

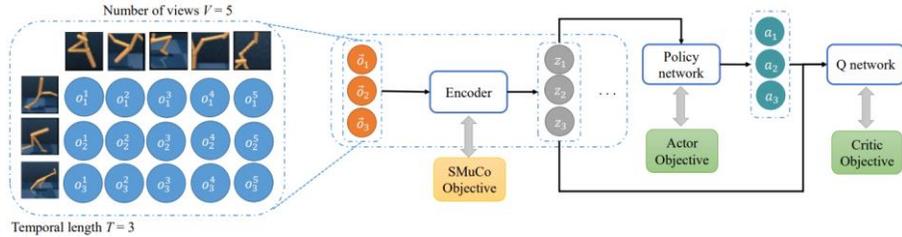
# SMuCo: Reinforcement Learning for Visual Control via Sequential Multi-view Total Correlation

Tong Cheng<sup>1</sup>, Hang Dong<sup>2</sup>, Lu Wang<sup>2</sup>, Bo Qiao<sup>2</sup>, Qingwei Lin<sup>2</sup>, Saravan Rajmohan<sup>3</sup>, Thomas Moscibroda<sup>4</sup>

Nanyang Technological University<sup>1</sup>, Microsoft AI<sup>2</sup>, Microsoft 365<sup>3</sup>, Microsoft Azure<sup>4</sup>

## Reinforcement Learning for Visual Control

- Importance: leveraging multi-view observation could enhance performance of agent for visual control task
- Challenge: 1) compress task-relevant information; 2) remove task-irrelevant information
- Solution: maximizing sequential multi-view total correlation



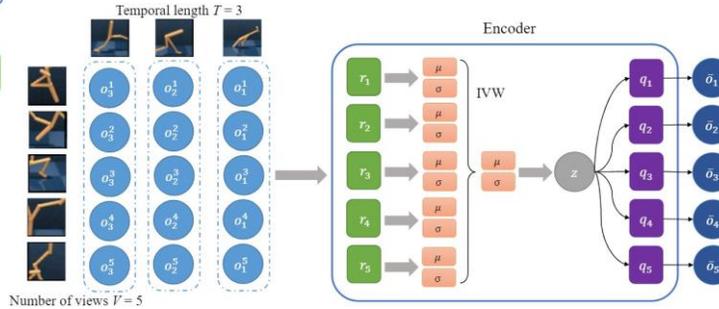
- Inverse Variance Weighted method (IVW)

$$\mathcal{L} = \mathcal{L}_{\text{REC}} + \mathcal{L}_{\text{LL}} + \mathcal{L}_{\text{TC}}$$

$$\mathcal{L}_{\text{REC}} = - \sum_{v=1}^V \sum_{t=1}^T H(O_t^v | Z_{t-1}, A_{t-1}),$$

$$\mathcal{L}_{\text{LL}} = - \sum_{v=1}^V \sum_{t=1}^T \mathbb{E}_{p_1} \ln q_{\psi}^v(o_t^v | z_t, z_{t-1}, a_{t-1}),$$

$$\mathcal{L}_{\text{TC}} = \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}_{p_2} [D_{\text{KL}}(p(z_t | o_s, \iota) \| r_{\phi}(z_t | \iota))],$$



## Sequential Multi-view Total Correlation

- Extend Multi-view Total Correlation into SMTC
- Derive Lower bound of SMTC

$$\text{SMTC}(\vec{O}_{1:T}; Z_{1:T} | A_{1:T}) = \sum_{v=1}^V I(O_{1:T}^v; Z_{1:T} | A_{1:T}) - I(\vec{O}_{1:T}; Z_{1:T} | A_{1:T})$$

$$\text{SMTC}(\vec{O}_{1:T}; Z_{1:T} | A_{1:T}) \geq \sum_{v=1}^V \sum_{t=1}^T [H(O_t^v | Z_{t-1}, A_{t-1}) + \mathbb{E}_{p(z_t, o_t^v | z_{t-1}, a_{t-1})}$$

$$\ln q_{\psi}^v(o_t^v | z_t, z_{t-1}, a_{t-1})] - \sum_{t=1}^T \sum_{s=1}^T \mathbb{E}_{p(\vec{o}_s)} [D_{\text{KL}}(p(z_t | o_s, \iota) \| r_{\phi}(z_t | \iota))].$$

## Training Encoder

- Formulate surrogate function
- Product of Expert (PoE)

## Experiment Results

Scores at 500k Steps	DrQ	RAD	DreamerV2	PI-SAC	SLAC	DRIBO	SMuCo
Cheetah, run	797 ± 116	880 ± 104	841 ± 57	802 ± 119	881 ± 116	864 ± 52	<b>1019 ± 107</b>
Walker, walk	930 ± 46	858 ± 82	966 ± 117	959 ± 103	930 ± 107	881 ± 90	<b>1036 ± 30</b>
Ball in cup, catch	958 ± 102	-9 ± 80	955 ± 82	963 ± 75	983 ± 98	<b>1006 ± 39</b>	853 ± 115
Finger, spin	738 ± 79	880 ± 32	366 ± 105	787 ± 112	947 ± 58	960 ± 53	<b>969 ± 52</b>
Acrobot, swingup	228 ± 50	163 ± 34	209 ± 106	246 ± 85	<b>256 ± 45</b>	242 ± 67	247 ± 51
Humanoid, run	470 ± 60	375 ± 117	436 ± 104	482 ± 37	453 ± 117	497 ± 40	<b>507 ± 105</b>
Hopper, hop	454 ± 89	357 ± 44	347 ± 52	431 ± 67	453 ± 41	<b>488 ± 38</b>	484 ± 33
Fish, swim	729 ± 90	546 ± 114	697 ± 109	694 ± 102	750 ± 101	<b>787 ± 38</b>	748 ± 114
Basic Manipulation (uh)	130 ± 30	310 ± 34	108 ± 39	168 ± 45	246 ± 34	84 ± 42	<b>377 ± 45</b>
Basic Manipulation (ud)	59 ± 34	<b>366 ± 32</b>	41 ± 38	198 ± 38	242 ± 42	93 ± 43	221 ± 41
Basic Manipulation (hd)	68 ± 45	105 ± 31	78 ± 33	173 ± 32	177 ± 33	70 ± 48	<b>358 ± 42</b>
Basic Manipulation (uhd)	86 ± 47	368 ± 36	62 ± 40	47 ± 33	136 ± 30	101 ± 38	<b>446 ± 34</b>

