

Big Models for Decision Making

Weinan Zhang Shanghai Jiao Tong University http://wnzhang.net

May 2022

Content

- Review the process of reinforcement learning
- Transformers in decision-making tasks
- Gato: a generalist agent
- Discussion and prospects

Overview of RL Training Paradigms

(a) online reinforcement learning



(b) off-policy reinforcement learning



(c) offline reinforcement learning



- Though both off-policy RL and offline RL evaluate the policy using the data sampled from a replay buffer, they are different.
- The key difference is whether the agent can interact with the environment while learning
- Offline RL techniques help deploy RL to realworld applications

Advantages of Offline RL

- Offline RL can help
 - 1. Pretrain an RL agent using an existing dataset
 - 2. Empirically evaluate RL algorithms based on their ability to exploit a fixed dataset of interactions
 - 3. Bridge the gap between academic interest in RL and real-world applications
- Offline RL makes RL more like supervised learning



Gulcehre, Caglar, et al. "RL unplugged: Benchmarks for offline reinforcement learning." arXiv preprint arXiv:2006.13888 (2020).

Key scientific problem

Extrapolation Error in Offline RL

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha(r + \gamma \max_{a'} Q(s',a'))$$

• Extrapolation error is introduced by the What if a' is an out-ofmismatch between the dataset and true distribution action? state-action visitation of the current policy.



Overview of Offline RL Methods



- The most severe problem that offline RL faces is the extrapolation error, i.e., the out-of-distribution problem: What if the agent performs unseen state-action?
- Implicit constraint methods directly learn the policy on the data support
 - with data instance selection or weighting

Behavior Cloning as Offline RL

 Behavior cloning, the simplest imitation learning method, requires no environment interaction





- Learning objective of BC Behavioral policy $\hat{\pi}^* = \arg \max_{\pi} \mathbb{E}_{\mathbf{s} \sim d_{\pi_E}} \mathbb{E}_{\mathbf{a} \sim \pi_E(\mathbf{s})}[\log \pi(\mathbf{a}|\mathbf{s})]$
- Obvious shortcomings of BC
 - 1. The policy upper bound is the behavioral policy
 - 2. Distribution shift

AWR: Advantage-Weighted Regression

Policy optimization objective

$$J(\pi) = \mathbb{E}_{\tau \sim p_{\pi}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^{t} r_{t} \right] = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{a \sim \pi(\mathbf{a}|\mathbf{s})} \left[r(\mathbf{s}, \mathbf{a}) \right]$$

Expected improvement

$$\begin{split} \eta(\pi) &= J(\pi) - J(\mu) \quad \text{[as derived in TRPO]} \\ &= \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[A^{\mu}(\mathbf{s}, \mathbf{a}) \right] = \mathbb{E}_{\mathbf{s} \sim d_{\pi}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} \left[\mathcal{R}_{\mathbf{s}, \mathbf{a}}^{\mu} - V^{\mu}(\mathbf{s}) \right] \end{split}$$

• Based on reward-weighted regression (RWR)

$$\pi_{k+1} = \arg \max_{\pi} \quad \mathbb{E}_{\mathbf{s} \sim d_{\pi_k}(\mathbf{s})} \mathbb{E}_{\mathbf{a} \sim \pi_k(\mathbf{a}|\mathbf{s})} \left[\log \pi(\mathbf{a}|\mathbf{s}) \exp\left(\frac{1}{\beta} \mathcal{R}_{\mathbf{s},\mathbf{a}}\right) \right]$$

return

 Regarded as solving a maximum likelihood problem that fits a new policy to samples collected under the current policy, where the likelihood is weighted by the exponentiated return.

Xue Bin Peng et al. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. arXiv preprint arXiv:1910.00177, 2019.

AWR: Advantage-Weighted Regression



• Performance of various algorithms on off-policy learning tasks with static datasets. AWR is able to learn policies that are comparable or better than the original demo policies.

Xue Bin Peng et al . Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. arXiv preprint arXiv:1910.00177, 2019.

BAIL: Best-Action Imitation Learning

- BAIL does not suffer from the extrapolation error, since it does not maximize over the actions space.
- Step 1: learn an upper envelope of state by solving a constrained optimization problem:

$$L^{K}(\phi) = \sum_{i=1}^{m} (V_{\phi}(s_{i}) - G_{i})^{2} \{ \mathbb{1}_{(V_{\phi}(s_{i}) \ge G_{i})} + K \cdot \mathbb{1}_{(V_{\phi}(s_{i}) < G_{i})} \} + \lambda \|w\|^{2}$$

- Step 2: Select actions satisfying G_i > x V(s_i) to perform simple imitation learning (BC). x is set such that 25% samples are selected.
- Step 3: perform supervised learning (BC) on selected data



X. Chen et al. BAIL: Best-action imitation learning for batch deep reinforcement learning. NeuIPS 2020.

BAIL: Best-Action Imitation Learning

Environment	BAIL	BCQ	BEAR	BC	MARWIL
$\sigma = 0.1$ Hopper B1	2173 ± 291	1219 ± 114	505 ± 285	626 ± 112	827 ± 220
$\sigma = 0.1$ Hopper B2	2078 ± 180	1178 ± 87	985 ± 3	579 ± 141	620 ± 336
$\sigma = 0.1$ Walker B1	1125 ± 113	576 ± 309	610 ± 212	514 ± 17	436 ± 24
$\sigma = 0.1$ Walker B2	3141 ± 300	2338 ± 388	2707 ± 425	1741 ± 239	1810 ± 200
$\sigma = 0.1 \text{ HC B1}$	5746 ± 29	5883 ± 43	0 ± 0	5546 ± 29	5573 ± 35
$\sigma = 0.1 \text{ HC B2}$	7212 ± 43	7562 ± 31	0 ± 0	6765 ± 108	6828 ± 111
$\sigma = 0.5$ Hopper B1	2054 ± 158	1145 ± 300	203 ± 42	919 ± 52	946 ± 103
$\sigma = 0.5$ Hopper B2	2623 ± 282	1823 ± 555	241 ± 239	694 ± 64	818 ± 112
$\sigma = 0.5$ Walker B1	2522 ± 51	1552 ± 455	1248 ± 181	2178 ± 178	2111 ± 52
$\sigma = 0.5$ Walker B2	3115 ± 133	2785 ± 123	2302 ± 630	2483 ± 94	2364 ± 228
$\sigma = 0.5 \text{ HC B1}$	1055 ± 9	1222 ± 38	924 ± 579	570 ± 35	512 ± 43
$\sigma = 0.5 \text{ HC B2}$	7173 ± 120	5807 ± 249	-114 ± 140	6545 ± 171	6668 ± 93
SAC HOPPER B1	3296 ± 105	2681 ± 438	1000 ± 110	2853 ± 318	2897 ± 227
SAC HOPPER B2	1831 ± 915	2134 ± 917	1139 ± 317	2240 ± 367	2063 ± 168
SAC WALKER B1	2455 ± 211	2408 ± 84	-3 ± 5	1674 ± 277	1484 ± 140
SAC WALKER B2	4767 ± 130	3794 ± 398	325 ± 75	2599 ± 145	2651 ± 268
SAC HC B1	10143 ± 77	8607 ± 473	7392 ± 257	8874 ± 221	9105 ± 90
SAC HC B2	10772 ± 59	10106 ± 134	7217 ± 273	9523 ± 164	9488 ± 136
SAC ANT B1	4284 ± 64	4042 ± 113	3452 ± 128	3986 ± 112	4033 ± 130
SAC ANT B2	4946 ± 148	4640 ± 76	3712 ± 236	4618 ± 111	4589 ± 130
SAC HUMANOID B1	3852 ± 430	1411 ± 250	0 ± 0	543 ± 378	589 ± 121
SAC HUMANOID B2	3565 ± 153	1221 ± 207	0 ± 0	1216 ± 826	1033 ± 257

X. Chen et al. BAIL: Best-action imitation learning for batch deep reinforcement learning. NeuIPS 2020.

Content

- Review the process of reinforcement learning
- Transformers in decision-making tasks
- Gato: a generalist agent
- Discussion and prospects

Upside-Down RL

O-value function

• Solve RL Problem with supervised learning methods?



Behavior Function

Key paper references

- Srivastava, Rupesh Kumar, Pranav Shyam, Filipe Mutz, Wojciech Jaśkowski, and Jürgen Schmidhuber. "Training agents using upside-down reinforcement learning." *arXiv preprint arXiv:1912.02877* (2019).
- Schmidhuber, Juergen. "Reinforcement Learning Upside Down: Don't Predict Rewards--Just Map Them to Actions." *arXiv preprint arXiv:1912.02875* (2019).

Upside-Down RL

Algorithms and experiments

Algorithm 1 Upside-Down Reinforcement Learning: High-level Description.

- 1: Initialize replay buffer with warm-up episodes using random actions
- 2: Initialize a behavior function
- 3: while stopping criteria is not reached do
- Improve the behavior function by training on replay buffer 4:
- Sample exploratory commands based on replay buffer 5:
- Generate episodes using Algorithm 2 and add to replay buffer 6:
- if evaluation required then 7:
- Evaluate current agent using Algorithm 2 8:
- end if 9:
- 10: end while

Algorithm 2 Generates an Episode using the Behavior Function.

Input: Initial command $c_0 = (d_0^r, d_0^h)$, Initial state s_0 , Behavior function $B(; \theta)$ **Output:** Episode data E 1: $E \leftarrow \emptyset$ 2: $t \leftarrow 0$ 3: while episode is not over do Compute $P(a_t|s_t, c_t) = B(s_t, c_t; \theta)$ 4: Execute $a_t \sim P(a_t|s_t, c_t)$ to obtain reward r_t and next state s_{t+1} from the environment 5: Append (s_t, a_t, r_t) to E 6: $s_t \leftarrow s_{t+1}$ // Update state 7: // Update desired reward $d_t^r \leftarrow d_t^r - r_t$ 8: $d_t^h \leftarrow d_t^h - 1$ // Update desired horizon 9: $c_t \leftarrow (d_t^r, d_t^h)$ 10: 11: $t \leftarrow t + 1$ 12: end while

A2C -100UDRL -1500.0 0.2 0.4 0.6 1.0 Environment Steps ×10

(a) On LunarLander-v2, TH is able to train agents that land the spacecraft, but is beaten by traditional RL algorithms.



(b) On TakeCover-v0, TH is able to consistently yield highperforming agents, while outperforming DQN and A2C.





(a) LunarLander-v2

(b) TakeCover-v0

Transformer for Sequence Modeling

- Encoder
 - Inputs: A sequence of vectors (words/states)
 - Outputs: A sequence of representation
- Decoder
 - Q: Start token + previous results -> Q
 - K & V: Outputs of encoder
 - Outputs: A sequence of symbol (labels/actions)



Attention in Transformer



Attention in RNN vs. Transformer

• Whether to use recent architecture or position embedding?



Here the same recent weights are used for multiple times of recent calculation, which introduce compounding error when decoding



Position embedding does not introduce such a compounding error when decoding

Decision Transformer

- Re-build RL task as a sequence prediction problem
- Causal transformer: GPT (i.e., decoder-only transformer)



Decision Transformer



- States, actions, and returns are fed into modality-specific linear embeddings and a positional episodic timestep encoding is added.
- Tokens are fed into a GPT architecture which predicts actions autoregressively using a causal self-attention mask

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
```

• Results comparing Decision Transformer to TD learning (CQL) and behavior cloning across Atari, OpenAI Gym, and Minigrid.



• Sampled (evaluation) returns accumulated by Decision Transformer when conditioned on the specified target (desired) returns.



Chen, Lili, et al. "Decision transformer: Reinforcement learning via sequence modeling." NeurIPS 2021.

• Results comparing Decision Transformer to TD learning (CQL) and behavior cloning across Atari, OpenAI Gym, and Minigrid.



• Sampled (evaluation) returns accumulated by Decision Transformer when conditioned on the specified target (desired) returns.



Chen, Lili, et al. "Decision transformer: Reinforcement learning via sequence modeling." NeurIPS 2021.

 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 ± 0.1	42.9	43.0	43.1	43.1	44.4			
Medium	Hopper	67.6 ± 1.0	65.9	65.2	65.3	63.9	58.0			
Medium	Walker	74.0 ± 1.4	78.8	80.9	78.8	77.3	79.2			
Medium	Reacher	51.2 ± 3.4	51.0	48.9	58.2	58.4	26.0			
Medium-Replay	HalfCheetah	36.6 ± 0.8	40.8	40.9	41.1	4.3	46.2			
Medium-Replay	Hopper	82.7 ± 7.0	70.6	58.6	31.0	27.6	48.6			
Medium-Replay	Walker	66.6 ± 3.0	70.4	67.8	67.2	36.9	26.7			
Medium-Replay	Reacher	18.0 ± 2.4	33.1	16.2	10.7	5.4	19.0			
Average		56.1	56.7	52.7	49.4	39.5	43.5			

OpenAl Gym D4RL

Atari Offline

Game	DT (Ours)	10%BC	25%BC	40%BC	100%BC
Breakout	267.5 ± 97.5	28.5 ± 8.2	73.5 ± 6.4	108.2 ± 67.5	138.9 ± 61.7
Qbert	15.1 ± 11.4	6.6 ± 1.7	16.0 ± 13.8	11.8 ± 5.8	$\bf 17.3 \pm 14.7$
Pong	106.1 ± 8.1	2.5 ± 0.2	13.3 ± 2.7	72.7 ± 13.3	85.2 ± 20.0
Seaquest	2.4 ± 0.7	1.1 ± 0.2	1.1 ± 0.2	1.6 ± 0.4	2.1 ± 0.3

Chen, Lili, et al. "Decision transformer: Reinforcement learning via sequence modeling." NeurIPS 2021.

- Credit assignment over long horizon
- A grid-based environment (Key-to-door) with a sequence of three phases:
 - 1. in the first phase, the agent is placed in a room with a key;
 - 2. then, the agent is placed in an empty room;
 - 3. and finally, the agent is placed in a room with a door.
- The agent receives a binary reward when reaching the door in the third phase, but only if it picked up the key in the first phase.

Dataset	DT (Ours)	CQL	BC	%BC	Random
1K Random Trajectories 10K Random Trajectories	71.8 % 94.6%	$13.1\%\ 13.3\%$	$1.4\% \\ 1.6\%$	$69.9\% \\ 95.1\%$	$3.1\%\ 3.1\%$

Key-to-door



TD learning (CQL) cannot effectively propagate Q-values over the long horizons involved and gets poor performance

Chen, Lili, et al. "Decision transformer: Reinforcement learning via sequence modeling." NeurIPS 2021.

Trajectory Transformer

- Offline RL as a sequence prediction problem
- Build a world model for RL task



• Discretize each dimension independently

$$\boldsymbol{\tau} = \left(\mathbf{s}_1, \mathbf{a}_1, r_1, \mathbf{s}_2, \mathbf{a}_2, r_2, \dots, \mathbf{s}_T, \mathbf{a}_T, r_T\right)$$

= $\left(\dots, s_t^1, s_t^2, \dots, s_t^N, a_t^1, a_t^2, \dots, a_t^M, r_t, \dots\right)$ $t = 1, \dots, T.$

• Training loss

$$\mathcal{L}(\tau) = \sum_{t=1}^{T} \left(\sum_{i=1}^{N} \log P_{\theta} \left(s_{t}^{i} \mid \mathbf{s}_{t}^{< i}, \boldsymbol{\tau}_{< t} \right) + \sum_{j=1}^{M} \log P_{\theta} \left(a_{t}^{j} \mid \mathbf{a}_{t}^{< j}, \mathbf{s}_{t}, \boldsymbol{\tau}_{< t} \right) + \log P_{\theta} \left(r_{t} \mid \mathbf{a}_{t}, \mathbf{s}_{t}, \boldsymbol{\tau}_{< t} \right) \right)$$

Trajectory Transformer

- Offline RL as a sequence prediction problem
- Build a world model for RL task



TT is similar with DT, but:

- DT is based-on UDRL, could be used as policies directly, while TT is a world model.
- Inputs all elements and predicts all (trajectories and rewards).
- Applicable in more cases, e.g., imitation learning, goal-conditioned reinforcement learning, and offline reinforcement learning.

Trajectory Transformer



Trajectory Transformer Decision Making

Taking actions via planning with beam search

Algorithm 1 Beam search

- 1: **Require** Input sequence x, vocabulary \mathcal{V} , sequence length T, beam width B
- 2: Initialize $Y_0 = \{ () \}$
- 3: for t = 1, ..., T do 4:
- $\begin{array}{c} \mathcal{C}_t \leftarrow \{\mathbf{y}_{t-1} \circ y \mid \mathbf{y}_{t-1} \in Y_{t-1} \text{ and } y \in \mathcal{V}\} \\ Y_t \leftarrow \operatorname{argmax} \log P_{\theta}(Y \mid \mathbf{x}) \end{array} // B \text{ most likely sequences from candidates} \end{array}$ 5:
 - $Y \subset \mathcal{C}_t, |Y| = B$
- 6: end for
- 7: **Return** argmax $\log P_{\theta}(\mathbf{y} \mid \mathbf{x})$ $\mathbf{v} \in Y_T$
- Planning in three tasks
 - Imitation learning
 - Setting $x = (\tau_{< t}, s_t)$
 - Goal-conditioned RL (s_T as goal)
 - Setting the input trajectory as $(\mathbf{s}_T, \boldsymbol{\tau}_{< t}, \boldsymbol{s}_t)$
 - Offline RL (seek for the highest return)
 - Keep the top reward + reward-to-go ($R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$) trajectories as Y_t

Beam width	maximum number of hypotheses retained during beam search	256
Planning horizon	number of transitions predicted by the model during	15
Vocabulary size	number of bins used for autoregressive discretization	100
Context size	number of input $(\mathbf{s}_t, \mathbf{a}_t, \mathbf{r}_t, R_t)$ transitions	5
$k_{ m obs}$	top- k tokens from which observations are sampled	1
$k_{ m act}$	top- k tokens from which actions	20

Trajectory Transformer Experiments

over 15 random seeds Dataset **Environment** BC **MBOP** BRAC CQL DT TT (uniform) **TT** (quantile) Med-Expert HalfCheetah 59.9105.941.9 91.6 86.8 40.8 ± 2.3 95.0 ± 0.2 Med-Expert 79.655.10.9105.4107.6 106.0 ± 0.28 110.0 ± 2.7 Hopper Med-Expert Walker2d 36.670.281.6 108.8 108.1 91.0 ± 2.8 101.9 ± 6.8 Medium HalfCheetah 43.144.646.344.042.6 44.0 ± 0.31 46.9 ± 0.4 Medium 63.9 48.831.3 58.567.6 67.4 ± 2.9 61.1 ± 3.6 Hopper Medium Walker2d 77.3 41.081.1 72.574.0 81.3 ± 2.1 79.0 ± 2.8 HalfCheetah 4.342.347.745.536.6 44.1 ± 0.9 41.9 ± 2.5 Med-Replay 82.7 99.4 ± 3.2 91.5 ± 3.6 Med-Replay Hopper 27.612.40.695.0Med-Replay Walker2d 36.99.7 0.977.266.6 79.4 ± 3.3 82.6 ± 6.9 47.7 47.8 36.9 77.6 74.7 72.6 78.9 Average

Offline RL performance comparison

TT performs on par with or better than the best prior offline reinforcement learning algorithms on D4RL locomotion (v2) tasks.



The Progress So Far

- Transformer
 - Currently most powerful sequential model, derived from NLP
 - Long-sequence modeling capabilities
 - Capable for large volume of parameters
 - Stronger generalization
- Upside-down RL
 - Theoretical foundation of DT
 - Solves reinforcement learning problems with supervised learning methods
- Decision Transformer
 - Based on UDRL, could be used as policies directly
 - Perform credit assignment directly via self-attention, bypass the need for bootstrapping for long term credit assignment
 - Avoids the need for discounting future rewards, which can induce undesirable shortsighted behaviors
 - Avoids value overestimation
- Trajectory Transformer
 - Similar with DT in implementation, but inputs all and predicts all
 - More like a world model
 - Applicable in more cases, e.g., imitation learning, goal-conditioned reinforcement learning, and offline reinforcement learning

Content

- Review the process of reinforcement learning
- Transformers in decision-making tasks
- Gato: a generalist agent
- Discussion and prospects



A Generalist Agent

Scott Reed^{*,†}, Konrad Żołna^{*}, Emilio Parisotto^{*}, Sergio Gómez Colmenarejo[†], Alexander Novikov, Gabriel Barth-Maron, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar and Nando de Freitas[†]

^{*}Equal contributions, [†]Equal senior contributions, All authors are affiliated with DeepMind

- From specialist agents in RL to generalist agents
- Prospects
 - Reduce the need for handcrafting policy models.
 - Generic models have tended to be better than domainspecific approaches.

Gato: A Generalist Agent

- Handle a variety of tasks with different types of data
- Multi-modal, multi-task, multi-embodiment: agents have different characteristics in different domains, e.g., shape, action spaces ...



Gato Model

- A single and large decoder-only transformer
- Purely supervised learning (in principle, also offline / online RL)

Hyperparameter	Gато 1.18В	364M	79M
TRANSFORMER BLOCKS	24	12	8
ATTENTION HEADS	16	12	24
LAYER WIDTH	2048	1536	768
Feedforward hidden size	8192	6144	3072
Key/value size	128	128	32
Shared embedding		True	
LAYER NORMALIZATION Activation Function	PF	RE-NORM GEGLU	1



Gato Model - Tokenization

Modality-specific tokenization

- Text is encoded via SentencePiece (Kudo and Richardson, 2018) with 32000 subwords into the integer range [0, 32000).
- Images are first transformed into sequences of non-overlapping 16×16 patches in raster order, as done in ViT (Dosovitskiy et al., 2020). Each pixel in the image patches is then normalized between [-1, 1] and divided by the square-root of the patch size (i.e. $\sqrt{16} = 4$).
- Discrete values, e.g. Atari button presses, are flattened into sequences of integers in row-major order. The tokenized result is a sequence of integers within the range of [0, 1024).
- Continuous values, e.g. proprioceptive inputs or joint torques, are first flattened into sequences of floating point values in row-major order. The values are mu-law encoded to the range [-1, 1] if not already there (see figure 13 for details), then discretized to 1024 uniform bins. The discrete integers are then shifted to the range of [32000, 33024).

After converting data into tokens, we use the following canonical sequence ordering.

- Text tokens in the same order as the raw input text.
- Image patch tokens in raster order.
- Tensors in row-major order.
- Nested structures in lexicographical order by key.
- Agent timesteps as observation tokens followed by a separator, then action tokens.
- Agent episodes as timesteps in time order.

Gato Model - Tokenization

- The data is sequenced as follows
 - Episodes are presented to the agent in order of time (timesteps).
 - **Timesteps** in turn are presented in the following order:
 - **Observations** ($[y_{1:k}, x_{1:m}, z_{1:n}]$) are ordered lexicographically by key, each item is sequenced as follows:
 - * Text tokens $(y_{1:k})$ are in the same order as the raw input text.
 - * Image patch tokens $(x_{1:m})$ are in raster order.
 - * Tensors $(z_{1:n})$ (such as discrete and continuous observations) are in row-major order.
 - **Separator** ('|'); a designated separator token is provided after observations.
 - Actions $(a_{1:A})$ are tokenized as discrete or continuous values and in row-major order.

A full sequence of tokens is thus given as the concatenation of data from T timesteps:

$$s_{1:L} = [[y_{1:k}^1, x_{1:m}^1, z_{1:n}^1, '|', a_{1:A}^1], \dots, [y_{1:k}^T, x_{1:m}^T, z_{1:n}^T, '|', a_{1:A}^T]],$$

where L = T(k + m + n + 1 + A) is the total number of tokens.

Details of Tokenization

- Text and discrete data are directly embedded
- Tensor data are mu-law encoded and discretized, then embedded
- Image patches are embedded via ResNet

Patch

Embed

(ResNet)

1.2

4.7

÷



ResNet embedding for image patches

Addition

 \mathbf{x}_{l+1}

 \mathbf{x}_l

GroupNorm

GeLU

Weight

GroupNorm

GeLU

Weight

Gato Model – Training

- Output
 - Only text tokens, discrete and continuous values, and actions.
 - No image tokens or observations are used as output targets.
- Training
 - Purely supervised learning
 - Minimize log-loss on outputs

$$\mathcal{L}(\theta, \mathcal{B}) = -\sum_{b=1}^{|\mathcal{B}|} \sum_{l=1}^{L} m(b, t) \log p_{\theta}(s_{l}^{(b)} | s_{1}^{(b)}, s_{2}^{(b)}, \dots, s_{l-1}^{(b)})$$

Masking function such at m(b,t)=1 for output tokens, 0 otherwise

Gato Model – Prompt

- Prompt to identify tasks
 - 'Prompt' means the tokenized sequences of partitions of trajectories
 - Training stage: 25% sequence in a batch are prepended a prompt sequence
 - from an episode generated by the same source agent on the same task
 - Half of the prompt are sampled from the end of the episode
 - Half of the prompt are uniformly sampled from the episode
- Evaluation stage: a successful demo of the desired task.



Gato Model – Inference Procedure

- 1. Start with prompt, such as a demonstration token sequence
- 2. For each timestep
 - Sample the action vector autoregressively for one token at a time
 - Once all tokens comprising the action vector have been sampled, decode the action by inverting the tokenization procedure
 - Deliver the action to the environment and get the new observation
- NOTE: The model always sees all previous observations and actions in its context window of 1024 tokens.



Gato Experiments - Datasets

Control onvironment	Tacks	Enicodos	Approx.	Sample	Vision / language dataset	Sample
Control environment	14585	Episodes	Tokens	Weight	VISION / Tanguage dataset	Weight
DM Lab	254	16.4M	194B	9.35%	MassiveText	6.7%
ALE Atari	51	63.4K	1.26B	9.5%	M3W	4%
ALE Atari Extended	28	28.4K	565M	10.0%	ALIGN	0.67%
Sokoban	1	27.2K	298M	1.33%	MS-COCO Captions	0.67%
BabyAI	46	4.61M	22.8B	9.06%	Conceptual Captions	0.67%
DM Control Suite	30	395K	22.5B	4.62%	LTIP	0.67%
DM Control Suite Pixels	28	485K	35.5B	7.07%	OKVQA	0.67%
DM Control Suite Random Small	26	10.6M	313B	3.04%	VQAV2	0.67%
DM Control Suite Random Large	26	26.1M	791B	3.04%	Total	14.7%
Meta-World	45	94.6K	3.39B	8.96%		II
Procgen Benchmark	16	1.6M	4.46B	5.34%		
RGB Stacking simulator	1	387K	24.4B	1.33%		
RGB Stacking real robot	1	15.7K	980M	1.33%		
Modular RL	38	843K	69.6B	8.23%		
DM Manipulation Playground	4	286K	6.58B	1.68%		
Playroom	1	829K	118B	1.33%		
Total	596	63M	1.5T	85.3%		

• Data filtering

return of episode *i*

• Expert return: $\max_{j \in [0,1,...,N-W]} \sum_{i=j}^{j+W-1} R_i / W$

 $W = \min(1000, 0.1 \times N)$

• Select the episodes $\{i\}$ with $R_i \ge 80\%$ of expert return

Gato Experiments – Simulated Control



- 450/604 tasks: performance over 50% expert score
- E.g., ALE Atari: achieve average human (or better) scores for 23 games
- Computing resource for final training
 - 16x16 TPUv3 cluster x 4 days \simeq 1600 V100 GPU days \simeq 480k RMB
 - Batch size 512, token sequence length 1024
- Prospects (as previously mentioned)
 - Reduce the need for handcrafting policy models.
 - Generic models have tended to be better than domain-specific approaches.

Gato Experiments – Robotics





- Robotics RGB Stacking Benchmark
 - Three plastic blocks colored red, green and blue with varying shapes
 - The goal is to stack red on blue, ignoring green
 - Observations: two 128×128 camera images, robot arm and gripper joint angles, the robot's end-effector pose
 - Two challenge settings
 - Skill mastery: training provided the 5 test object triplets it is later tested on
 - Skill generalization: training not provided 5 test object triplets

Gato Experiments – Robotics

Real robot Skill Mastery results. Gato is competitive with the filtered BC baseline.

Agent	Group 1	GROUP 2	Group 3	Group 4	GROUP 5	Average
Gato	58%	57.6%	78.5 %	89 %	95.1 %	75.6 %
BC-IMP (Lee et al., 2021)	75.6 %	60.8%	70.8%	87.8%	78.3%	74.6%

Gato real robot Skill Generalization results. In addition to performing hundreds of other tasks, Gato also stacks competitively with the comparable published baseline.

Agent	GROUP 1	GROUP 2	GROUP 3	Group 4	GROUP 5	Average
Gato	24.5 %	33%	50.5 %	76.5%	66.5 %	50.2 %
BC-IMP (Lee et al., 2021)	23%	39.3 %	39.3%	77.5 %	66%	49%

 The standard "stack red on blue" task tested in the Skill Generalization benchmark.



• Results: robot successfully stacking the red object on the blue object, and the data does not include the object shapes in the test set

Gato Experiments – Robotics

 The novel "stack blue on green" task that demonstrates Gato's out of distribution adaptation to perceptual variations.



5000

79M

• Additionally adding simulated demonstrations of the stack blue on green task to the fine-tuning dataset improved performance.



Gato Experiments – Image Captioning



The colorful ceramic toys are on Man standing in the street the living room floor.

a living room with three different color deposits on the floor

a room with a long red rug a tv and some pictures



wearing a suit and tie.

A man in a blue suit with a white bow tie and black shoes.

A man with a hat in his hand looking at the camera



A bearded man is holding a plate of food.

Man holding up a banana to take a picture of it.

a man smiles while holding up a slice of cake



a group of people that is next to a big horse

A tan horse holding a piece of cloth lying on the ground.

Two horses are laying on their side on the dirt.



Man biting a kite while standing on a construction site

a big truck in the middle of a road

A truck with a kite painted on the back is parked by rocks.



a white horse with a blue and silver bridle

A white horse with blue and gold chains.

A horse is being shown behind a wall.



a couple of people are out in the A baseball player pitching a ball ocean

A surfer riding a wave in the ocean.

A surfer with a wet suit riding a wave.



on top of a baseball field.

A man throwing a baseball at a pitcher on a baseball field.

A baseball player at bat and a catcher in the dirt during a baseball game



Pistachios on top of a bowl with coffee on the side.

A bowl and a glass of liquid sits on a table.

A white plate filled with a banana bread next to a cup of coffee.



A group of children eating pizza at a table.

Two boys having pizza for lunch with their friends.

The boys are eating pizza together at the table.

The first three captions sampled using temperature 0.9, without cherry-picking.

Gato Experiments – Dialogues



Dialogues with Gato when it is prompted to be a chat bot. Usually Gato replies with a relevant response, but is
often superficial or factually incorrect, which could likely be improved with further scaling.

Gato Experiments

- Model size scaling laws results on indistribution overall performance
- Fine-tuning is very similar to pretraining with minor changes, such as different learning rate schedule (still supervised learning)
- Few-shot performance, ablating over various pretraining settings on 364M parameter variants of Gato





Summary of Big Decision-Making Models

- Main idea: leverage the power of large sequence models (transformers) to generate good actions given the conditions
- Decision transformer
 - A policy outputs the action conditioning on state, action, return-to-go pre-sequence
 - Offline supervised learning
- Trajectory transformer
 - A world model to autoregressively generate the next state, action, reward conditioning on state, action, reward pre-sequence
 - Take actions via planning with beam search (or RL policies)
- Both DT and TT are specialist agents
- Gato a generalist agent
 - Multi-task, multi-model, multi-embodiment
 - Prompt to identify tasks
 - Purely supervised learning



Content

- Review the process of reinforcement learning
- Transformers in decision-making tasks
- Gato: a generalist agent
- Discussion and prospects

Discussion and Prospects

- Now it seems supervised learning has become comparable with reinforcement learning on decision-making tasks
 - Transformer: model complex sequential patterns and deal with compounding error problems
 - Huge dataset: sufficient data to train the transformer
 - Offline training: do not require the agent with interact with the environment during training
- Near-future research on this direction
 - Environment model: build better world models
 - Data: how to collect the data for more effective training
 - Goal: better goal representation in decision transformers
 - Transfer: zero-shot transfer
 - MARL: multi-agent decision-making tasks

Discussion and Prospects



$$\hat{\pi}^* = \arg \max_{\pi} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim d_{\pi_D}} [m(\mathbf{s}, \mathbf{a}) \log \pi(\mathbf{a} | \mathbf{s})]$$

- In traditional RL, m(s, a) is the value function corresponds to $Q^{\pi}(s, a)$
 - Learned by TD propagation
 - Need the agent to interact with environment
- In big DM models, m(s, a) is a masking (or weighting) function
 - Just predefined by rule
 - No need to interact with environment

Thank You! Questions?





SHANGHAI JIAO TONG UNIVERSITY

Weinan Zhang

Associate Professor

APEX Data & Knowledge Management Lab

Department of Computer Science and Engineering

Shanghai Jiao Tong University

http://wnzhang.net