# Reinforcement Learning (RL)

## Goal

Design **agent** that steers an **environment** to maximize a **reward**



**Agent**      **Environment**

## Applications



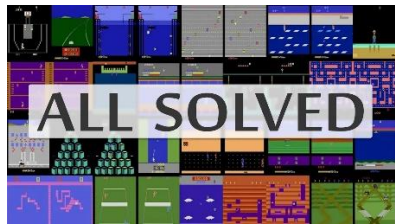| Computer gaming | Playing Go | Autonomous driving | Medical treatment | Manufacturing optimization |

# Learning via Trial & Error

Learning an agent typically entails **trial & error** in environment



Feasible in some applications; **prohibitively costly/dangerous** in others



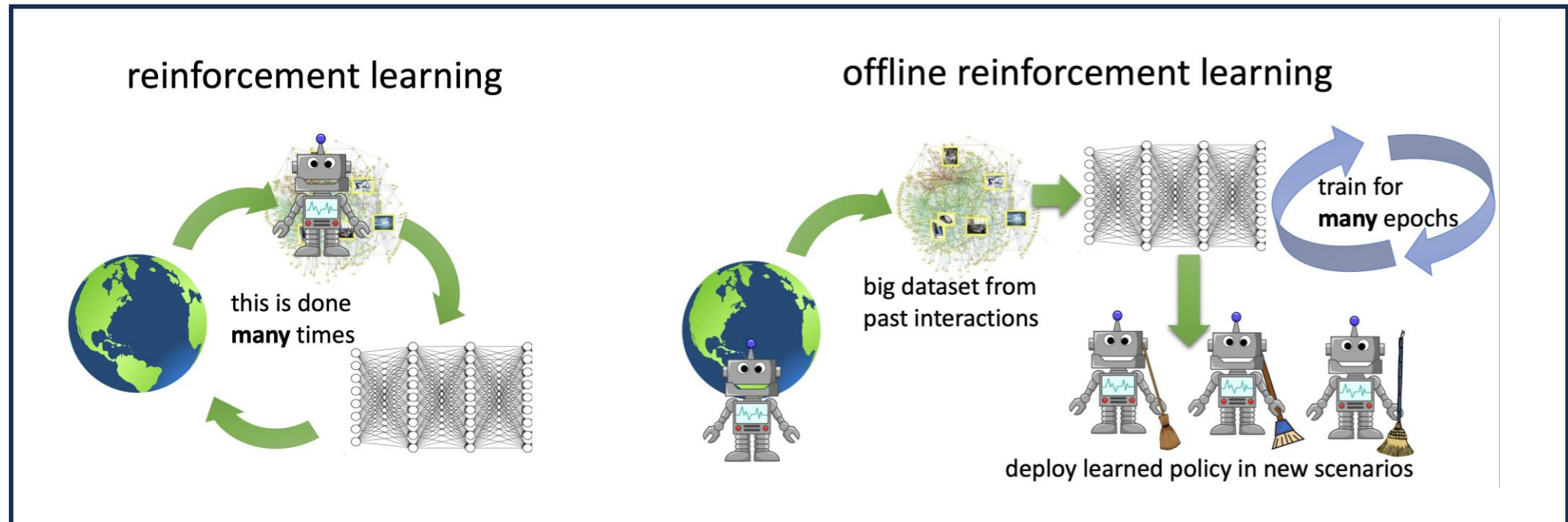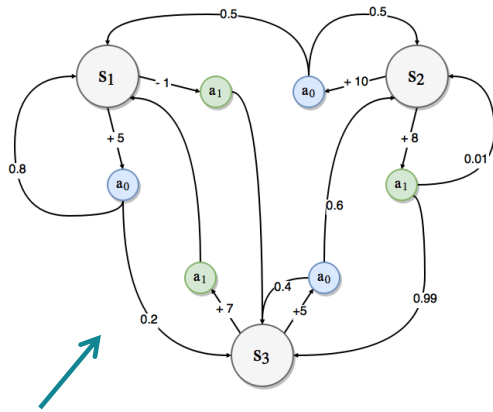| Computer gaming | Playing Go | Autonomous driving | Medical treatment | Manufacturing optimization |

# Offline RL

## Goal

Learn an agent without trial & error in environment

# Conventional Offline RL Methods

Designed for **Markov Decision Process** (**MDP**) environments

**MDP environment**

**Conventional offline RL methods**

*Value-based methods*

|  | A₁ | A₂ | ... | Aₘ |
|---|---|---|---|---|
| $S_1$ | Q(S₁, A₁) | Q(S₁, A₂) |  | Q(S₁, Aₘ) |
| $S_2$ | Q(S₂, A₁) | Q(S₂, A₂) |  | Q(S₂, Aₘ) |
| ⋮ |  |  | ⋱ | ⋮ |
| $S_N$ | Q(Sₙ, A₁) | Q(Sₙ, A₂) | ... | Q(Sₙ, Aₘ) |

$$Q(s, a) = r(s, a) + \gamma \max_a Q(s', a)$$

*Policy-based methods*

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}\left[\sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t) \hat{R}(\tau, t)\right]$$

MDP environment is **fully observable**: its observations reveal its full state

**Challenge:** many real-world environments are not fully observable

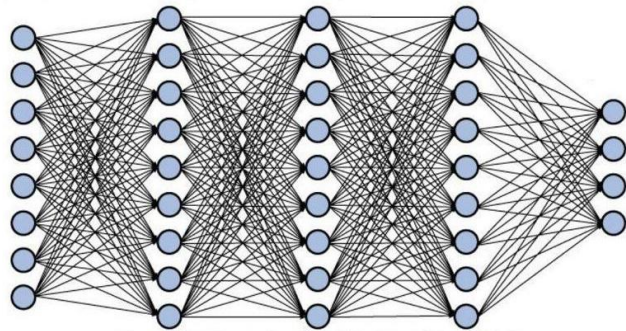Autonomous driving ❌

Medical treatment ❌

Manufacturing optimization ❌

# An Appeal to Supervised Learning

In supervised learning:

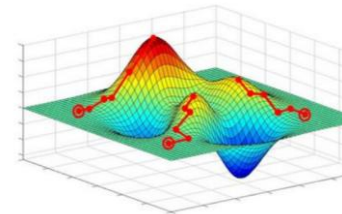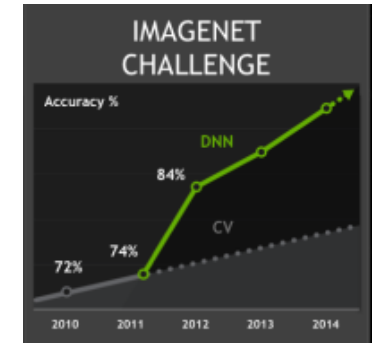**Overparameterized neural networks** (**NNs**) trained by **gradient descent** (**GD**) led to a breakthrough



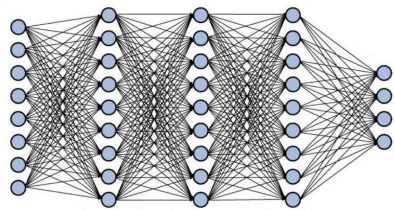**# of learned weights**

>>

**Training set size**

**GD training**

**Breakthrough results**

**Q:** Can a similar approach be taken in offline RL?

# An Approach to Offline RL Inspired by Supervised Learning

**Step 1** **Learn Environment Model**

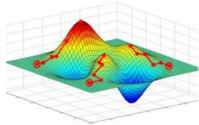Overparameterized NN trained by GD over pre-recorded data

| Time | Action | Observation |
|------|--------|-------------|
| $t$ | $a_t$ | $o_t$ |
| $t+1$ | $a_{t+1}$ | $o_{t+1}$ |
| $t+2$ | $a_{t+1}$ | $o_{t+1}$ |
| ... | ... | ... |

**Overparam NN**

**Pre-recorded data**

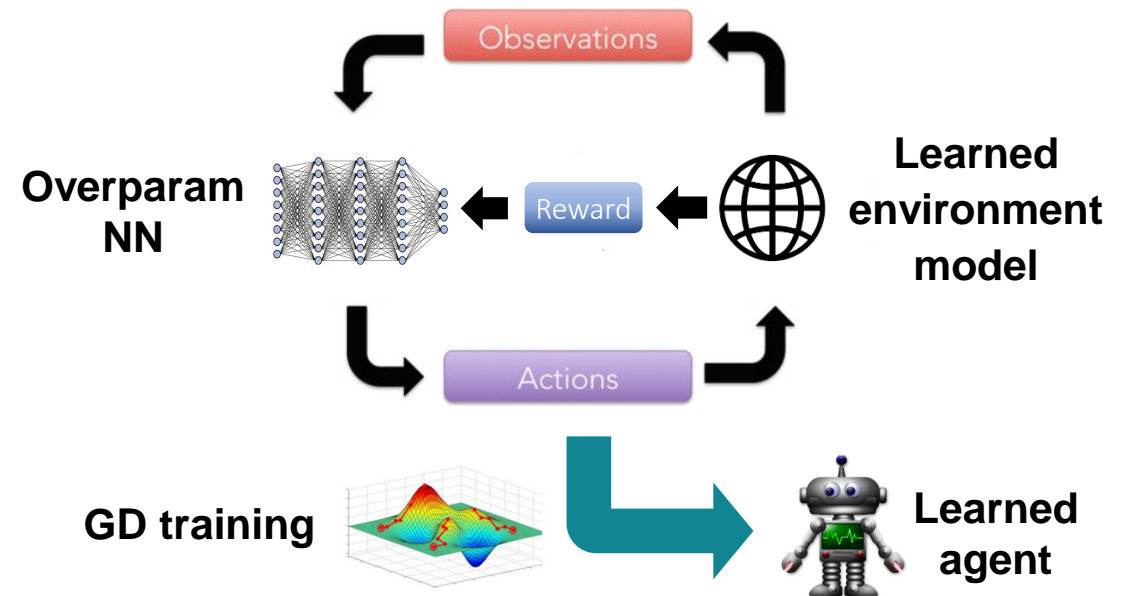**GD training**

**Learned environment model**

**Step 2** **Learn Agent**

Overparameterized NN trained by GD over learned environment model

Observations

**Overparam NN**

Reward

**Learned environment model**

Actions

**GD training**

**Learned agent**

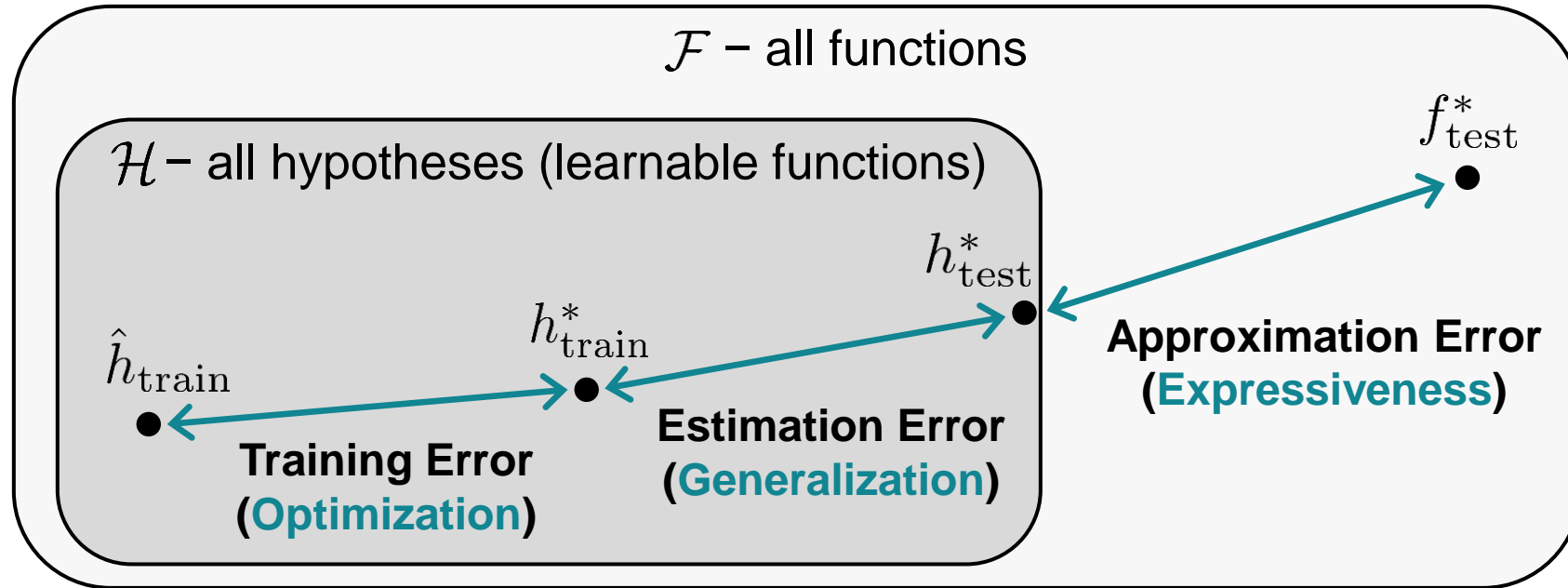**Q:** Can this approach work well enough in critical environments?

Medical treatment

Manufacturing optimization

# Three Pillars of Statistical Learning:
# Expressiveness, Generalization and Optimization
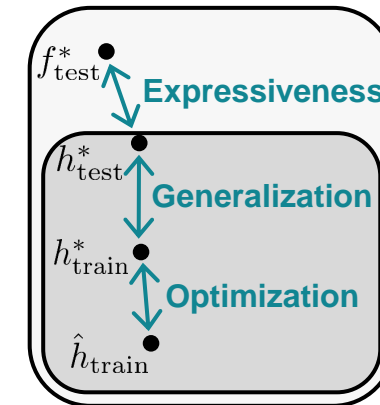


$f_{\text{test}}^*$ – optimal function (minimizer of test error over $\mathcal{F}$)

$h_{\text{test}}^*$ – optimal hypothesis (minimizer of test error over $\mathcal{H}$)

$h_{\text{train}}^*$ – empirically optimal hypothesis (minimizer of train error over $\mathcal{H}$)

$\hat{h}_{\text{train}}$ – returned hypotheiss

# Three Pillars in Supervised Learning



**Overparam NN**  **Training Set**  **GD training**  **Breakthrough results**

Various theoretical guarantees:

| Expressiveness | Generalization | Optimization |
|---|---|---|
| [Telgarsky 15'] | [Lampinen and Ganguli 19'] | [Saxe et al. 14'] |
| [Eldan and Shamir 15'] | [Arora et al. 19'] | [Bartlett el al. 18'] |
| [Cohen et al. 16'] | [Advani and Saxe 20'] | [Arora et al. 18'] |
| [Raghu et al. '16] | [Chizat and Bach 20'] | [Arora et al. 19'] |
| [Levine el al. 17'] | [Razin and Cohen 20'] | [Ji and Telgarsky 20'] |
| [Razin et al. 22'] | [Razin et al. 21] | [Elkabetz and Cohen 21'] |
| ... | ... | ... |

# Three Pillars in Offline RL



Significant challenges:

o **Expressiveness**: capacity of NN arch to reach low test error is **highly obscured** by **dynamics**

o **Generalization**: test distribution can **vastly differ** from train (**distribution shift**)

o **Optimization**: train loss is **extremely complex** (GD faces **instability**, **vanishing gradients**, etc.)

# Three Pillars in Offline RL (cont.)

Nascent theory gives positive indications:

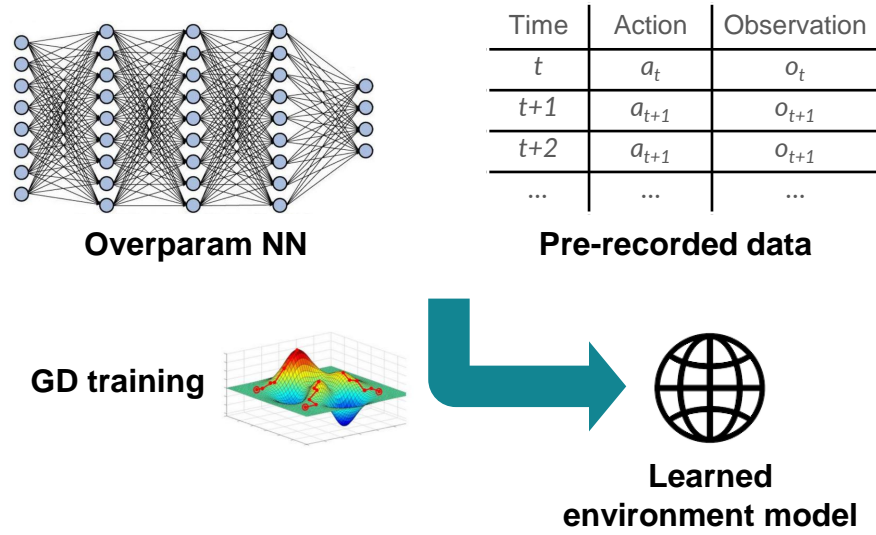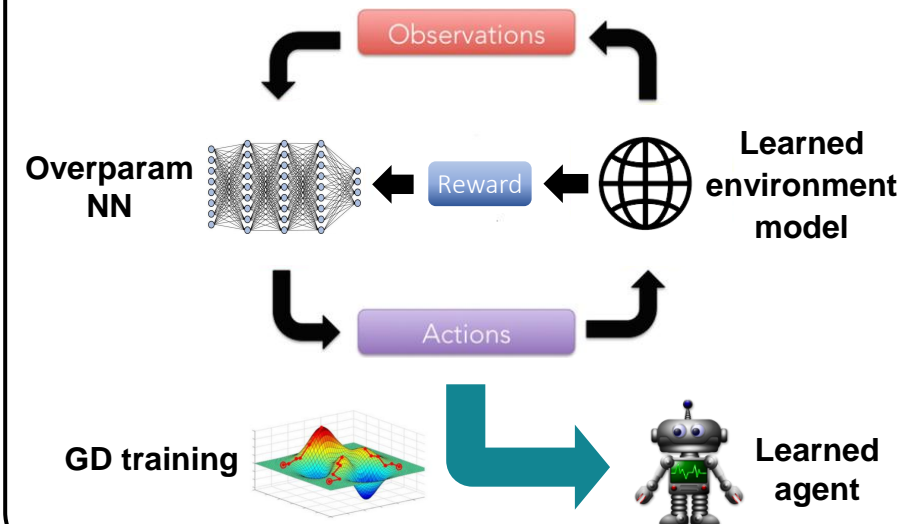| | |
|---|---|
| **On the Implicit Bias of Gradient Descent for Temporal Extrapolation**<br><br>Cohen-Karlik + Ben David + *C* + Globerson<br>*AISTATS 2022* | **Learning Low Dimensional State Spaces with Overparameterized Recurrent Neural Nets**<br><br>Cohen-Karlik + Menuhin-Gruman + Giryes + *C* + Globerson<br>*ICLR 2023* |
| **Implicit Bias of Policy Gradient in Linear Quadratic Control: Extrapolation to Unseen Initial States**<br><br>Razin* + Alexander* + Cohen-Karlik + Giryes + Globerson + *C*<br>*ICML 2024* | **Provable Benefits of Complex Parameterizations for Structured State Space Models**<br><br>Ran-Milo + Lumbroso + Cohen-Karlik + Giryes + Globerson + *C*<br>*NeurIPS 2024* |
| **The Implicit Bias of Structured State Space Models Can Be Poisoned with Clean Labels**<br><br>Slutzky* + Alexander* + Razin + *C*<br>*Under Review 2025* | **Implicit Bias of Neural Networks for Control: A Tendency for Safety (tentative)**<br><br>Slutzky + Alexander + Nagel + *C*<br>*Work in Progress 2025* |

# Offline RL in the Wild?



**Step 1** Learn Environment Model

Overparam NN

| Time | Action | Observation |
|------|--------|-------------|
| $t$ | $a_t$ | $o_t$ |
| $t+1$ | $a_{t+1}$ | $o_{t+1}$ |
| $t+2$ | $a_{t+1}$ | $o_{t+1}$ |
| ... | ... | ... |

**Pre-recorded data**

GD training

**Learned environment model**

**Step 2** Learn Agent

Overparam NN

Observations

Reward

Learned environment model

Actions

GD training

**Learned agent**

*in the wild*

**Q:** Can this approach work well enough in critical environments?

Medical treatment

Manufacturing optimization

# Case Study I: Medical Treatment

## Machine Learning for Mechanical Ventilation Control

Daniel Suo[*†], Naman Agarwal[*], Wenhan Xia[*†], Xinyi Chen[*†], Udaya Ghai[*†], Alexander Yu[*], Paula Gradu[*], Karan Singh[*†], Cyril Zhang[*†], Edgar Minasyan[*†], Julienne LaChance[†], Tom Zajdel[†], Manuel Schottdorf[†], Daniel Cohen[†], Elad Hazan[*†]
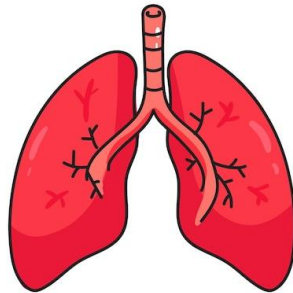
### Abstract

Mechanical ventilation is one of the most widely used therapies in the ICU. However, despite ventilation, a form of assist-control ventilation, evidence suggests that a combination of high peak pressure and high tidal volume can lead to tissue injury in

**Step 1** **Learn Environment Model**
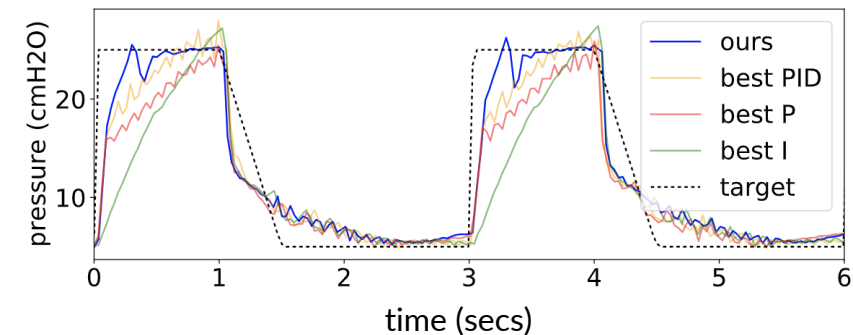
Use pre-recorded data for learning an NN lungs model

**Step 2** **Learn Agent**

Use learned lungs model for learning an NN mechanical ventilator controller
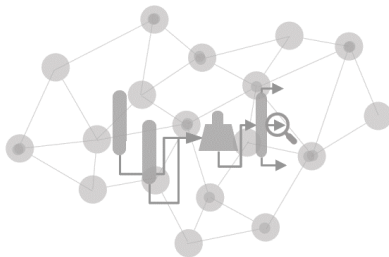


Critical environment? ✅   In the wild? ❌

# Case Study II: Manufacturing Optimization



## Step 1 Learn Environment Model
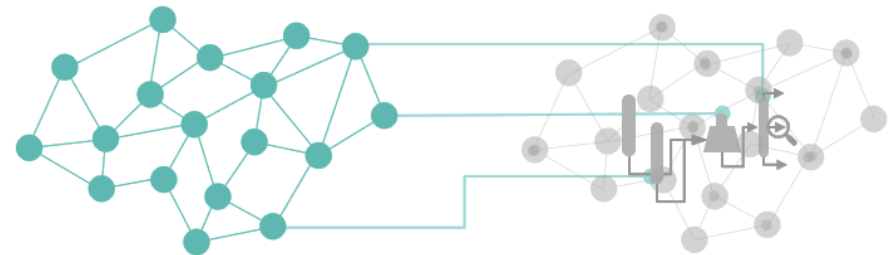
Use pre-recorded data for learning an NN plant model



## Step 2 Learn Agent

Use learned plant model for learning an NN controller



Critical environment? ✅

In the wild? ✅

# Case Study II: Manufacturing Optimization (cont.)

**70+** Applications

**30+** Process Plants

**6+** Years of Model Engagement

**15-30%** Reduced Natural Gas Usage

**1-3%** Yield Improvement

Optimized by **IMUBIT**

Chevron

ExxonMobil

BAZAN GROUP

HUNT REFINING COMPANY

Big West Oil LLC

Oxbow

FLINT HILLS resources

HF Sinclair

preem

INVISTA

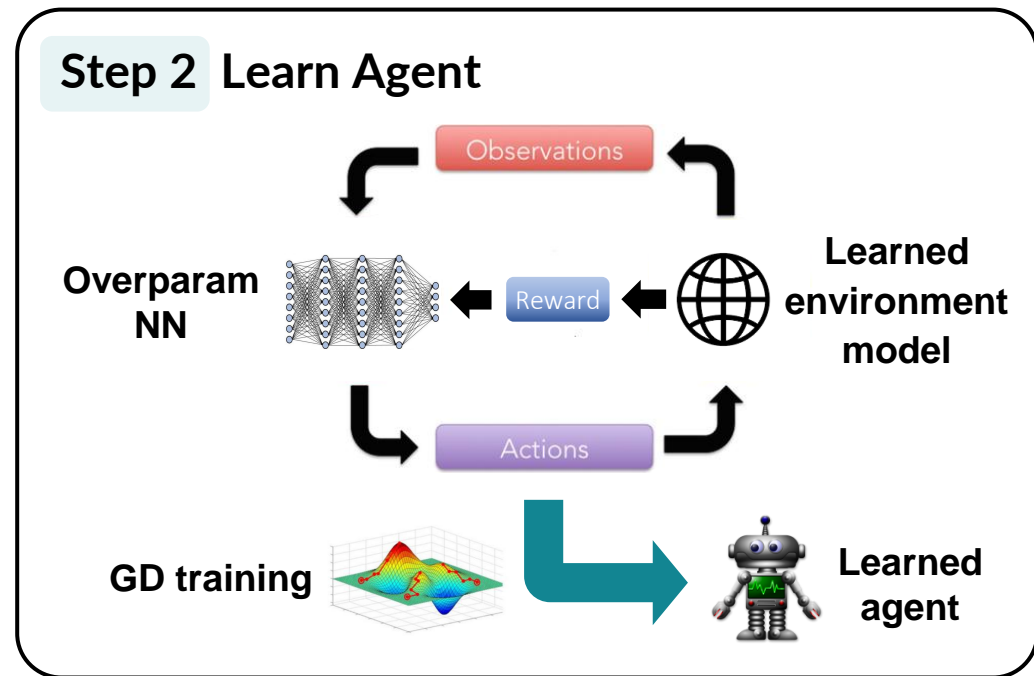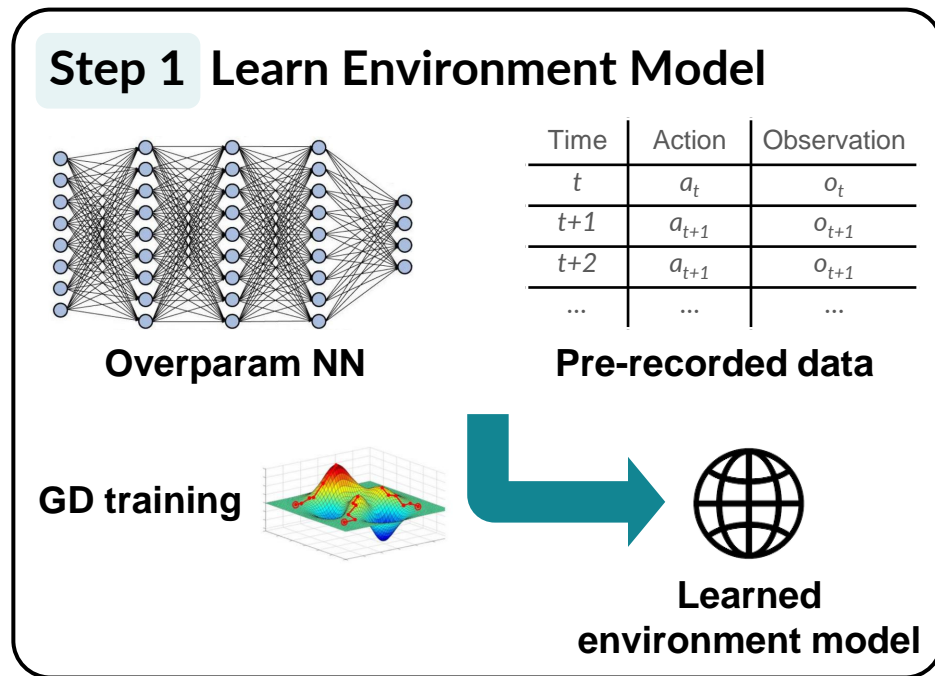Monroe Energy

ASH GROVE A CRH COMPANY

CITGO

Eagle Materials

# Conclusion

In critical applications, **trial & error is prohibitively costly/dangerous** ➡ **RL** must be **offline**

Supervised learning success of **overparameterized NNs** trained by **GD** inspires offline RL approach :



Nascent **theory supports** the approach

Approach **successfully demonstrated** in critical application **in the wild**!

# Perspective

Practical progress in **AI is currently driven by trial & error**



**Less suitable for critical applications**

Medical
treatment



Manufacturing
optimization



For AI to proliferate in critical applications, **theory may be necessary**

# Thank You!