

---

# Decision Transformer: Reinforcement Learning via Sequence Modeling

---

**Lili Chen<sup>\*,1</sup>, Kevin Lu<sup>\*,1</sup>, Aravind Rajeswaran<sup>2</sup>, Kimin Lee<sup>1</sup>,  
Aditya Grover<sup>2</sup>, Michael Laskin<sup>1</sup>, Pieter Abbeel<sup>1</sup>, Aravind Srinivas<sup>†,1</sup>, Igor Mordatch<sup>†,3</sup>**

<sup>\*</sup>equal contribution   <sup>†</sup>equal advising

<sup>1</sup>UC Berkeley   <sup>2</sup>Facebook AI Research   <sup>3</sup>Google Brain

{lilichen, kzl}@berkeley.edu

CS 6784 Advanced Topics in Machine Learning: ML for Feedback Systems

Paper Presentation

Menghan Xu

Nov. 4, 2025

# Problem Statement

**Treat sequential decision making as RL and solve it via Transformers**

## **What we Do?**

- how to formalize?
- how to apply Transformer?

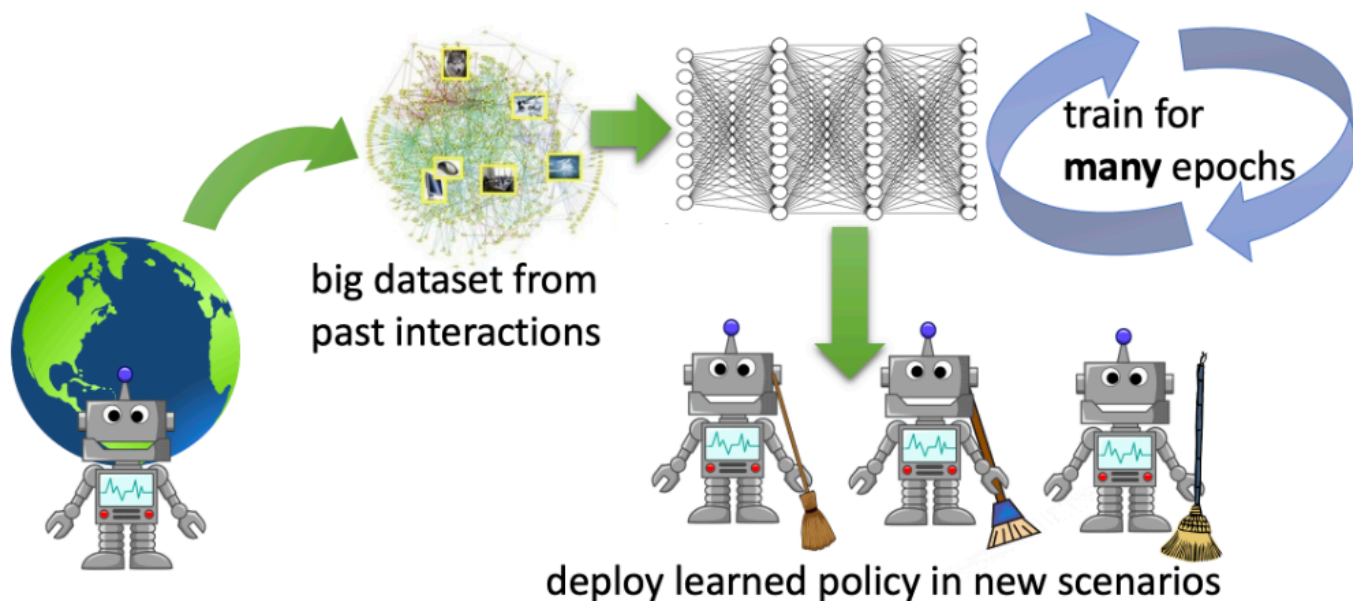
## **Why we do it?**

- does it actually work?
- why can this approach succeed?

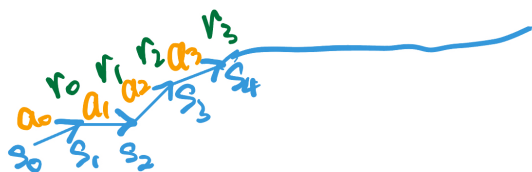
# What we do?

## Problem set up - offline RL

offline reinforcement learning



trajectory  $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$



Markov Decision Process (MDP)

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, \mathcal{R})$$

- $s_t \in \mathcal{S}$ : state
  - $a_t \in \mathcal{A}$ : action
  - $P(s'|s, a)$ : transition dynamics
  - $r_t = \mathcal{R}(s_t, a_t)$ : reward
- Don't learn this (i.e. model free)

$$\tau = (s_0, a_0, r_0, \dots, s_T, a_T, r_T) \quad (\text{trajectory})$$

$$R = \sum_{t=0}^T r_t \quad (\text{total return})$$

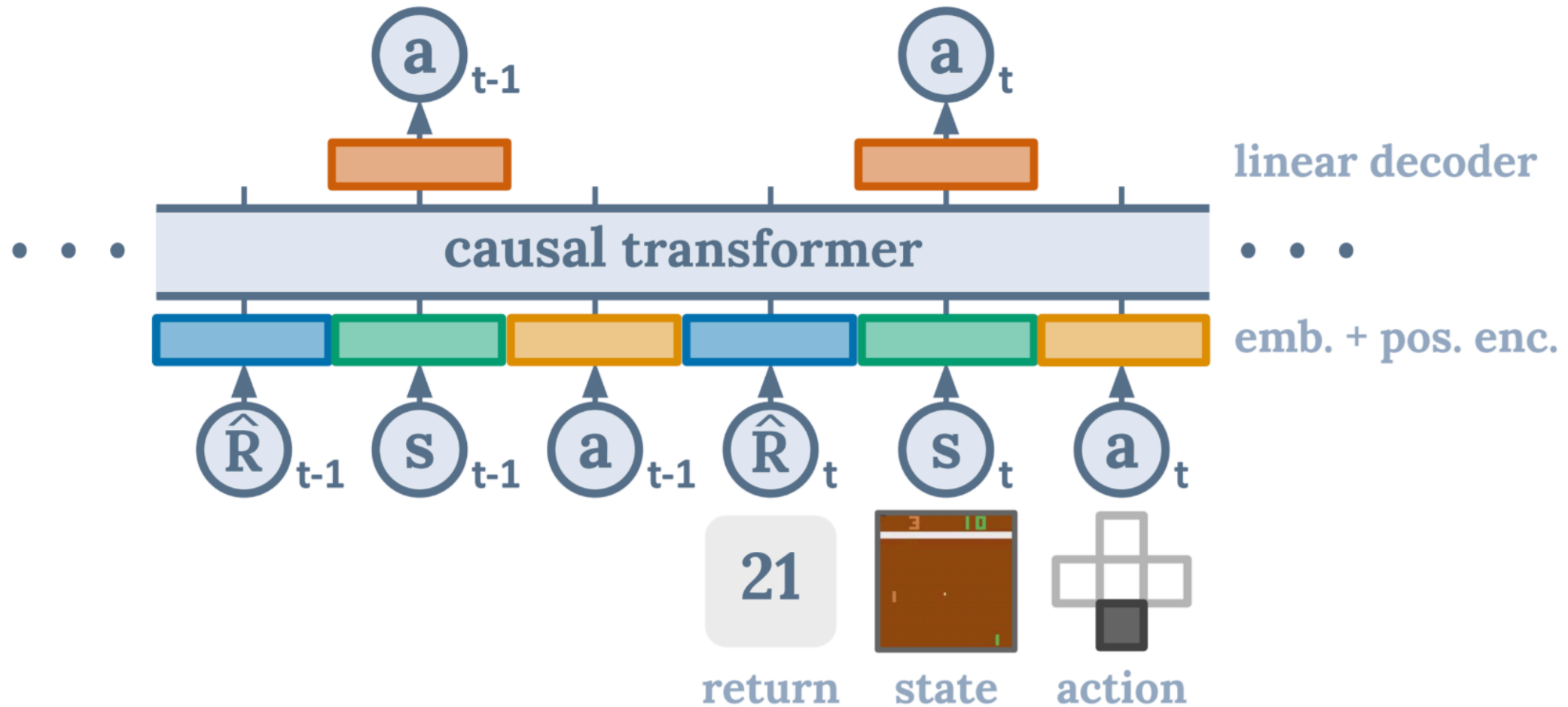
$$\hat{R}_t = \sum_{t'=t}^T r_{t'} \quad (\text{return-to-go})$$

Goal:

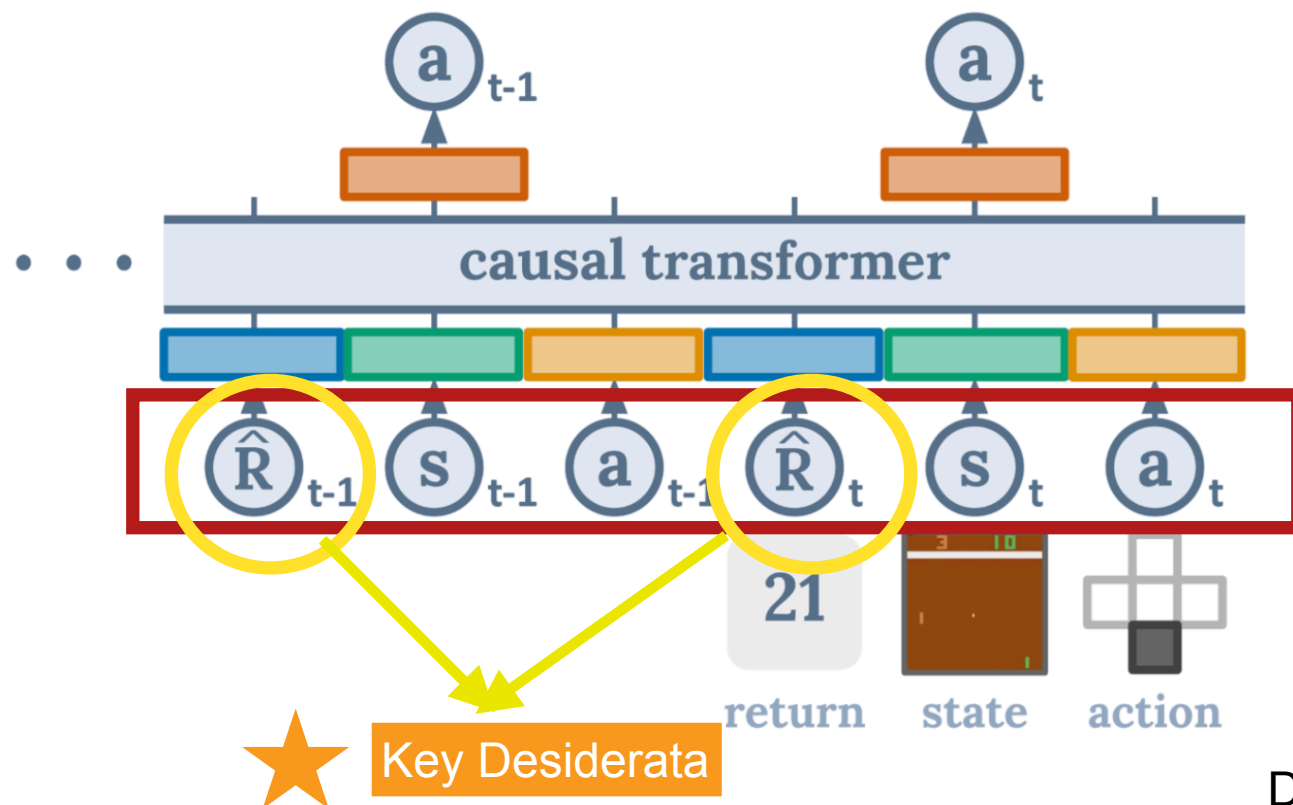
$$\max_{\pi} \mathbb{E} \left[ \sum_{t=0}^T r_t \right]$$

# What we do?

## Overall Architecture



# What we do?



Trajectory representation:

$$\hat{R}_t = \sum_{t'=t}^T r_{t'} \quad (\text{return-to-go})$$

$$\tau = (\hat{R}_1, s_1, a_1, \hat{R}_2, s_2, a_2, \dots, \hat{R}_T, s_T, a_T)$$

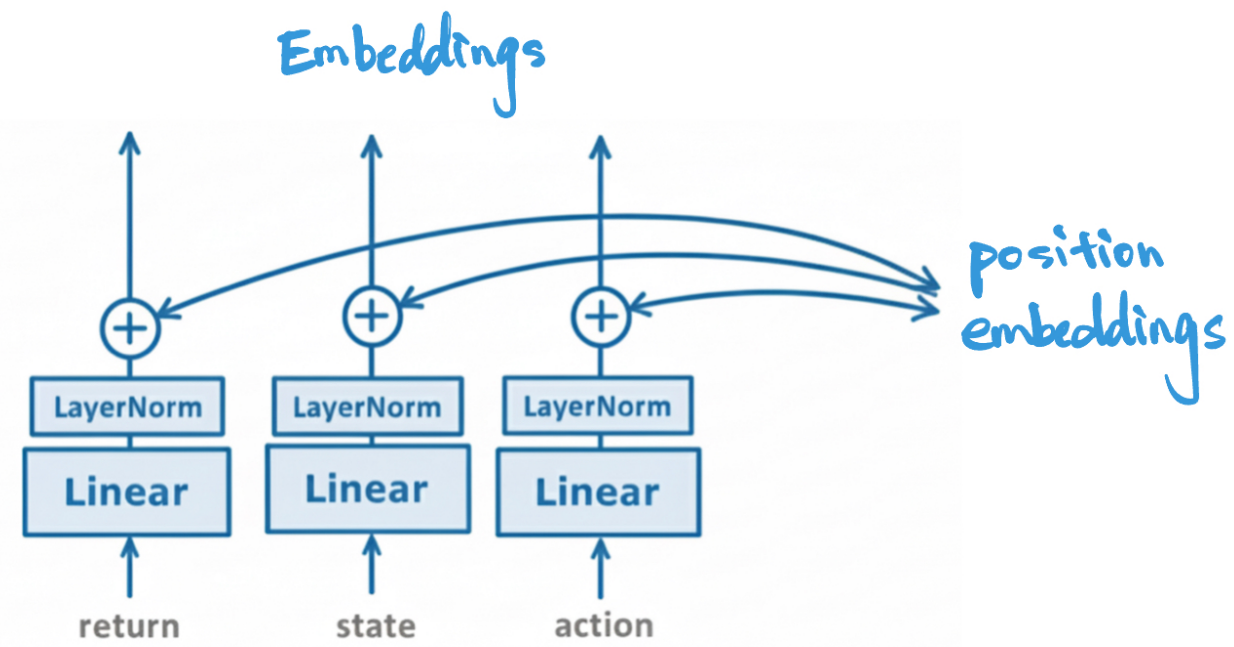
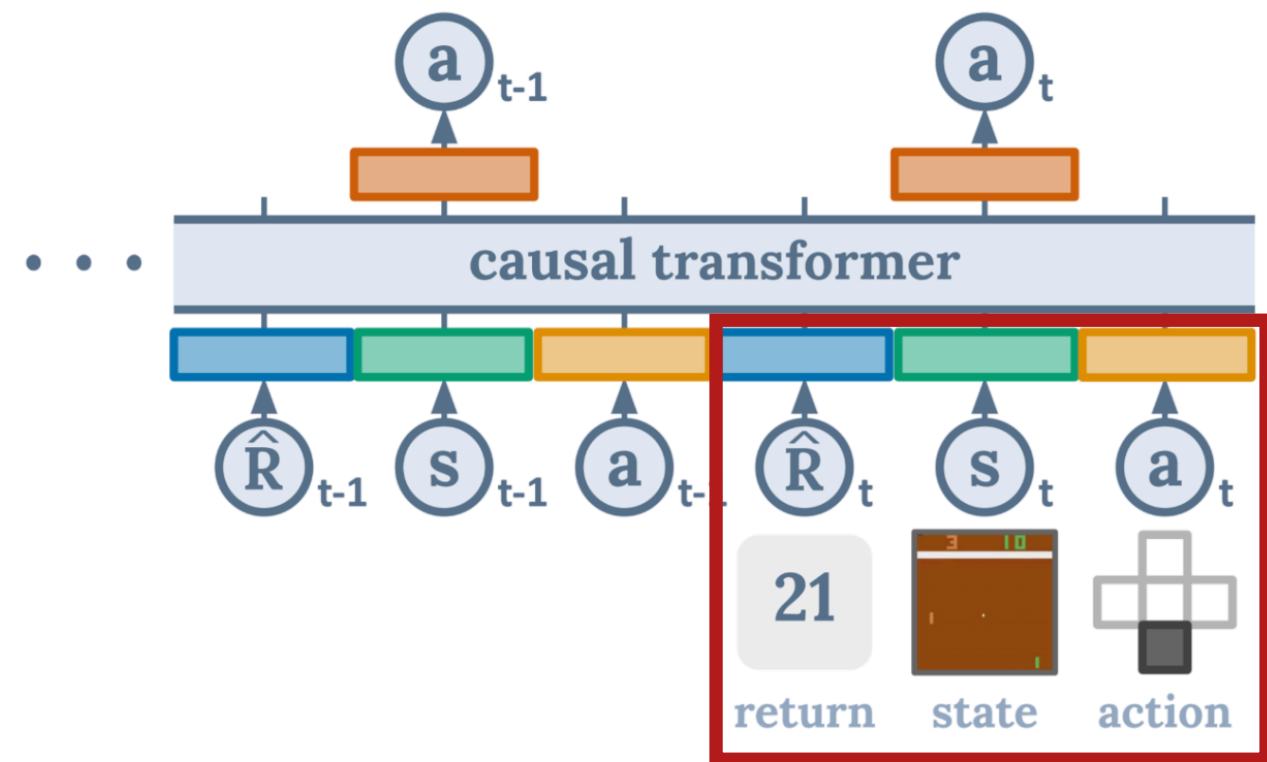
Desired performance

Initial state

At test time

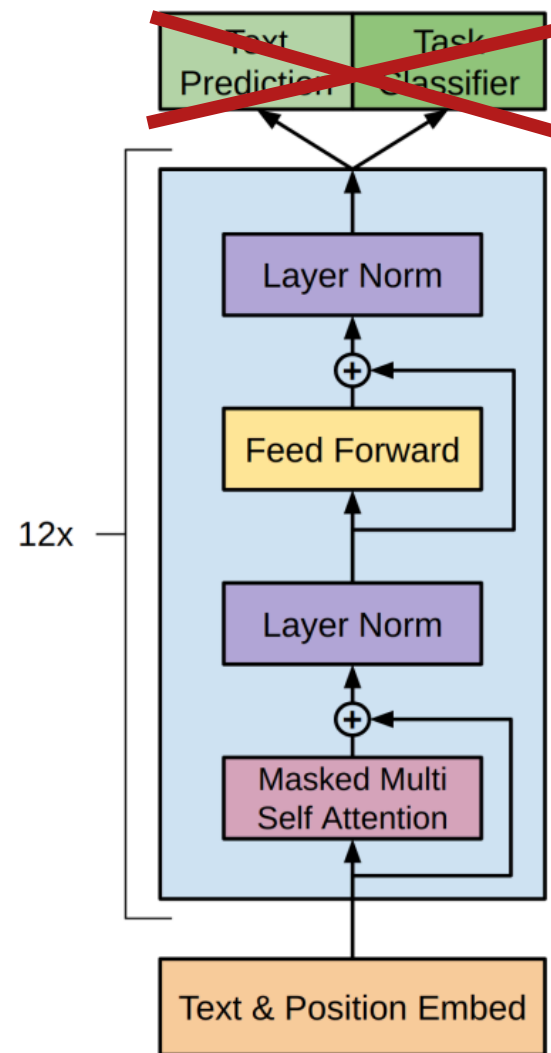
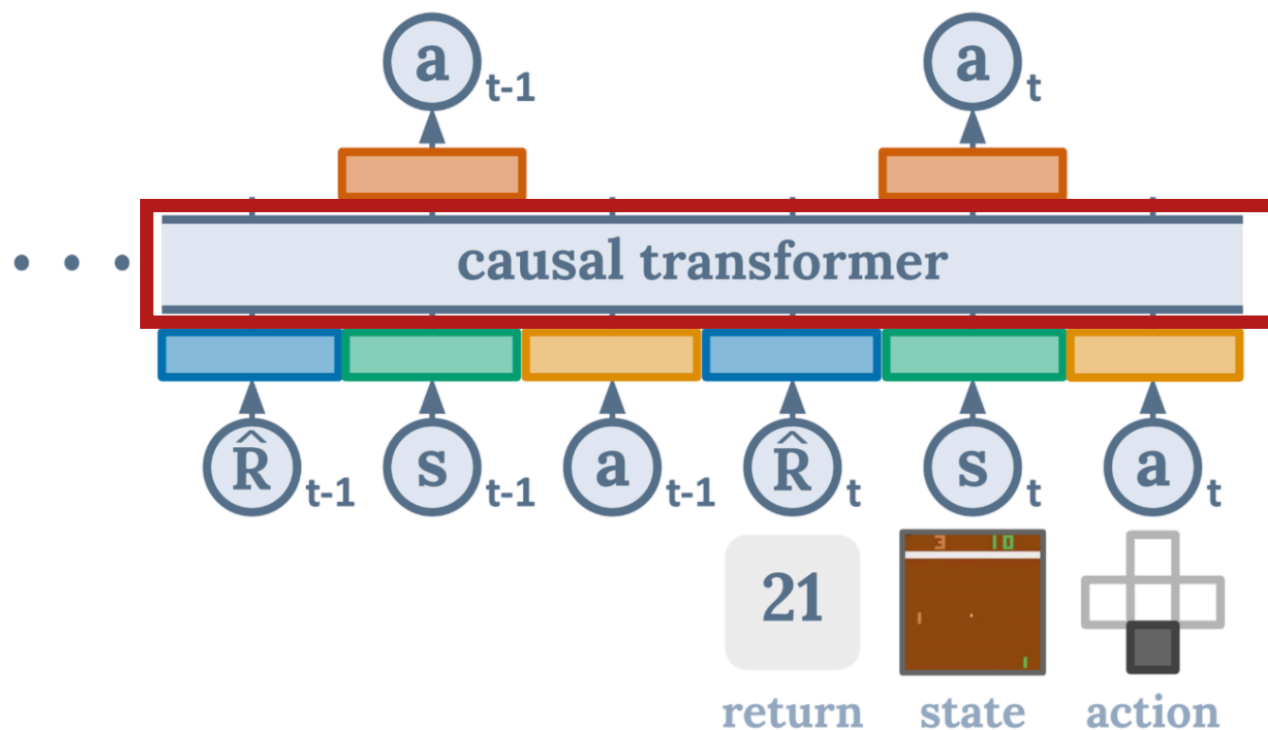
$$\hat{R}_{t+1} = \hat{R}_t - r_t$$

# What we do?

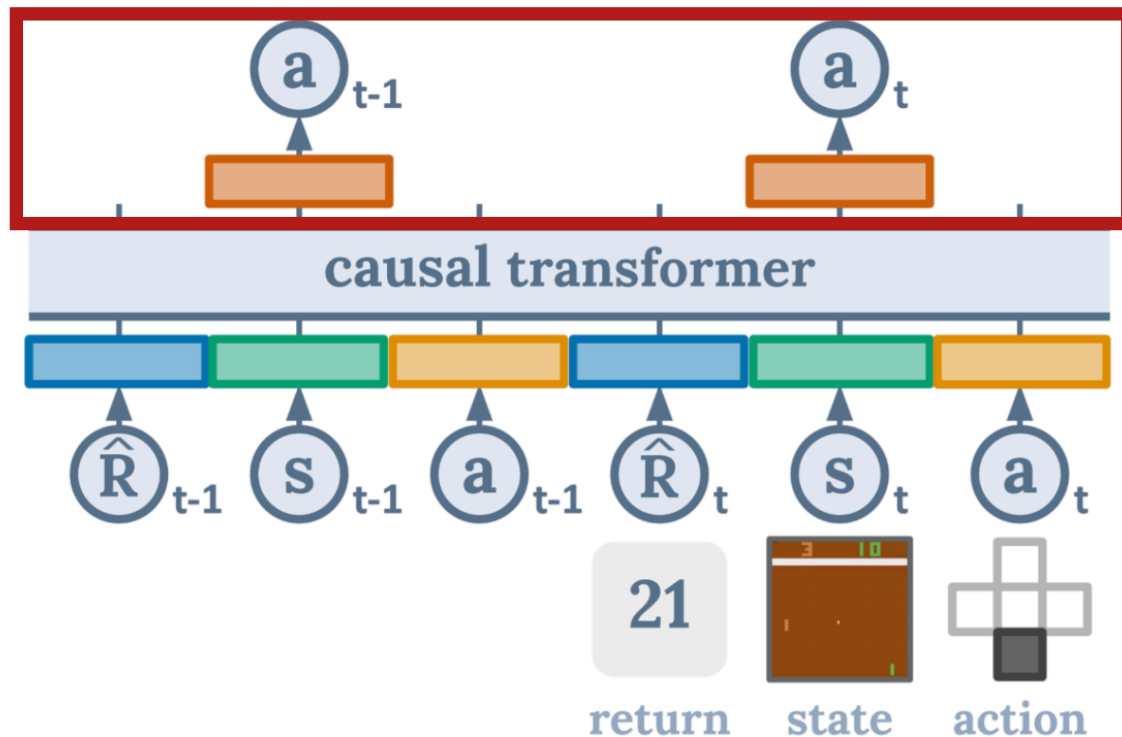


# What we do?

GPT



# What we do?



Only select hidden states corresponding to **action tokens**

$$h_t = \text{Transformer}(\hat{R}_{0:t}, s_{0:t}, a_{0:t-1})_t$$

Continuous actions: MSE loss

$$\mathcal{L}_{\text{cont}} = \frac{1}{K} \sum_{t=1}^K \|a_t - \hat{a}_t\|_2^2$$

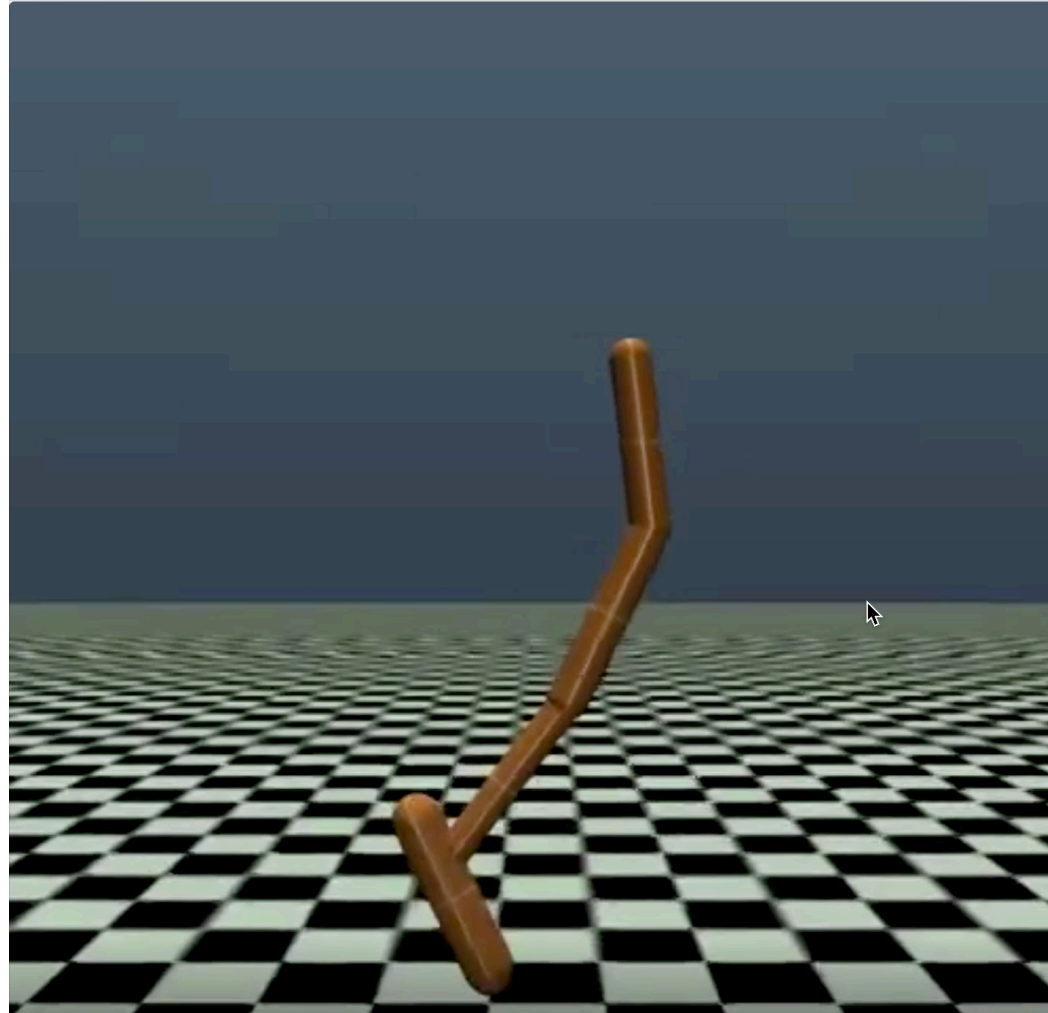
Discrete actions: cross-entropy loss

$$\mathcal{L}_{\text{disc}} = \frac{1}{K} \sum_{t=1}^K -\log \pi_{\theta}(a_t | h_t)$$



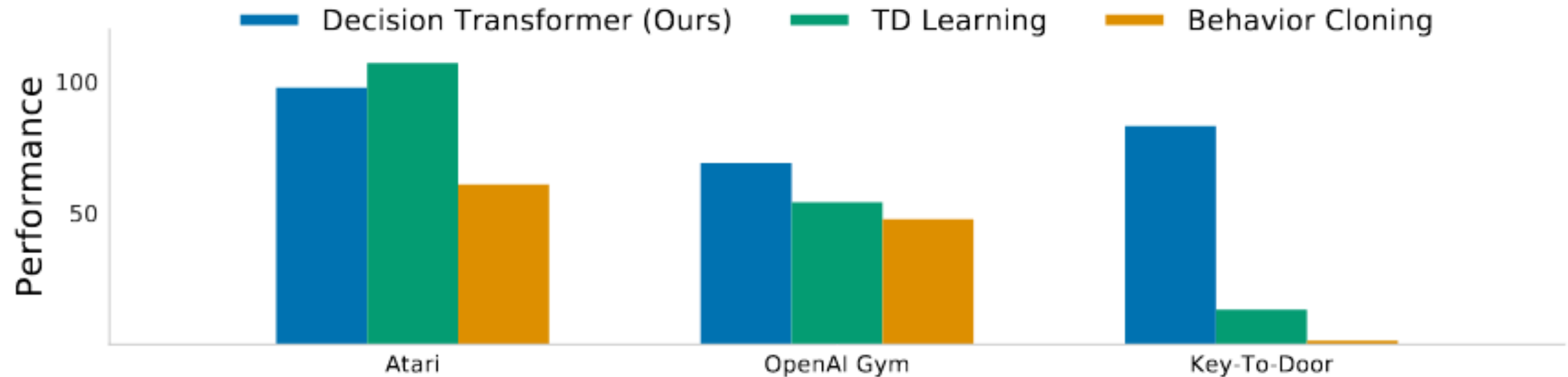
# Why we do it?

## Experiments



# Why we do it?

## Experiments



# Why we do it?

## Experiments

### TD Learning

Atari

Game	DT (Ours)	COL	OR-DON	REM	BC
Breakout	<b>267.5 ± 97.5</b>	211.1	17.1	8.9	138.9 ± 61.7
Qbert	15.4 ± 11.4	<b>104.2</b>	0.0	0.0	17.3 ± 14.7
Pong	106.1 ± 8.1	<b>111.9</b>	18.0	0.5	85.2 ± 20.0
Seaquest	<b>2.5 ± 0.4</b>	1.7	0.4	0.7	2.1 ± 0.3

OpenAI Gym

Dataset	Environment	DT (Ours)	CQL	BEAR	BRAC-v	AWR	BC
Medium-Expert	HalfCheetah	<b>86.8 ± 1.3</b>	62.4	53.4	41.9	52.7	59.9
Medium-Expert	Hopper	107.6 ± 1.8	<b>111.0</b>	96.3	0.8	27.1	79.6
Medium-Expert	Walker	<b>108.1 ± 0.2</b>	98.7	40.1	81.6	53.8	36.6
Medium-Expert	Reacher	<b>89.1 ± 1.3</b>	30.6	-	-	-	73.3
Medium	HalfCheetah	42.6 ± 0.1	44.4	41.7	<b>46.3</b>	37.4	43.1
Medium	Hopper	<b>67.6 ± 1.0</b>	58.0	52.1	31.1	35.9	63.9
Medium	Walker	74.0 ± 1.4	79.2	59.1	<b>81.1</b>	17.4	77.3
Medium	Reacher	<b>51.2 ± 3.4</b>	26.0	-	-	-	<b>48.9</b>
Medium-Replay	HalfCheetah	36.6 ± 0.8	46.2	38.6	<b>47.7</b>	40.3	4.3
Medium-Replay	Hopper	<b>82.7 ± 7.0</b>	48.6	33.7	0.6	28.4	27.6
Medium-Replay	Walker	<b>66.6 ± 3.0</b>	26.7	19.2	0.9	15.5	36.9
Medium-Replay	Reacher	<b>18.0 ± 2.4</b>	<b>19.0</b>	-	-	-	5.4
Average (Without Reacher)		<b>74.7</b>	63.9	48.2	36.9	34.3	46.4
Average (All Settings)		<b>69.2</b>	54.2	-	-	-	47.7



## Why we do it?

### Does Decision Transformer perform behavior cloning on a subset of the data?

Dataset	Environment	DT (Ours)	10%BC	25%BC	40%BC	100%BC	CQL
Medium	HalfCheetah	42.6 $\pm$ 0.1	42.9	43.0	43.1	43.1	<b>44.4</b>
Medium	Hopper	<b>67.6 <math>\pm</math> 1.0</b>	65.9	65.2	65.3	63.9	58.0
Medium	Walker	74.0 $\pm$ 1.4	78.8	<b>80.9</b>	78.8	77.3	79.2
Medium	Reacher	51.2 $\pm$ 3.4	51.0	48.9	58.2	<b>58.4</b>	26.0
Medium-Replay	HalfCheetah	36.6 $\pm$ 0.8	40.8	40.9	41.1	4.3	<b>46.2</b>
Medium-Replay	Hopper	<b>82.7 <math>\pm</math> 7.0</b>	70.6	58.6	31.0	27.6	48.6
Medium-Replay	Walker	66.6 $\pm$ 3.0	<b>70.4</b>	67.8	67.2	36.9	26.7
Medium-Replay	Reacher	18.0 $\pm$ 2.4	<b>33.1</b>	16.2	10.7	5.4	19.0
Average		56.1	<b>56.7</b>	52.7	49.4	39.5	43.5

## Why we do it?

### Does Decision Transformer perform behavior cloning on a subset of the data?

low data regime - Atari

use 1% of a replay buffer as the dataset, so then %BC is weak

Game	DT (Ours)	10%BC	25%BC	40%BC	100%BC
Breakout	<b>267.5 <math>\pm</math> 97.5</b>	28.5 $\pm$ 8.2	73.5 $\pm$ 6.4	108.2 $\pm$ 67.5	138.9 $\pm$ 61.7
Qbert	15.4 $\pm$ 11.4	6.6 $\pm$ 1.7	16.0 $\pm$ 13.8	11.8 $\pm$ 5.8	<b>17.3 <math>\pm</math> 14.7</b>
Pong	<b>106.1 <math>\pm</math> 8.1</b>	2.5 $\pm$ 0.2	13.3 $\pm$ 2.7	72.7 $\pm$ 13.3	85.2 $\pm$ 20.0
Seaquest	<b>2.5 <math>\pm</math> 0.4</b>	1.1 $\pm$ 0.2	1.1 $\pm$ 0.2	1.6 $\pm$ 0.4	2.1 $\pm$ 0.3

Takeaway: DT is not just cloning the best trajectories

## Why we do it?

## How well does Decision Transformer model the distribution of returns?

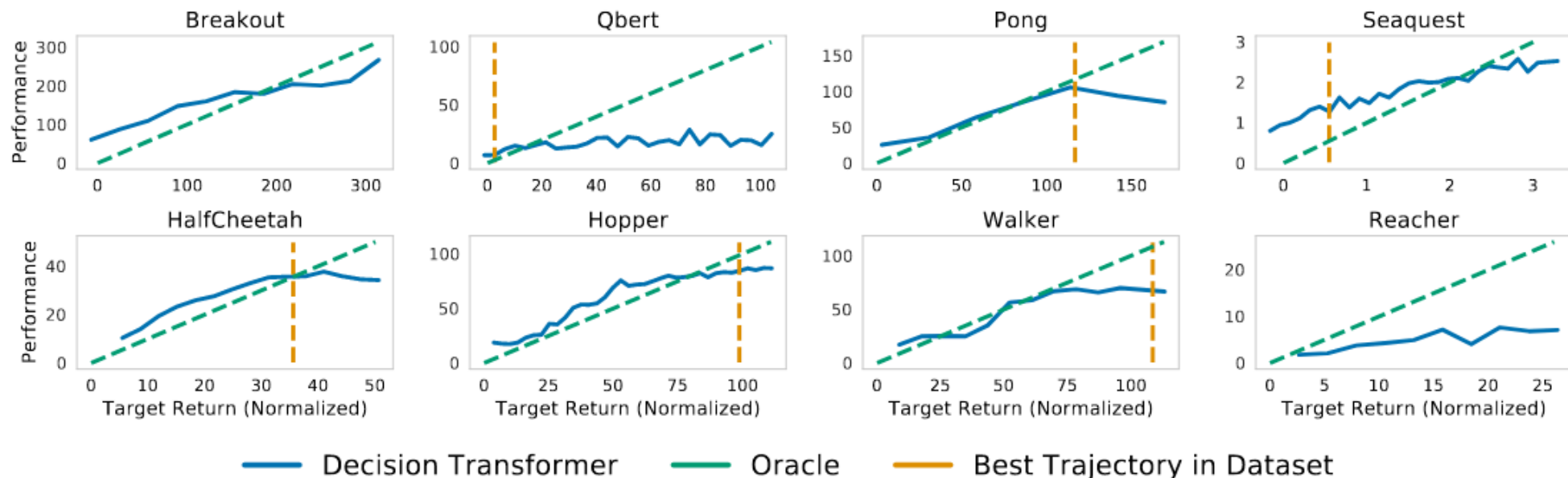


Figure 4: Sampled (evaluation) returns accumulated by Decision Transformer when conditioned on the specified target (desired) returns. **Top:** Atari. **Bottom:** D4RL medium-replay datasets.

**Takeaway:** DT successfully conditions on return-to-go and can adjust behavior across return levels, sometimes even extrapolating beyond the dataset

## Why we do it?

## What is the benefit of using a longer context length?

Game	DT (Ours)	DT with no context ( $K = 1$ )
Breakout	<b><math>267.5 \pm 97.5</math></b>	$73.9 \pm 10$
Qbert	<b><math>15.1 \pm 11.4</math></b>	$13.6 \pm 11.3$
Pong	<b><math>106.1 \pm 8.1</math></b>	$2.5 \pm 0.2$
Seaquest	<b><math>2.5 \pm 0.4</math></b>	$0.6 \pm 0.1$

Table 5: Ablation on context length. Decision Transformer (DT) performs better when using a longer context length ( $K = 50$  for Pong,  $K = 30$  for others).

Takeaway: DT really uses the past

# Why we do it?

## In sparse reward setting

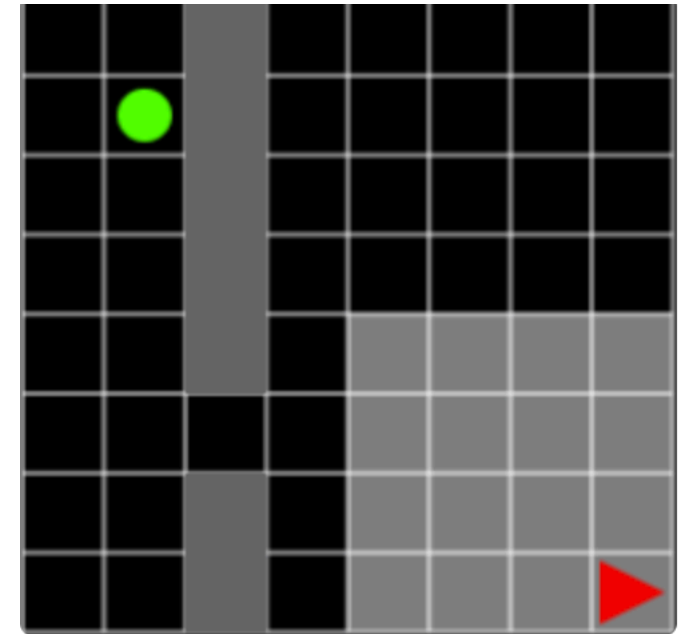
a variant to the key to door

3 phases:

1. the agent is placed in a room with a key
2. the agent is placed in an empty room
3. the agent is placed in a room with a door

3 problems in this setting:

1. Does DT perform effective long-term credit assignment? **YES**
2. Can DT be accurate critics in sparse reward settings? **YES**
3. Does DT perform well in sparse reward settings? **YES**





# Why we do it? In sparse reward setting

## Does Decision Transformer perform effective long-term credit assignment? YES!

Dataset	DT (Ours)	CQL	BC	%BC	Random
1K Random Trajectories	<b>71.8%</b>	13.1%	1.4%	69.9%	3.1%
10K Random Trajectories	94.6%	13.3%	1.6%	<b>95.1%</b>	3.1%

Table 6: Success rate for Key-to-Door environment. Methods using hindsight (Decision Transformer, %BC) can learn successful policies, while TD learning struggles to perform credit assignment.

# Why we do it? In sparse reward setting

## Can transformers be accurate critics? YES!

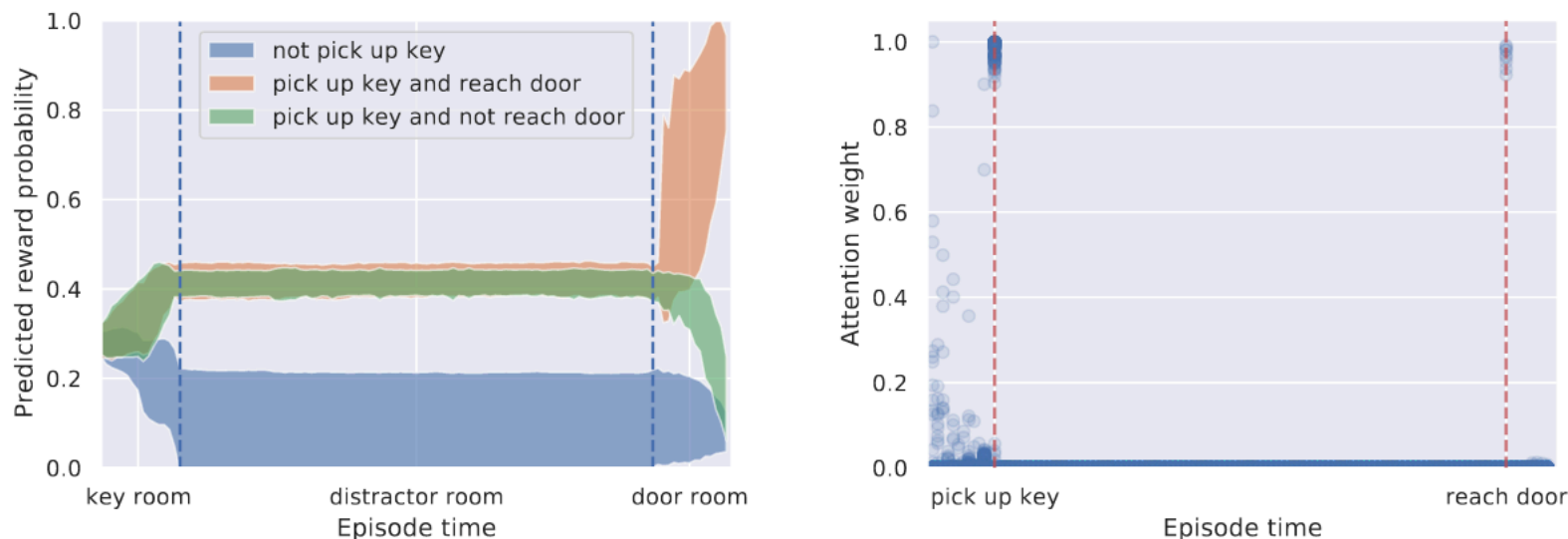


Figure 5: **Left:** Averages of running return probabilities predicted by the transformer model for three types of episode outcomes. **Right:** Transformer attention weights from all timesteps superimposed for a particular successful episode. The model attends to steps near pivotal events in the episode, such as picking up the key and reaching the door.

# Why we do it? In sparse reward setting

## Does Decision Transformer perform well? YES!

Dataset	Environment	Delayed (Sparse)		Agnostic		Original (Dense)	
		DT (Ours)	CQL	BC	%BC	DT (Ours)	CQL
Medium-Expert	Hopper	<b><math>107.3 \pm 3.5</math></b>	9.0	59.9	102.6	107.6	111.0
Medium	Hopper	$60.7 \pm 4.5$	5.2	63.9	<b>65.9</b>	67.6	58.0
Medium-Replay	Hopper	<b><math>78.5 \pm 3.7</math></b>	2.0	27.6	70.6	82.7	48.6

Table 7: Results for D4RL datasets with delayed (sparse) reward. Decision Transformer (DT) and imitation learning are minimally affected by the removal of dense rewards, while CQL fails.

## Why we do it?

### Why does Decision Transformer avoid the need for value pessimism or behavior regularization?

Traditional offline RL (e.g., TD-learning) learns an **approximate value function** and **optimizes a policy** on top of it

But optimizing an imperfect value function can **amplify errors** → **unstable policy improvement**

So prior methods require **policy regularization / conservatism** (e.g., CQL) to stay close to the dataset and avoid value overestimation

**Decision Transformer does NOT optimize a value function**

- It treats RL as **sequence modeling**, using supervised learning
- Therefore, it **does not need explicit regularization** to remain stable

# Why we do it?

## How can Decision Transformer benefit online RL regimes?

### **Strong offline behavior prior**

Learns diverse high-quality trajectories from offline data

→ good starting policy for online learning

### **Likelihood-based training transfers well**

Sequence modeling objective has shown smooth offline→online adaptation

### **Acts as a “memory engine”**

Remembers good behaviors and can regenerate them during online training

### **Supports broad behavior exploration**

When combined with exploration methods (e.g., Go-Explore),  
enables generating diverse and high-return trajectories

# POST DT work

- Online Decision Transformer <https://arxiv.org/abs/2202.05607>
- Value-Guided/Q-Guided DT <https://arxiv.org/abs/2209.03993>
- Diffussion policies/Multimodal Agents <https://arxiv.org/abs/2409.00588>

# Question?



# References

- <https://arxiv.org/abs/2106.01345>
- <https://sites.google.com/berkeley.edu/decision-transformer>
- <https://openreview.net/forum?id=a7APmM4B9d>
- <https://bair.berkeley.edu/blog/2020/12/07/offline/>
- <https://greentec.github.io/reinforcement-learning-fifth-en/>
- <https://arxiv.org/pdf/2006.04779>
- [https://cdn.openai.com/research-covers/language-unsupervised/language\\_understanding\\_paper.pdf](https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf)
- <https://github.com/souradipp76/SeqModRL?tab=readme-ov-file>

