

# DECISION TRANSFORMER: REINFORCEMENT LEARNING VIA SEQUENCE MODELING

Youssef Fathi,  
CS 885: Reinforcement Learning  
Winter 2022

Presented to: Prof. Pascal Poupart



# Outline

- Introduction
- Background
- Methodology
- Evaluation
- Discussion
- Conclusion

# INTRODUCTION



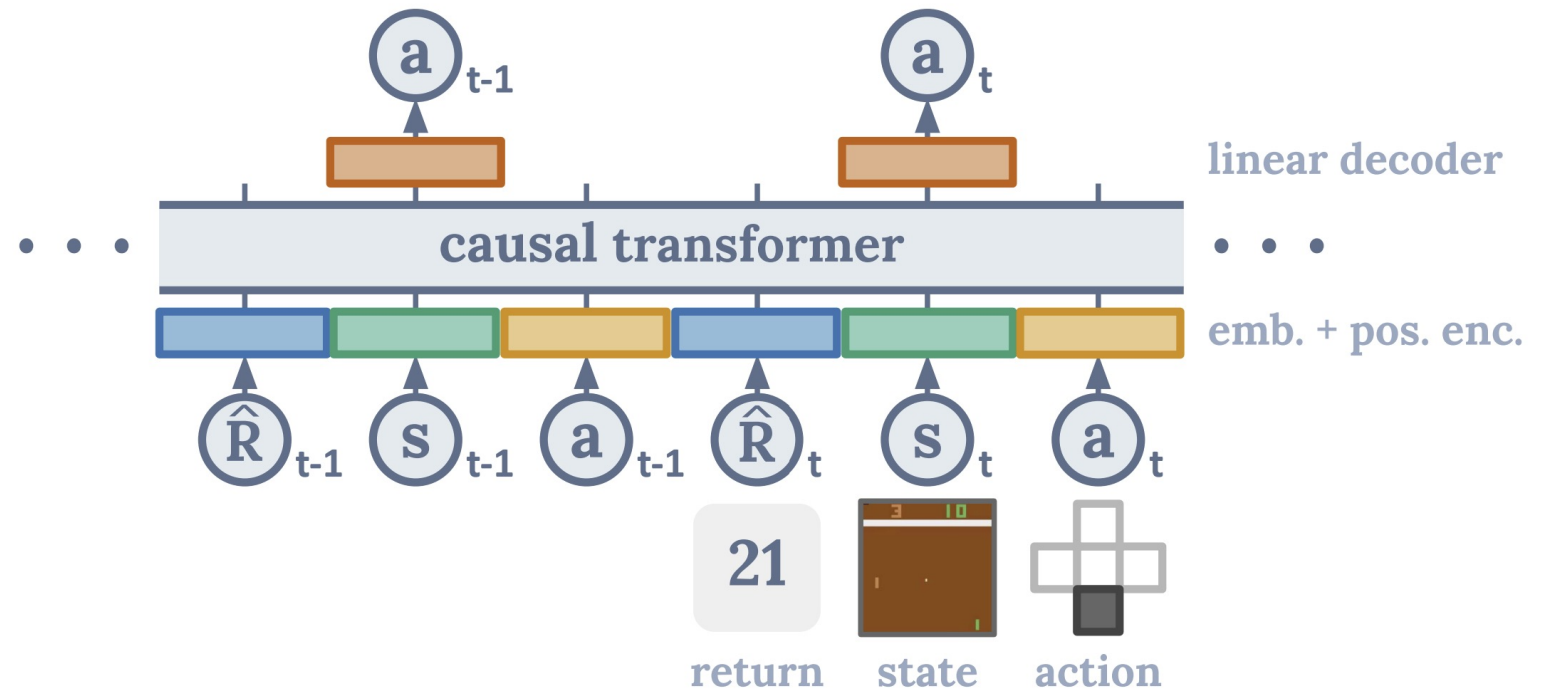
# Introduction

**Supervised RL**

**Model-Free**

**Offline RL**

**Sequence Modelling**



# BACKGROUND



# Background: Markov Decision Process (MDP)

## Model-based

**States:**  $S$ .

**Actions:**  $A$ .

**Transition Model:**  $P(s_t | s_{t-1}, a_{t-1})$

**Reward Model:**  $R(s_t, a_t)$

**Discount Factor:**  $0 \leq \gamma \leq 1$

**Horizon:**  $h$

## Model-free

**States:**  $S$ .

**Actions:**  $A$ .

~~**Transition Model:**  $P(s_t | s_{t-1}, a_{t-1})$~~

~~**Reward Model:**  $R(s_t, a_t)$~~

**Discount Factor:**  $0 \leq \gamma \leq 1$

**Horizon:**  $h$

# Background: Online Reinforcement Learning (RL)

**Target:** Maximize expected sum of discounted rewards.

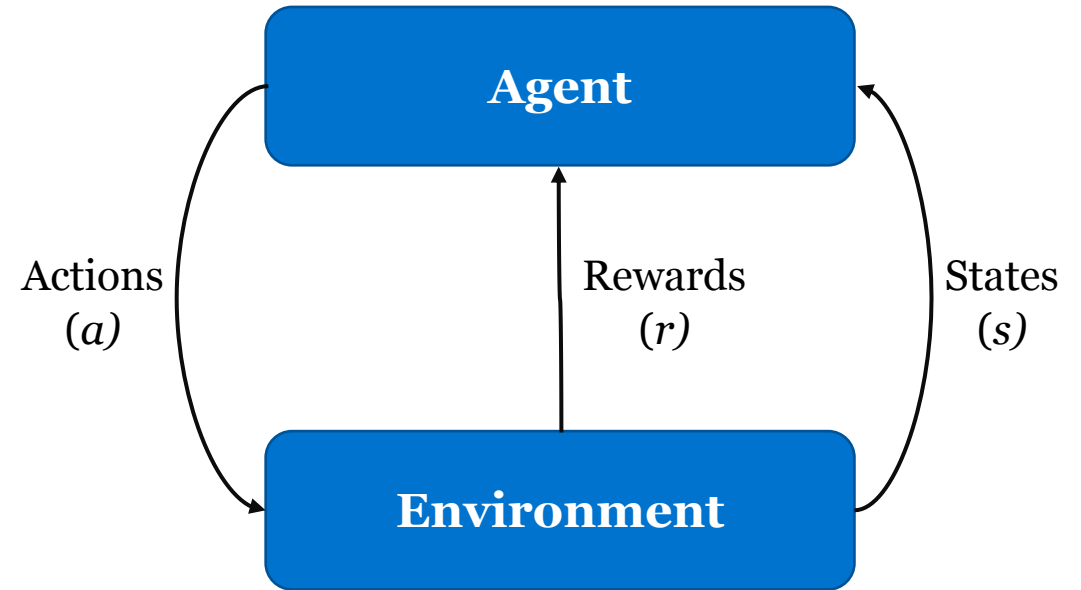
$$E[R] = \sum_i \gamma^i r_i$$

**Method:** Using temporal difference in bellman backups to estimate optimal value function [1,2].

$$Q_i(s, a) = Q_{i-1}(s, a) + \alpha(r_i + \gamma Q_{i-1}(s', \pi(s')) - Q_{i-1}(s, a))$$

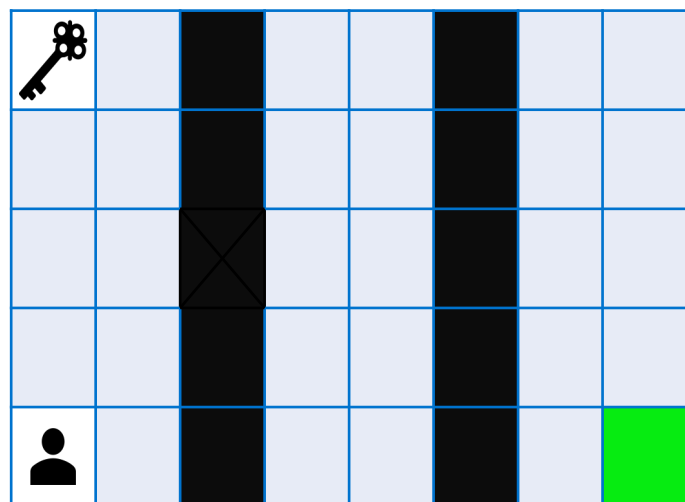
## Challenges:

- High cost/risk of interaction with environment
- Credit assignment problem due to discounting in bellman backups.

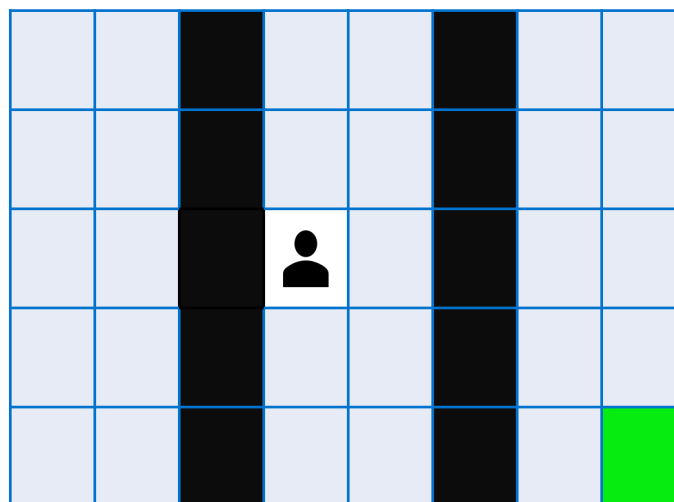


$$\tau = (s_0, a_0, r_0) \rightarrow (s_1, a_1, r_1) \rightarrow (s_2, a_2, r_2) \rightarrow \dots$$

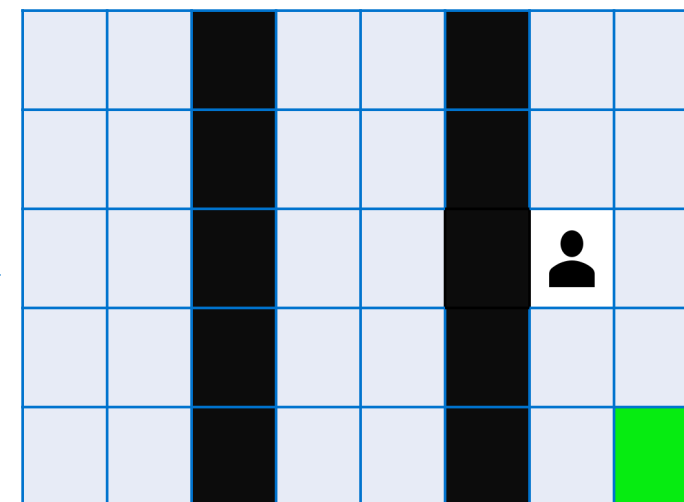
# Background: Credit Assignment Problem (Key-to-Door Env. [27])



*Phase 1*



*Phase 2*



*Phase 3*

**Problem:** Assigning delayed rewards to their originating actions.

**Possible Solution:** State association [16,17,18,19]



# Background: Offline (Batch) Reinforcement Learning

- **Objective:** Learning from a fixed dataset without further interactions with the environment.

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_n, a_n, r_n)$$

*Pre-generated*

- **Popular Examples:** DDPG [4].
- **Main Challenge:** Distribution Shift (Extrapolation Error) [6].

$$Q_i(s, a) = Q_{i-1}(s, a) + \alpha(r_i + \gamma Q_{i-1}(s', \pi(s')) - Q_{i-1}(s, a))$$
$$\pi(s') = \operatorname{argmax}_a Q(s', a)$$



$\pi(s')$  chooses rarely visited  $(s', a')$

- **Solutions Proposed:**
  - Constrain policy action space [6,7]
  - Incorporate value pessimism [6,8]

# Background: Supervised RL

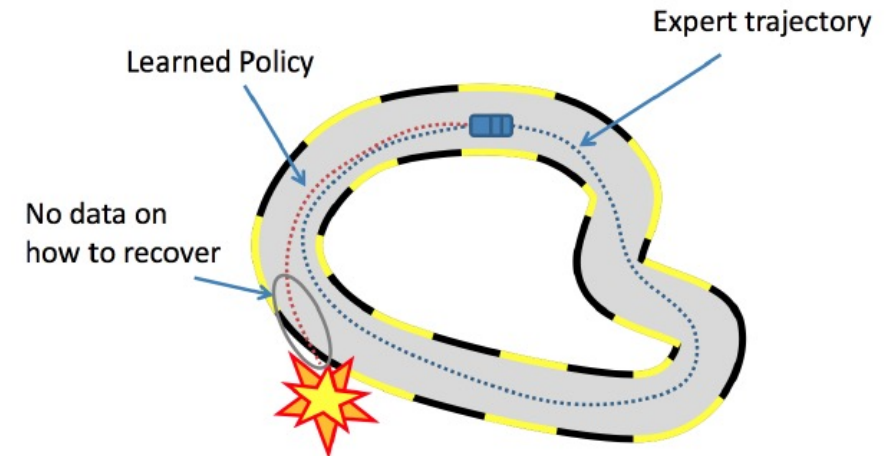
**Imitation Learning:** imitating the behaviour observed in existing trajectories.

- Behavioural Cloning (Basic Version) : using supervised losses to map existing states to actions with no regards to rewards. [11]

$$f(s) = a, \quad (s, a) \in \{(s_1, a_1), (s_2, a_1), \dots, (s_n, a_n)\}$$

## Drawbacks

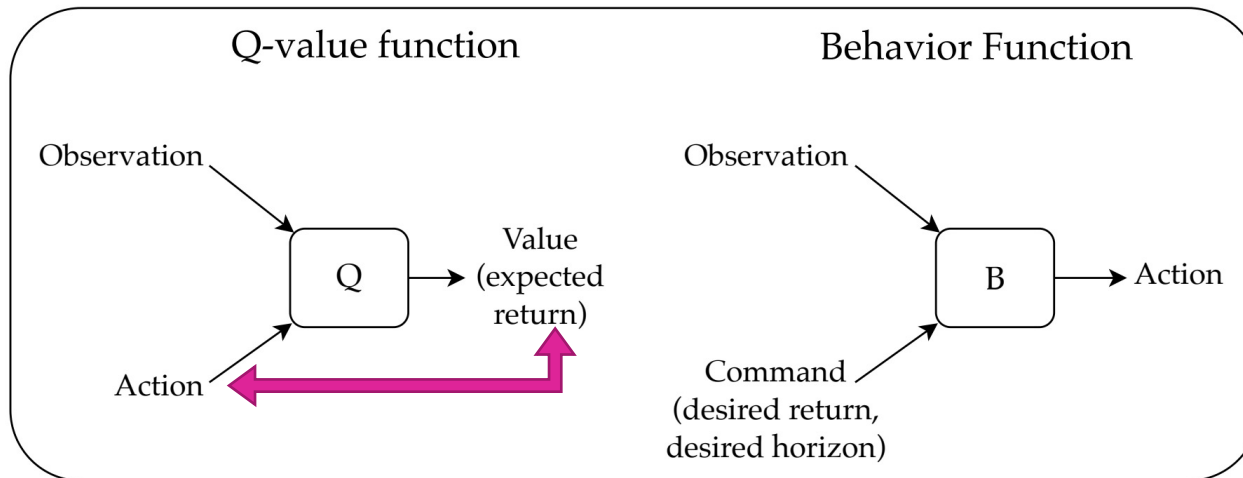
- Impossible to generalize to new scenarios.
- Requires large amount of optimal (expert) actions in the trajectories
- Assumes state-action pairs are i.i.d.



Source: <http://web.stanford.edu/class/cs234/slides/lecture7.pdf>

# Background: Supervised RL

**Upside-Down RL [12]:** trains agents to follow commands such as “obtain so much total reward in so much time.”



## Variants

Kumar et.al  
[13]

- Fully Offline RL
- Reward Conditioning

Ghosh et al.  
[14]

- State Conditioning

Paster et al.  
[15]

- Online RL
- LSTM with State Conditioning

## Supervised Objective

$$B = \operatorname{argmin}_B \sum_{t_1, t_2 \in \tau} L(B(a_{t_1}, s_{t_1}, d^r, d^h), a_{t_2})$$

# Background: Attention

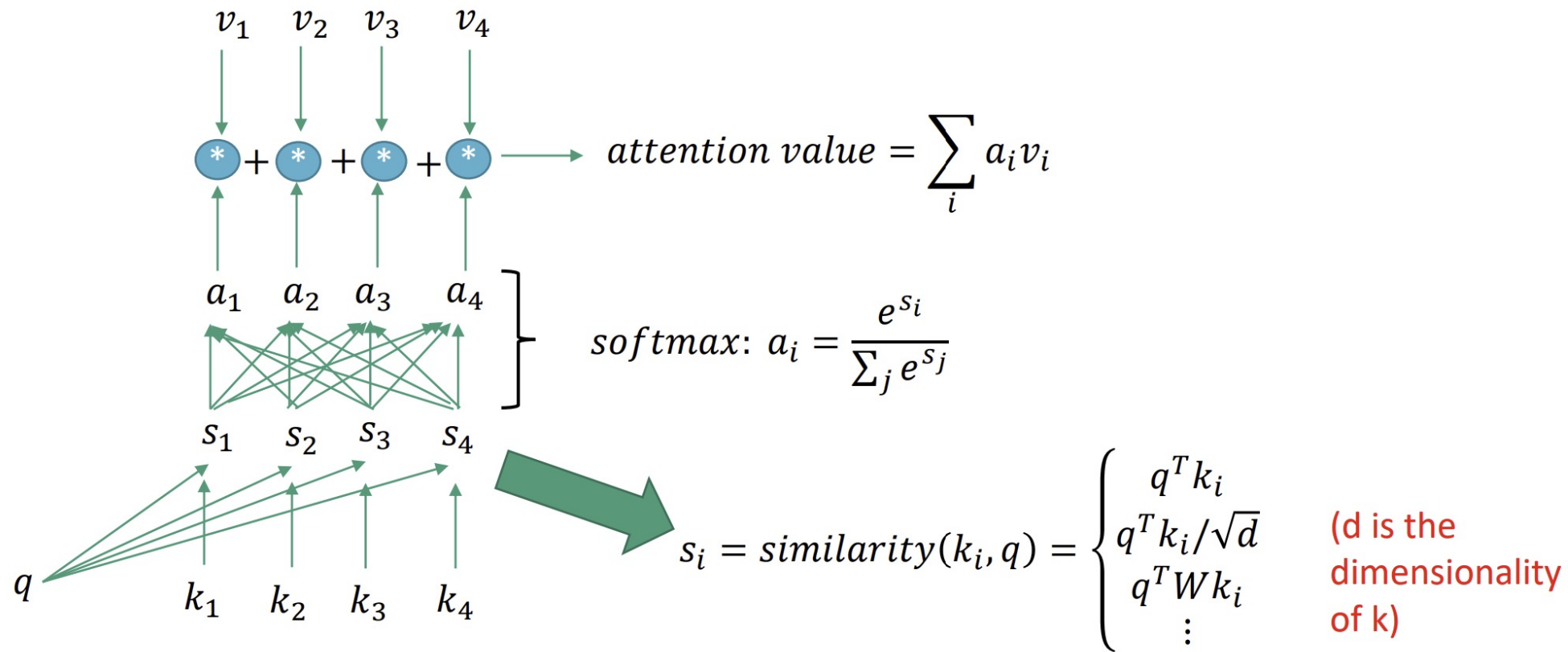
“The cat drank the milk because **it** was hungry.”

“The cat drank the milk because **it** was sweet.”



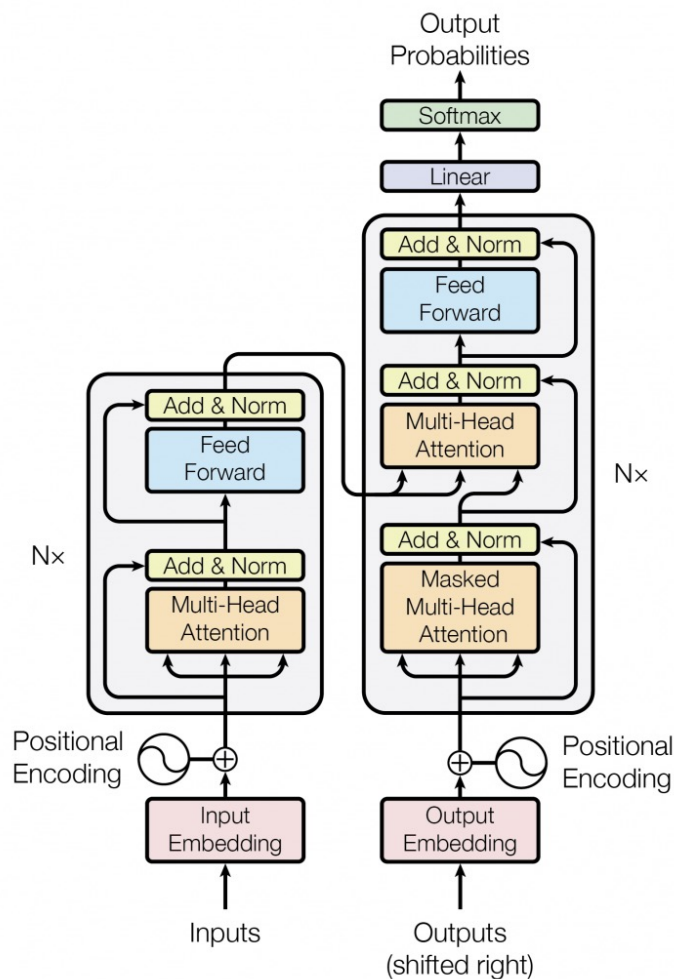
Credit: <https://towardsdatascience.com/transformers-explained-visually-part-1-overview-of-functionality-95a6dd460452>

# Background: Attention

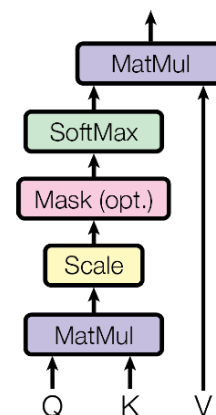


Credit: STAT940, Prof. Ali Ghodsi, University of Waterloo

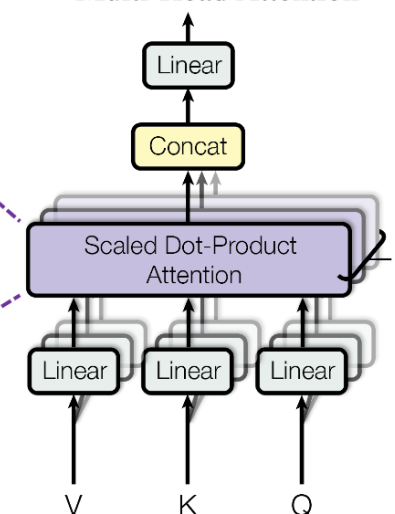
# Background: Sequence Modeling with Transformers



Scaled Dot-Product Attention



Multi-Head Attention



Credit: <https://data-science-blog.com/blog/2021/04/07/multi-head-attention-mechanism/>

# Background: Transformers

## Original [20]

$$z_i = \sum_{j=1}^n \text{softmax} \left( \{ \langle q_i, k_{j'} \rangle \}_{j'=1}^n \right) \cdot v_j$$

“The cat drank the milk because it was hungry.”

## GPT [21]

$$z_i = \sum_{j=1}^i \text{softmax} \left( \{ \langle q_i, k_{j'} \rangle \}_{j'=1}^i \right) \cdot v_j$$

“The cat drank the milk because it was hungry.”

# METHODOLOGY





# Methodology: Decision Transformer [28] Overview

Upside-down RL [13]

Model-Free RL

Offline RL

Sequence Modelling  
using GPT

Implicit Credit  
Assignment



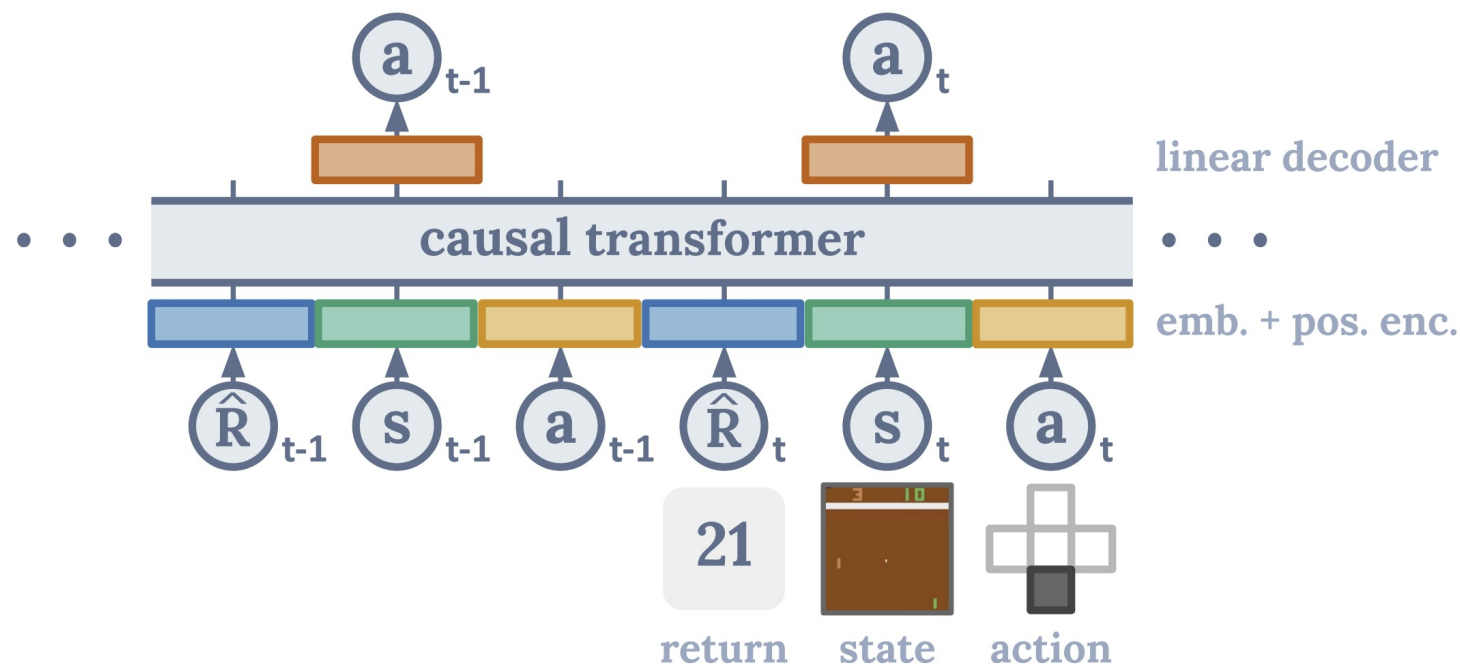
No Distribution Shift



No expert  
demonstrations



Match or Exceed S.O.A.



# Methodology: Input Setup

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

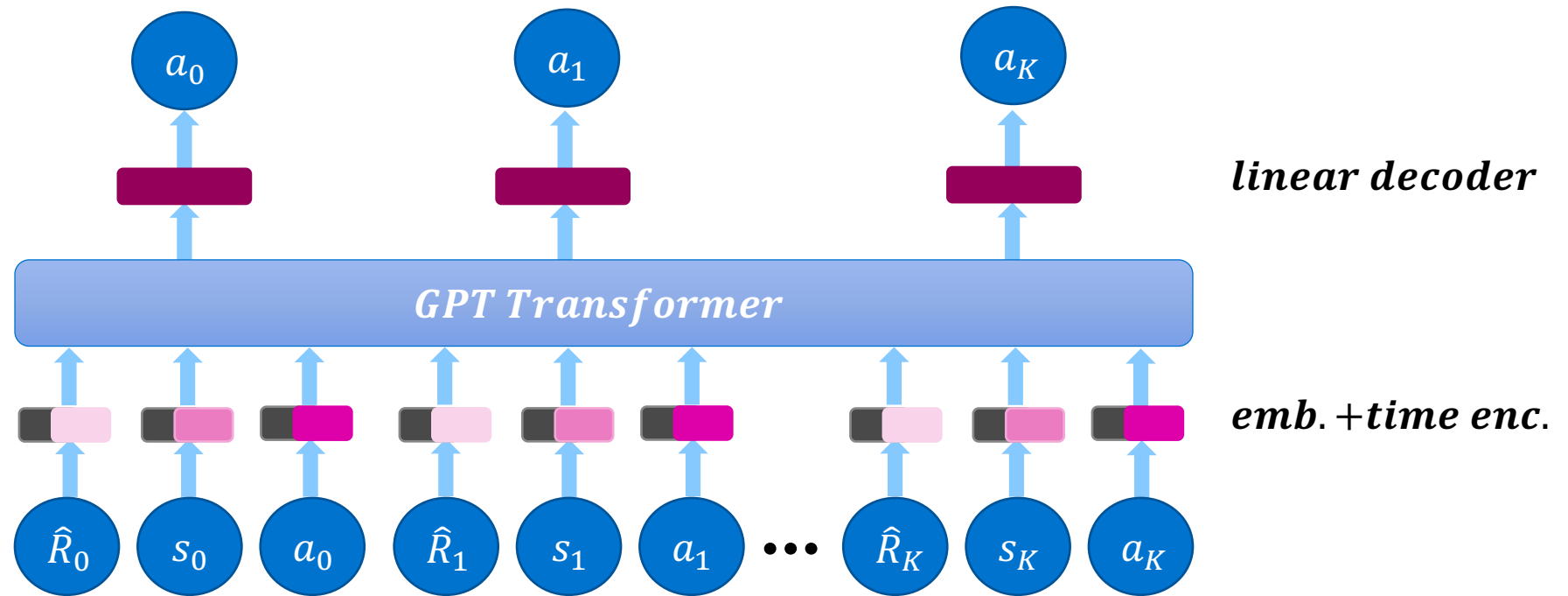


$$\tau = (\hat{R}_0, s_0, a_0, \hat{R}_1, s_1, a_1, \dots, \hat{R}_T, s_T, a_T)$$

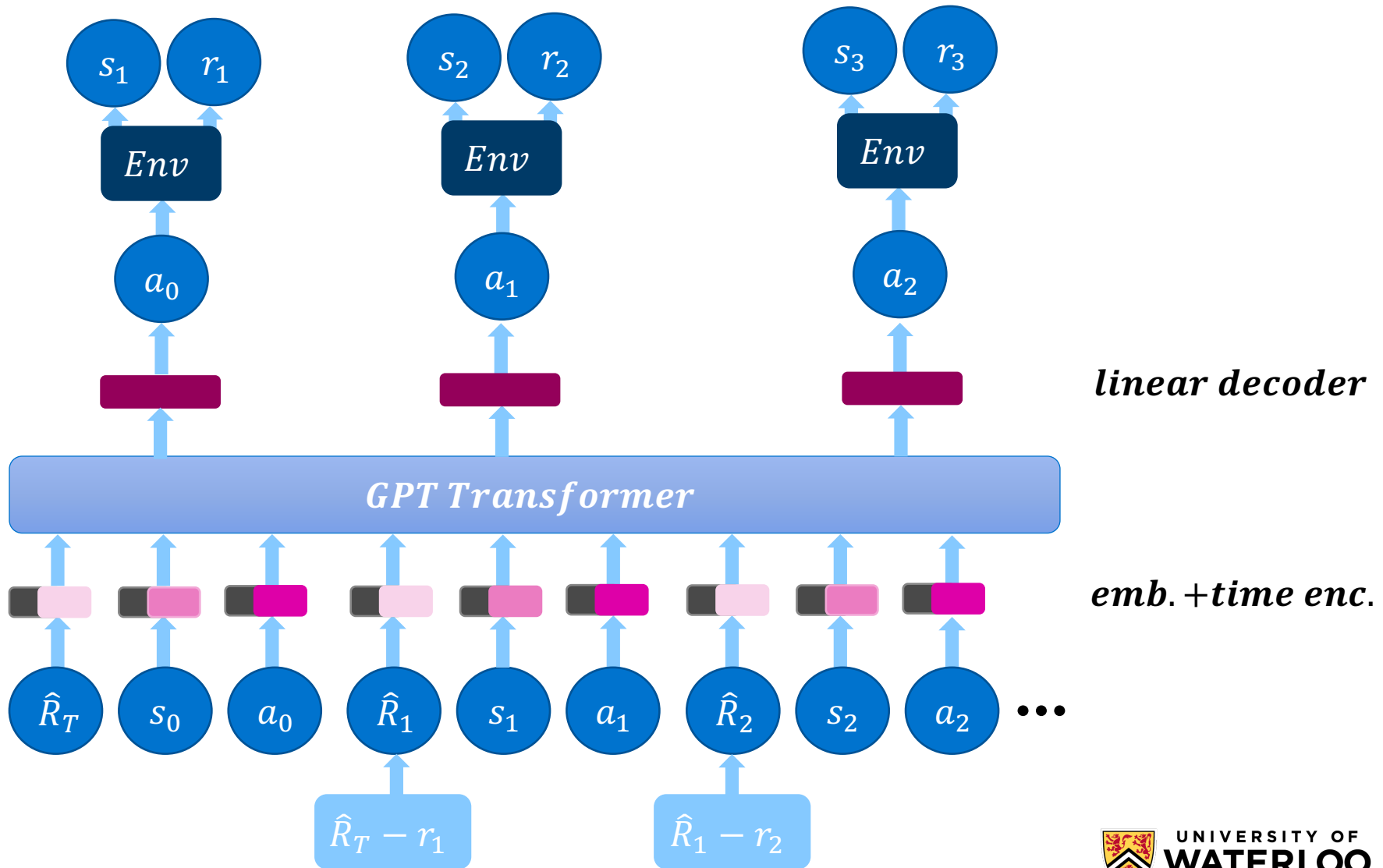
$$\hat{R}_t = \sum_{t'=t}^T r_{t'}$$

***Rewards-to-go***

# Methodology: Training Pipeline



# Methodology: Inference Pipeline



# Methodology: Psuedo-Code

---

## Algorithm 1 Decision Transformer Pseudocode (for continuous actions)

---

```
# R, s, a, t: returns-to-go, states, actions, or timesteps
# transformer: transformer with causal masking (GPT)
# embed_s, embed_a, embed_R: linear embedding layers
# embed_t: learned episode positional embedding
# pred_a: linear action prediction layer

# main model
def DecisionTransformer(R, s, a, t):
    # compute embeddings for tokens
    pos_embedding = embed_t(t) # per-timestep (note: not per-token)
    s_embedding = embed_s(s) + pos_embedding
    a_embedding = embed_a(a) + pos_embedding
    R_embedding = embed_R(R) + pos_embedding

    # interleave tokens as (R_1, s_1, a_1, ..., R_K, s_K)
    input_embeds = stack(R_embedding, s_embedding, a_embedding)

    # use transformer to get hidden states
    hidden_states = transformer(input_embeds=input_embeds)

    # select hidden states for action prediction tokens
    a_hidden = unstack(hidden_states).actions

    # predict action
    return pred_a(a_hidden)

# training loop
for (R, s, a, t) in dataloader: # dims: (batch_size, K, dim)
    a_preds = DecisionTransformer(R, s, a, t)
    loss = mean((a_preds - a)**2) # L2 loss for continuous actions
    optimizer.zero_grad(); loss.backward(); optimizer.step()

# evaluation loop
target_return = 1 # for instance, expert-level return
R, s, a, t, done = [target_return], [env.reset()], [], [1], False
while not done: # autoregressive generation/sampling
    # sample next action
    action = DecisionTransformer(R, s, a, t)[-1] # for cts actions
    new_s, r, done, _ = env.step(action)

    # append new tokens to sequence
    R = R + [R[-1] - r] # decrement returns-to-go with reward
    s, a, t = s + [new_s], a + [action], t + [len(R)]
    R, s, a, t = R[-K:], ... # only keep context length of K
```

---

# EXPERIMENTS & RESULTS



# Experiments: Atari Benchmark

## Baselines

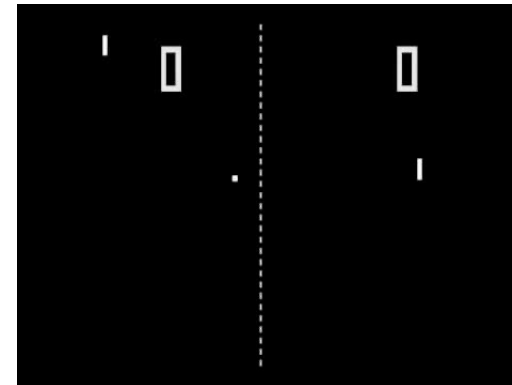
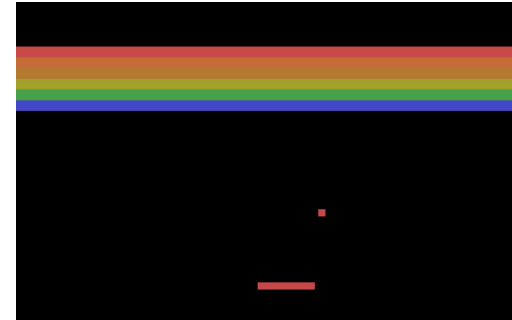
- CQL [22]
- REM [23]
- QE-DQN [24]
- BC (New)

## Games

- Breakout
- Qbert
- Pong (K=50)
- Seaquest

## Challenges

- Visual Inputs
- Long-term credit assignment



# Experiments: D4RL [3] Benchmark

## Baselines

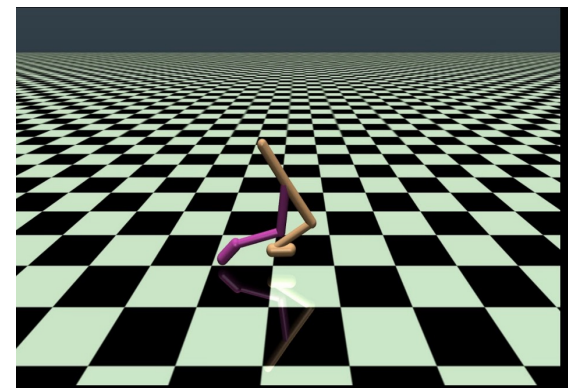
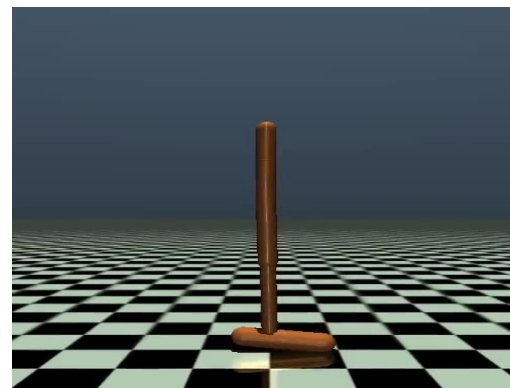
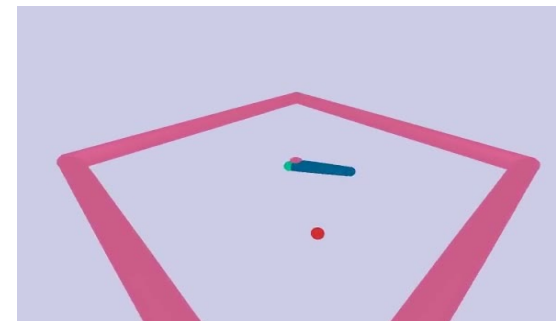
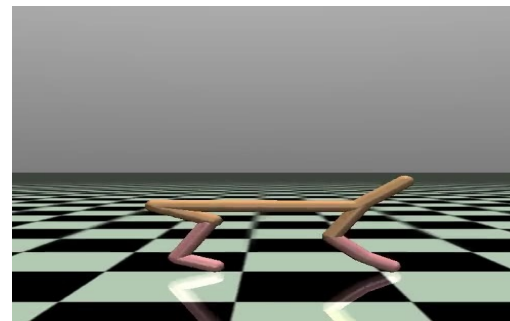
- CQL [22]
- BEAR [25]
- BRAC [26]
- AWR [5]
- BC (New)

## Games

- HalfCheetah
- Hopper
- Walker
- Reacher (New)

## Dataset Settings

- Medium
- Medium-Replay
- Medium-Expert





# Results: Atari Benchmark

Game	DT (Ours)	CQL	QR-DQN	REM	BC
Breakout	<b>267.5 <math>\pm</math> 97.5</b>	211.1	17.1	8.9	138.9 $\pm$ 61.7
Qbert	15.4 $\pm$ 11.4	<b>104.2</b>	0.0	0.0	17.3 $\pm$ 14.7
Pong	106.1 $\pm$ 8.1	<b>111.9</b>	18.0	0.5	85.2 $\pm$ 20.0
Seaquest	<b>2.5 <math>\pm</math> 0.4</b>	1.7	0.4	0.7	2.1 $\pm$ 0.3

# Results: D4RL [3] Benchmark

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

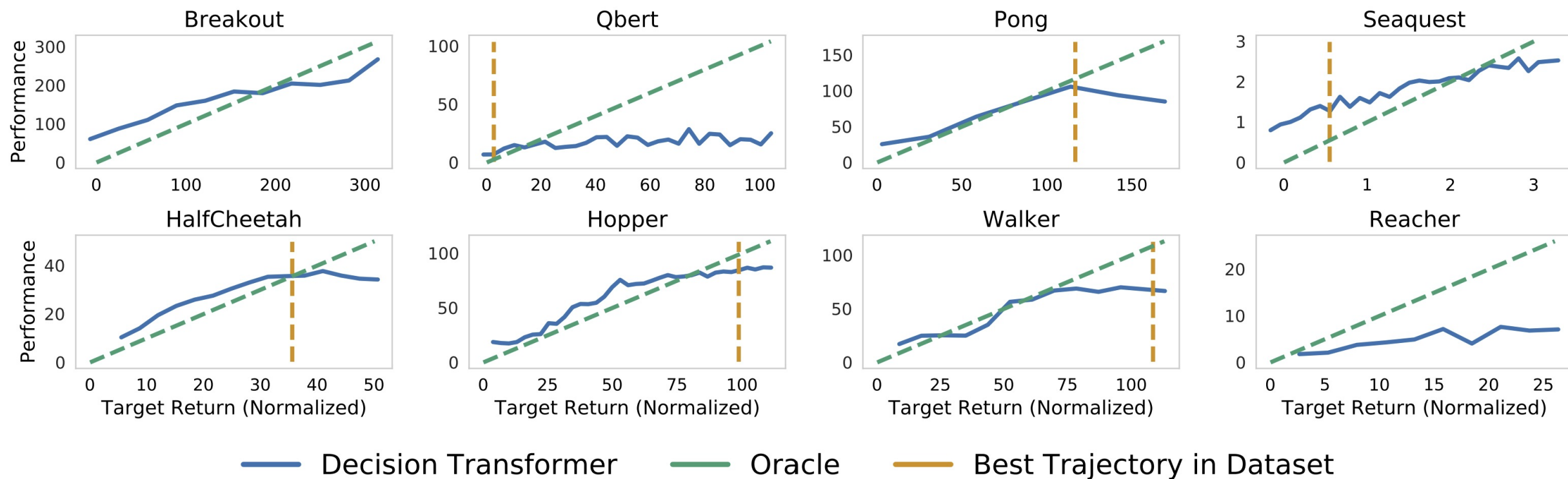
# DISCUSSION



# Q1: Does Decision Transformer perform behavior cloning on a subset of the data?

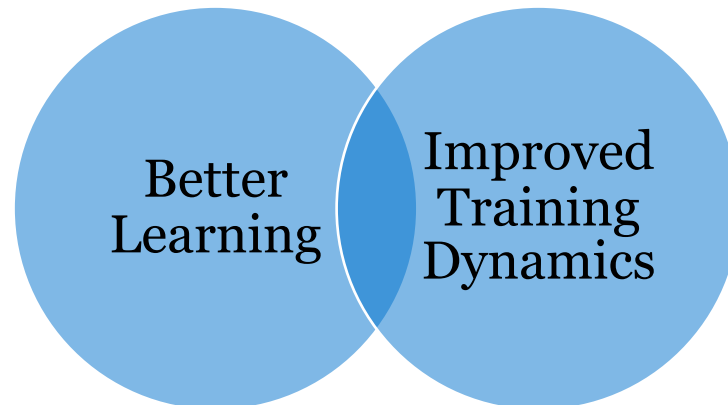
<i>Large Dataset</i>	Dataset	Environment	DT (Ours)	10%BC	25%BC	40%BC	100%BC	CQL
	Medium	HalfCheetah	$42.6 \pm 0.1$	42.9	43.0	43.1	43.1	<b>44.4</b>
	Medium	Hopper	<b><math>67.6 \pm 1.0</math></b>	65.9	65.2	65.3	63.9	58.0
	Medium	Walker	$74.0 \pm 1.4$	78.8	<b>80.9</b>	78.8	77.3	79.2
	Medium	Reacher	$51.2 \pm 3.4$	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><math>82.7 \pm 7.0</math></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
<i>Small Dataset</i>	Game	DT (Ours)	10%BC	25%BC	40%BC	100%BC		
	Breakout	<b><math>267.5 \pm 97.5</math></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><math>17.3 \pm 14.7</math></b>		
	Pong	<b><math>106.1 \pm 8.1</math></b>	$2.5 \pm 0.2$	$13.3 \pm 2.7$	$72.7 \pm 13.3$	$85.2 \pm 20.0$		
	Seaquest	<b><math>2.5 \pm 0.4</math></b>	$1.1 \pm 0.2$	$1.1 \pm 0.2$	$1.6 \pm 0.4$	$2.1 \pm 0.3$		

## Q2: How well does Decision Transformer model the distribution of returns?

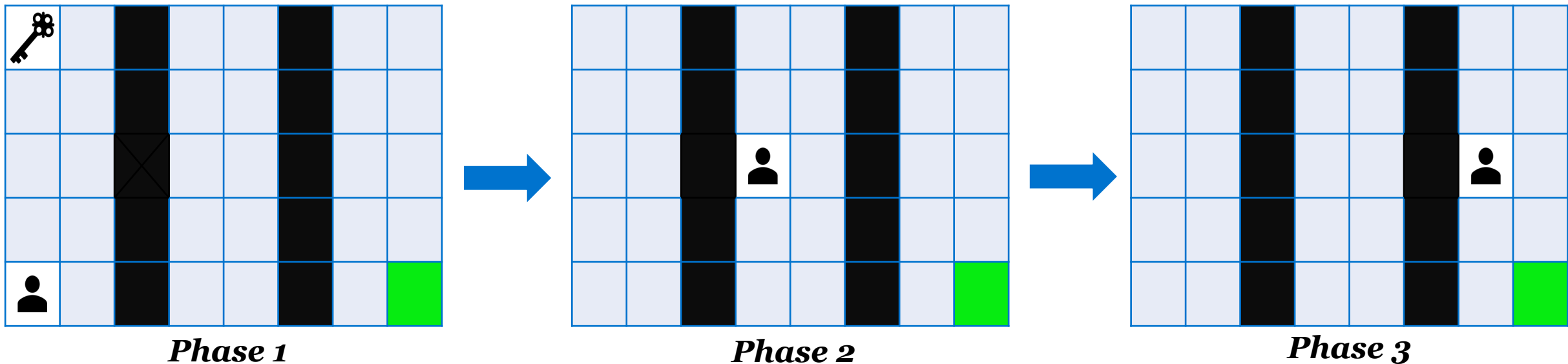


### Q3: What is the benefit of using a longer context length?

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

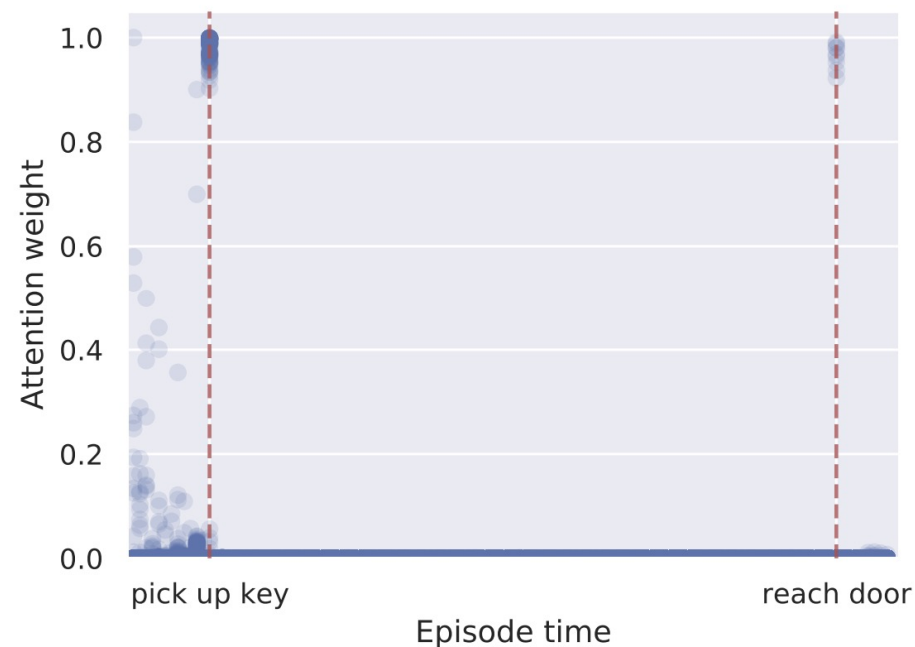
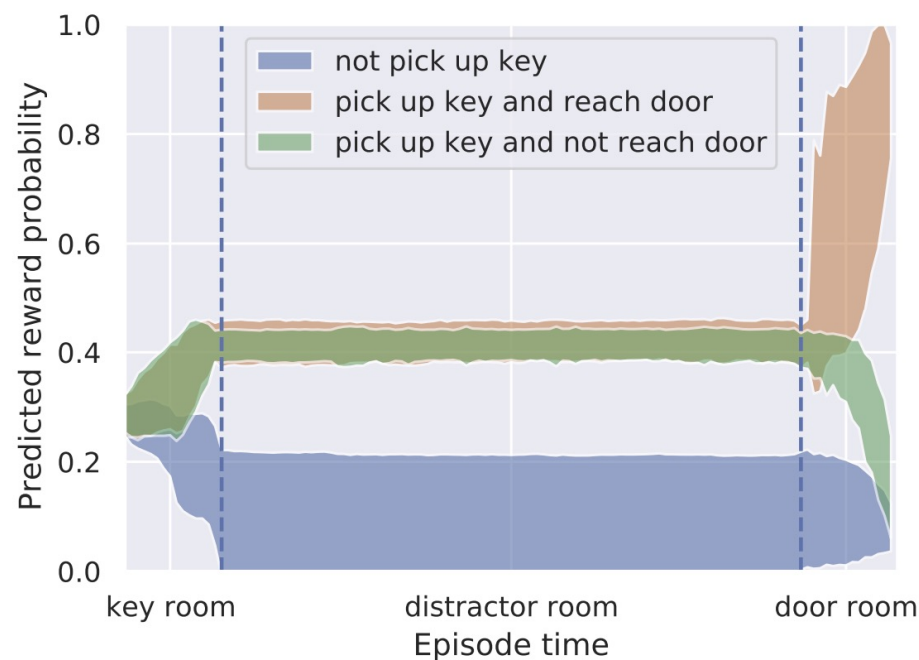


# Q4: Does Decision Transformer perform effective long-term credit assignment?



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	<b>94.6%</b>	13.3%	1.6%	<b>95.1%</b>	3.1%

# Q5: Can transformers be accurate critics in sparse reward settings?





## Q6: Does Decision Transformer perform well in sparse reward settings?



Dataset	Environment	Delayed (Sparse)		Agnostic		Original (Dense)	
		DT (Ours)	CQL	BC	%BC	DT (Ours)	CQL
Medium-Expert	Hopper	<b>107.3 <math>\pm</math> 3.5</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>78.5 <math>\pm</math> 3.7</b>	2.0	27.6	70.6	82.7	48.6

# Extra Observations

- No regularization or value pessimism needed
- Implicit representation of the value function
- Decision Transformer can benefit sample-efficient online regimes
- Can act as a strong model for behaviour generation

# CONCLUSION



# Conclusion

- Effective model-free supervised offline RL algorithm using sequence modelling.
- No reliance on any of the traditional RL concepts.
- Solves credit assignment and distribution shift problems seen in other RL algorithms.
- Match or surpass offline model-based RL state-of-the-art methods.

# Future Work

- 1 Use larger transformer models
- 2 Conditioning on return distributions instead of discrete returns
- 3 Model the state evolution using the transformer model to be an alternative for model-based RL.
- 4 Understand the errors made by transformers for risks in real-world settings..

# Limitations

Dependency on Context Length



Computational Time



Prior Knowledge on rewards



Loss of theoretical guarantees



# References

- [1] Christopher Watkins. Learning from delayed rewards. 01 1989
- [2] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602, 2013.
- [3] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. arXiv preprint arXiv:2004.07219, 2020.
- [4] Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N.M., Erez, T., Tassa, Y., Silver, D., & Wierstra, D. (2016). Continuous control with deep reinforcement learning. CoRR, abs/1509.02971.
- [5] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. arXiv preprint arXiv:1910.00177, 2019
- [6] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In International Conference on Machine Learning, 2019.
- [7] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. In Advances in Neural Information Processing Systems, 2019.
- [8] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In Advances in Neural Information Processing Systems, 2020.
- [9] Oh, J., Y. Guo, S. Singh, and H. Lee 2018. Self-Imitation Learning. arXiv:1806.05635 [cs, stat].
- [10] Arjona-Medina, J. A., M. Gillhofer, M. Widrich, T. Unterthiner, J. Brandstetter, and S. Hochreiter 2019. Rudder: Return decomposition for delayed rewards. In Advances in Neural Information Processing Systems, Pp. 13544–13555.
- [11] Pomerleau, D.A. (1988). ALVINN: An Autonomous Land Vehicle in a Neural Network. NIPS.

# References

- [12] Rupesh Kumar Srivastava, Pranav Shyam, Filipe Mutz, Wojciech Jaskowski, and Jürgen Schmidhuber. Training agents using upside-down reinforcement learning. arXiv preprint arXiv:1912.02877, 2019.
- [13] Aviral Kumar, Xue Bin Peng, and Sergey Levine. Reward-conditioned policies. arXiv preprint arXiv:1912.13465, 2019.
- [14] Dibya Ghosh, Abhishek Gupta, Justin Fu, Ashwin Reddy, Coline Devin, Benjamin Eysenbach, and Sergey Levine. Learning to reach goals without reinforcement learning. arXiv preprint arXiv:1912.06088, 2019.
- [15] Keiran Paster, Sheila A McIlraith, and Jimmy Ba. Planning from pixels using inverse dynamics models. arXiv preprint arXiv:2012.02419, 2020.
- [16] Johan Ferret, Raphaël Marinier, Matthieu Geist, and Olivier Pietquin. Self-attentional credit assignment for transfer in reinforcement learning. arXiv preprint arXiv:1907.08027, 2019.
- [17] Anna Harutyunyan, Will Dabney, Thomas Mesnard, Mohammad Azar, Bilal Piot, Nicolas Heess, Hado van Hasselt, Greg Wayne, Satinder Singh, Doina Precup, et al. Hindsight credit assignment. arXiv preprint arXiv:1912.02503, 2019.
- [18] Jose A Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brandstetter, and Sepp Hochreiter. Rudder: Return decomposition for delayed rewards. arXiv preprint arXiv:1806.07857, 2018
- [19] Chia-Chun Hung, Timothy Lillicrap, Josh Abramson, Yan Wu, Mehdi Mirza, Federico Carnevale, Arun Ahuja, and Greg Wayne. Optimizing agent behavior over long time scales by transporting value. Nature communications, 10(1):1–12, 2019.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems, 2017.
- [21] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- [22] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In Advances in Neural Information Processing Systems, 2020.
- [23] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In International Conference on Machine Learning, 2020.
- [24] Will Dabney, Mark Rowland, Marc Bellemare, and Rémi Munos. Distributional reinforcement learning with quantile regression. In Conference on Artificial Intelligence, 2018.



# References

- [25] Aviral Kumar, Justin Fu, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. arXiv preprint arXiv:1906.00949, 2019.
- [26] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. arXiv preprint arXiv:1911.11361, 2019.
- [27] Thomas Mesnard, Théophane Weber, Fabio Viola, Shantanu Thakoor, Alaa Saade, Anna Harutyunyan, Will Dabney, Tom Stepleton, Nicolas Heess, Arthur Guez, et al. Counterfactual credit assignment in model-free reinforcement learning. arXiv preprint arXiv:2011.09464, 2020.
- [28] Chen, L., Lu, K., Rajeswaran, A., Lee, K., Grover, A., Laskin, M., Abbeel, P., Srinivas, A., & Mordatch, I. (2021). Decision Transformer: Reinforcement Learning via Sequence Modeling. ArXiv, abs/2106.01345.

UNIVERSITY OF  
**WATERLOO**



**FACULTY OF MATHEMATICS**

<https://arxiv.org/pdf/2106.01345.pdf>

