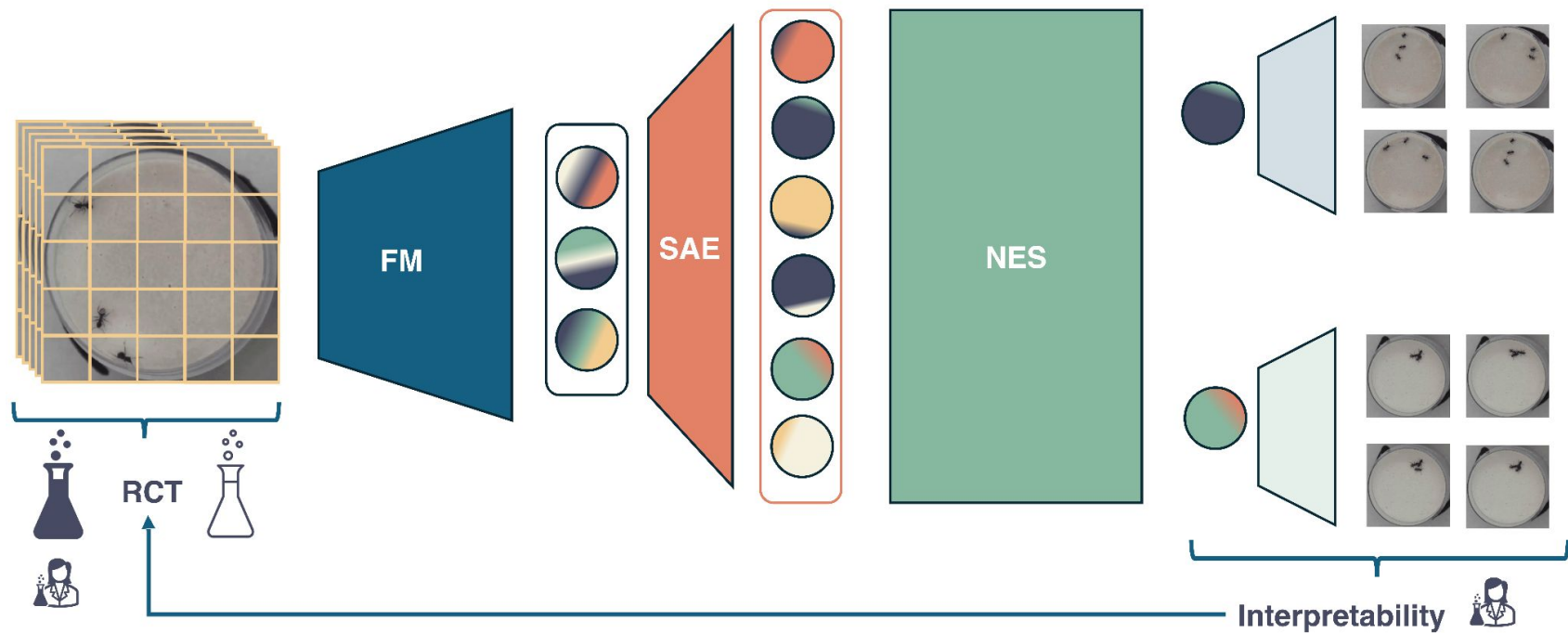


Exploratory causal inference

Exploratory causal inference

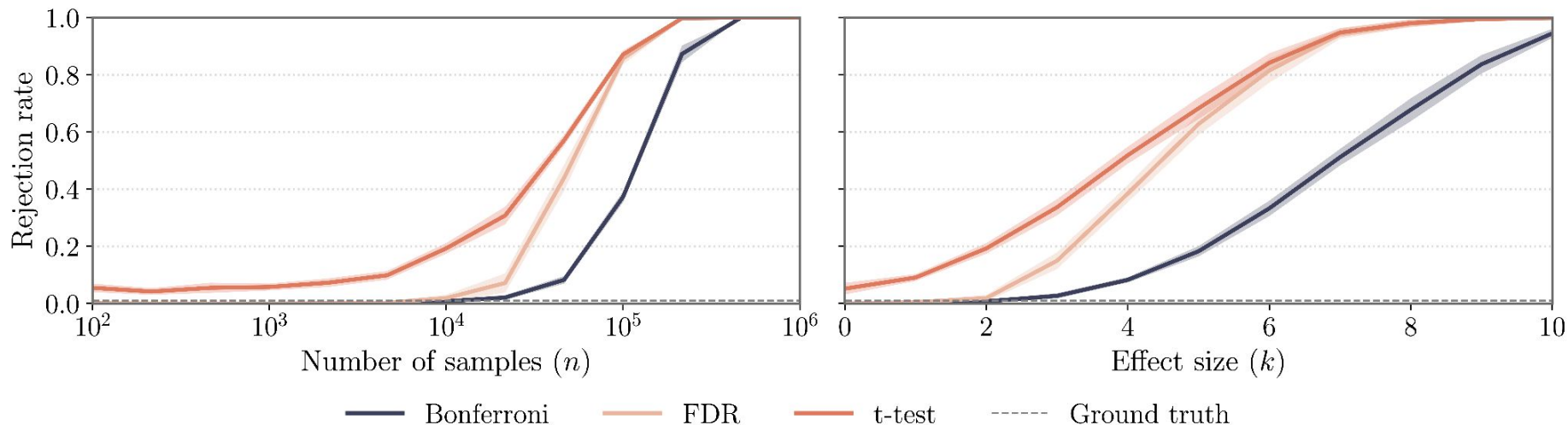


Deep networks as data-driven measurement devices



Hypothesis testing challenges: the paradox of ECI

NES: Consistent recursive testing procedure that corrects the entanglement from previously discovered hypotheses.



Idea: With powerful tests (strong effects or large sample sizes), all correlations are statistically significant. Small entanglement \rightarrow all neurons are individual effects.


Results on synthetic trials

Most activated images for Neuron 38



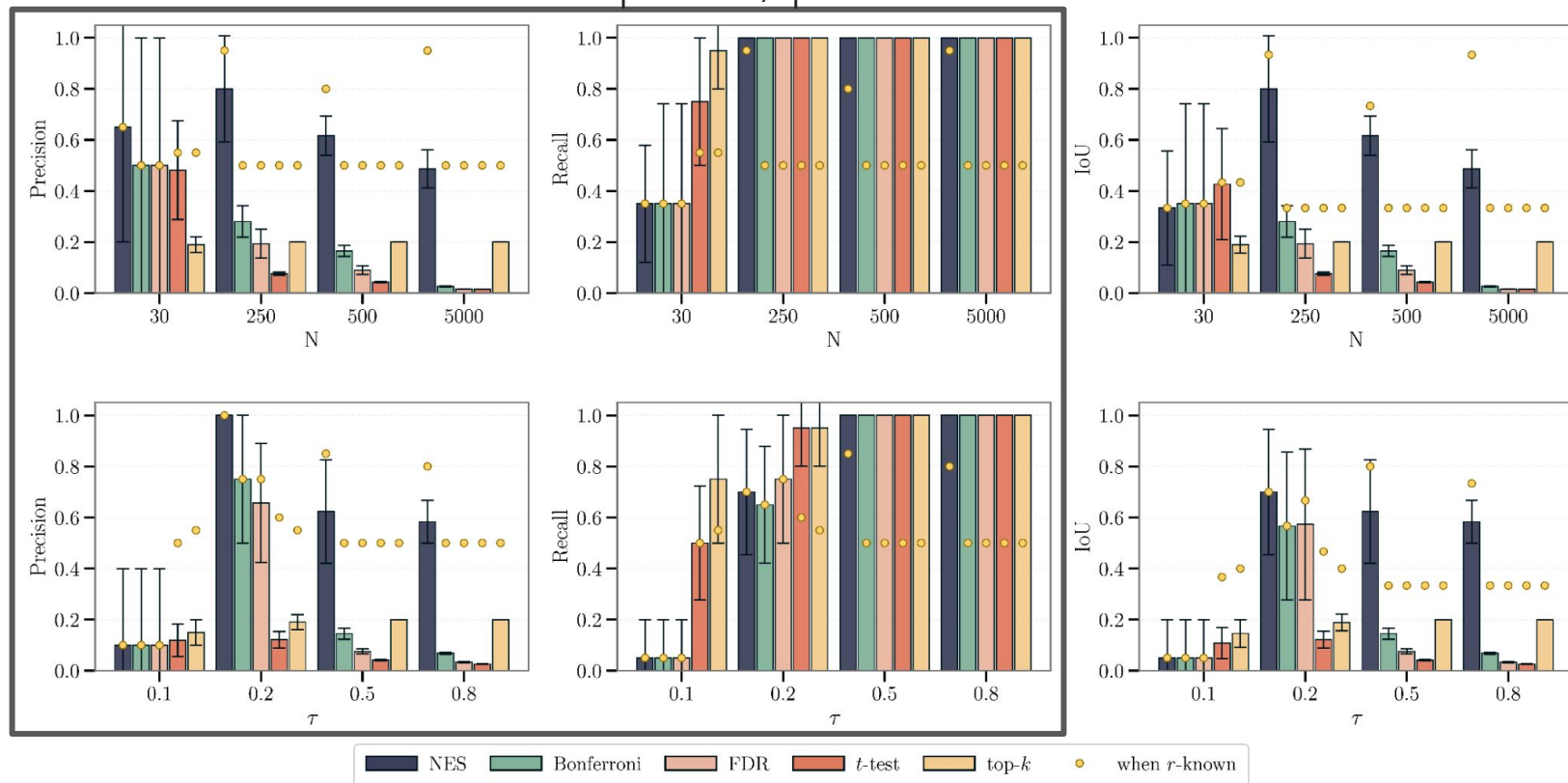
Most activated images for Neuron 6051



- | | |
|--|--|
|  stronger activation (left panel) |  stronger activation (right panel) |
|  weaker activation (left panel) |  weaker activation (right panel) |

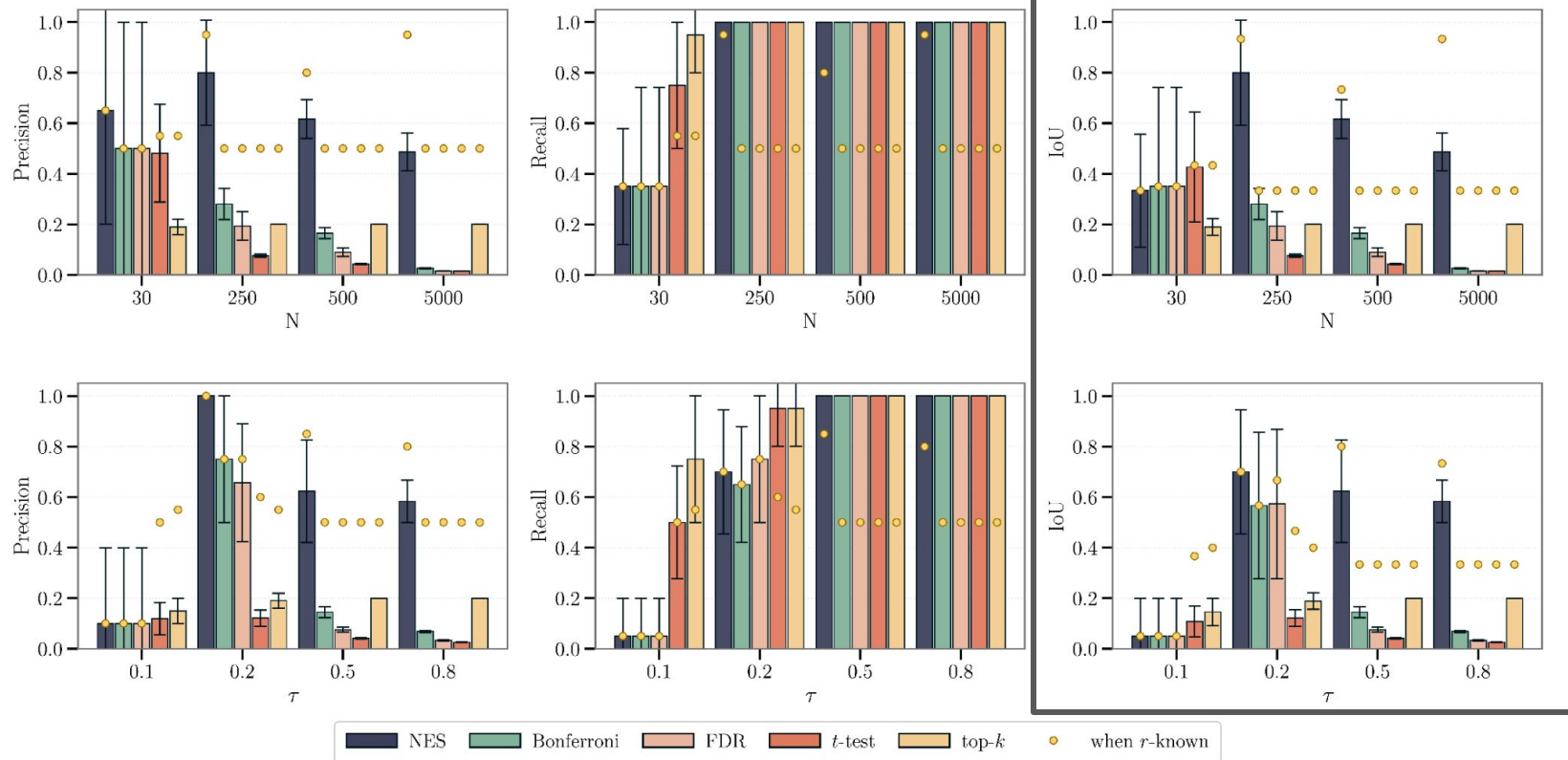
Results on synthetic trials

Paradox of ECI: Precision collapses w/ powerful tests



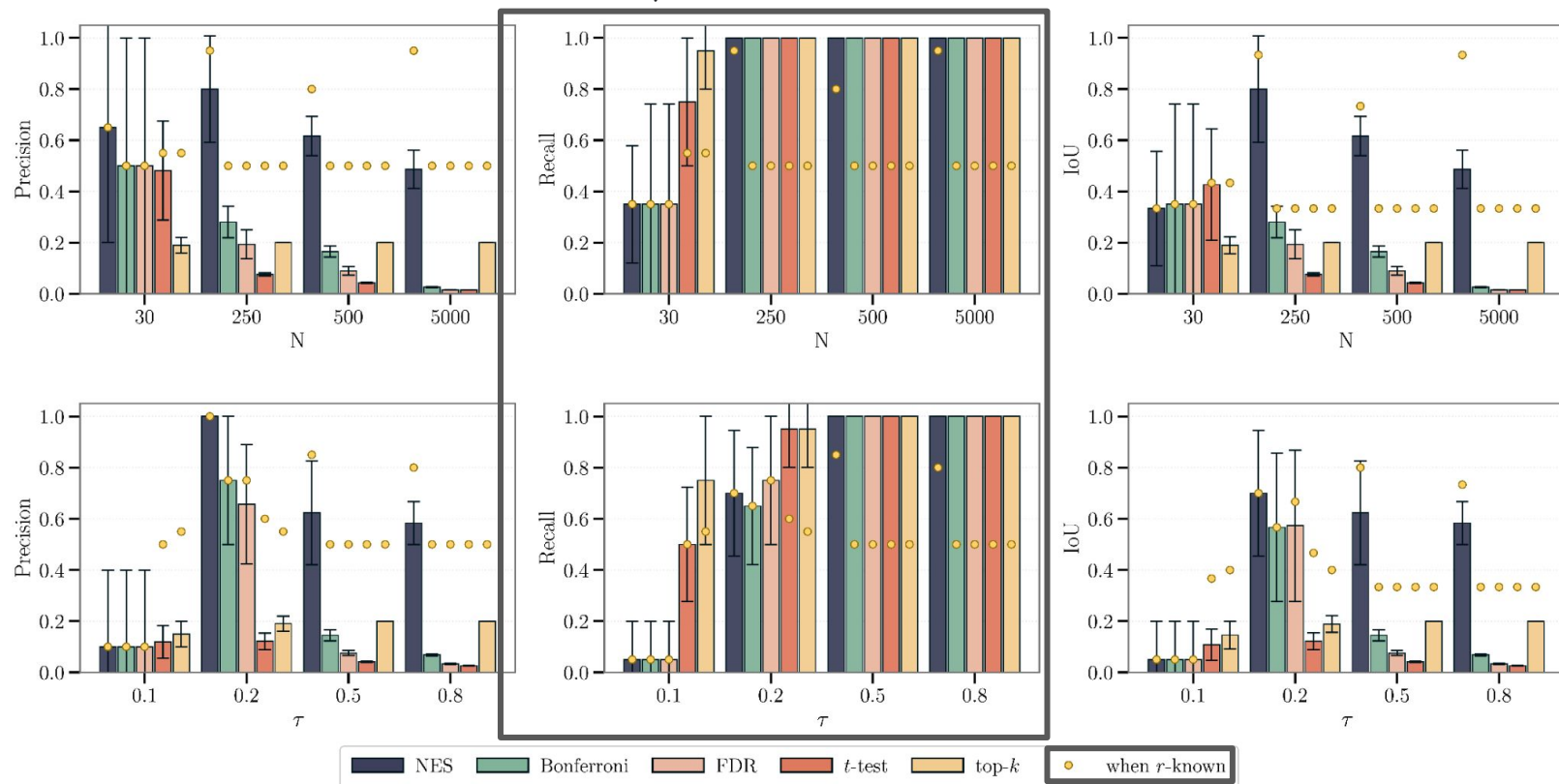
Results on synthetic trials

NES is robust and works well



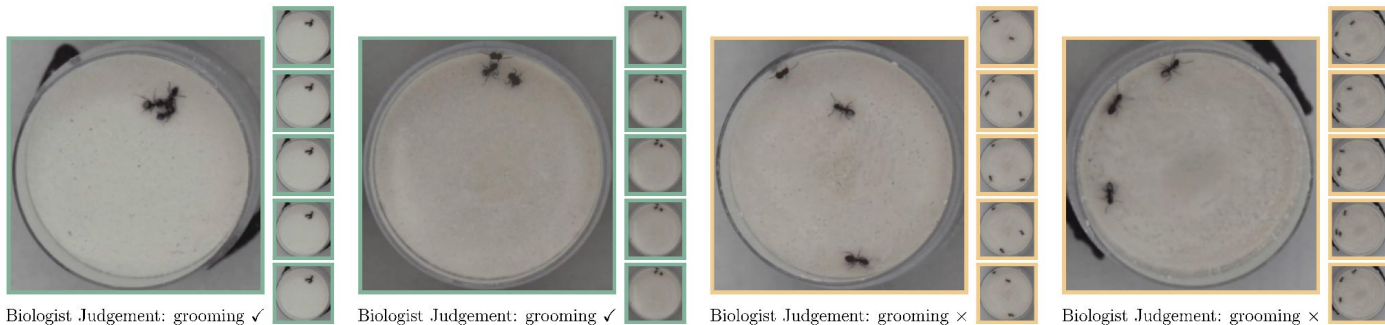
Results on synthetic trials

If r is known, NES does not miss effects

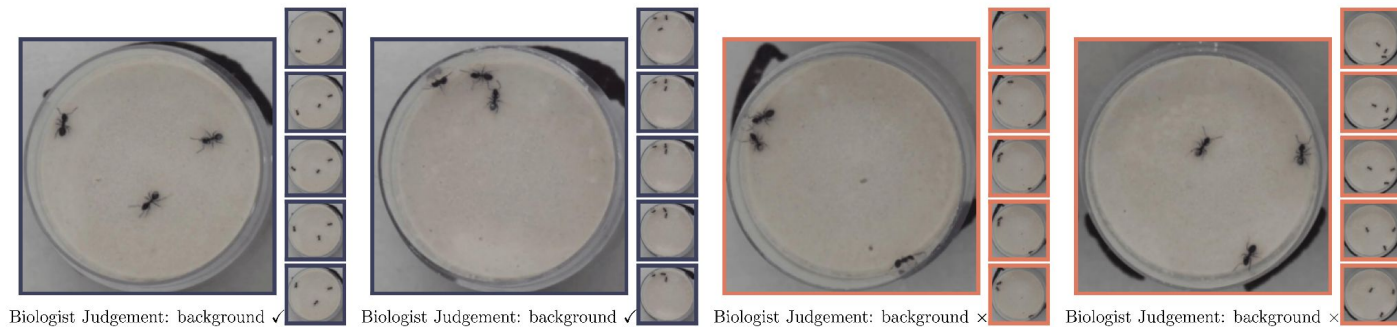


Results on ISTAnt

Qualitative Interpretation for Neuron 394



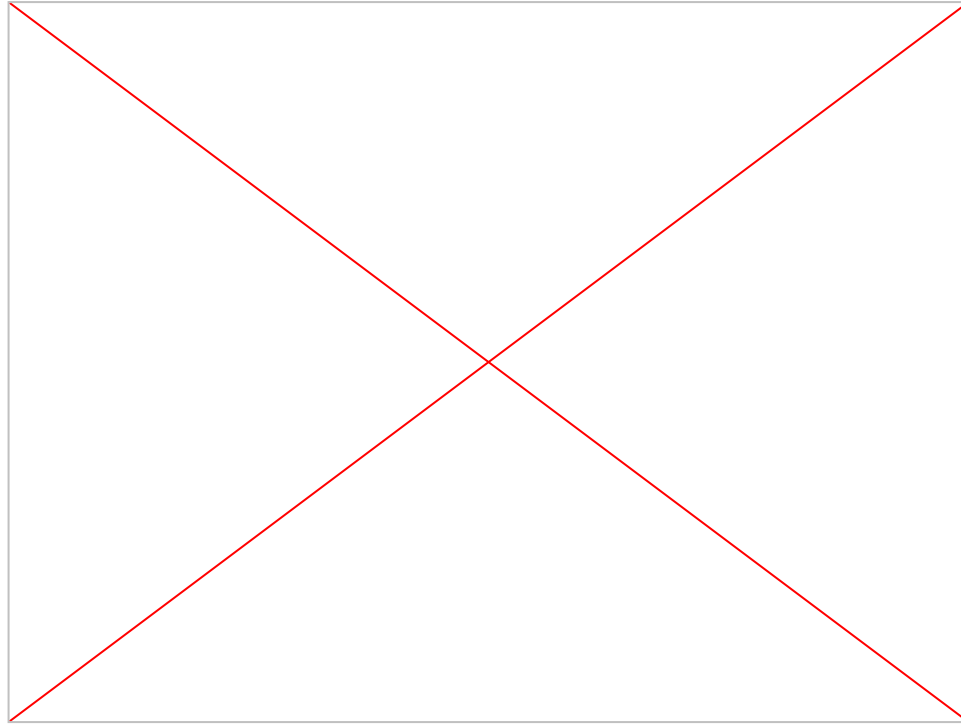
Qualitative Interpretation for Neuron 550



Max-Activating (Neuron 394) Non-Activating (Neuron 394) Max-Activating (Neuron 550) Non-Activating (Neuron 550)

Beyond “standard” causal models

Temporal dynamics



Dynamical systems as causal models

System of coupled differential equations *modeling physical mechanisms* responsible for time evolution

$$\frac{dx}{dt} = f(x), \quad x \in \mathbb{R}^d$$

Future states are “caused” by immediate past

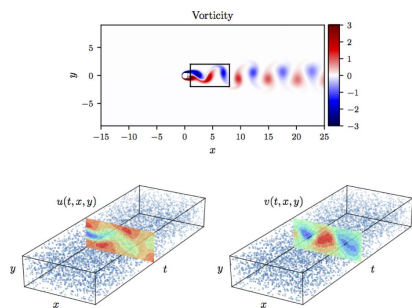
$$x(t + dt) = x(t) + dt \cdot f(x(t))$$

Unclear to which extent these can be learned from data for non-linear systems

For the equation to have causal meaning, f must be a causal mechanism (i.e., interventions are well defined and align with physical experiments)

ML methodologies

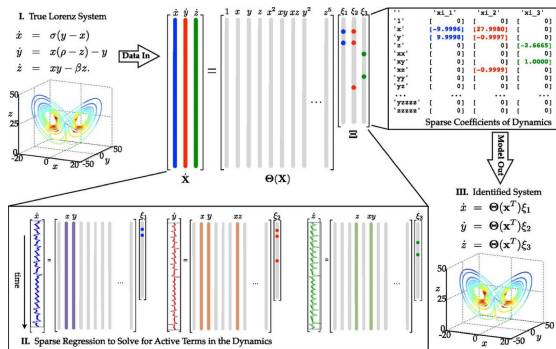
PINNs



Raissi et al., 2019

- Constrain evolution prediction to **follow physics**
- Embeds physics **implicitly via loss**

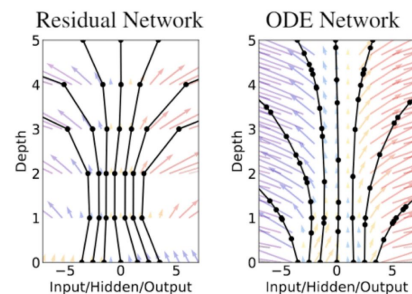
SINDy



Brunton et al., 2016

- **Discover governing equation** regressing time derivatives onto basis functions
- **No learned representations**
- **Reconstruct time derivatives**

NeuralODE

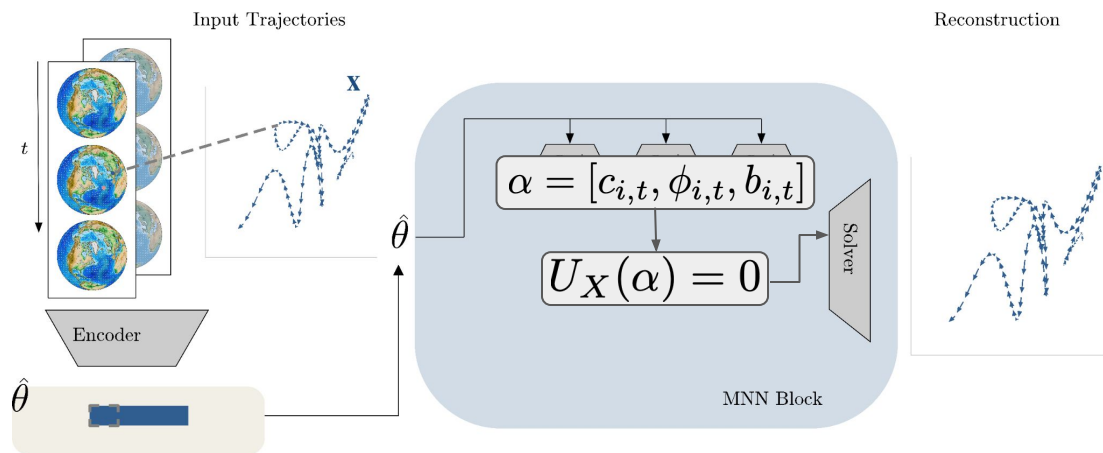


Chen et al., 2018

- Embed **network in solver** to predict next states
- Train by **reconstruction error**

Model Architecture

Parameterize equation as a combination of **learnable** coefficients parameterized by deep neural networks. Trained fully end-to-end.

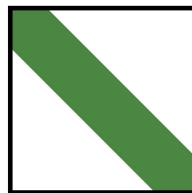


$$U_X : \underbrace{\sum_i c_i(t, X) u^{(i)}}_{\text{linear terms}} + \underbrace{\sum_k g_k(t, u, u', \dots; \phi(t, X))}_{\text{nonlinear terms}} = b(t, X)$$

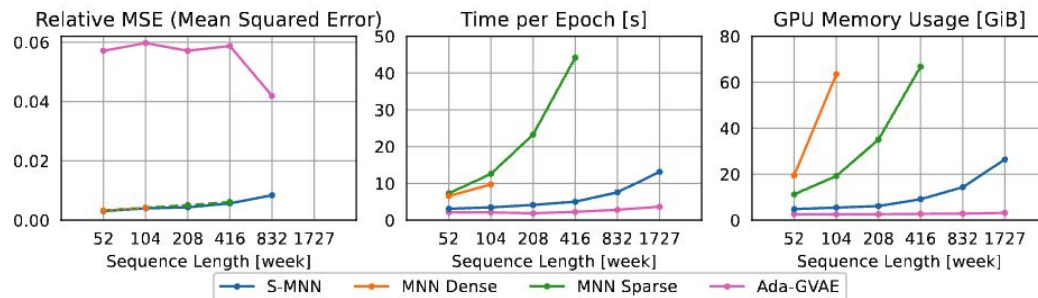
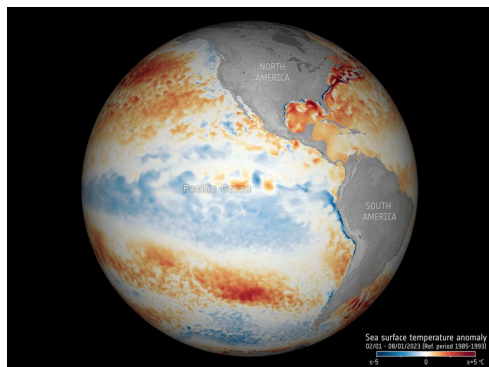
Direct least squares solution

	Time Complexity	Space Complexity
Dense Solver	$O(T^3)$	$O(T^2)$
Sparse Solver	$O(T^2)$	$O(T^2)$

- Linear system $\sum_i c_{i,t} u^{(i)} = b_t \rightarrow Az = b$
$$z = (A^T A)^{-1} A^T b$$
$$z = M^{-1} \beta$$
- M is a banded symmetric matrix \rightarrow solver with linear time and space complexity
- In both cases, we also need to derive the backward gradients for GPU implementation

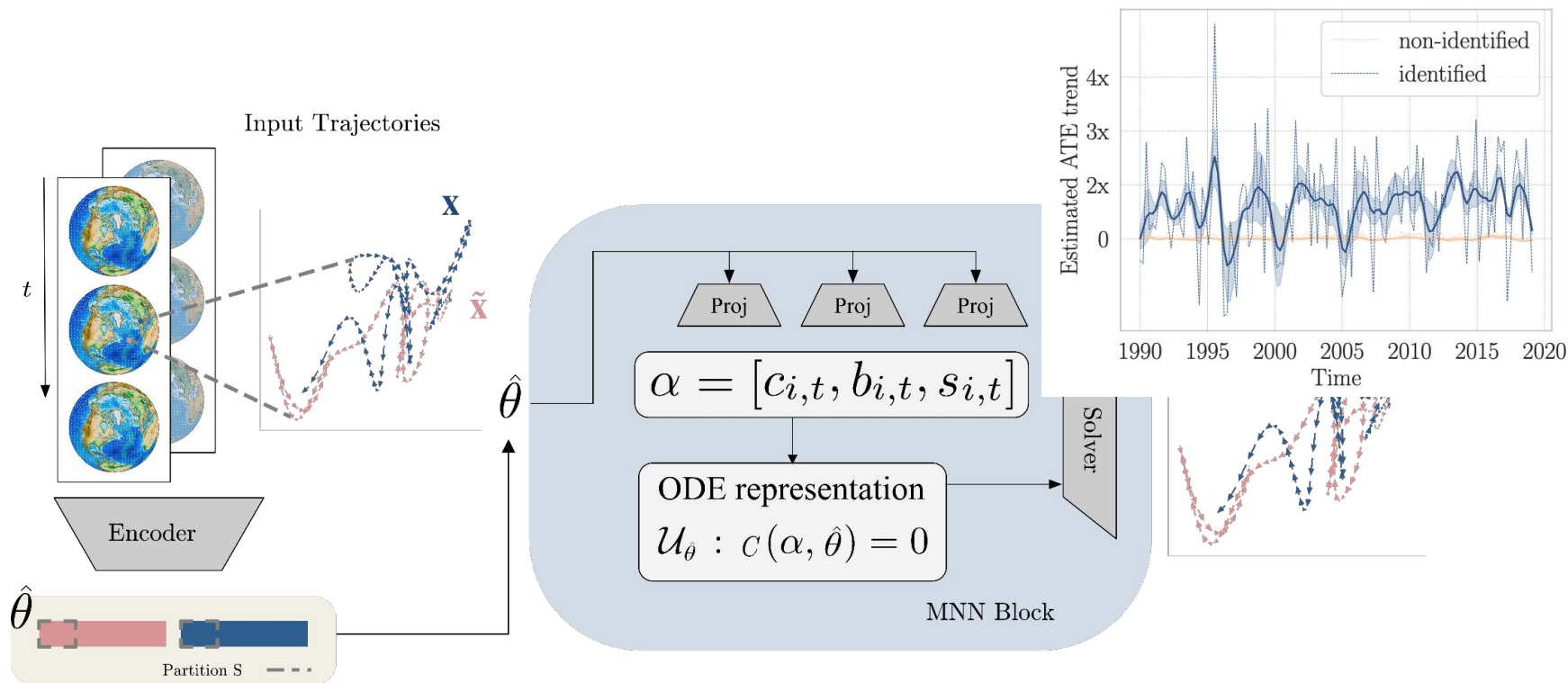


Long term sea surface temperature forecasting



Implication: Better solver immediately improves scalability without affecting accuracy

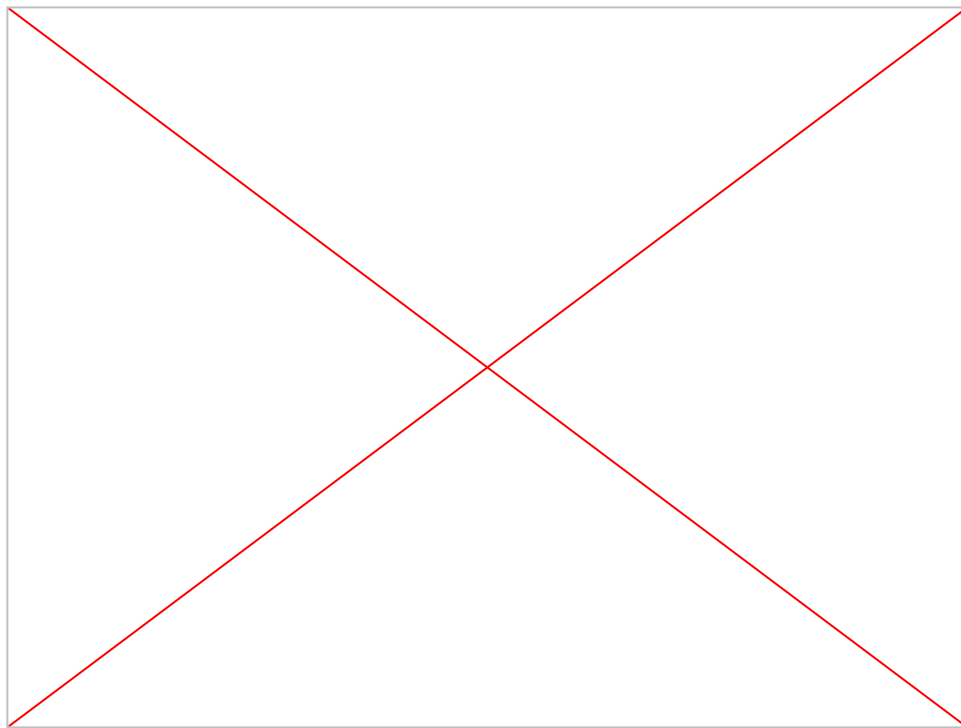
CRL recipes also extend to dynamical systems



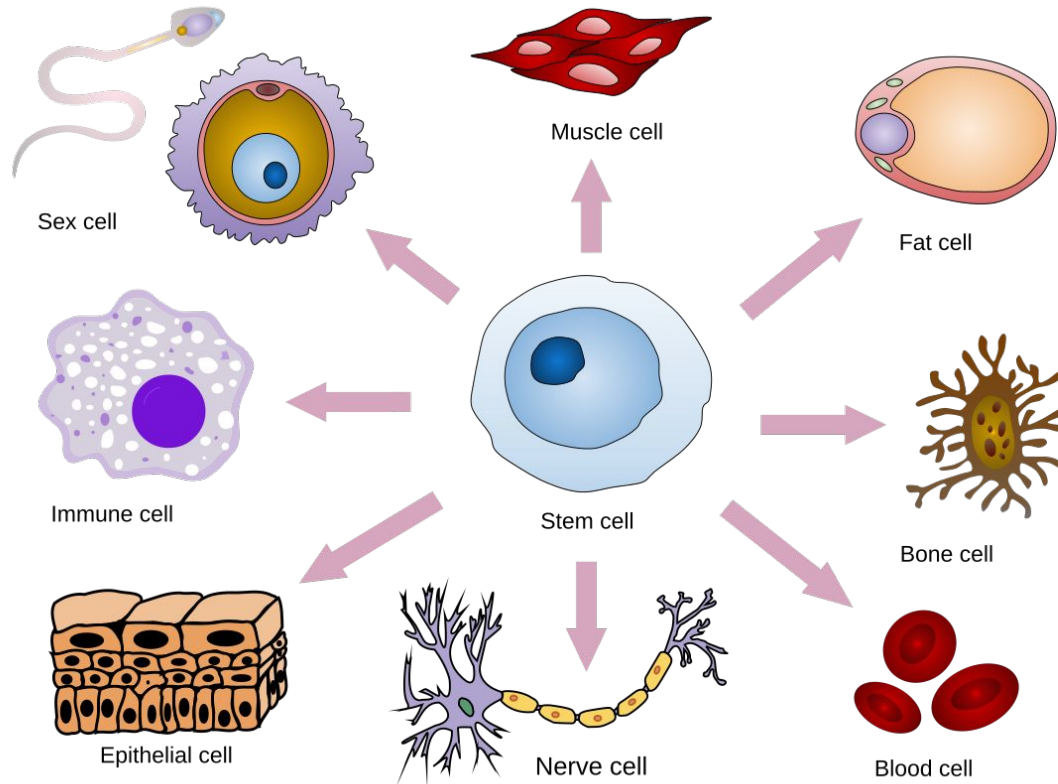
"Marrying Causal Representation Learning with Dynamical Systems for Science", Yao, Muller, L; NeurIPS 2024

"The arctic has warmed nearly four times faster than the globe since 1979", Rantanen et al., Nature Communications Earth & Environment 2022.

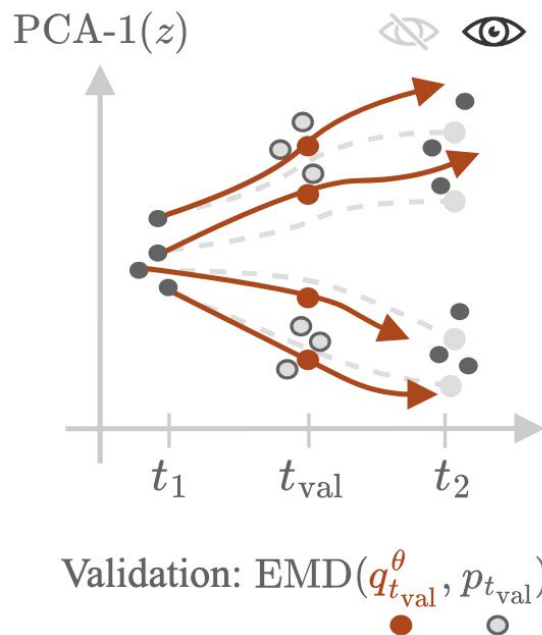
PDE extension: Ginzburg-Landau Reaction Diffusion



Application: Modeling cell differentiation

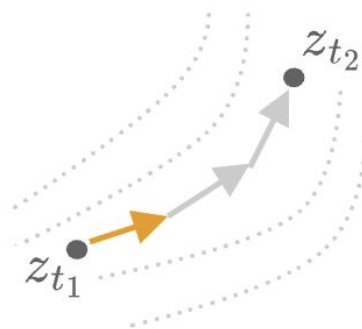


Application: Modeling cell differentiation



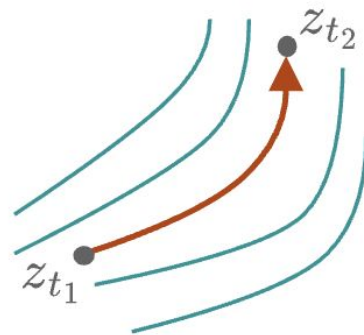
NODE / FM predicts

a velocity
 $\dot{z} = f_\theta(z_{t_1}, t_1)$



Cell-MNN predicts

a linear ODE
 $\dot{z} = A_\theta(z_{t_1}, t_1) z$



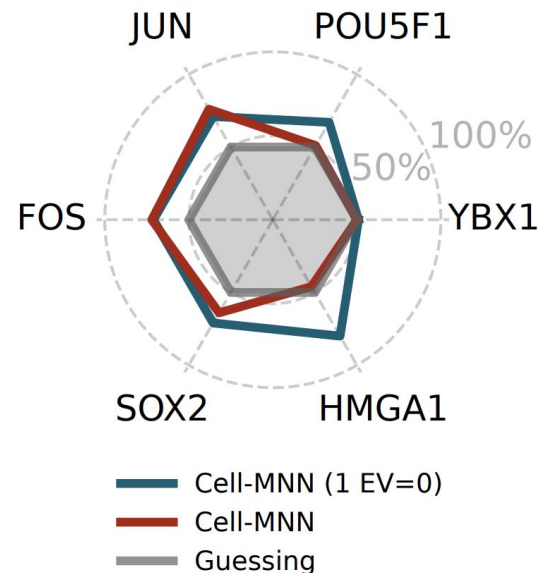
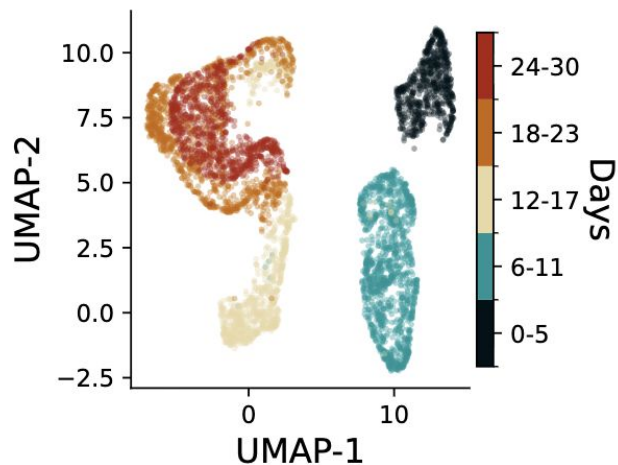
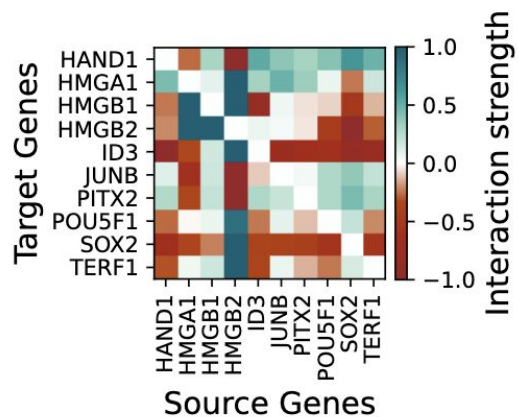
analytical solve

$$z_{t_2} = \exp[A_\theta(z_{t_1}, t_1) \cdot (t_2 - t_1)] z_{t_1}$$

Performance

Method	Cite	EB	Multi	Average ↓
TrajectoryNet [50]	—	0.848	—	—
WLF-UOT [39]	—	0.800 ± 0.002	—	—
NLSB [27]	—	0.777 ± 0.021	—	—
SB-CFM [52]	1.067 ± 0.107	1.221 ± 0.380	1.129 ± 0.363	1.139 ± 0.077
[SF] ² M-Sink [53]	1.054 ± 0.087	1.198 ± 0.342	1.098 ± 0.308	1.117 ± 0.074
[SF] ² M-Geo [53]	1.017 ± 0.104	0.879 ± 0.148	1.255 ± 0.179	1.050 ± 0.190
I-CFM [52]	0.965 ± 0.111	0.872 ± 0.087	1.085 ± 0.099	0.974 ± 0.107
DSB [14]	0.965 ± 0.111	0.862 ± 0.023	1.079 ± 0.117	0.969 ± 0.109
I-MFM [24]	0.916 ± 0.124	0.822 ± 0.042	1.053 ± 0.095	0.930 ± 0.116
[SF] ² M-Exact [53]	0.920 ± 0.049	0.793 ± 0.066	0.933 ± 0.054	0.882 ± 0.077
OT-CFM [52]	0.882 ± 0.058	0.790 ± 0.068	0.937 ± 0.054	0.870 ± 0.074
DeepRUOT [57]*	0.845 ± 0.167	0.776 ± 0.079	0.919 ± 0.090	0.846 ± 0.071
OT-Interpolate*	0.821 ± 0.004	0.749 ± 0.019	0.830 ± 0.053	0.800 ± 0.044
OT-MFM [24]	0.724 ± 0.070	0.713 ± 0.039	0.890 ± 0.123	0.776 ± 0.099
Cell-MNN (ours)*	0.791 ± 0.022	0.690 ± 0.073	0.742 ± 0.100	0.741 ± 0.050

Discovering GRN



Concluding remarks

Discussion

- ML has a great opportunity in powering causal analysis with key applications in scientific discovery but causal questions are subtle
- Chasing predictive accuracy may not lead to more accurate causal conclusions. In AI4Science, scientific questions should be part of the benchmark (especially if they are causal)
- Seeing the hidden world requires assumptions that the statistical causality language can describe
- **CRL:** Representation that makes it easier/ possible to extract causal information with some downstream estimator

Thanks!



Institute of
Science and
Technology
Austria



You?



**I'm hiring PhD
students and
PostDocs!**

Causal
Learning and