EFSTRATIOS GAVVES

# GROUNDING FOUNDATION MODELS IN REALITY
## PHYSICS- & CAUSALITY-INFORMED WORLD MODELS
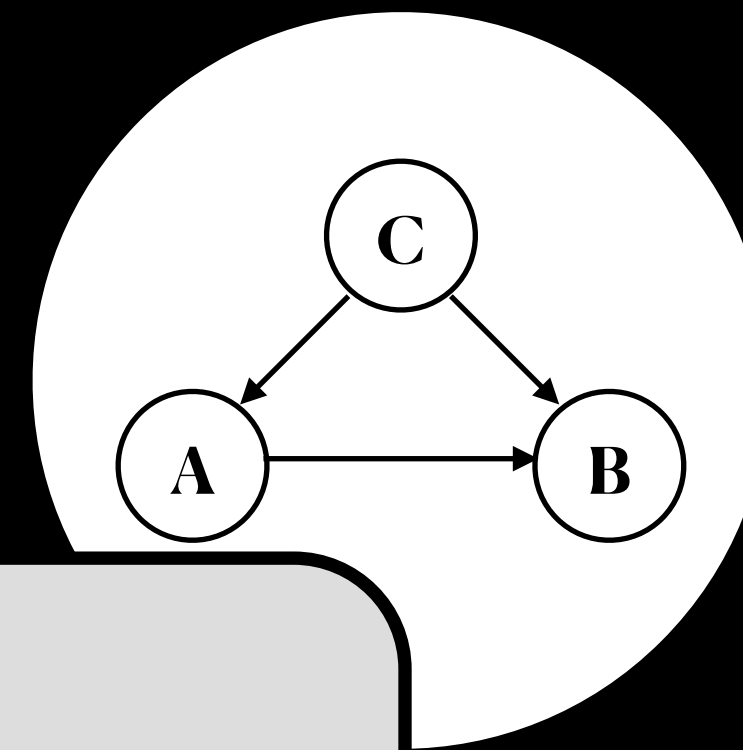
# WHO AM I?

- Associate Professor Deep Learning

- ERC StG & NWO VIDI

- ELLIS Scholar

- Co-director of ICAI QUVA & POP-AART

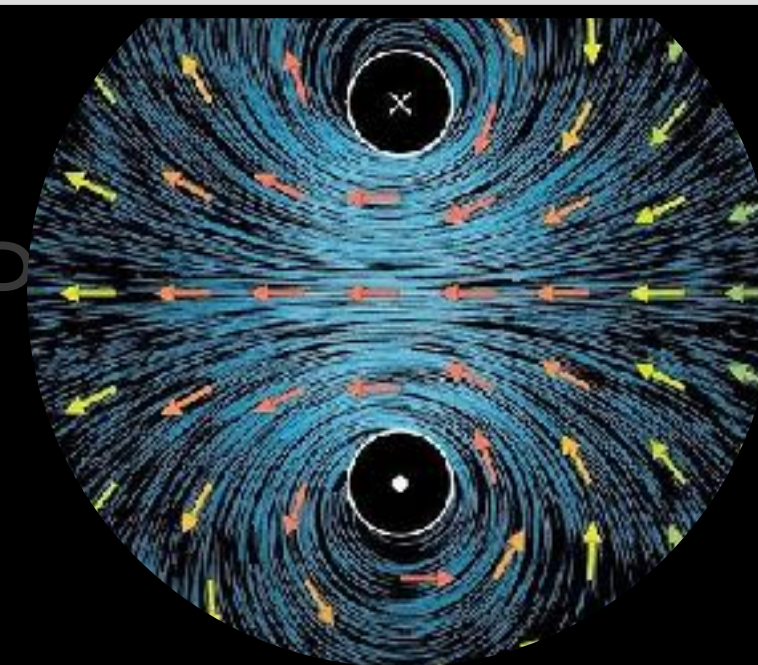- Co-founder of Ellogon.AI and LYDS Partners

# WHO AM I?

**Robot Learning**

**Causal Representation Learning**

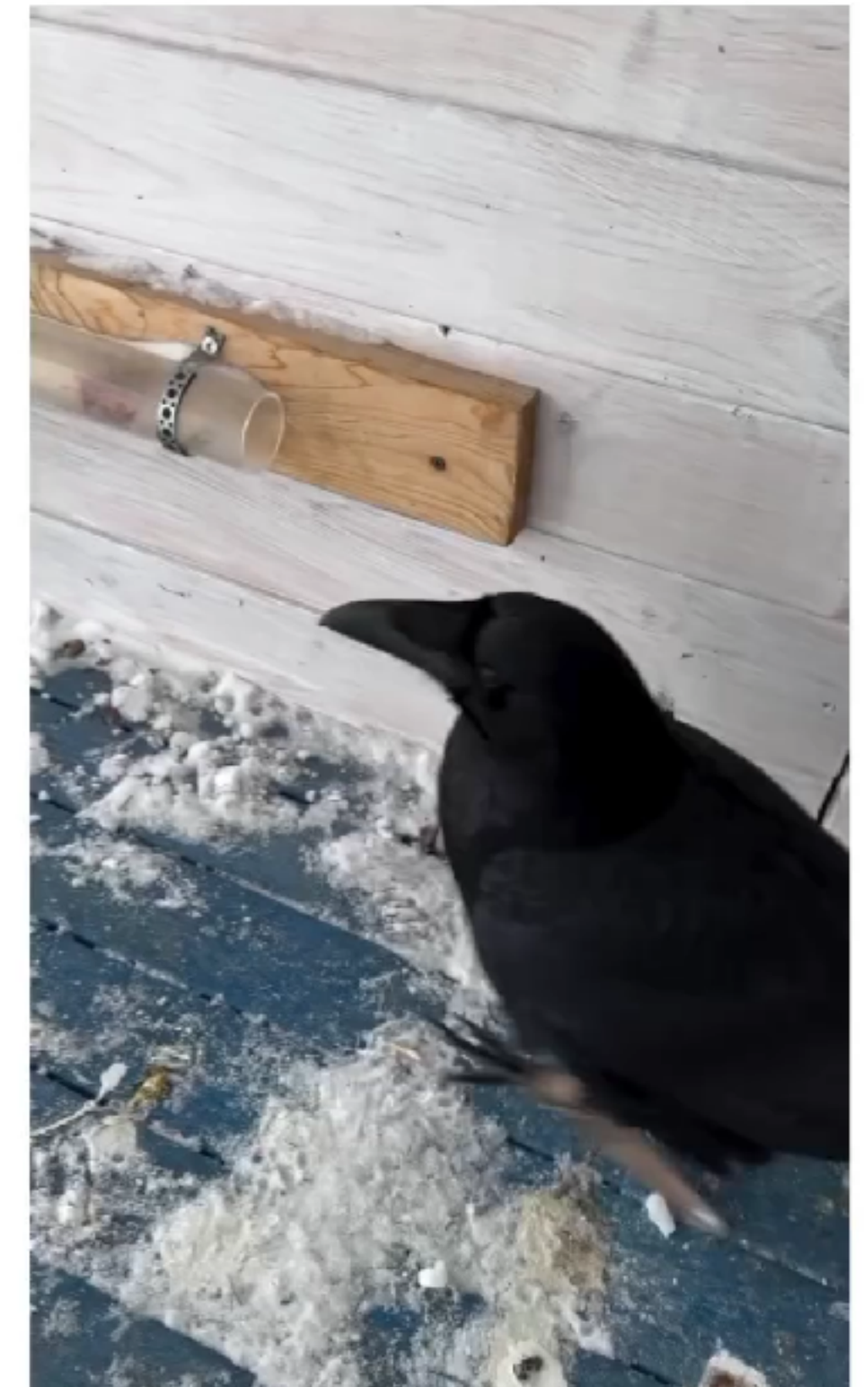- Associate Profe...                  earning

- ERC StG & NWO VIDI

- ELLIS Scholar

- Co-director of ICAI QUVA & POP

- Co-founder of Ellogon.AI

**Dynamical Deep Learning**

**AI & Science**

# CROW VS FRONTIER MODELS

QUO VADIS AI?

# FROM SEEING TO
## → INDUCING PHYSICS & CAUSALITY
## → SPATIAL MEMORY & NAVIGATION
## → REASONING
## → PLANNING
## ⇒ INTERACTING WITH REALITY

# GROUNDING ROBOT WORLD MODELS



**Dynamics/Physics inductive biases**

**Causal inductive biases**

**Embodied inductive biases**

**Open World**

**AI OS**

DREAM TO MANIPULATE

# COMPOSITIONAL MANIPULATION WORLD MODELS

https://dreamtomanipulate.github.io/DreMa/

w. L. Barcellona, A. Zadaianchuk, D. Allegro, S. Papa, S. Ghidoni

# OUR INDUCTIVE BIASES

EXPLICIT GROUNDING == PHOTOREALISTIC RECON@OBJECT-CENTRIC+ EXPLICIT PHYSICS
⇒ COMPOSITIONAL MANIPULATION WORLD MODELS ⇒ IMAGINATION
⇒↑ ROBOT IMITATION LEARNING

# PHOTOREALISTIC RECONSTRUCTION …

- Gaussian Splatting is like 'sparse 3D pixels'

- Real-time rendering

- High-quality

- Good depth rendering

- Fast training

- Explicitly grounded representation



Novel view



Mesh

- Zero-shot object localization

→ Grounded SAM, DEVA, or our VISA[1] & LV-VIS[2]

- Prompts: "`object`" & "`table`"

- Segment and group objects across views



RAM-Grounded-SAM

Only Image Inputs

Automatically Generate Tag, Box, and Mask Annotations

"armchair, blanket, lamp, carpet, couch, dog, floor, furniture, gray, green, living room, picture frame, pillow, plant, room, sit, stool, wood floor"

[1] Wang, Gavves et al., Towards Open-Vocabulary Video Instance Segmentation, ICCV 2023
[2] Yan, Wang, Gavves et al., VISA: Reasoning Video Object Segmentation via Large Language Models, ECCV 2024

# … WITH "SELF-AWARENESS" …

- The robot is also an object asset

- Articulated but with known joints

- Movement given by the URDF file, no need to estimate

# ... AND MANIPULABLE WITH PYBULLET ...

- Integrate explicit physics engines

- Manipulate object assets by exerting forces on them

- PyBullet requires mesh grids

- Convert Gaussian Splats to Meshes

- "Decode" effect of manipulations with Gaussian Splatting

- Play seen trajectories

- And render it from any angle

- "Re-imagining" past experiences

# ... AND IMAGINE NEW ONES ...

- Since we have a 'digital twin'

- That reconstructs photorealistically

- Understands physics

- And can be intervened with

- We can "imagine" novel trajectories

Figure 3: The effect of equivariant translation, equivariant rotation, and the object rotation transformations. Top row: start of demonstration. Bottom row: target of demonstration.

- Single demonstration per variant

- Excluding tasks with articulation

- 60-110 imaginations generated

- Works even with pure imagination

Table 1: Comparison of PerAct (Shridhar et al., 2023) trained on original demonstrations to DREMA trained on only imagination demonstrations and the combination of both. The table reports the mean ± std and maximum success rate over 5 test runs.

| | Single-task | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | close jar | | place shape | | slide block | | avg single-task | |
| | mean ± std | max | mean ± std | max | mean ± std | max | mean | max |
| PerAct (Original data) | 38.4 ± 0.80 | 40 | 6.4 ± 1.50 | 8 | 48.4 ± 3.20 | 50 | 31.1 | 32.7 |
| DREMA (Imagined data) | 41.2 ± 2.40 | 46 | 9.6 ± 1.50 | 12 | 54.4 ± 2.15 | 62 | 35.1 | 40.0 |
| DREMA (All data) | **51.2 ± 1.60** | **54** | **11.2 ± 1.60** | **12** | **62.0 ± 2.19** | **66** | **41.5** | **44.0** |

| | Multi-task | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | close jar | | place shape | | slide block | | avg multi-task | |
| | mean ± std | max | mean ± std | max | mean ± std | max | mean | max |
| PerAct (Original data) | 26.0 ± 3.10 | 28 | 7.2 ± 1.60 | 10 | 34.0 ± 5.06 | 38 | 22.4 | 25.3 |
| DREMA (Imagined data) | 28.0 ± 3.35 | 32 | **18.0 ± 2.83** | **22** | 48.0 ± 1.79 | 50 | 31.3 | 34.7 |
| DREMA (All data) | **46.0 ± 3.58** | **52** | 6.4 ± 3.20 | 12 | **54.0 ± 2.19** | **58** | **35.5** | **40.7** |

| | close jar | place cups | place shape | put groceries | slide blocks | stack block | AVERAGE |
|---|---|---|---|---|---|---|---|
| RVT2 original four cameras | 62 | 0.4 | 2.4 | 8.8 | 41.2 | 3.2 | 19.67 |
| RVT2 original three cameras (val results) | 17.5 | 2.5 | 7.5 | 5 | 50 | 2.5 | 14.17 |
| RVT2 drema three cameras (val results) | 25 | 0 | 15 | 10 | 50 | 10 | 18.33 |
| PERACT original | 26 | 1.6 | 7.2 | 1.2 | 34 | 3.6 | 12.27 |
| PERACT drema | 46 | 0.4 | 6.4 | 1.6 | 54 | 0.4 | 18.13 |

**Table 3: Localization errors.**

| Task | Error (m) |
|---|---|
| pick block | 0.010 |
| pick shape | 0.050 |
| push | 0.049 |
| Average | 0.038 |

**Table 4: In- and out-of-distribution evaluation with real robots.**

| | pick block | | pick shape | | push | |
|---|---|---|---|---|---|---|
| | In distr. | OOD | In distr. | OOD | In distr. | OOD |
| PerAct | 55 | 50 | 30 | 10 | 40 | 10 |
| DREMA (All) | **90** | **90** | **35** | **30** | **80** | **60** |

Figure 5: Original (top) and imagined demonstration (bottom) after a 90° rotation transformation.

# MORE IMAGINED AUGMENTATIONS HELP

- Simple transformations for now

- Main challenge: make sure demonstrations are valid

- Lower efficacy but better scalability

Table 2: Performance of DREMA trained on single-task demonstrations from different types of transformations. The table reports the mean ± std and maximum success rate over 5 test runs.

| | close jar | | place shape | | slide block | |
|---|---|---|---|---|---|---|
| | mean ± std | max | mean ± std | max | mean ± std | max |
| Replay | $10.0 \pm 3.10$ | 16 | $1.2 \pm 0.98$ | 2 | $26.0 \pm 0.00$ | 26 |
| Object Rotation | $25.2 \pm 0.98$ | 26 | $10.4 \pm 2.65$ | **14** | $42.0 \pm 0.00$ | 42 |
| Roto-translation | $41.2 \pm 2.40$ | 44 | $10.0 \pm 2.53$ | **14** | $50.4 \pm 1.50$ | 52 |
| DREMA (*All data*) | $\mathbf{51.2 \pm 1.60}$ | **54** | $\mathbf{11.2 \pm 1.60}$ | 12 | $\mathbf{62.0 \pm 2.19}$ | **66** |

Equivariant transformations are complementary …



Figure 4: Imagined demonstrations keep improving imitation learning even with increasing number of original data.

… and scale

19

# WHAT MAKES IT NOT A WORLD MODEL?

- Not end-to-end yet

- Dynamics are not learned

- In fact, not much learning (from us) yet

- BUT, the idea is not to return to feature engineering

- Start with grounding explicitly

- Then everything else neural



UniSim

World Model Learning

Learning Value and Actor Networks

Environment Interaction

A1 Quadruped Walking

UR5 Multi-Object Visual Pick Place

XArm Visual Pick and Place

Sphero Ollie Visual Navigation

# LEARNING & GENERALIZATION

- For one, add representations to object assets

- Learn (fine-tune) on observed trajectories

- Learn with differentiable physics

- Learn with partial observability

- Learn with manipulations

- Learn dynamics beyond physics (eg, causality, third-person actors, theory of mind, …)

- Learn with dynamic scenes

# BEYOND IMITATION: LEARNING RL POLICIES

- Why stop with imitation learning?

- Grounded imagination & reasoning opens up lots of exciting possibilities

- Maximizing future reward in future imagination?

- …

# SKY IS THE LIMIT

- Articulated objects (identifying joints)

- Deformable objects (with complex physics engines)

- Closed feedback loop ('eye-hand coordination')

- Physical parameter identification (friction, etc)

- Safety and interpretability

- Stochasticity (many possible futures)

CAUSAL INDUCTIVE BIASES

LIPPE ET AL, ICLR, ICML, UAI, 2022-23

Unknown

Input image 1 | Input image 2 | Generated Output | Latents from image 2

Microwave Active
Stove (front-left)

Manipulating Image 1 by turning on the Microwave and the front-left Stove. Note the egg staying uncooked despite the stove being turned on, which the model has never seen in training and shows BISCUIT's ability to perform novel interventions.

# TEMPORAL CAUSAL REPRESENTATION LEARNING



https://phlippe.github.io/

[1] B Schölkopf, F Locatello, S Bauer, N R Ke, N Kalchbrenner, A Goyal, Y Bengio, Towards Causal Representation Learning, Proceedings of the IEEE, 2021

# TEMPORAL CAUSAL REPRESENTATION LEARNING

**Environment**

NN

**Representation**

**Agent**

**Action**

## Representation Learning Tasks

What are the causal variables of the environment?

How do they interact with each other?

How can the robot agent intervene on causal variables?

**https://phlippe.github.io/**

[1] B Schölkopf, F Locatello, S Bauer, N R Ke, N Kalchbrenner, A Goyal, Y Bengio, Towards Causal Representation Learning, Proceedings of the IEEE, 2021

# CAUSAL LEARNING FROM BINARY INTERACTIONS

- Many interactions are binary

(Turn lights on/off, open/close door …)

- Learn latents to reflect change

- Provable & manipulable causal factors

- By integrating probabilistic causal priors



🍪 **BISCUIT: Causal Representation Learning from Binary Interactions**
We aim to learn the *causal variables* by interactions with a dynamical environment.

https://phlippe.github.io/BISCUIT/

# CAUSAL LEARNING FROM BINARY INTERACTIONS

- Causal model: a temporal DAG

- Observed: images $X^t$ and "regime" variables $R^t$

- Latents: causal $C^{t-1}$ and interactions $I^t = f(R^t, C^{t-1})$



https://phlippe.github.io/BISCUIT/

# BINARY INTERACTIONS FOR IDENTIFIABILITY

- Assumption #1: interactions described by a binary variable

- Assumption #2: distinct interaction patterns → interactions not functions of other interactions

- Assumption #3: mechanisms vary sufficiently with interactions or time

Time step t+1

Time step t

No interaction (observational)

Interaction (interventional)

A. **(Dynamics Variability)** *Each variable's log-likelihood difference is twice differentiable and not always zero:*

$$\forall C_i^t, \exists C^{t-1} : \frac{\partial^2 \Delta(C_i^t | C^{t-1})}{\partial (C_i^t)^2} \neq 0;$$

B. **(Time Variability)** *For any $C^t \in \mathcal{C}$, there exist $K+1$ different values of $C^{t-1}$ denoted with $c^1, ..., c^{K+1} \in \mathcal{C}$, for which the vectors $v_1, ..., v_K \in \mathbb{R}^{K+1}$ with*

$$v_i = \left[ \frac{\partial \Delta(C_i^t | C^{t-1} = c^1)}{\partial C_i^t} \quad \cdots \quad \frac{\partial \Delta(C_i^t | C^{t-1} = c^{K+1})}{\partial C_i^t} \right]^T$$

*are linearly independent.*

https://phlippe.github.io/BISCUIT/

# OPTIMIZING WITH VARIATIONAL INFERENCE



https://phlippe.github.io/BISCUIT/

# TWO-STAGE LEARNING IN COMPLEX SETTINGS



**https://phlippe.github.io/BISCUIT/**

CAUSALWORLD TRI-FINGER

# DISCOVERING INTERACTIONS

Example Sequence    Ground Truth    $\beta$-VAE    BISCUIT – AE + NF

$R^2$ scores (diag ↑/sep ↓)

| Models | CausalWorld |
|--------|-------------|
| iVAE (Khemakhem et al., 2020a) | 0.28 / 0.00 |
| LEAP (Yao et al., 2022b) | 0.30 / 0.00 |
| DMS (Lachapelle et al., 2022b) | 0.32 / 0.00 |
| BISCUIT-NF (Ours) | **0.97** / 0.01 |

https://phlippe.github.io/BISCUIT/

ITHOR OBJECTS AS CAUSAL VARIABLES

# CAUSE-AND-EFFECT DRIVES SEMANTICS (?)

Input Image  Learned Interactions  Combined Image

$R^2$ scores (diag ↑/sep ↓)

| Models | iTHOR |
|---|---|
| iVAE (Khemakhem et al., 2020a) | 0.48 / 0.35 |
| LEAP (Yao et al., 2022b) | 0.63 / 0.45 |
| DMS (Lachapelle et al., 2022b) | 0.61 / 0.40 |
| BISCUIT-NF (Ours) | **0.96 / 0.15** |

https://phlippe.github.io/BISCUIT/

ELEMENTS OF "IMAGINATION"

# HACKING THE SIMULATOR

Input Image 1　　Input Image 2　　Generated Output

Scaling-up with foundation vision models towards autonomous learning?

https://phlippe.github.io/BISCUIT/

# LIMITATIONS



BISCUIT: Causal Representation Learning from Binary Interactions
We aim to learn the *causal variables* by interactions with a dynamical environment.

- Sufficient intervention data

- Works with temporal data only

- Assumes binary interactions

https://phlippe.github.io/BISCUIT/

# IS IT REALLY CAUSAL?

- Patterns are often correlations

- Cause-and-effect is a strong (albeit sometimes biased) framework to learn

- If we go past the chicken-egg problem

- Power of causal representations is in autonomousness and controllability (imho)



"Essentially, all models are wrong, but some are useful."
George E. P. Box

# SCALE & ROBOT LEARNING

- Scale up to many environments → reuse semantics

- LLMs for guidance and sample efficiency?

- Ideally, Gaussian Splats for *de novo interactive* environments and scaling-up

- System IIa: First causal principles for novel problem-solving & Causal World Models

- System IIb: Safe & human-robot-aligned planning

https://phlippe.github.io/BISCUIT/

LIU ET AL, ICML 2023, ONGOING
AUZINA ET AL, NEURIPS 2023
PERVEZ ET AL, ONGOING

# DYNAMICS INDUCTIVE BIASES

# SWITCHING DYNAMICS IN INTERACTING SYSTEMS



- In many settings in perception and sciences, we have systems of multiple objects

- These objects may interact (or not) with higher-order complex & switching temporal dynamics[1,2]

- Finding dynamical patterns is often critical

- Generalization of temporal clustering

[1] Z Gharamani, G Hinton, Variational learning for switching state-space models, NeurIPS, 2020
[2] A Ansari, K Benidis, R Kurle, A Turkmen, H So, A Smola, B Wang, T Januschowski, NeurIPS, 2022

# GRAPH SWITCHING DYNAMICAL SYSTEMS



- Switching Dynamical Systems focus on finding out when objects behave differently

- Key idea #1: Scale up by NN function approximation to amortize pairwise transition dynamics between multiple objects and dynamic behaviors

- Key idea #2: Graph NNs and message passing  and VI for dynamic interactions between objects

- Divide and Conquer: Breaking complex dynamics into switching between simpler systems

**https://github.com/yongtuoliu/Graph-Switching-Dynamical-Systems**

Latent temporal graph
$N$ objects

$K$ motion modes

$e_t^{m \to n}$

$v_t^m = \{z_t^m, x_t^m, c_t^m, y_t^m\}$

LIU, MAGLIACANE, KOFINAS, GAVVES, ONGOING

# EQUATION DISCOVERY+SWITCHING DYNAMICS

- Hypothesis #1: Symbolic learning critical for extrapolation

- Hypothesis #2: Disentangled representation learning critical for generalization[1,2]

- Hypothesis #3: Graph learning critical for interacting dynamics

Hybrid-SINDy            Latent dynamics as governing equations

[1] J von Kügelgen*, Y Sharma*, L Gresele*, W Brendel, B Schölkopf, M Besserve, F Locatello, Self-Supervised Learning with Data Augmentations Provably Isolates Content from Style, NeurIPS, 2021
[2] I Auzina, C Yildiz, S Magliacane, M Bethge, E Gavves, Moduleated Neural ODEs, NeurIPS, 2023

# EQUATION DISCOVERY + SWITCHING DYNAMICS

- #H1: Link GRASS dynamics modes to parametric basis functions $\Theta(y_t)$

- #H2: Learn $k$-th latent dynamics $\dot{y}_t^m \approx \Theta(y_t)w_m$ modulated by $w_m$ in shared NN

- Think of SINDy-governed latent spaces

SINDy



The dynamic mode $z_k$ indexes candidate basis $\Theta_k \cdot w_k$

# EQUATION DISCOVERY+SWITCHING DYNAMICS

- So far, single-object dynamics

- #H3: Graph message passing for interacting dynamics

- Approximate inference for edges, exact inference for discrete variables



Probabilistic model

Latent temporal graph
$N$ objects
$K$ motion modes

$$\ddot{y}_t^m = \Theta(y_t)w_m$$

$$e_t^{m \to n}$$

$$v_t^m = \{z_t^m, c_t^m, y_t^m\}$$

# EQUATION DISCOVERY+SWITCHING DYNAMICS

## Scientific data



**Table 1.** Segmentation results on Mass-spring Hopper dataset.

| Method | NMI ↑ | ARI ↑ | Accuracy ↑ | $F_1$ ↑ |
|---|---|---|---|---|
| Hybrid-SINDy | 0.426 | 0.383 | 0.705 | 0.691 |
| AMORE (ours) | **0.934** | **0.970** | **0.993** | **0.994** |

**Table 2.** Forecasting results of Location/Velocity on the Mass-spring Hopper dataset.

| Method | NMAE ↓ | NRMSE ↓ |
|---|---|---|
| LLMTime | 0.120 / 0.320 | 0.430 / 0.500 |
| SVI | 0.063 / 0.070 | 0.140 / 0.250 |
| AMORE (ours) | **0.009 / 0.037** | **0.024 / 0.056** |

**Table 7.** Forecasting results of in terms of NMAE / NRMSE on ODE-driven Particle dataset.

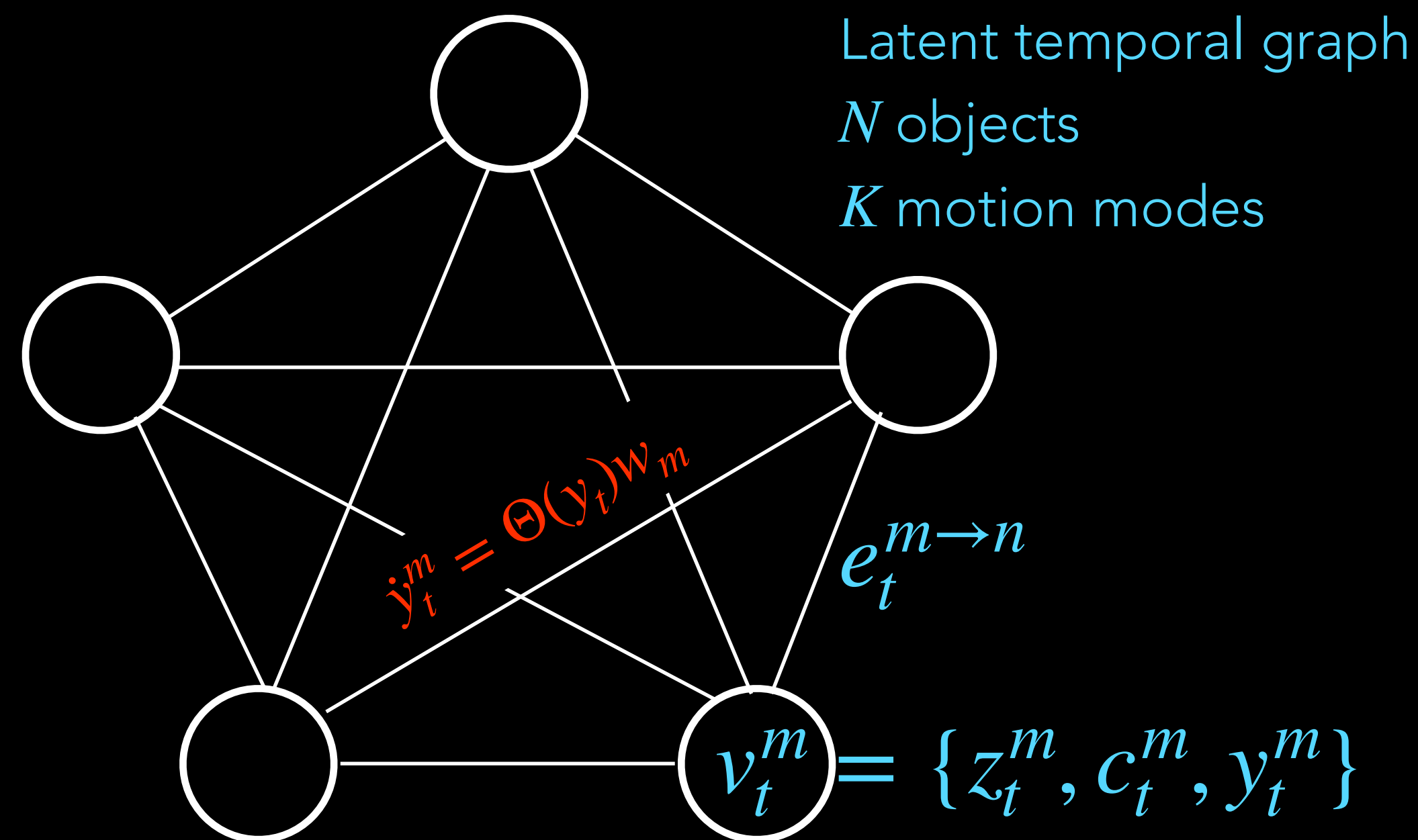| Method | One-step | Multi-step |
|---|---|---|
| LLMTime | 0.335 / 0.438 | 0.370 / 0.473 |
| SVI | 0.319 / 0.432 | 0.346 / 0.465 |
| GRASS | 0.151 / 0.224 | 0.193 / 0.270 |
| AMORE (ours) | 0.184 / 0.265 | 0.217 / 0.302 |
| AMORE-MIO (ours) | **0.146 / 0.217** | **0.186 / 0.259** |

**Table 5.** Forecasting results on non-hybrid dynamical systems. Results are shown in $\log_{10}$(NRMSE) where lower is better.

| System | LLMTime | SVI | AMORE (ours) |
|---|---|---|---|
| Coupled linear | -0.39 | -1.13 | **-1.18** |
| Cubic oscillator | -0.45 | -1.02 | **-1.06** |
| Lorenz'63 | -0.41 | **-1.27** | -1.23 |
| Hopf bifurcation | -0.32 | -0.94 | **-1.03** |
| Selkov glycolysis | -0.68 | **-1.55** | -1.49 |
| Duffing oscillator | -0.53 | -1.12 | **-1.17** |

**Table 6.** Segmentation results on ODE-driven Particle Dataset.

| Method | NMI ↑ | ARI ↑ | Accuracy ↑ | $F_1$ ↑ |
|---|---|---|---|---|
| Hybrid-SINDy | 0.205 | 0.192 | 0.414 | 0.407 |
| AMORE (ours) | 0.418 | 0.405 | 0.692 | 0.684 |
| AMORE-MIO (ours) | **0.453** | **0.442** | **0.741** | **0.735** |

**Table 10.** Analyses on robustness to different orders of polynomial as candidate basis functions on Mass-spring Hopper dataset.

| Polynomial order | 2 | | 3 | | 5 | |
|---|---|---|---|---|---|---|
| | NMI↑ | RER↓ | NMI↑ | RER↓ | NMI↑ | RER↓ |
| Hybrid-SINDy | 0.426 | $7.5e^{-3}$ | 0.384 | $8.1e^{-3}$ | 0.316 | $9.7e^{-3}$ |
| AMORE (ours) | **0.934** | **$2.1e^{-4}$** | **0.936** | **$2.3e^{-4}$** | **0.933** | **$2.8e^{-4}$** |

## Perception data



Forward:
$$\dot{x} = 5.18 - 0.82y + 0.17xy$$
$$\dot{y} = 0.42 + 0.16xy$$

Backward:
$$\dot{x} = -4.39 - 0.66y + 0.21x^2$$
$$\dot{y} = -0.72 + 0.03y^2$$

$$\dot{x} = 4.73 - 1.04x^2y + 0.28xy$$
$$\dot{y} = 1.27 - 0.90xy + 0.07y^2$$

$$\dot{x} = 3.13 - 1.26x^2y + 0.33x^3$$
$$\dot{y} = 2.51 + 0.84y^3 - 0.14xy$$

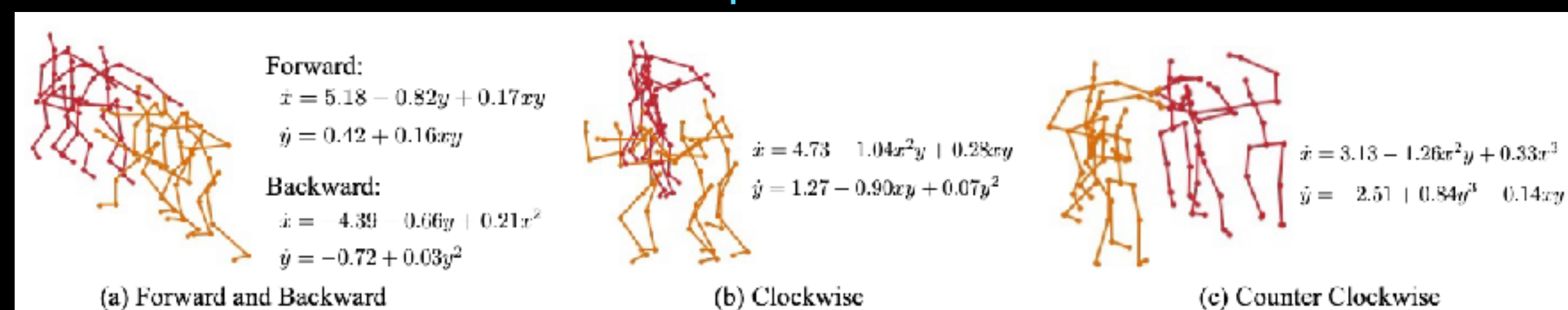(a) Forward and Backward     (b) Clockwise     (c) Counter Clockwise

Figure 5. Discovered equations on the Salsa-dancing dataset. Locations $(x, y)$ of the hip joints are used as observations.
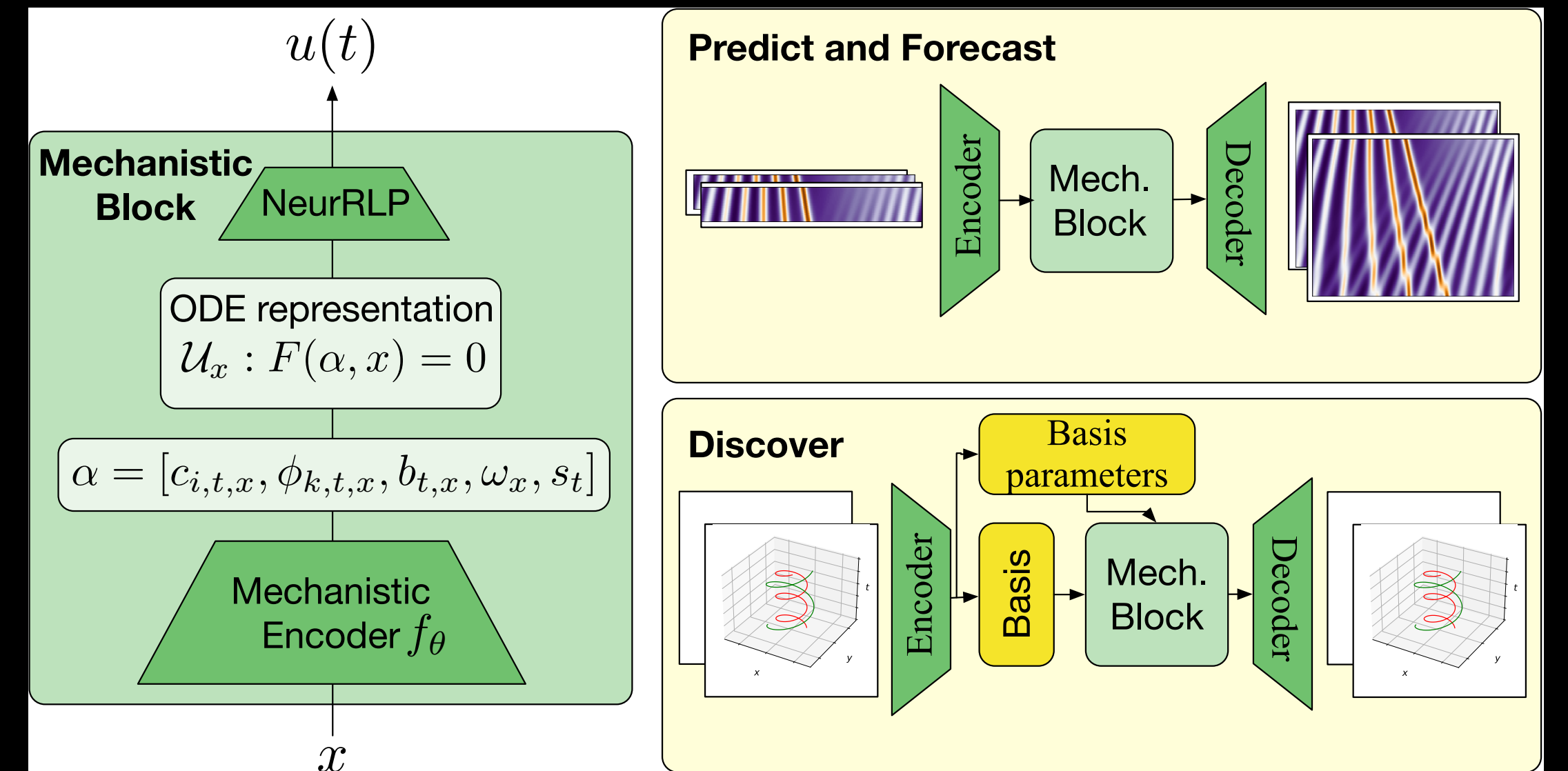
**Table 8.** Segmentation results on the Salsa-dancing dataset.

| Method | NMI ↑ | ARI ↑ | Accuracy ↑ | $F_1$ ↑ |
|---|---|---|---|---|
| Hybrid-SINDy | 0.102 | 0.097 | 0.325 | 0.309 |
| AMORE (ours) | 0.167 | 0.173 | 0.565 | 0.518 |
| AMORE-MIO (ours) | **0.179** | **0.182** | **0.583** | **0.531** |

**Table 9.** Forecasting results in terms of NMAE / NRMSE on the Salsa-dancing dataset.

| Method | One-step | Multi-step |
|---|---|---|
| LLMTime | 0.402 / 0.452 | 0.449 / 0.480 |
| SVI | 0.384 / 0.441 | 0.423 / 0.465 |
| GRASS | 0.285 / 0.344 | 0.313 / 0.359 |
| AMORE (ours) | 0.291 / 0.361 | 0.334 / 0.373 |
| AMORE-MIO (ours) | **0.272 / 0.335** | **0.301 / 0.352** |

# MECHANISTIC NEURAL NETWORKS

- Neural Networks built on data-driven numerical representations

- Uninterpretable → unfit for scientific exploration and analysis

- Mechanisms would be a great alternative but cannot easily learn from data[1]

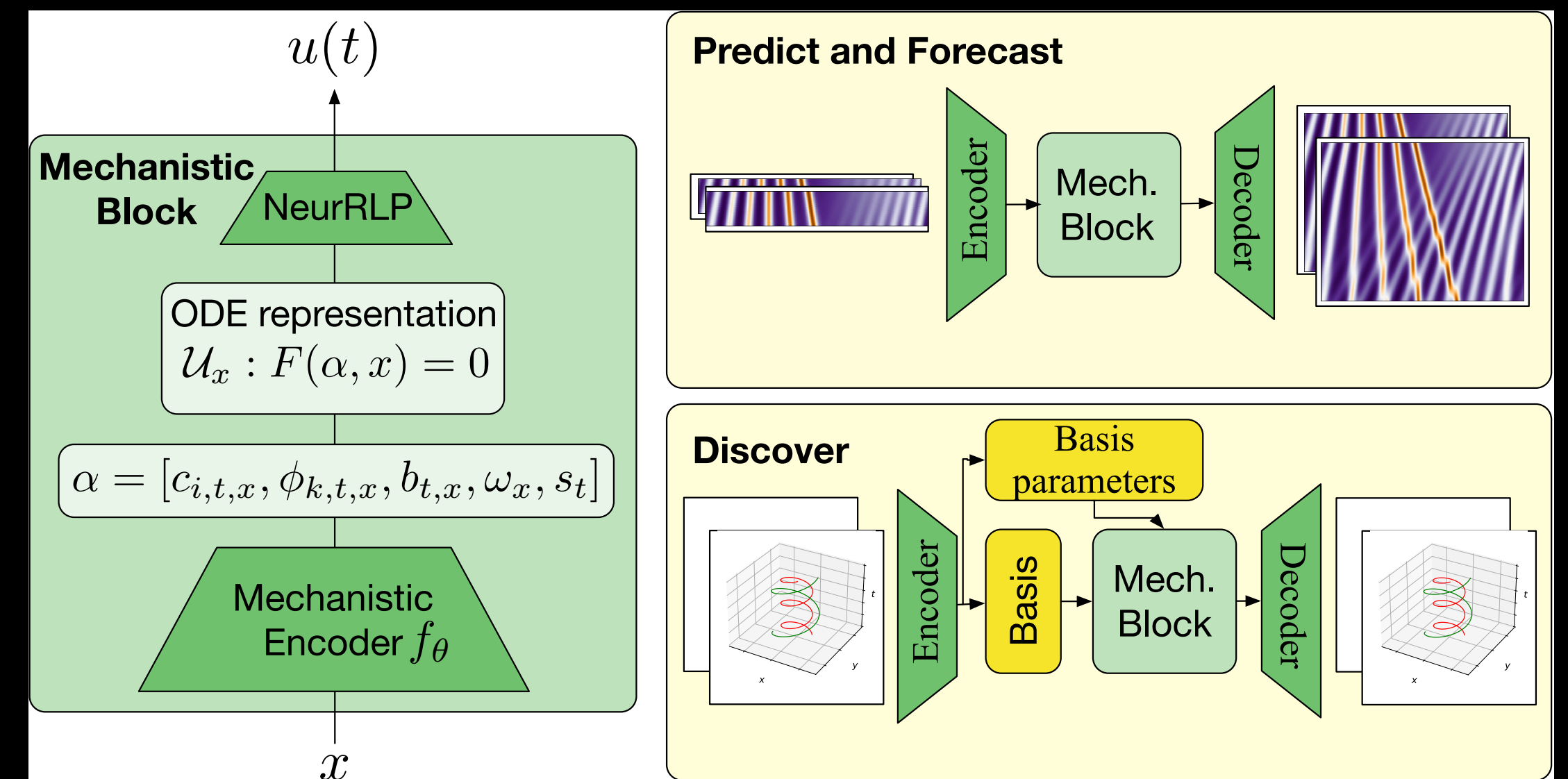- Mechanistic NN modules built on governing equations as representations



$$F(\alpha, x) = \sum_i^d c_i(t; x)u^{(i)} + \sum_k^r \phi_k(t; x)g_k(t, u, u', \dots) - b(t; x)$$

**https://github.com/alpz/mech-nn**

[1] B Schölkopf, F Locatello, S Bauer, N R Ke, N Kalchbrenner, A Goyal, Y Bengio, Towards Causal Representation Learning, Proceedings of the IEEE, 2021

# KEY IDEA

- Define general family of ODEs as governing mechanisms

- Mechanistic NN simultaneously

  - learns the governing ODE explicitly

  - generate new ODEs that explain input

  - solves the ODEs

- Forward pass through the ODE

- Backward pass requires custom, NN-native solver



$$F(\alpha, x) = \sum_i^d c_i(t;x)u^{(i)} + \sum_k^r \phi_k(t;x)g_k(t, u, u', \dots) - b(t;x)$$

**https://github.com/alpz/mech-nn**

# THE MODEL

- Traditional ODE solvers suboptimal: Hard to parallelize, no learned step sizes

- Young showed that Linear ODEs can be solved as Linear Programs[1]

- Continuous ODE: $\displaystyle\sum_i^d c_i(t;x)u^{(i)} + \sum_k^r \phi_k(t;x)g_k(t,u,u',\dots) - b(t;x) = 0$

- Discretize it: $\displaystyle\sum_i^d c_{i,t}u^{(i)} + \sum_k^r \phi_{k,t}g_k(t,u_t,u_t',\dots) - b_t = 0$

- Set up the linear program and solve for $u, u^{(i)}, i = 1,2,\dots$

https://github.com/alpz/mech-nn

[1] J Young, Linear Programming Applied to Linear Differential Equations, Lawrence Berkeley National Laboratory, 1961

# NEURAL RELAXED LP SOLVER

- Define ODE coefficients etc as LP variables

- But LP solvers are not neural network friendly

  - Solutions not differentiable wrt parameters

  - Specialized solvers cannot parallelize easily

  - Constraint matrices too large for solvers

- Relaxing LP to QP and GPU solve KKT conditions

- Similar error bounds like the Euler solver and much faster

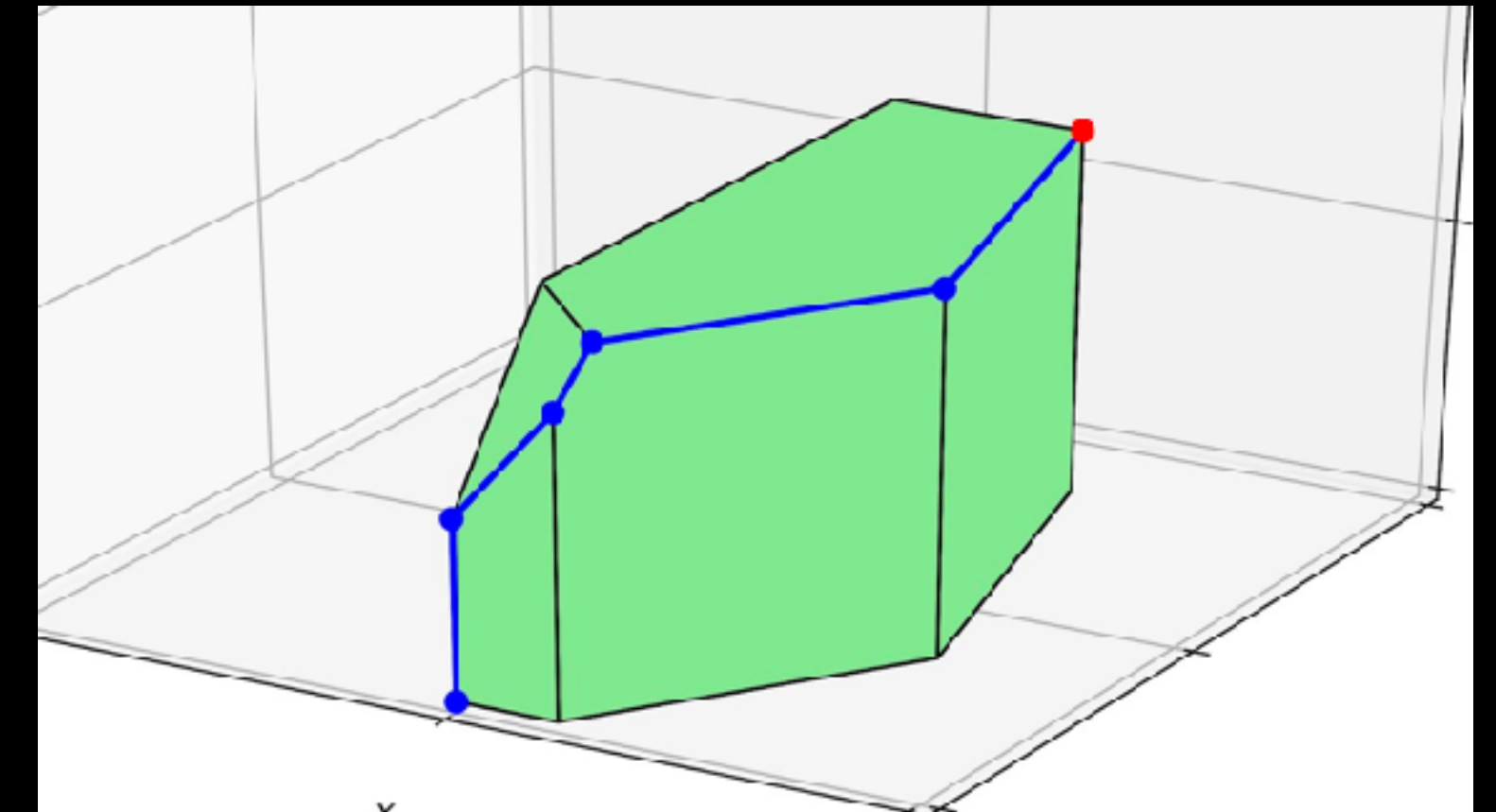- For nonlinear ODEs, the nonlinear terms pushed to the NN loss function



Table 2: Comparing the NeuRLP solver with the RK4 solver with a step size of 0.1 on fitting noisy sinusoidal waves of 300 and 1000 steps. Showing MSE loss and time.

| Steps | QP (seconds) | RK4 (seconds) | QP Loss | RK4 Loss |
|-------|--------------|---------------|---------|----------|
| 40 | 1.52 | 28.06 | 11.4 | 29.3 |
| 100 | 1.61 | 64.57 | 27.9 | 35.6 |
| 300 | 1.76 | 211.52 | 52 | 96.8 |
| 500 | 2.12 | 359.7 | 128 | 301 |
| 1000 | 3.68 | 666.69 | 292 | 589 |

https://github.com/alpz/mech-nn
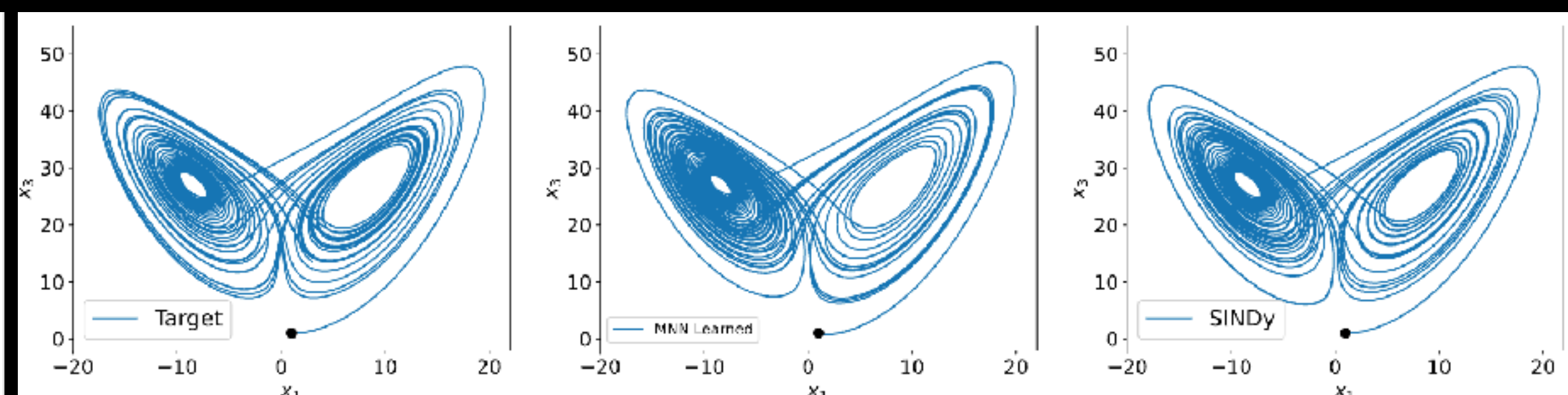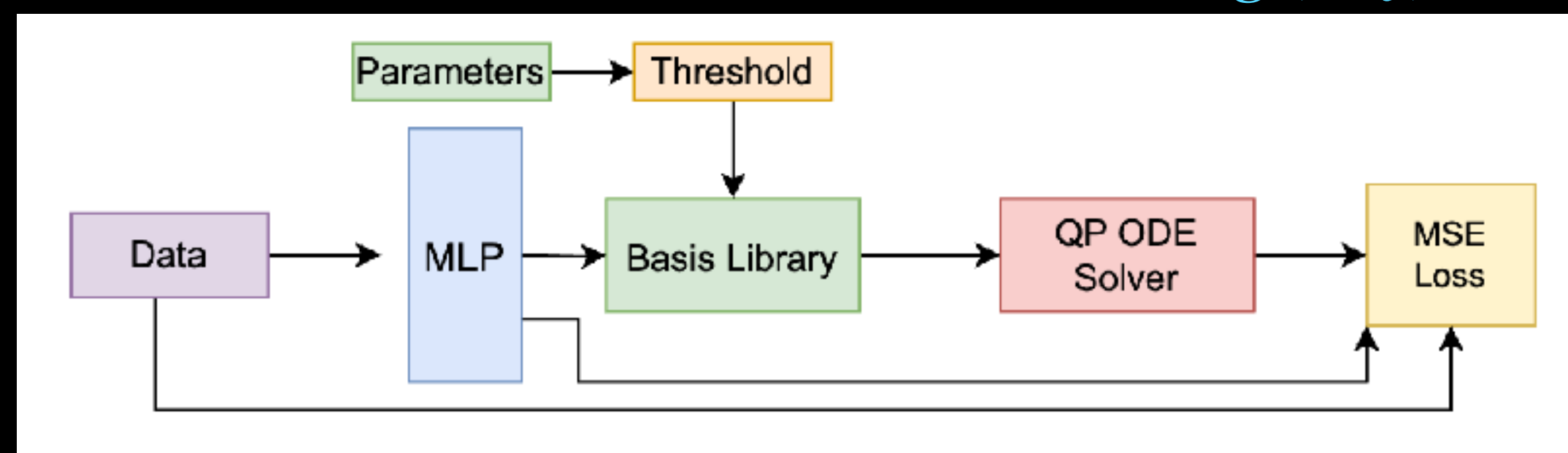
# MECHANISTIC NEURAL NETWORKS

- Weaving in governing equations to an NN impacts lots of scientific applications

- With the same framework, we outperform per task specialized methods

| | Neural ODE,UDE Chen et al. (2018) Rackauckas et al. (2020) | SINDy Brunton et al. (2016) | Neural Operators Li et al. (2020c) | Mech. NN |
|---|---|---|---|---|
| Linear discovery | – | ✓ | – | ✓ |
| Nonlinear discovery | – | – | – | ✓ |
| Physical parameters | ✓ | ✓ | – | ✓ |
| Forecasting | ✓ | – | ✓ | ✓ |
| Interpretability | – | ✓ | – | ✓ |

https://github.com/alpz/mech-nn

# MECH NN FOR DISCOVERY

Solve ODEs of the form: $u' = g(\Theta \xi)$

Chaotic Lorenz



Nonlinear functions

Rational functions



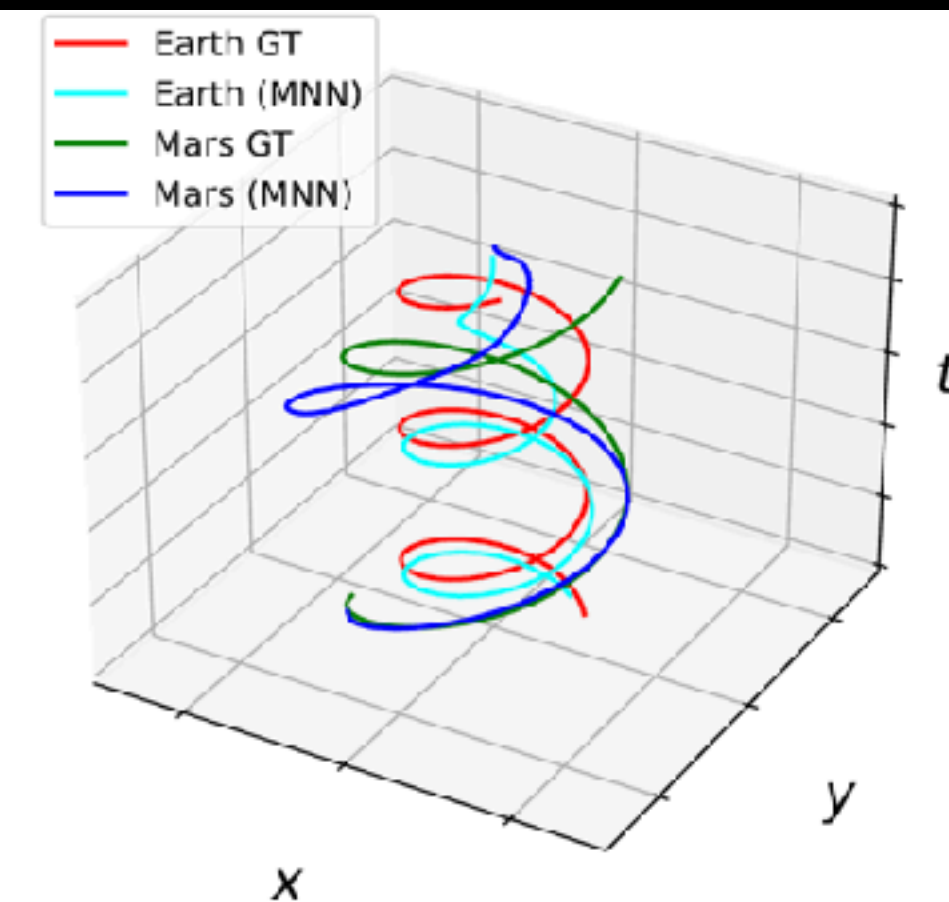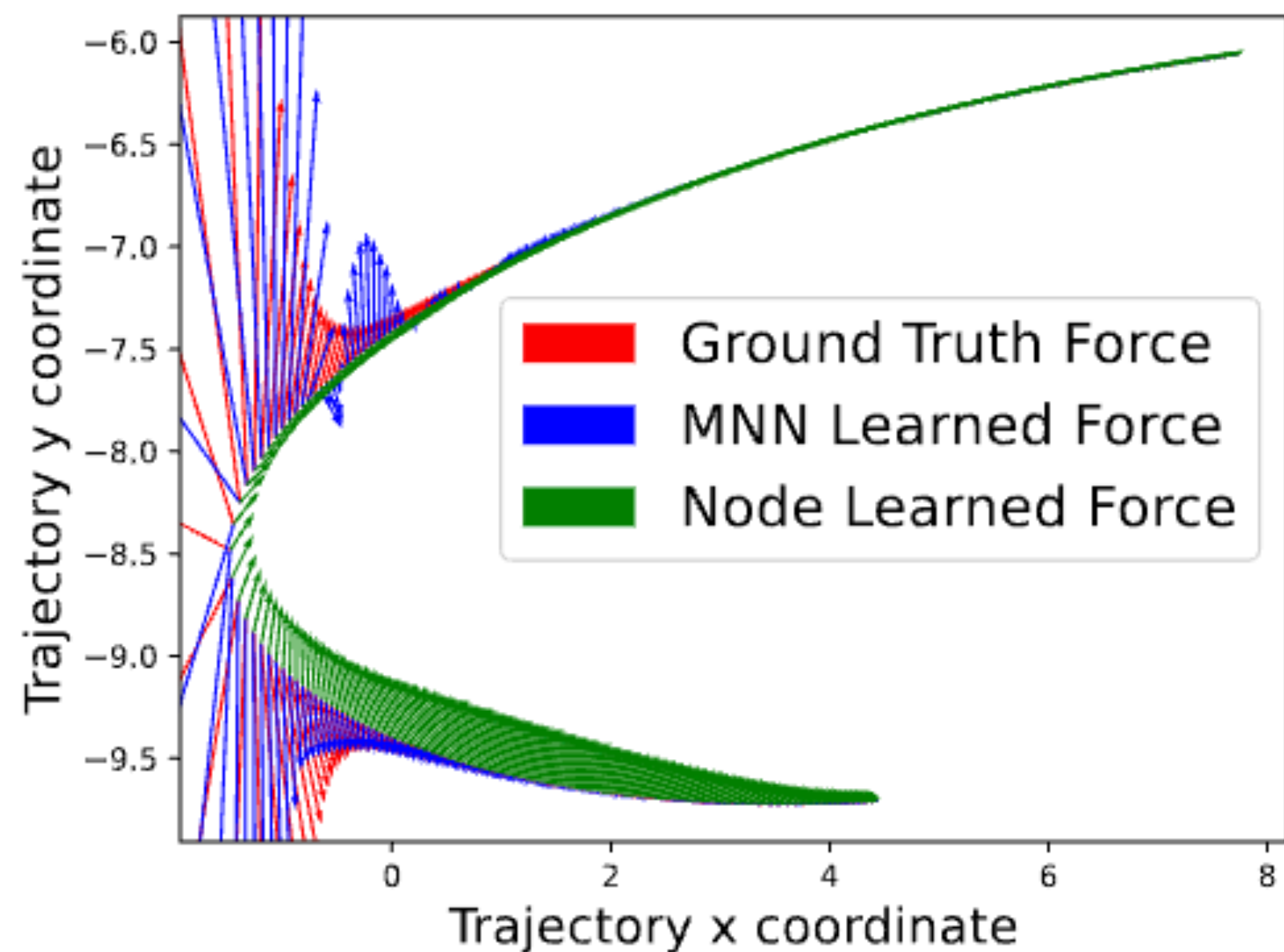$$\frac{dx}{dt} = \tanh(-2x + y)$$
$$\frac{dy}{dt} = \tanh(x + y)$$

Figure 2: Learned ODE vector fields for MNN and SINDy with non-linear tanh function of basis combination and training and test trajectories. Ground truth equation is on the right.
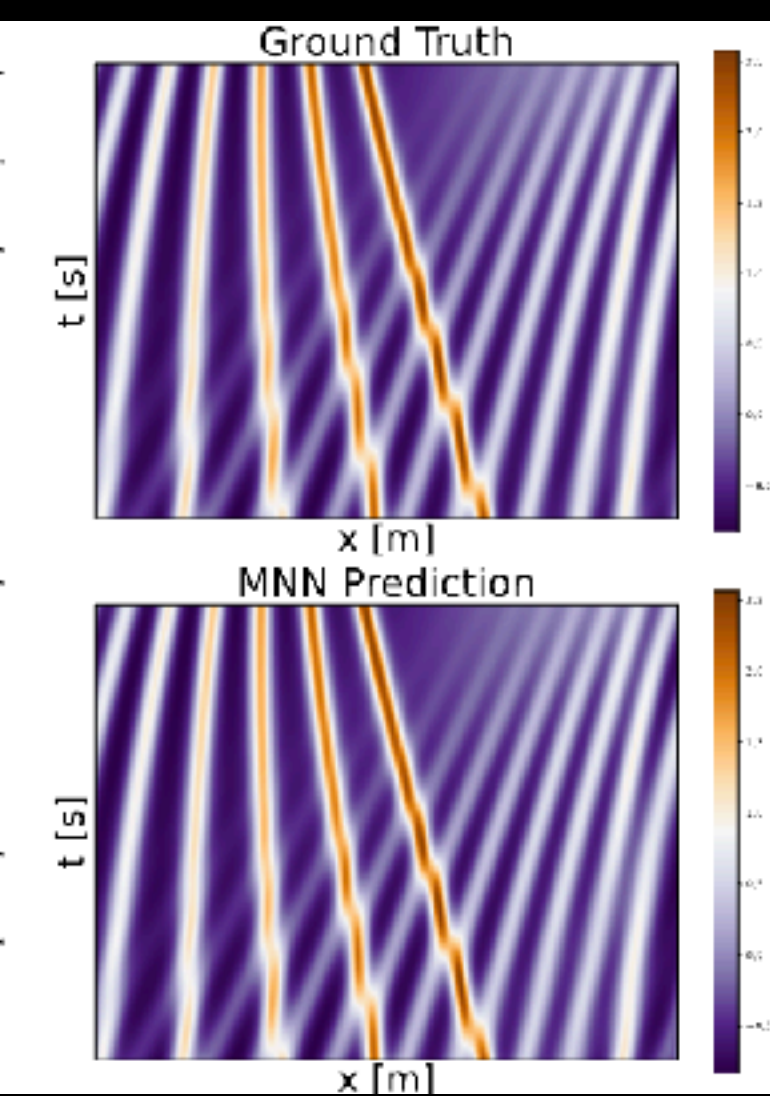
$$\frac{dx}{dt} = \frac{-2x + y}{1 + x^2}$$
$$\frac{dy}{dt} = \frac{x + y}{1 + y^2},$$

# MORE APPLICATIONS

Forecasting
JPL Horizon planetary ephemerides



| Method | Eval. MSE |
|--------|-----------|
| ANODE | 0.0470 |
| NODE | 0.0485 |
| SONODE | 12.200 |
| MNN | 0.0034 |

| Method | Force MSE ↓ | Cosine sim. ↑ | Mass Ratio GT=2 |
|--------|-----------|-------------|-----------------|
| SONODE | 879 | -0.26 | 2.11 |
| MNN | 345 | 0.85 | 2.02 |

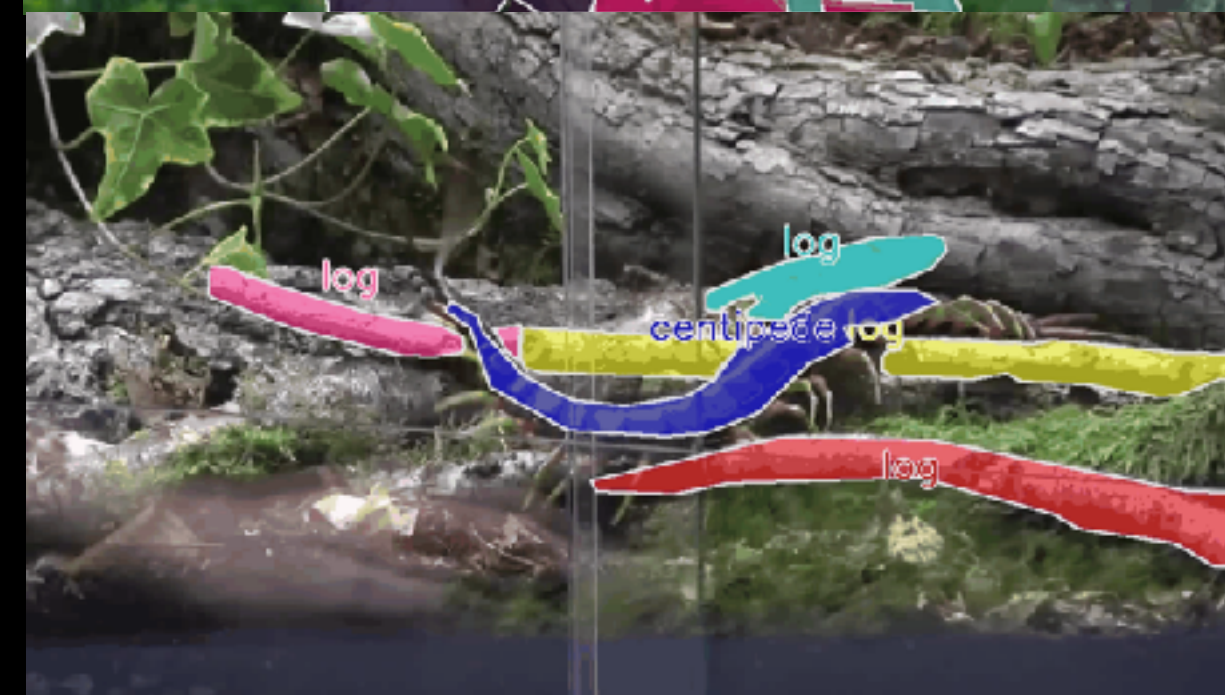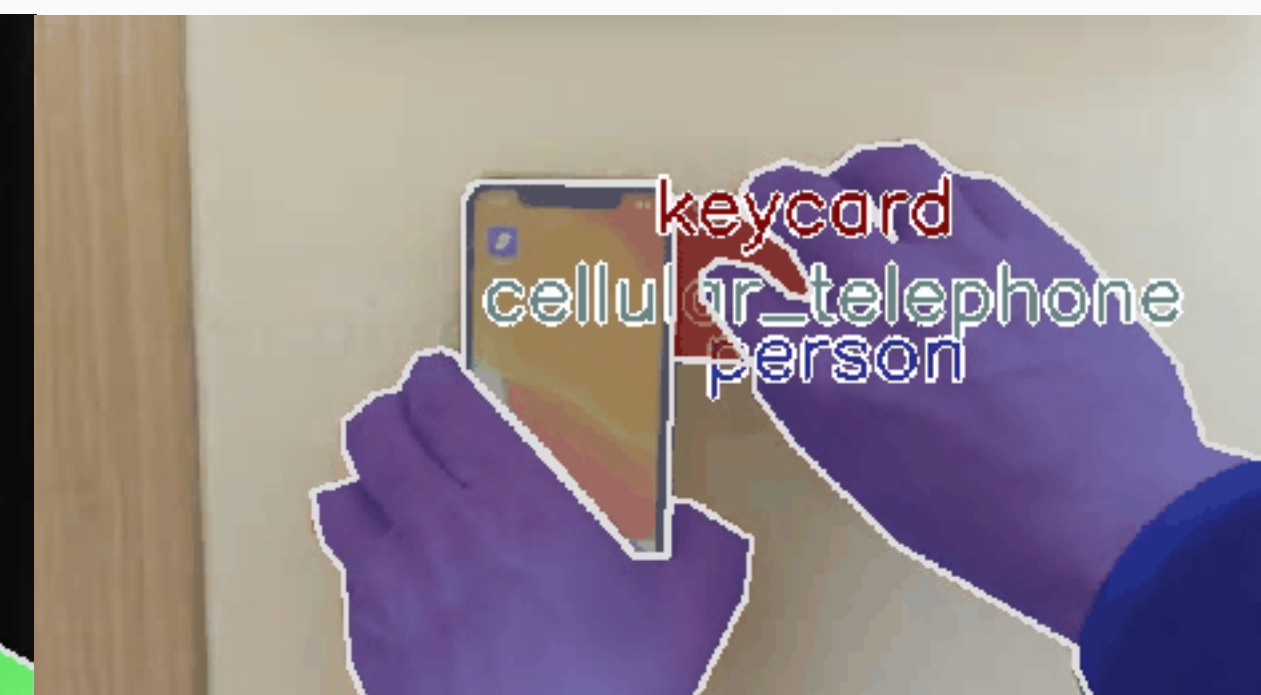| Method | RMSE | |
|--------|------|------|
| | N=512 | N=256 |
| ResNet | 0.0223 | 0.0392 |
| ResNet-LPSDA-1 | 0.0200 | 0.0284 |
| ResNet-LPSDA-2 | 0.0111 | 0.0185 |
| ResNet-LPSDA-3 | 0.0155 | 0.0269 |
| ResNet-LPSDA-4 | 0.0113 | 0.0184 |
| FNO | 0.0276 | 0.0407 |
| FNO-LPSDA | 0.0055 | 0.0132 |
| FNO-AR | 0.0030 | 0.0058 |
| FNO-AR-LPSDA | 0.0010 | 0.0037 |
| Mechanistic NN (50 sec) | 0.0039 | 0.0086 |

Discovering physical parameters
Mass ratio & force distribution 2-body problem

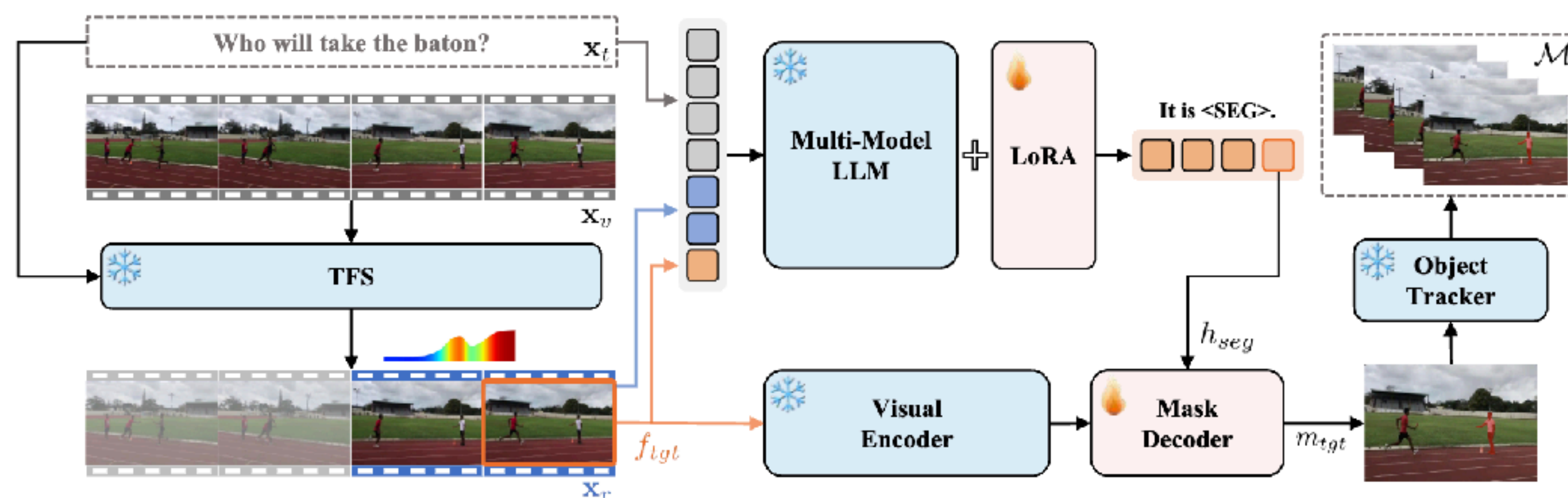Neural PDE Solving
1d KdV and 2d Darcy Flow

WANG ET AL, CVPR, ICCV, 2022-24

# OPEN-WORLD INDUCTIVE BIASES

# REASONING VIDEO INSTANCE SEGMENTATION

- Segment, classify, track by 'commonsense reasoning' w.r.t. the pixels and world knowledge

- Key idea: Scale-up sequence length in the input and relay to LLMs for 'reasoning'

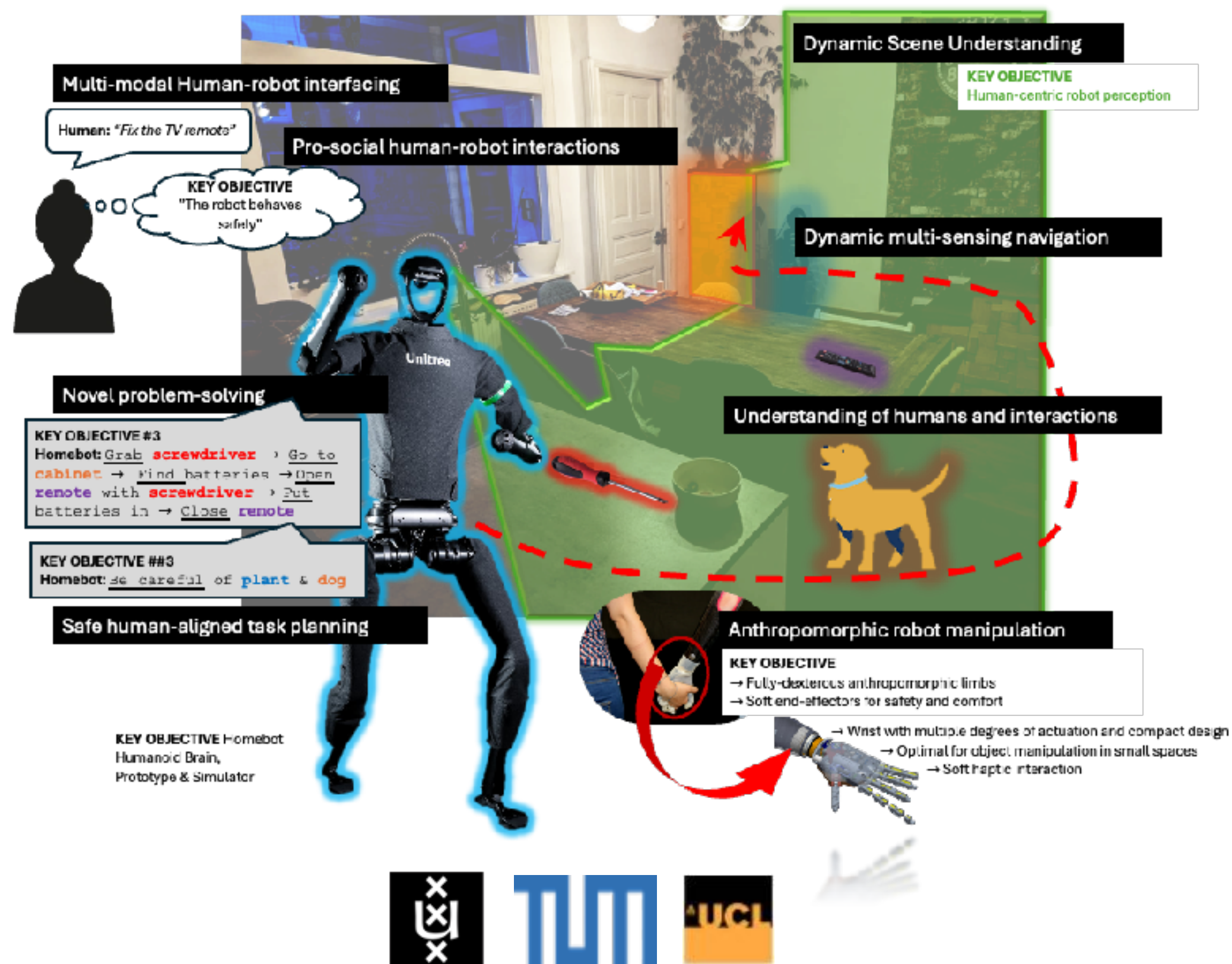- Fantastic if we could pair this with formal reasoning



*"Which ball should first be hit according to the rules?"*
*"Which ball is the target of this shot?"*
*"If this shot goes in, which ball is most likely to be hit next?"*

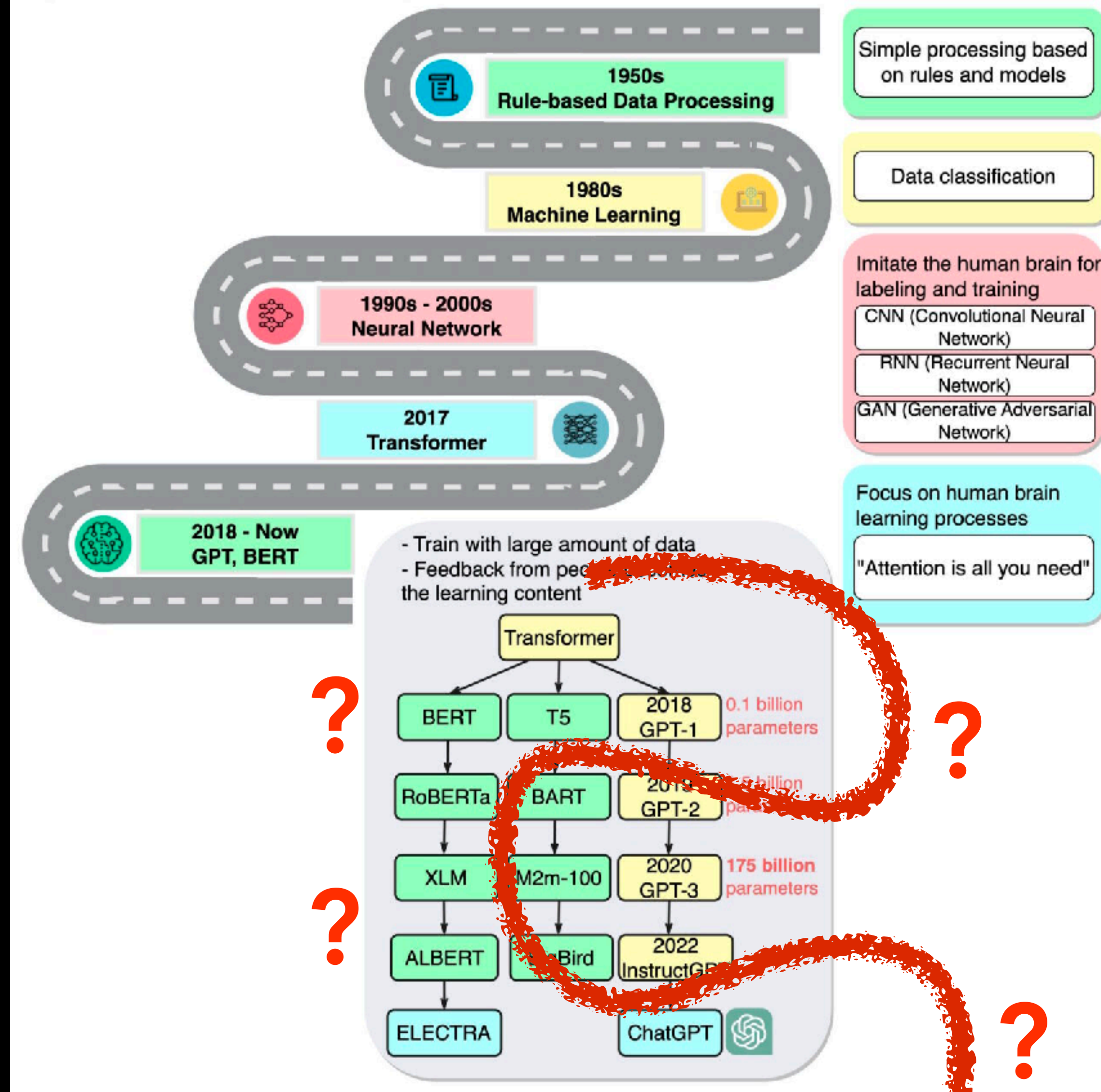noneACADEMIC AI IN THE YEARS OF CHATGPT

# TO WRAP UP

# CONCLUSION



- Not yet end-to-end training, but all parts can be made differentiable

- "Physics & causal grounding": great start toward fully neural robot world models

- And lots of exciting possibilities towards generalization and extrapolation

- We are in good company of Fei Fei Li: https://github.com/cremebrule/digital-cousins

# C H E E R S !

egavves@uva.nl
@egavves

Dynamics
Causal Learning
Open-World
Robots & Systems