# CtrlFormer: Learning Transferable State Representation for Visual Control via Transformer
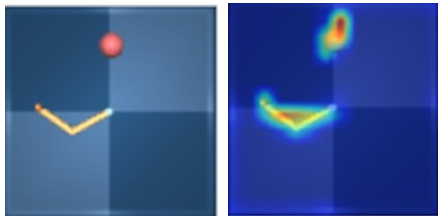
Yao Mu, Shoufa Chen, Mingyu Ding, Jianyu Chen, Runjian Chen, Ping Luo

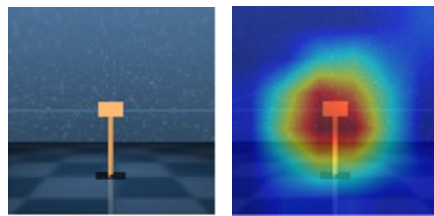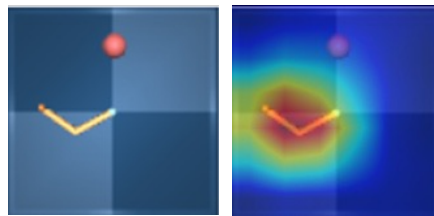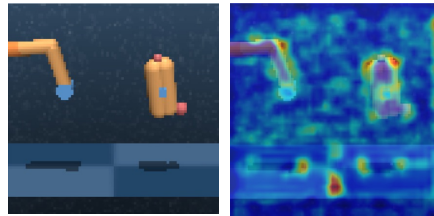The University of Hong Kong, Tsinghua University

2022.06.27

# 1 Motivation


Visualization of the attention of task specific visual representation

- Task specific
- High sample efficiency for RL learning
- Difficult to transfer across tasks

Each task has its own focus


Visualization of the attention of Pretrained ResNet

- Task-independent
- Easy to transfer across tasks
- Low sample efficiency for RL learning

Can we learn the visual representation mechanism like a human, which can capture the characteristic of every task and can be easily transferred to a new task?


Human behavior learning

learn to learn tasks

quickly learn new task

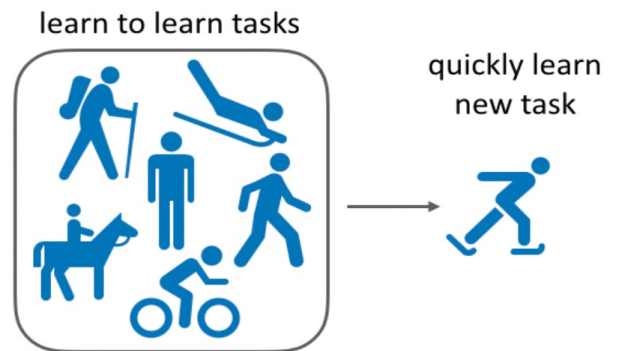- CtrlFormer jointly learns self-attention mechanisms between visual tokens and policy tokens among different control tasks, where multitask representation can be learned and transferred without catastrophic forgetting.
- We carefully design a contrastive reinforcement learning paradigm to train CtrlFormer, enabling it to achieve high sample efficiency, which is important in control problems.
- Extensive experiments show that CtrlFormer outperforms previous works in terms of both transferability and sample efficiency without catastrophic forgetting.



(a) Maintainability

(b) Transferability

Effect of CtrlFormer

Visual tokens          Policy tokens

Each task has its own policy token

# Overall framework of CtrlFormer



- Policy token is a learnable variable that learns a context for its task during the learning process

- Task-related information can be extracted by computing the attention of policy tokens and other tokens.

Explicitly model the attention mechanism between the new task and the old task and input images thus enabling fast transfer of knowledge learned from the old task to the new one

- Contrastive Co-training to improve the sample efficiency
- Reduce the number of parameters to be learned by multi-stage Pooling

Thus, the input of the transformer is

$$\mathbf{z}_{\ell_0} = \left[ \mathbf{x}_{\mathrm{con}}; \mathbf{x}_\pi^1; \ldots; \mathbf{x}_\pi^K; \mathbf{x}_p^1; \cdots; \mathbf{x}_p^N \right] + \mathbf{E}_{\mathrm{pos}}$$

$$\mathbf{z}'_{\ell_j} = \mathrm{MHSA}\left( \mathrm{LN}\left( \mathbf{z}_{\ell_{j-1}} \right) \right) + \mathbf{z}_{\ell_{j-1}}$$

$$\mathbf{z}_{\ell_j} = \mathrm{MLP}\left( \mathrm{LN}\left( z'_{\ell_j} \right) \right) + \mathbf{z}'_{\ell_j}$$
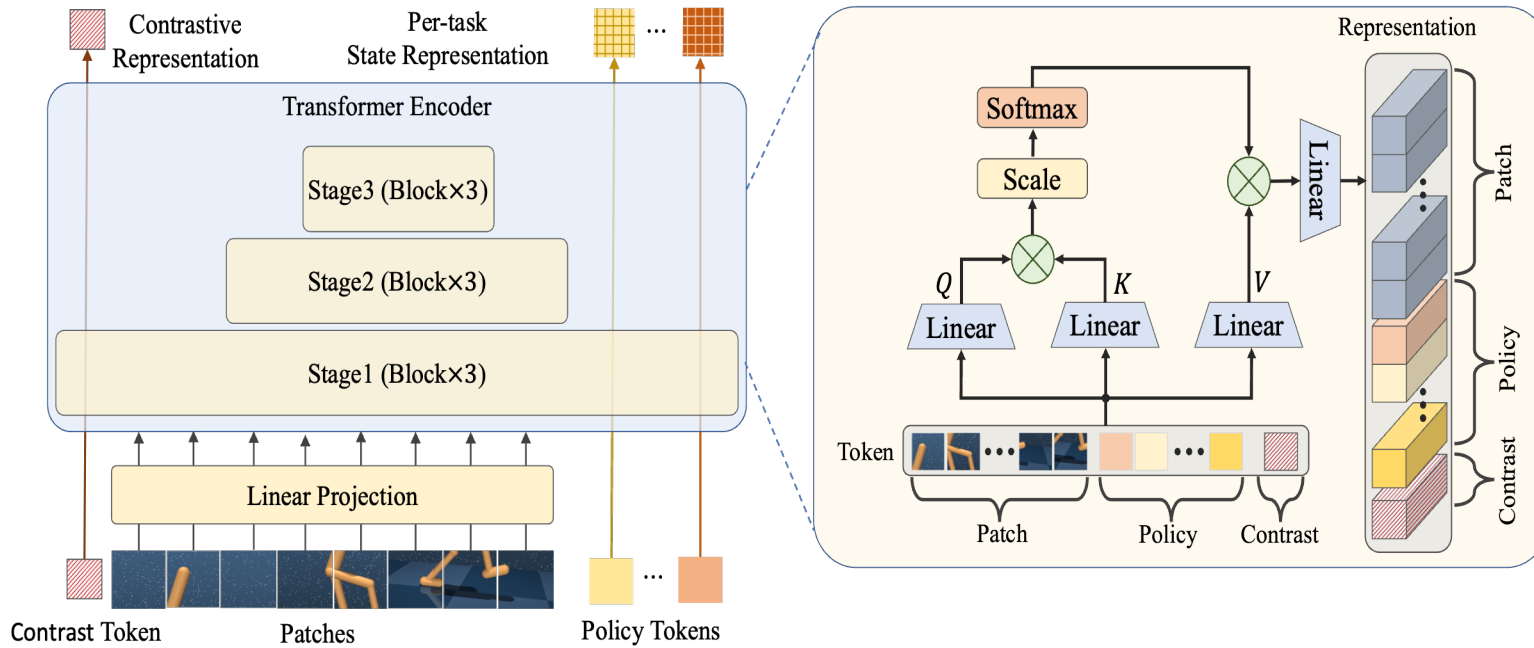
$$\mathbf{q}_\pi = \mathbf{z}_\pi \mathbf{W}^q, \mathbf{k} = \mathbf{z} \mathbf{W}^k$$

$$\mathbf{z} = \left[ \mathbf{z}_{con}; \mathbf{z}_\pi^1; \ldots; \mathbf{z}_\pi^{K+1}; \mathbf{z}_p^1; \ldots; \mathbf{z}_p^N \right]$$

$$\mathbf{z}_\pi = \left[ \mathbf{z}_\pi^1; \ldots; \mathbf{z}_\pi^{K+1} \right]$$

# Visualization



Visualization of CtrlFormer



Visualization of Pretrained ResNet



(a) DrQ

(b) CtrlFormer

Comparison of the attention map change before and after the transferring

| Method | Learn from scratch 100k | 500k | Retest after new task fine-tune |
|---|---|---|---|
| DrQ | $549_{\pm36}$ | $\mathbf{854_{\pm22}}$ | $373_{\pm24}$ |
| Dreamer | $326_{\pm27}$ | $762_{\pm27}$ | $704_{\pm33}$ |
| Resnet+SAC | $192_{\pm19}$ | $357_{\pm85}$ | $357_{\pm85}$ |
| CtrlFormer | $\mathbf{759_{\pm48}}$ | $846_{\pm25}$ | $\mathbf{842_{\pm22}}$ |

(a) **Left:** Learn old task in **Cartpole** (`swingup`)

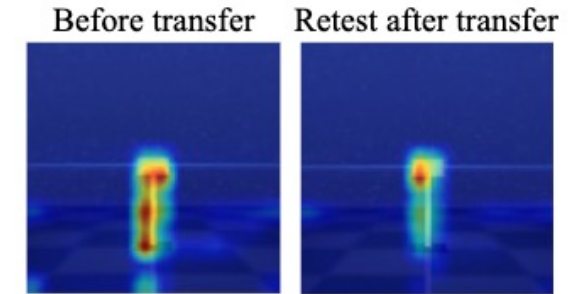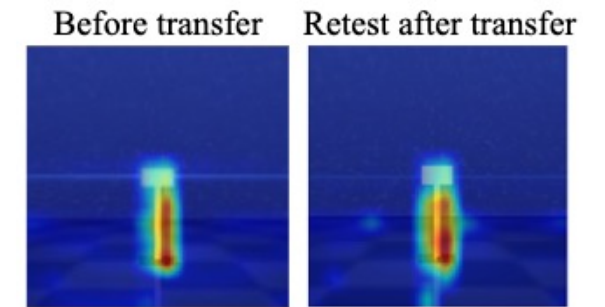| Method | Learn from scratch 100k | 500k | Learn with transfer 100k | 500k |
|---|---|---|---|---|
| DrQ | 0 | $505_{\pm335}$ | 0 | $75.5_{\pm41}$ |
| Dreamer | $8_{\pm4}$ | $376_{\pm214}$ | 0 | $589_{\pm122}$ |
| Resnet+SAC | 0 | 0 | 0 | 0 |
| CtrlFormer | 0 | $\mathbf{671_{\pm81}}$ | $\mathbf{769_{\pm34}}$ | $\mathbf{804_{\pm26}}$ |

**Right:** Transfer to new task **Cartpole** (`swingup-sparse`)

| Method | Scratch (previous) 500 k | Transfer (new task) 100 k | 500 k | Retest (previous) 500 k |
|---|---|---|---|---|
| DrQ | $\mathbf{971_{\pm27}}$ | $283_{\pm121}$ | $332_{\pm96}$ | $124_{\pm22}$ |
| Resnet+SAC | $382_{\pm299}$ | $298_{\pm17}$ | $300_{\pm29}$ | $382_{\pm299}$ |
| CtrlFormer | $918_{\pm33}$ | $\mathbf{299_{\pm38}}$ | $\mathbf{547_{\pm56}}$ | $\mathbf{889_{\pm34}}$ |

(a) Transfer from **Reacher**(`easy`) to **Finger**(`turn-easy`)

| Method | Learn from scratch 100k | 500k | Retest after new task fine-tune |
|---|---|---|---|
| DrQ | $\mathbf{346_{\pm33}}$ | $448_{\pm65}$ | $300_{\pm42}$ |
| Dreamer | $25_{\pm18}$ | $245_{\pm159}$ | $182_{\pm34}$ |
| Resnet+SAC | $298_{\pm17}$ | $300_{\pm29}$ | $300_{\pm29}$ |
| CtrlFormer | $281_{\pm67}$ | $\mathbf{493_{\pm35}}$ | $\mathbf{475_{\pm43}}$ |

(b) **Left:** Learn old task in **Finger** (`turn-easy`)

| Method | Learn from scratch 100k | 500k | Learn with transfer 100k | 500k |
|---|---|---|---|---|
| DrQ | $8_{\pm24}$ | $274_{\pm137}$ | $133_{\pm26}$ | $455_{\pm34}$ |
| Dreamer | 0.0 | $17_{\pm9}$ | 0.0 | $38_{\pm18}$ |
| Resnet+SAC | 0.0 | $17_{\pm10}$ | 0.0 | $17_{\pm10}$ |
| CtrlFormer | $197_{\pm78}$ | $\mathbf{344_{\pm47}}$ | $\mathbf{294_{\pm37}}$ | $\mathbf{569_{\pm32}}$ |

**Right:** Transfer to new task **Finger** (`turn-hard`)

| Method | Scratch (previous) 500 k | Transfer (new task) 100 k | 500 k | Retest (previous) 500 k |
|---|---|---|---|---|
| DrQ | $\mathbf{448_{\pm65}}$ | $203_{\pm87}$ | $693_{\pm282}$ | $184_{\pm57}$ |
| Resnet+SAC | $300_{\pm29}$ | $322_{\pm285}$ | $382_{\pm299}$ | $300_{\pm29}$ |
| CtrlFormer | $424_{\pm35}$ | $\mathbf{416_{\pm117}}$ | $\mathbf{770_{\pm71}}$ | $\mathbf{409_{\pm31}}$ |

(b) Transfer from **Finger**(`turn-easy`) to **Reacher**(`easy`)

| Method | Learn from scratch 100k | 500k | Retest after new task fine-tune |
|---|---|---|---|
| DrQ | $558_{\pm38}$ | $971_{\pm27}$ | $243_{\pm52}$ |
| Dreamer | $314_{\pm155}$ | $793_{\pm164}$ | $485_{\pm67}$ |
| Resnet+SAC | $322_{\pm285}$ | $382_{\pm299}$ | $382_{\pm299}$ |
| CtrlFormer | $\mathbf{642_{\pm42}}$ | $\mathbf{973_{\pm53}}$ | $\mathbf{906_{\pm31}}$ |

(c) **Left:** Learning old task in **Reacher** (`easy`)

| Method | Learn from scratch 100k | 500k | Learn with transfer 100k | 500k |
|---|---|---|---|---|
| DrQ | $194_{\pm84}$ | $616_{\pm274}$ | $96_{\pm43}$ | $524_{\pm68}$ |
| Dreamer | $13_{\pm32}$ | $115_{\pm98}$ | $63_{\pm07}$ | $148_{\pm12}$ |
| Resnet+SAC | $26_{\pm4}$ | $31.3_{\pm12}$ | $26_{\pm4}$ | $31_{\pm12}$ |
| CtrlFormer | $104 \pm 48$ | $\mathbf{548_{\pm131}}$ | $\mathbf{147_{\pm44}}$ | $\mathbf{657_{\pm68}}$ |

**Right:** Transfer to new task **Reacher** (`hard`)

## Transfer across multiple tasks

| Method | Task 0 → Task 1 → Task 2 → Task 3 | | | |
|---|---|---|---|---|
| Scratch (100k) | $967_{\pm27}$ | $869_{\pm61}$ | $759_{\pm48}$ | 0 |
| Train together (100k) | $433_{\pm23}$ | $143_{\pm34}$ | $310_{\pm41}$ | 0 |
| CtrlFormer (100k) | $\mathbf{967_{\pm27}}$ | $\mathbf{981_{\pm29}}$ | $\mathbf{988_{\pm36}}$ | $\mathbf{853_{\pm69}}$ |
| Scratch (500k) | $995_{\pm18}$ | $949_{\pm44}$ | $846_{\pm25}$ | $671_{\pm81}$ |
| Train together (500k) | $947_{\pm32}$ | $942_{\pm53}$ | $632_{\pm44}$ | $40_{\pm15}$ |
| CtrlFormer (500k) | $\mathbf{995_{\pm18}}$ | $\mathbf{1000_{\pm0}}$ | $\mathbf{992_{\pm26}}$ | $\mathbf{878_{\pm64}}$ |

*Table 4.* **Performance comparison with a series tasks.**

| Method | Learn from scratch 100k | 500k | Retest after new task fine-tune |
|---|---|---|---|
| DrQ | $875_{\pm76}$ | $\mathbf{973_{\pm65}}$ | $698_{\pm57}$ |
| Dreamer | $583_{\pm21}$ | $974_{\pm31}$ | $912_{\pm19}$ |
| Resnet+SAC | $177_{\pm32}$ | $190_{\pm24}$ | $190_{\pm24}$ |
| CtrlFormer | $\mathbf{877_{\pm42}}$ | $954_{\pm38}$ | $\mathbf{950_{\pm42}}$ |

(d) **Left:** Learning old task in **Walker** (`stand`)

| Method | Learn from scratch 100k | 500k | Learn with transfer 100k | 500k |
|---|---|---|---|---|
| DrQ | $504_{\pm191}$ | $\mathbf{947_{\pm101}}$ | $321_{\pm54}$ | $947_{\pm36}$ |
| Dreamer | $277_{\pm12}$ | $897_{\pm49}$ | $851_{\pm44}$ | $949_{\pm22}$ |
| Resnet+SAC | $63_{\pm7}$ | $148_{\pm12}$ | $63_{\pm7}$ | $148_{\pm12}$ |
| CtrlFormer | $\mathbf{593_{\pm52}}$ | $903_{\pm43}$ | $\mathbf{857_{\pm47}}$ | $\mathbf{959_{\pm42}}$ |

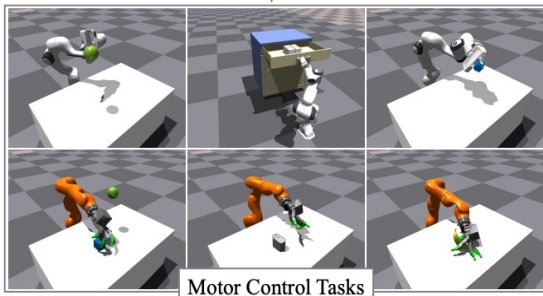**Right:** Transfer to new task **Walker** (`walk`)

# Future works

1. Pretrain the CtrlFormer with the unlabeled data from wild
2. Replace the frame stacking with better temporal modeling



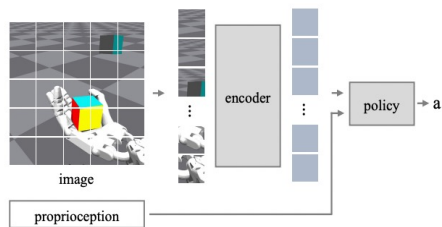Masked visual pre-training for motor control

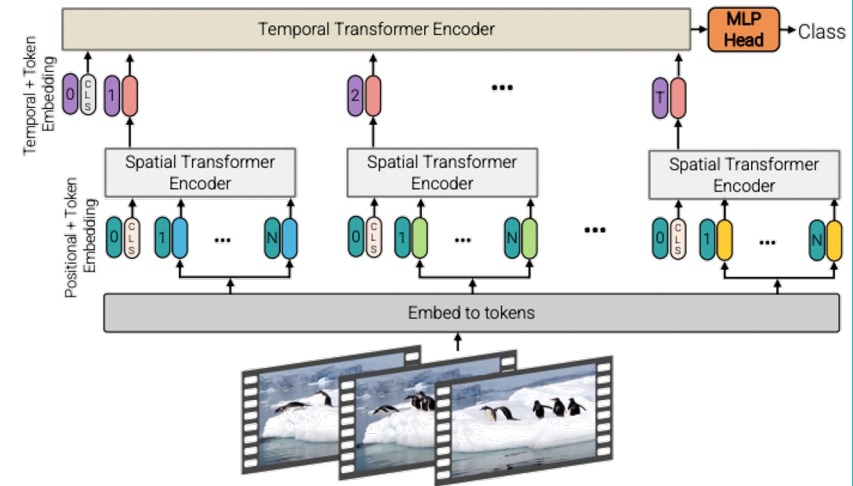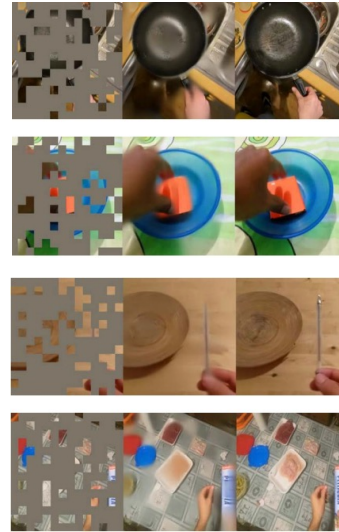

Temporal Spatial Transformer

Xiao, Tete, et al. "Masked visual pre-training for motor control." *arXiv preprint arXiv:2203.06173* (2022).

Thanks for Listening! (Q&A)

HKU
MMLAB