

Input

State:



+

History:

$[s_{t-T+1}, s_{t-T+2}, \dots, s_t]$

+

Prompt

Contrastive Learning

Negative Contrast

(-)

Positive Contrast

(+)

Hierarchical Predictive Contrastive Correction

Trajectory-Level

z^{traj}

Align

Sinkhorn

Prompt

Trajectory

z^{traj}

Subgoal-Level

Predict

$MLP_{f_{\theta}^{sub}}$

Compare

\hat{z}^{traj}

Cross-Level Correction

z^{sub}

Align

Score-field

Prompt

Trajectory

Action-Level

Predict

$MLP_{f_{\theta}^{act}}$

Compare

\hat{z}^{sub}

Cross-Level Correction

z^{act}

Align

Score-field

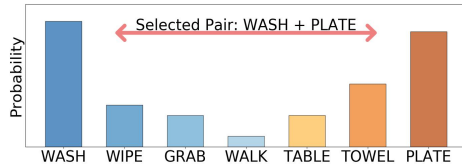
Prompt

Trajectory

Total Loss

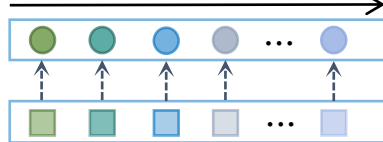
$emb_{grads} = \partial L / \partial \theta$

Execution Network



MLP

Execution Gradient ($\nabla_{\theta\pi}$)



Embedding Gradient (∇_z)

Transformer

Input:

$State_{seq}[T_{seq}, image_{emb} | text_{emb}]$