

# Transforming CoPhy

---

Alexandre St-Aubin

The University of British Columbia

CPSC 532Y – Causal Machine Learning

# Problem Definition

---

# Background: CoPhy

- **Origin:** Baradel et al. (Facebook AI, SFU), ICLR 2020.
- **Objective:** Teach a model to reason about counterfactual scenarios in physics.
- **Contributions:** A testing benchmark (CoPhy) and a model (CoPhyNet).

## CoPHY: COUNTERFACTUAL LEARNING OF PHYSICAL DYNAMICS

Fabien Baradel<sup>1</sup> Natalia Neverova<sup>2</sup> Julien Mille<sup>3</sup> Greg Mori<sup>4</sup> Christian Wolf<sup>1,5</sup>

<sup>1</sup>Université Lyon, INSA Lyon, CNRS, LIRIS, Villeurbanne, France

<sup>2</sup>Facebook AI Research, Paris, France

<sup>3</sup>Laboratoire d'Informatique de l'Univ. de Tours, INSA Centre Val de Loire, Blois, France

<sup>4</sup>Simon Fraser University and Borealis AI, Vancouver, Canada

<sup>5</sup>Inria, Chroma group, CITI Laboratory, Villeurbanne, France

### ABSTRACT

Understanding causes and effects in mechanical systems is an essential component of reasoning in the physical world. This work poses a new problem of counterfactual learning of object mechanics from visual input. We develop the CoPhy benchmark to assess the capacity of the state-of-the-art models for causal physical reasoning in a synthetic 3D environment and propose a model for learning the physical dynamics in a counterfactual setting. Having observed a mechanical experiment that involves, for example, a falling tower of blocks, a set of bouncing balls or colliding objects, we learn to predict how its outcome is affected by an arbitrary intervention on its initial conditions, such as displacing one of the objects in the scene. The alternative future is predicted given the altered past and a latent representation of the confounders learned by the model in an end-to-end fashion with no supervision of confounders. We compare against feedforward video prediction baselines and show how observing alternative experiences allows the network to capture latent physical properties of the environment, which results in significantly more accurate predictions at the level of super human performance.

# Project Objectives

**Core Goal:** Modernize the CoPhyNet architecture by replacing the sequential RNN encoder with a **Transformer**.

1. **Accelerate Training Efficiency**

Leverage the parallel nature transformers to process time-series data faster.

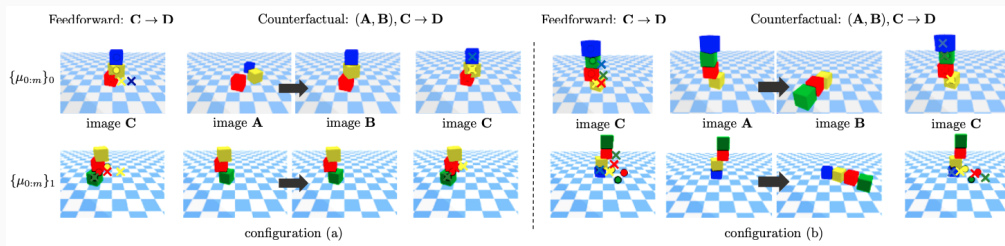
2. **Enhance Physical Reasoning**

Improve the estimation of latent confounders ( $U$ ).

3. **Benchmark Performance**

Validate improvements in stability classification and trajectory MSE on the CoPhy benchmark.

# Problem Formulation



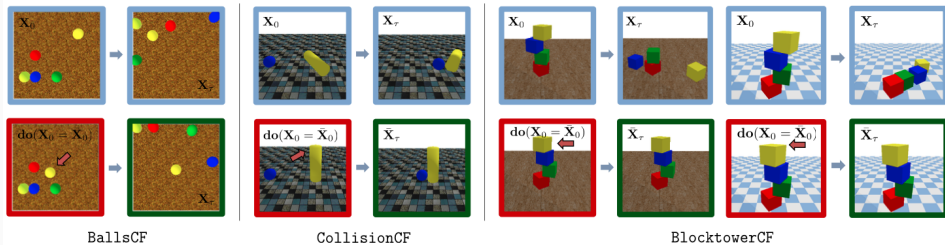
## The Objective

Given a triplet  $\{A, B, C\}$ :

- $A := X_0$  (Initial State)
- $B := \{X_1, \dots, X_\tau\}$  (Observed Outcome)
- $C := \overline{X}_0$  (Intervened Initial State)

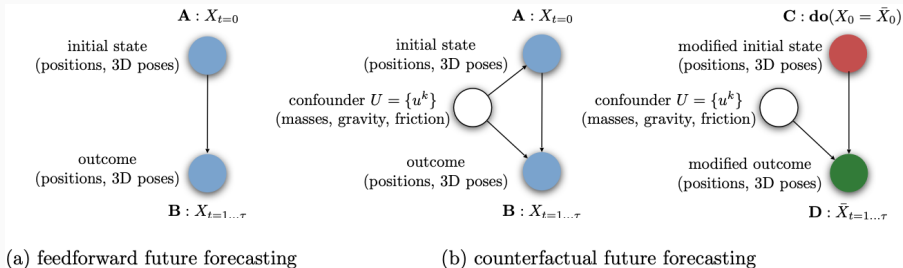
**Goal:** Predict the counterfactual outcome  $D := \{\overline{X}_1, \dots, \overline{X}_\tau\}$ .

# The Tasks



► Watch Demo Video (YouTube)

# Structural Causal Model (SCM)



The physical laws  $U$  determine valid initial states  $X_0$  and outcomes.

# The CoPhy Guarantee: Anti-Shortcut

## The Constraint

*"We enforce existence of at least two different confounder configurations resulting in significantly different trajectories."*

### Example: The Ambiguous Tower

- **Scenario A (High Friction):** The tower stands.
- **Scenario B (Low Friction):** The exact same tower collapses.

*Visually identical  $X_0 \rightarrow$  Divergent Futures.*

**Implication:** Simple extrapolation fails. Estimating  $U$  is **necessary**.



## Why does this matter?

- **Robotics (CausalCF):** Use latent confounder representation to train RL agents.
- **Robustness:** Makes agents robust to unseen changes in the environment (e.g., slippery floors, heavier payloads).
- **Trustworthiness:** Agents make decisions based on physical understanding, not pixel statistics.

# CoPhyNet

---

# The Pearl Framework (2009)

The theoretical foundation for the 3-step process:

## **Theorem 7.1.7**

*Given model  $\langle M, P(u) \rangle$ , the conditional probability  $P(B_A | e)$  of a counterfactual sentence “If it were A then B,” given evidence  $e$ , can be evaluated using the following three steps.*

1. **Abduction** – Update  $P(u)$  by the evidence  $e$  to obtain  $P(u | e)$ .
2. **Action** – Modify  $M$  by the action  $do(A)$ , where  $A$  is the antecedent of the counterfactual, to obtain the submodel  $M_A$ .
3. **Prediction** – Use the modified model  $\langle M_A, P(u | e) \rangle$  to compute the probability of  $B$ , the consequence of the counterfactual.

# The 3-Step Process in CoPhyNet

1. **Abduction (Infer  $U$ ):** Use observed data  $A = X_0, B = X_{1:\tau}$  to compute the latent representation  $U$ .

# The 3-Step Process in CoPhyNet

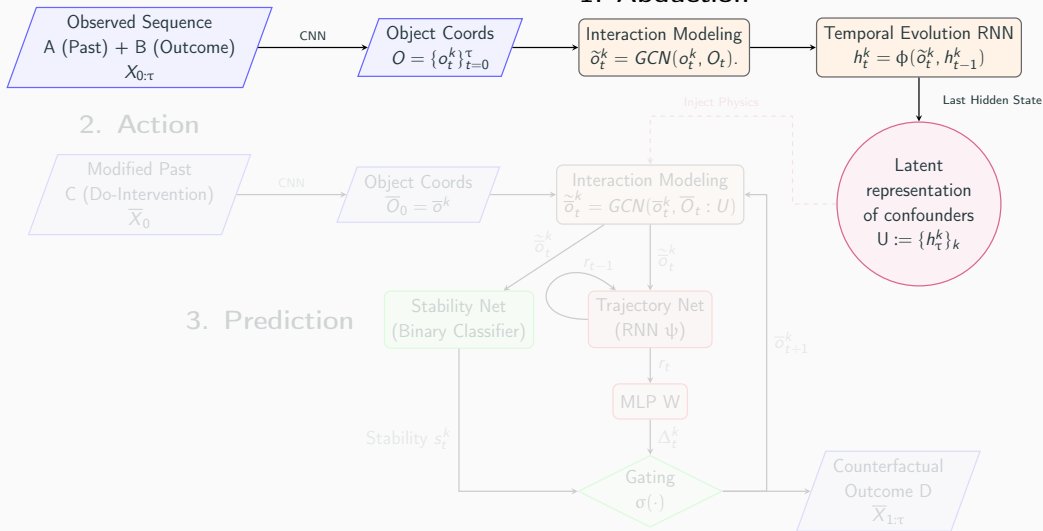
1. **Abduction (Infer  $U$ ):** Use observed data  $A = X_0, B = X_{1:\tau}$  to compute the latent representation  $U$ .
2. **Action (Intervene):** Update the causal model. Keep identified confounders  $U$ , but apply  $do(X_0 = \bar{X}_0)$ .

# The 3-Step Process in CoPhyNet

1. **Abduction (Infer  $U$ ):** Use observed data  $A = X_0, B = X_{1:\tau}$  to compute the latent representation  $U$ .
2. **Action (Intervene):** Update the causal model. Keep identified confounders  $U$ , but apply  $do(X_0 = \bar{X}_0)$ .
3. **Prediction (Simulate):** Compute counterfactual outcome  $D = \bar{X}_{1:\tau}$  using the modified causal model.

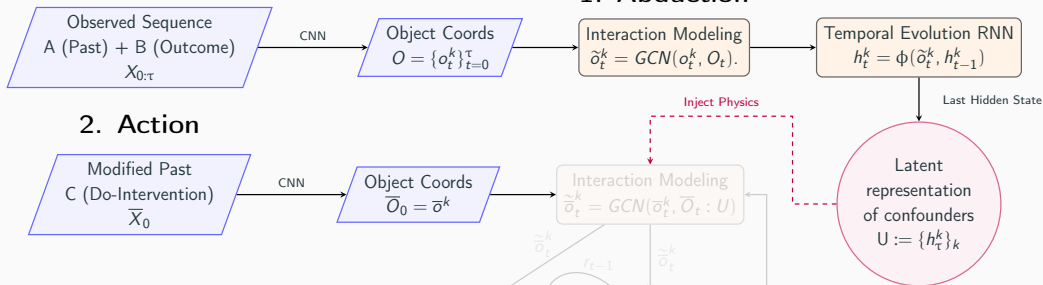
# The 3-Step Process in Architecture

## 1. Abduction

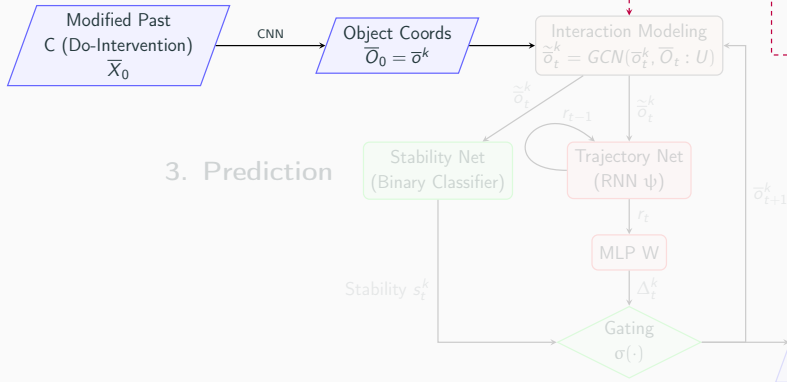


# The 3-Step Process in Architecture

## 1. Abduction



## 2. Action

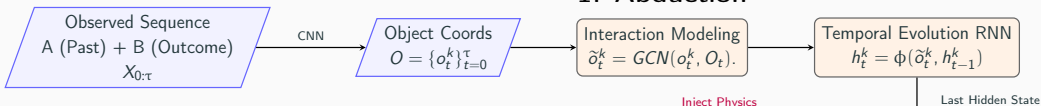


## 3. Prediction

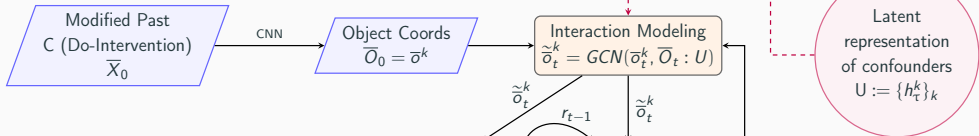


# The 3-Step Process in Architecture

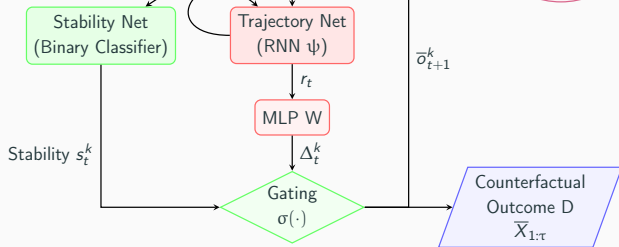
## 1. Abduction



## 2. Action



## 3. Prediction



# Loss Function: "Whether" vs "How"

Separating **Stability** (Whether-causation) from **Trajectory** (How-causation):

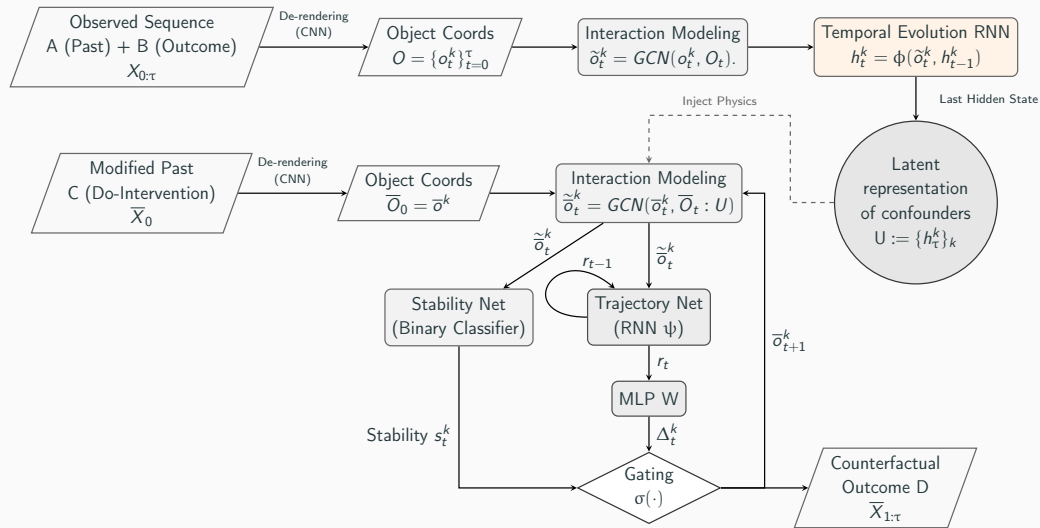
$$\mathcal{L}_{e2e} = \gamma \cdot \sum_{k=1}^K \underbrace{\mathcal{L}_{ce}(s^k, s^{k*})}_{\text{Whether (Classification)}} + \alpha \cdot \sum_{t=0}^{\tau} \left[ \sum_{k=1}^K \underbrace{\mathcal{L}_{mse}(\bar{o}_t^k, \bar{o}_t^{k*})}_{\text{How (Regression)}} \right]$$

Inspired by: Gerstenberg et al. (2015). "How, whether, why: Causal judgments as counterfactual contrasts."

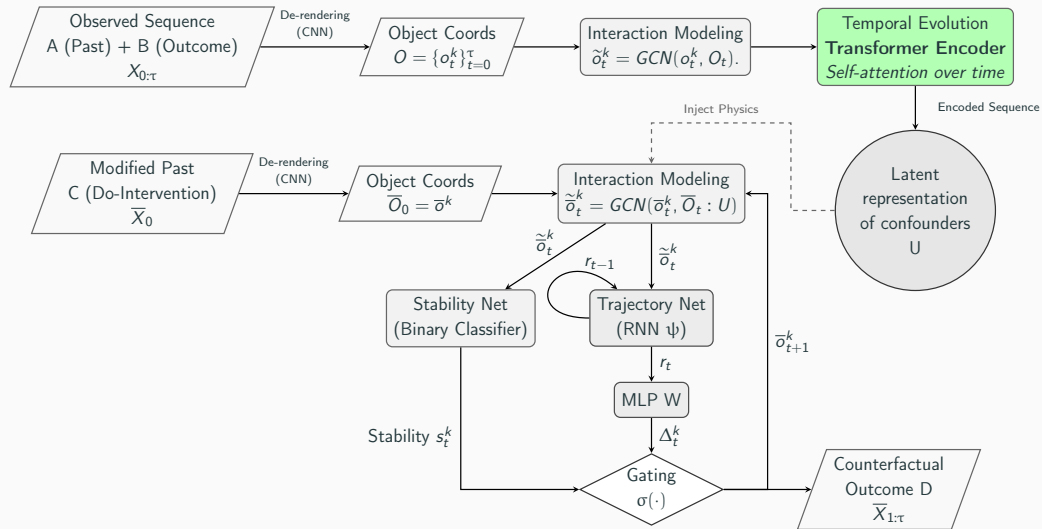
# Transforming CoPhyNet

---

# Original Architecture

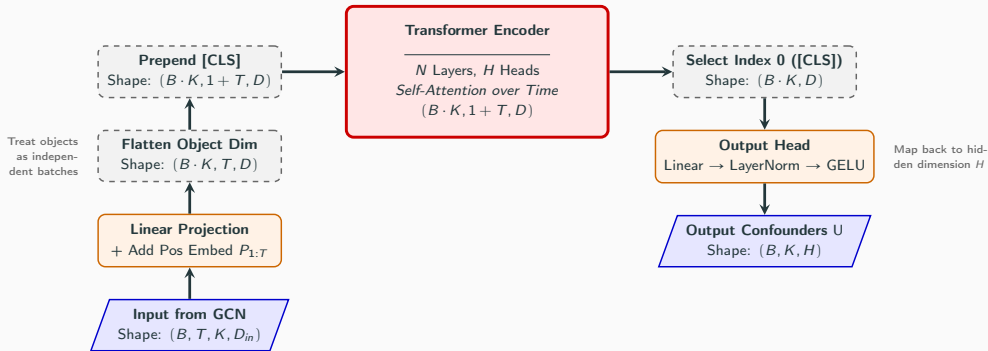


# New Architecture



# Proposal: Transformer Encoder

Replacing the RNN with a Transformer:



# Hypothesis: RNN vs. Transformer

Feature	RNN (GRU)	Transformer
<b>Time View</b>	Sequential	<b>Simultaneous</b>
<b>Info Flow</b>	Bottleneck at $h_t$	Direct access to history
<b>Event Detection</b>	Hard (Vanishing Gradients)	<b>Easy</b> (Self-Attention)
<b>Training</b>	Slow (Sequential)	Fast (Parallelizable)

# Experiments

---



# Experimental Setup

## Dataset & Task

- **Data:** CoPhy Ba11sCF (3 Objects, Normal Split)
- **Goal:** Predict counterfactual stability and 3D trajectory.

## Comparison Framework

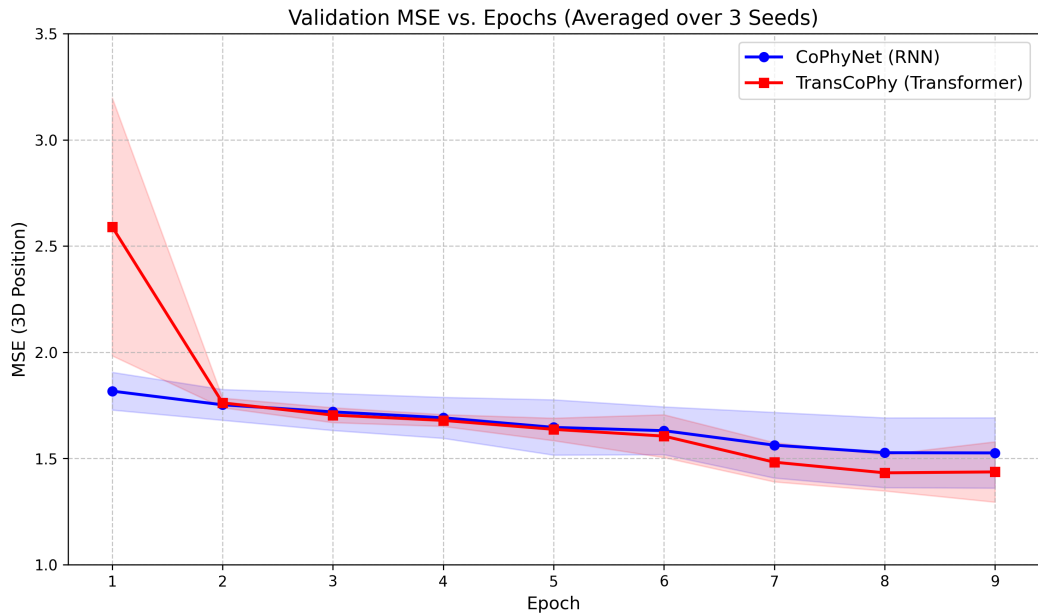
- **Models:** Baseline CoPhyNet vs. TransCoPhy.
- **Metrics:** Mean Squared Error (MSE) on 3D positions and computational cost (Training Time).

Both models use identical batch sizes and hardware. *Note: Learning rates and warmup schedules were tuned specifically for each architecture's stability requirements.*

# Results

Model	CoPhyNet (RNN)	TransCoPhy (Ours)
epoch train time	165s	135s (−18%)

# Results



# Conclusion

---

# Future Work & Conclusion

- **Current Limitation:** Benchmark restricted to 6 seconds @ 5fps.
- **Scaling Up:** Newer physics datasets use 25fps+ and longer horizons.
- **The Transformer Advantage:** RNNs struggle with long sequences (vanishing gradients). Transformers should show drastic improvements on longer, higher-frequency simulations.

**Thank You.**