

# Hybrid Reinforcement Learning with Expert State Sequences

Xiaoxiao Guo, Shiyu Chang, Mo Yu,  
Gerald Tesauro, Murray Campbell

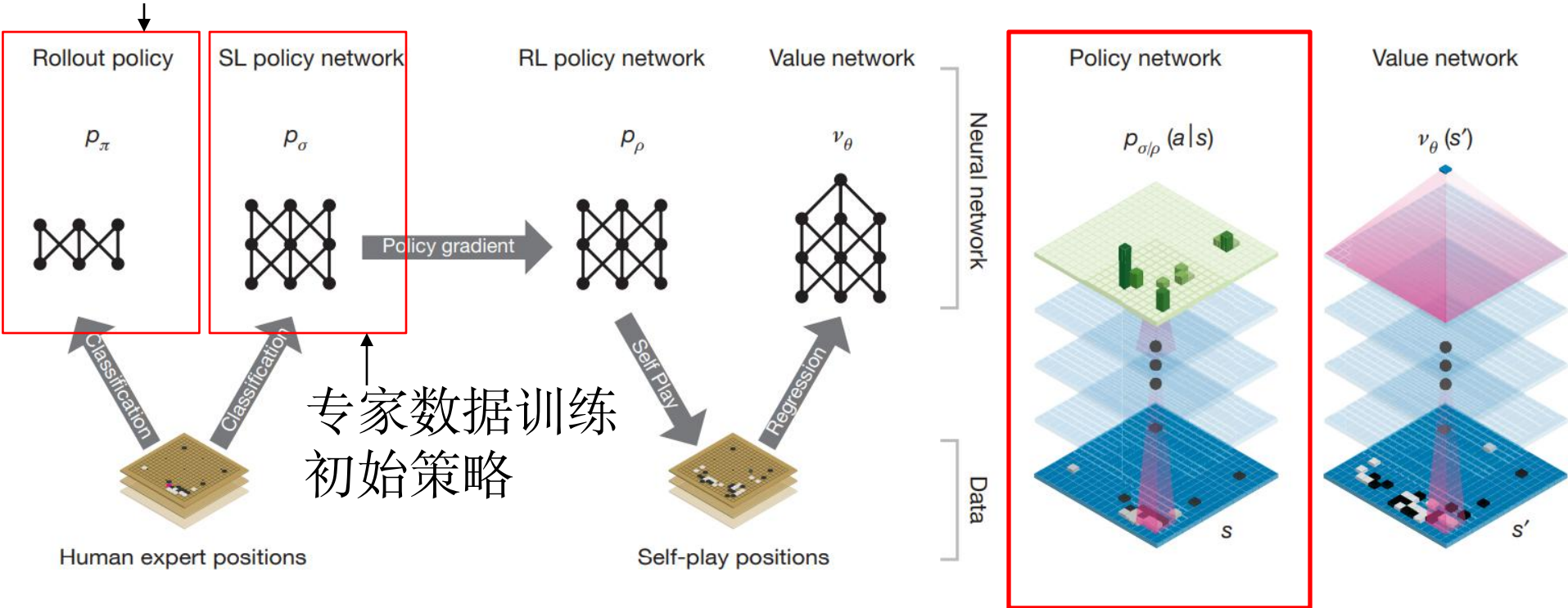
IBM Research AI

Presented by Quan He



# Alpha-Go

## 专家数据训练快速规划



专家数据训练  
初始策略

可以被视为扩展版的  
Behavior-Cloning

# Alpha-Go

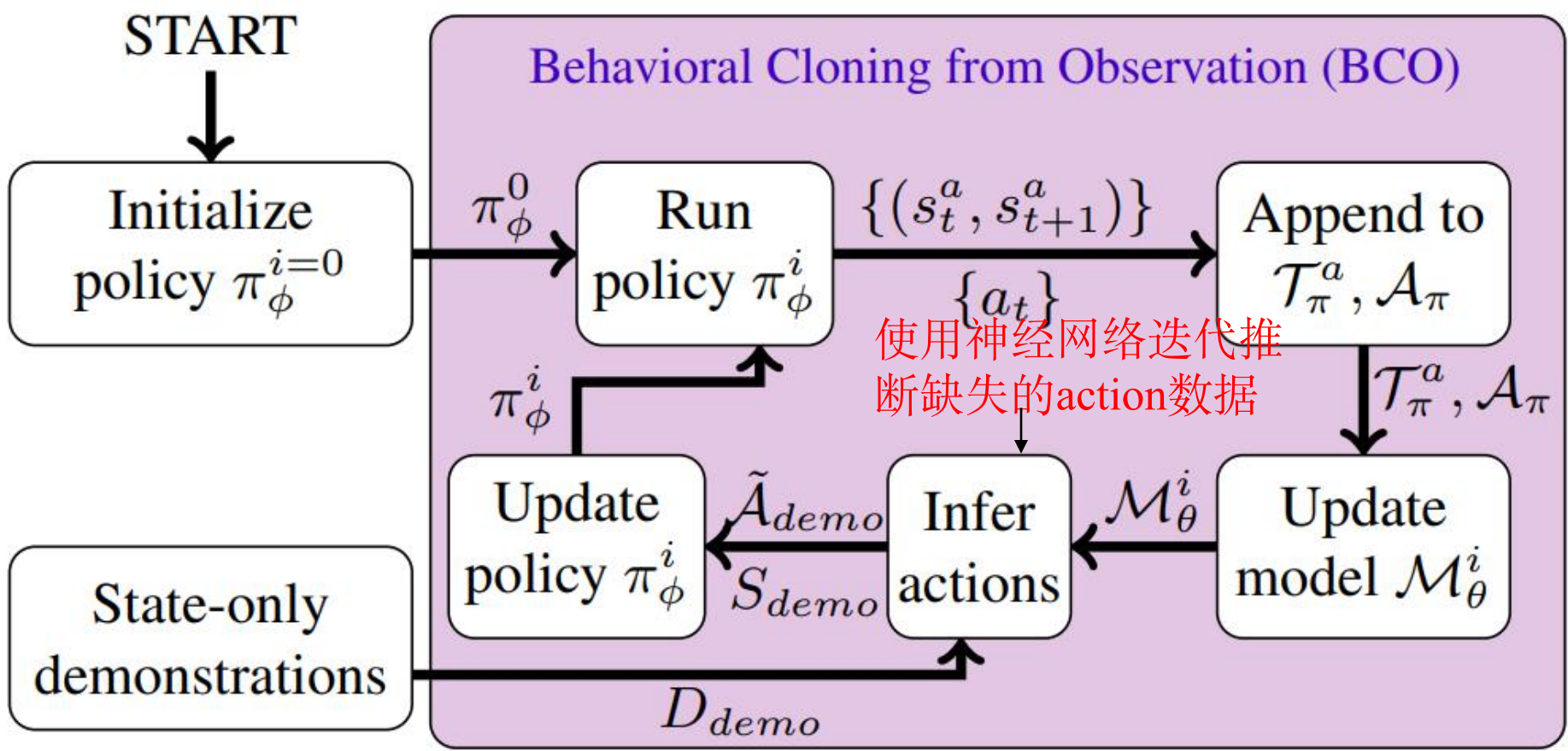
---

在Alpha-Go中，专家样本为棋谱（连续的state 轨迹），可直接推出action（下了哪步棋），从而直接计算 $P_{\sigma}(s, a)$ 进行Behavior Cloning

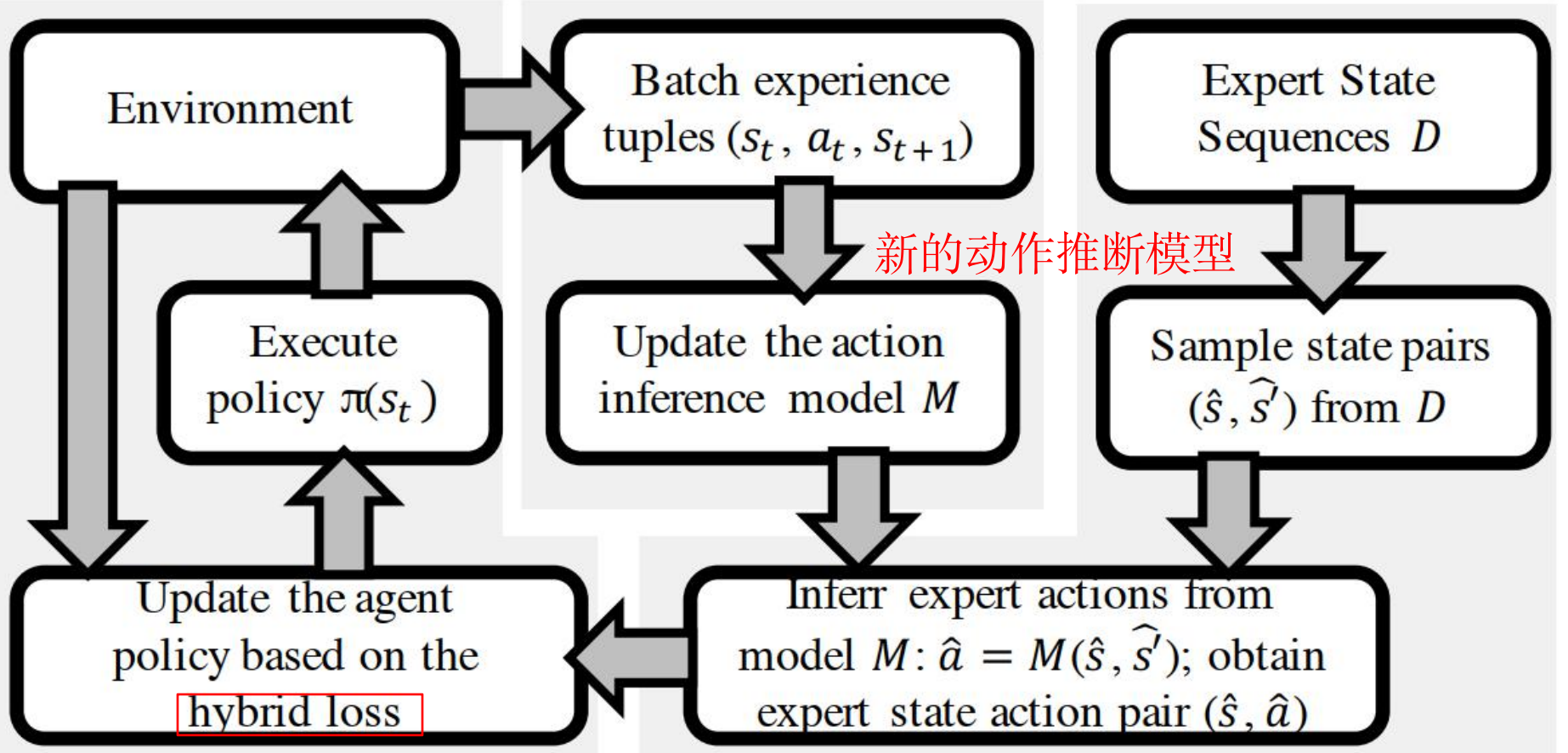
在其他环境中（例如机器人运动），专家样本通常为连续的状态，但未必能立即获得action信息——需要推断专家所使用的action

# State-only BC

## Framework



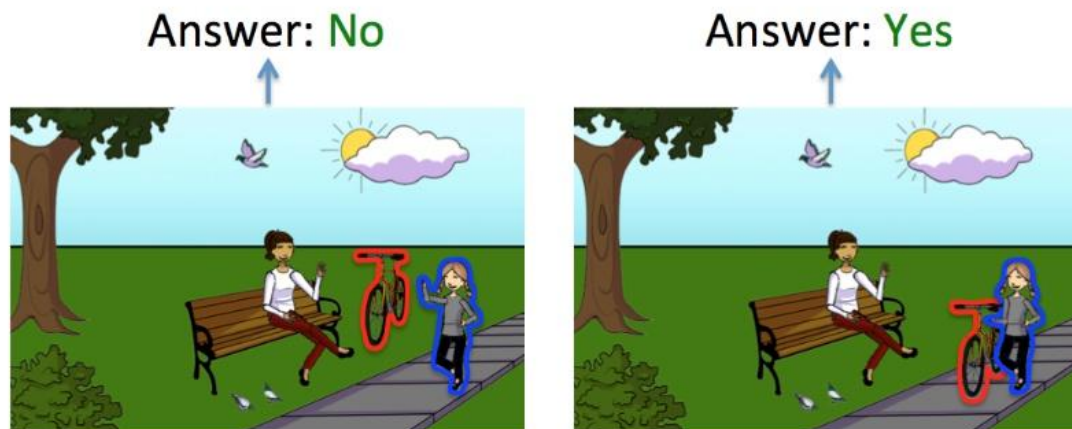
# Hybrid Framework



行为克隆和强化学习结合



# Visual Question Answering



complementary scenes


Tuple: <girl, walking, bike>

Question: Is the girl walking the bike?

已知Image, Question, 求Answer: 涉及模态融合, 三元组预测

# Visual Question Generation

---

INPUT		OUTPUT
Image	Expected answer category	Generated questions
	<b>object</b>	What is the person throwing? What is on top of the person's head?
	<b>attribute</b>	What material are those white pants made out of? What color is the frisbee?
	<b>relationship</b>	What is the person on the right doing with the frisbee? Did the person on the left throw or catch the frisbee?

已知Image,text,求question

# MUTAN: a bilinear model

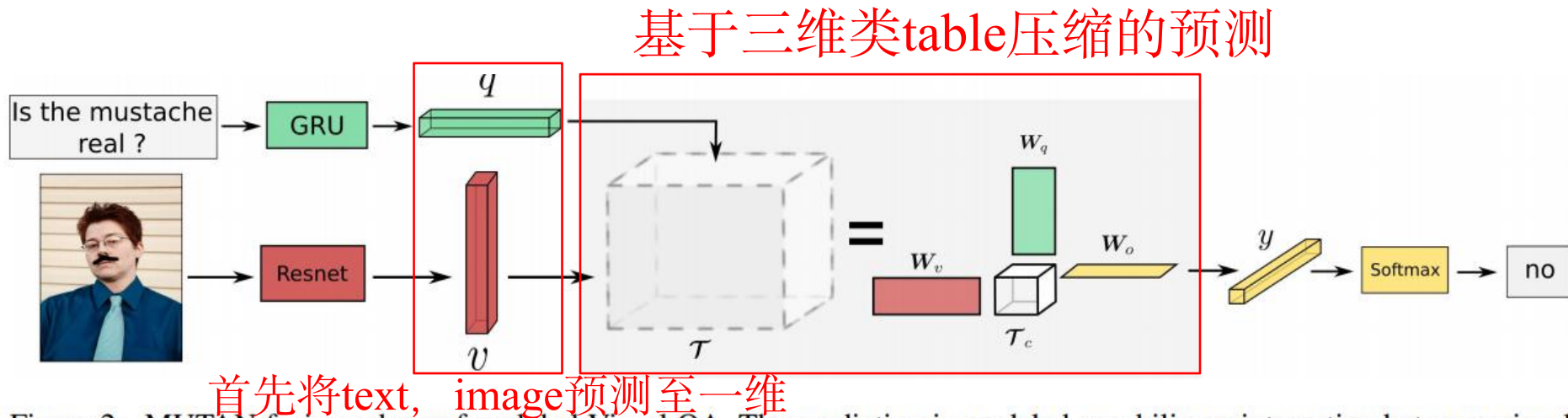


Figure 2: MUTAN fusion scheme for global Visual QA. The prediction is modeled as a bilinear interaction between visual and linguistic features, parametrized by the tensor  $\mathcal{T}$ . In MUTAN, we factorise the tensor  $\mathcal{T}$  using a Tucker decomposition, resulting in an architecture with three intra-modal matrices  $W_q$ ,  $W_v$  and  $W_o$ , and a smaller tensor  $\mathcal{T}_c$ . The complexity of  $\mathcal{T}_c$  is controlled *via* a structured sparsity constraint on the slice matrices of the tensor.

$\mathcal{T}$ : 3-way tensor  $\mathcal{T} \in \mathbb{R}^{d_q * d_v * d_a}$

即基于三维类tabular的方法预测特定answer的条件概率



# MUTAN: a bilinear model

$\times_n$ : mode-n product (张量的模式-n积)

$$(\mathbb{T} \times_n \mathbb{W})_{i_1, \dots, j, \dots, i_k} = \sum_{i_n=1}^{d_n} \mathbb{T}_{i_1, \dots, i_n, \dots, i_k} \mathbb{W}_{n,j}$$

相当于改变第n维的长度

对答案的预测  $\square$   $\mathbf{a}$ :

$$\square \mathbf{a} = (\mathbb{T} \times_1 \mathbf{q}) \times_2 \mathbf{v}$$

$\mathbf{q}$ 为GRU提取question向量,  $\mathbf{v}$ 为Resnet提取图片信息向量

例: 若 $\mathbf{1}(\mathbf{q}), \mathbf{1}(\mathbf{v})$ 为one-hot 向量(在 $\mathbf{q}, \mathbf{v}$ 处为1),  $\square$   $\mathbf{a}$ 向量为

$\mathbb{T}[\mathbf{q}, \mathbf{v}, :]$

# MUTAN: a bilinear model

---

tucker分解(高阶奇异值分解):

$$\mathbb{T} = ((\mathbb{T}_c \times_1 W_q) \times_2 W_v) \times_3 W_a$$

则

$$a = ((\mathbb{T}_c \times_1 q^T W_q) \times_2 v^T W_v) \times_3 W_a$$

令  $q^T W_q = q$ ,  $v^T W_v = v$ ,  $a^T W_a = a$ , 有

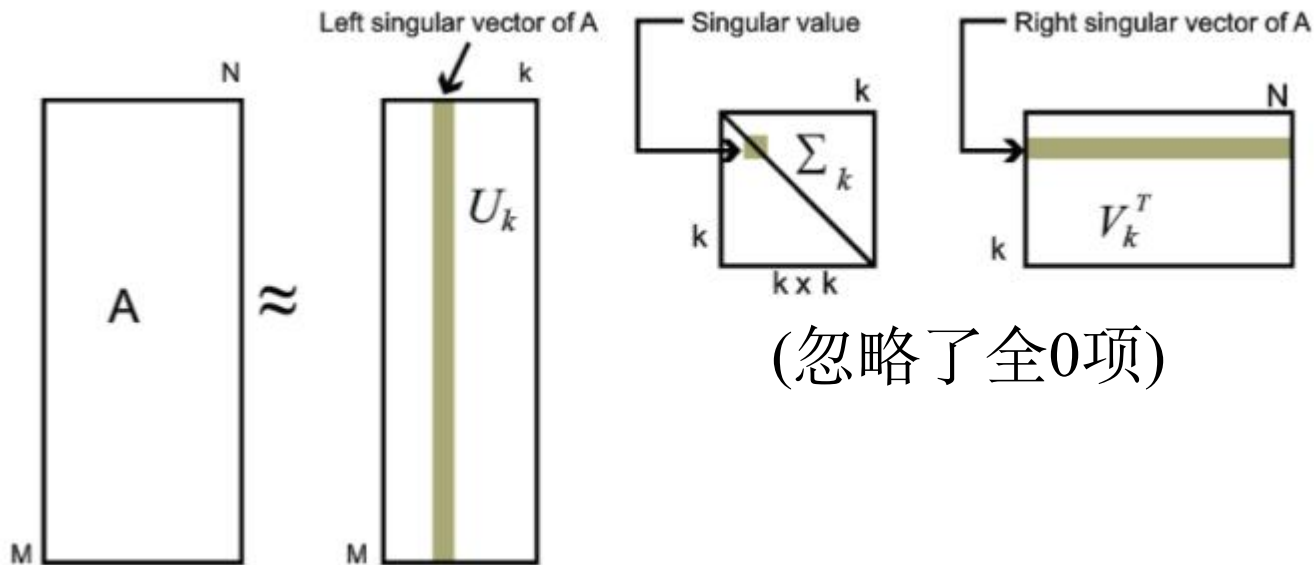
$$a = (\mathbb{T}_c \times_1 q) \times_2 v$$

$$a[k] = q^T \mathbb{T}_c[:, :, k] v \quad (1)$$

# MUTAN: a bilinear model

对  $T_c[:, :, k]$  降维(奇异值分解):

若  $T_c[:, :, k]$  为  $m \times n$  矩阵, 则  $T_c[:, :, k] = U \Sigma V^T$ , 其中  $U$  为  $m \times m$  矩阵,  $V$  为  $n \times n$  矩阵,  $\Sigma$  为  $m \times n$  对角矩阵. Rank 与  $T_c[:, :, k]$  相同



Ben-Younes, et al. MUTAN: Multimodal Tucker Fusion for Visual Question Answering, ICCV 2017

# MUTAN: a bilinear model

对  $\mathcal{T}_c[:, :, k]$  取前  $R$  个特征值, 有  $\mathcal{T}_c[:, :, k]$  为  $R$  个列向量与  $R$  个行向量的外积之和:

$$\mathcal{T}_c[:, :, k] = \sum_{r=1}^R M_r[:, k] \otimes N_r[:, k]^T$$

代入前(1)式得 双线性 模型 MUTAN:

$$a = \sum_{r=1}^R (q^T M_r) * (v^T N_r)$$

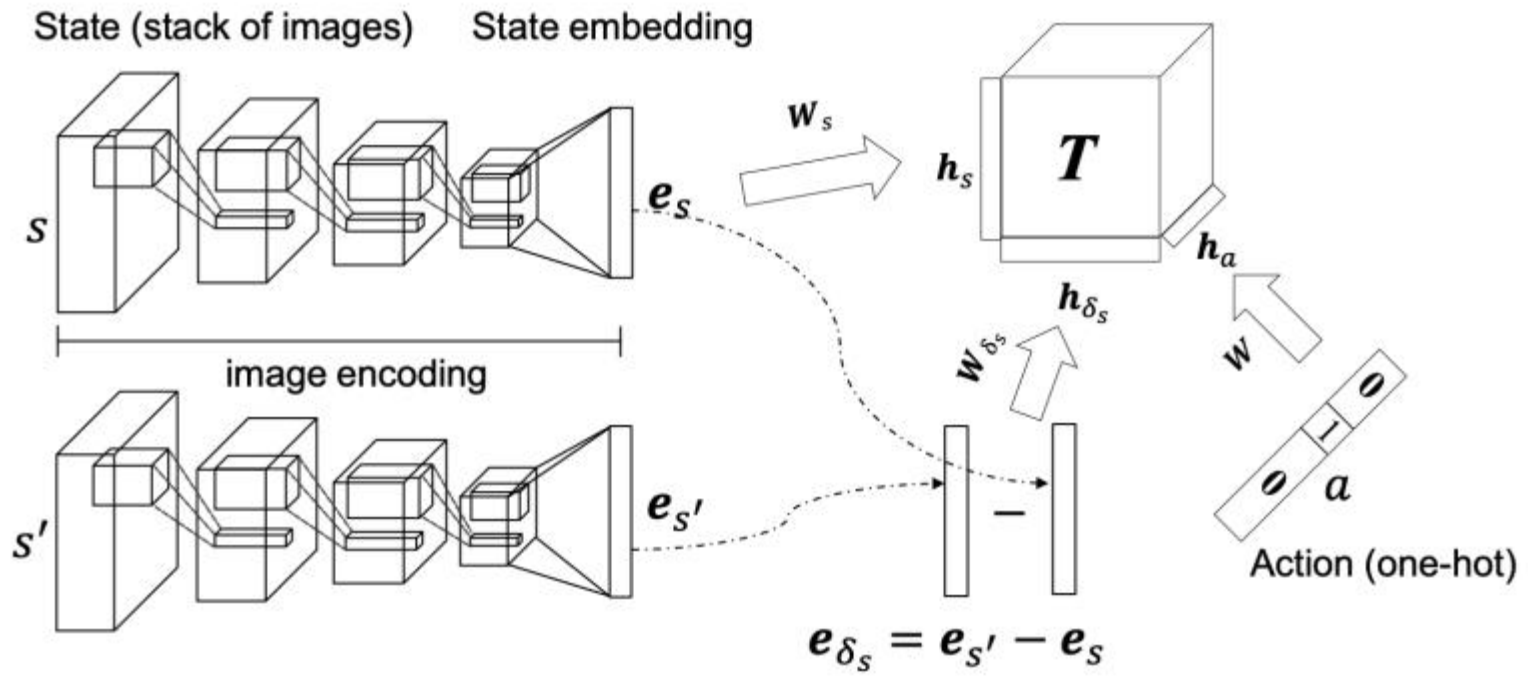
其中 \* 为逐元素点乘

一种理解:  $a_r[k] = (q \text{ similar to } M_r[:, k]) \text{ AND } (v \text{ similar to } N_r[:, k])$

Ben-Younes, et al. MUTAN: Multinomial Tucker Fusion for Visual Question Answering, ICCV 2017

$$a[k] = a_1[k] \text{ OR } \dots \text{ OR } a_R[k]$$

# tensor-based inference model



用  $\delta_s$  代替  $s'$  后，使用类MUTAN模型进行预测



# tensor-based inference model

$$\begin{aligned} \square_{\delta_s} &= \mathbb{T} \times_1 h_s \times_2 h_a \\ h_a &= \mathbb{T} \times_1 h_s \times_3 h_{\delta_s} \end{aligned}$$

使用双线性模型:

$$\square_a = \sum_{r=1}^R (h_s M_r) * (h_{\delta_s} N_r)$$

假设在相同的state下, embedding所得 $h_a$ 与 $h_{\delta_s}$ 有相似性, 引入对偶性简化 (借鉴了VQA与VQG对偶的思路):

$$\begin{aligned} \mathbb{T}[:, :, k] &= \mathbb{T}[:, k, :], \quad d_{\delta_s} = d_a = d \\ \square_{\delta_s} &= \sum_{r=1}^R (h_s M_r) * (h_a N_r) \end{aligned}$$

# tensor-based inference model

动作推断:  $P^i(\cdot|s, s') = \text{SoftMax}(W_a h_a + b_a)$

状态推断:  $P^i(\cdot|s, a) = \text{SoftMax}((W_{\delta_s} h_{\delta_s} + W_s h_s) + b_{s'})$

推断模型学习目标:

$$\mathcal{L}^{\text{dual-model}} = \mathbb{E}_{(s, a, s')} \left[ -\log P^f(s'|s, a) - \log P^i(a|s, s') + \|\mathbf{h}_a - \hat{\mathbf{h}}_a\|_1 + \|\mathbf{h}_{\delta_s} - \hat{\mathbf{h}}_{\delta_s}\|_1 \right]$$

推断结果:

$$\mathcal{M}(\hat{s}, \hat{s}') = \arg \max_a P^i(a|\hat{s}, \hat{s}')$$

# Hybrid BC and RL

---

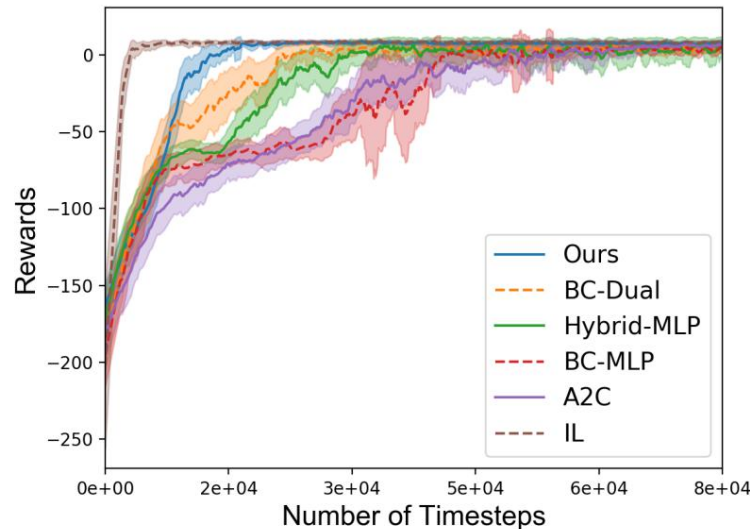
与RL (A2C) 结合的策略训练目标:

$$\begin{aligned} \mathcal{U}^{\text{hybrid}}(\theta) = & \mathbb{E}_{s,a} \left[ A(s) \log \pi(a|s; \theta) + \alpha \mathcal{H}(\pi(.|s)) \right] \\ & + \mathbb{E}_{(\hat{s}, \hat{s}') \sim \rho(\mathcal{D})} \left[ \log \pi(\mathcal{M}(\hat{s}, \hat{s}') | \hat{s}; \theta) \right] \end{aligned}$$

其中A(s)为优势值函数

$$A(s_t) = \sum_{n=1}^N \gamma^{n-1} r_{t+n} + \gamma^N V(s_{t+N}) - V(s_t)$$

# Experiments on Taxi Domain



BC-Dual: 仅做模仿学习，使用本文 **tensor-based** 推断模型

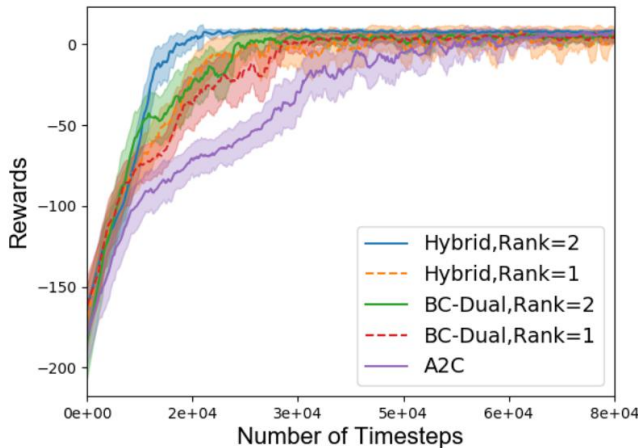
Hybrid-MLP: 结合 BC 与 RL，使用神经网络推断模型

BC-MLP: 仅做模仿学习，使用神经网络推断模型

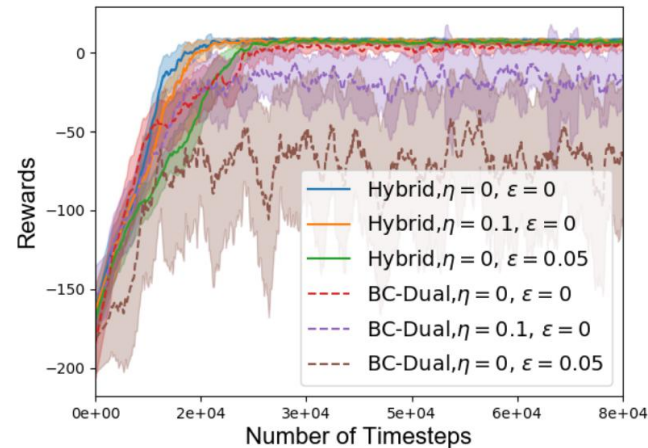
A2C: 纯强化学习

IL: 有 action 信息的模仿学习（相当于上界）

# Experiments on Taxi Domain



(b)



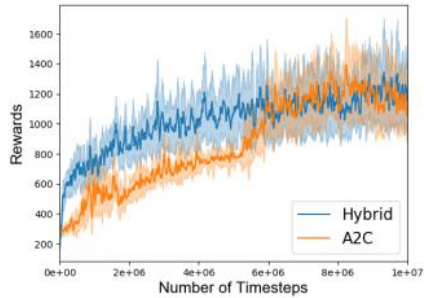
(c)

左图：推断模型rank参数的影响

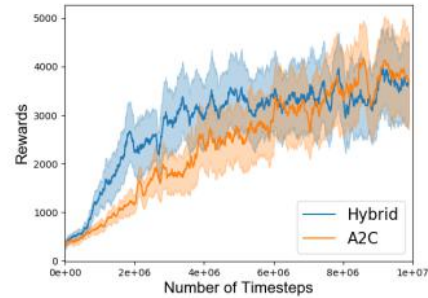
右图：引入RL对模型robust能力的影响（ $\eta$ 表示专家样本在state上的缺失率， $\epsilon$ 表示专家样本的错误率）



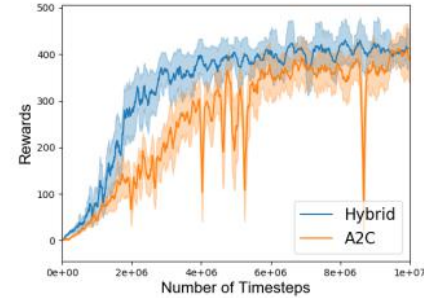
# Experiments on Atari



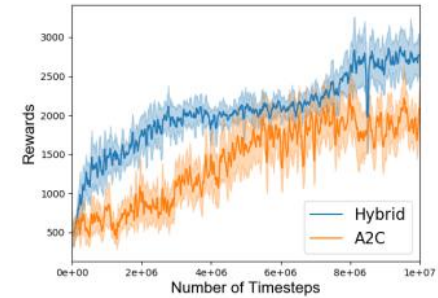
(a) Alien



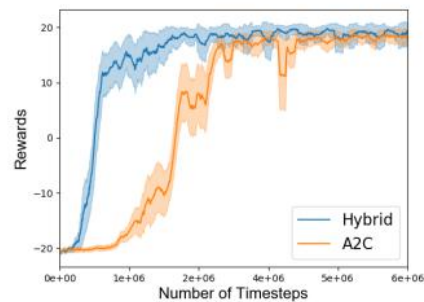
(b) BeamRider



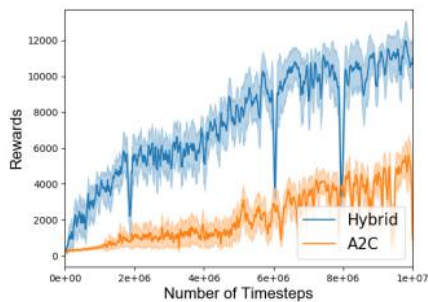
(c) Breakout



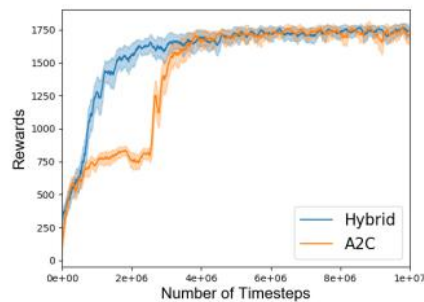
(d) MsPacman



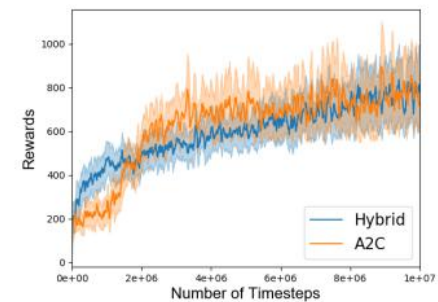
(e) Pong



(f) Qbert



(g) Seaquest



(h) SpaceInvaders

**Thanks!**

