

---

# Spatiotemporal Predictive Pre-training for Robotic Motor Control

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Robotic motor control necessitates the ability to predict the dynamics of envi-  
2        ronments and interaction objects. However, advanced self-supervised pre-trained  
3        visual representations (PVRs) in robotic motor control, leveraging large-scale  
4        egocentric videos, often focus solely on learning the static content features of  
5        sampled image frames. This neglects the crucial temporal motion clues in human  
6        video data, which implicitly contain key knowledge about sequential interacting  
7        and manipulating with the environments and objects. In this paper, we present  
8        a simple yet effective robotic motor control visual pre-training framework that  
9        jointly performs spatiotemporal prediction with dual decoders, utilizing large-scale  
10       video data, termed as **STP**. STP adheres to two key designs in a multi-task learning  
11       manner. First, we perform spatial prediction on the masked current frame for  
12       learning content features. Second, we utilize the future frame with an extremely  
13       high masking ratio as a condition, based on the masked current frame, to conduct  
14       temporal prediction of future frame for capturing motion features. This asymmetric  
15       masking and decoder architecture design is very efficient, ensuring that our  
16       representation focusing on motion information while capturing spatial details. We  
17       carry out the largest-scale BC evaluation of PVRs for robotic motor control to date,  
18       which encompasses 21 tasks within a real-world Franka robot arm and 5 simulated  
19       environments. Extensive experiments demonstrate the effectiveness of STP as well  
20       as unleash its generality and data efficiency by further post-pre-training and hybrid  
21       pre-training. Our code and weights will be released for further applications.

## 22 1 Introduction

23       In NLP and CV, adapting pre-trained foundation models from large-scale data to various downstream  
24       tasks has seen great success. For example, pre-trained visual representations using self-supervised [38,  
25       15, 67, 2, 93] or weakly-supervised [71, 25, 55] methods exhibit strong generalization ability for  
26       visual understanding. However, in robot learning, due to data scarcity and homogeneity, some  
27       groundbreaking methods [53, 1] resort to training from scratch only using domain-specific data.  
28       Recently, inspired by the success of transfer learning in CV, many works [69, 73, 65, 58, 59, 19] have  
29       explored developing a pre-trained visual representation (PVR) using large-scale out-of-domain data  
30       for various robotic motor control tasks. Currently, one successful paradigm [73, 99, 59, 19] is to use  
31       large-scale egocentric video datasets [29] and train vanilla vision transformers (ViT) [22] based on  
32       MAE [38], which exhibits excellent learning efficiency and generalization ability for learning policy  
33       from raw pixel. Among them, the Ego4D [29] dataset offers numerous first-person human-object  
34       interaction scenes and good motion clues. We argue that although learning static spatial structure  
35       priors from task-relevant pre-training data sources is crucial, designing a more relevant self-supervised  
36       proxy task for motor control should not be overlooked. Therefore, in this paper, we aim to develop a  
37       more relevant self-supervised proxy task for robotic motor control representation learning.

38 Robotic motor control typically requires fine-grained spatial localization and relatively dense se-  
39 mantics. With its ability to effectively capture low-level geometry and space structure, MAE [38]  
40 pre-training excels at this task. However, is dense spatial content sufficient for robotic motor control?  
41 Some neuroscientific studies [50, 21, 88] suggest the brain’s different areas or cells show special-  
42 ization. Some are dedicated to processing the information of temporal object motion, while others  
43 focus on static spatial details. Their combination results in subjective pattern perception. Inspired by  
44 this finding, we hypothesize that an effective robotic motor control pre-training proxy task should  
45 require joint learning of spatial content features and temporal motion features. However, current  
46 methods [73, 59, 19] use MAE pre-training with image frames from human videos, capturing only  
47 static content features. They overlook the temporal motion clues in human videos, which implicitly  
48 contain key knowledge about sequential interaction with environment and manipulation of objects.  
49 Therefore, we aim to bridge this gap by incorporating these motion clues into our proxy task.

50 Based on the analysis above, the most critical challenge is the absence of action annotations in human  
51 video data for modeling object motion. To model interaction and manipulation actions from actionless  
52 video data, we propose to implicitly capture them by predicting future frame pixels based on current  
53 frame. However, predicting the future frame without any conditions could contain high uncertainty  
54 and be extremely difficult. Therefore, we propose to use the future frame with an extremely high  
55 masking ratio as a prompt condition, specifically 95%, which serves to reveal some behavior and  
56 dynamic priors, i.e. what to do and how to do it. In the experiments section, we will further explore  
57 different condition alternatives, including language narration and their combination. Additionally,  
58 directly and simply executing temporal prediction could lead the model to overlook static spatial  
59 details, and it is also not efficient enough. Therefore, another technical contribution of STP is to jointly  
60 perform spatial prediction by masking the current frame with 75% masking ratio. In summary, we  
61 present **STP**, a multi-task self-supervised pre-training framework through spatiotemporal predictive  
62 learning. Our STP asymmetrically mask the current frame and future frame from a video clip, using  
63 a spatial decoder to conduct spatial prediction for content learning and a temporal decoder to conduct  
64 temporal prediction for motion learning. This asymmetric masking and decoder architecture design  
65 ensures that our pre-trained encoder focusing on motion information while capturing spatial details.

66 Subsequently, we establish our evaluation scheme. Currently, how to adapt pre-trained visual  
67 representations for robotic motor control still remain an open question. Considering the expensive  
68 cost of robot data collection or exploration, we employ a data-efficient paradigm of few-shot behavior  
69 cloning by learning from demonstrations (Lfd). To demonstrate the generalization ability of visual  
70 representation, our primary evaluation scheme involves freezing the visual encoder during policy  
71 training. Additionally, considering that fine-tuning ViT with few demonstrations might lead to  
72 overfitting and masked modeling exhibits excellent data efficiency [86, 102, 52] in domain-in data,  
73 we further follow the post-pre-training [7, 93, 59] paradigm to perform STP pre-training with task-  
74 specific data to achieve better results. It is noteworthy that different tasks do not share representation  
75 in this setting. Finally, we conduct the largest-scale BC evaluation of PVRs for robotic motor control  
76 to date to demonstrate the effectiveness of STP, which encompasses 21 tasks ( 2 real-world tasks and  
77 19 simulation tasks across 5 environments). These simulation tasks are derived from the union of  
78 manipulation and locomotion tasks from prior works [65, 59].

79 We make the following **four contributions**: (1) We present STP, a *self-supervised* visual pre-  
80 training framework for robotic motor control, which jointly conducts spatiotemporal prediction with  
81 *asymmetric masking and decoder architecture design* for content and motion features learning. (2)  
82 We further expand STP by performing hybrid pre-training with ImageNet-MAE and post-pre-training  
83 with task-specific data, unleashing its *generality* and *data efficiency*. (3) To our best knowledge, we  
84 conduct the *largest-scale BC evaluation* of PVRs for robotic motor control to date to demonstrate the  
85 effectiveness of STP. (4) Our experiments yield some insightful observations. In temporal prediction,  
86 language does not significantly enhance performance. Instead, *single-modality self-supervised*  
87 *paradigm* achieves the best results. This finding is highly encouraging for self-supervised robotic  
88 motor control representation learning. Moreover, in the few-shot BC setting, naively scaling up model  
89 size does not necessarily lead to improved outcomes. Finally, incorporating *more diverse* data and  
90 *domain-in* data into the pre-training can further enhance performance.

## 91 2 Related Work

92 **Pre-trained Visual Representation Learning.** Large-scale visual representation pre-training are  
93 continually empowering computer vision. The primary supervised learning methods include learning

94 image recognition [40, 87] from ImageNet [20] and learning multi-modal alignment [71] from image-  
 95 text pairs. Currently, self-supervised learning methods are enjoying significant popularity, primarily  
 96 falling into two main categories. The first category utilizes contrastive learning [39, 15, 14] technique  
 97 or joint-embedding architecture [13] to learn view-invariance. The second category performs masked  
 98 modeling [7, 38, 100, 95, 4, 2] and predict the pixel or representation of invisible parts in space. In  
 99 addition, some methods [106, 67, 8] have also proposed to combine different optimization objectives  
 100 in a multi-task learning manner. Recently pre-trained visual representation learning for robotic motor  
 101 control have been rapidly developing [69, 65, 73, 99, 58, 57, 46, 59, 19]. These methods cover different  
 102 backbones (ResNet [40], ViT [22]), different policy learning methods (reinforcement learning [99],  
 103 behavior cloning [69, 65, 59], reward function [58] and task specification [42]), different adaptation  
 104 schemes (linear probing [69, 65, 46, 59], fine-tuning [19] and designing adapters [78, 56]), and  
 105 different evaluation environments (diverse simulation benchmarks). At present, it is still unclear how  
 106 these factors collectively influence the performance. In this paper, we choose scalable vanilla vision  
 107 transformer [22] as our backbone and data-efficient few-shot behavior cloning paradigm to conduct  
 108 policy learning, while ensuring the backbone is frozen during policy training.

109 **Temporal Predictive Learning.** Early works once explored representation learning through future  
 110 prediction, encompassing image [61], video [35, 80] and audio [66]. VideoMAE [86, 93] extend  
 111 MAE [38] to 3D video architecture. Recently TrackMAE [17] and SiamMAE [33] predict the  
 112 masked future frame based on unmasked current frame, leading to a better capture of temporal  
 113 correspondence and achieving outstanding performance in object tracking and segmentation tasks. In  
 114 robot learning, predicting future visual states primarily serves as a transition dynamic model such as  
 115 World Models [62, 77] and Dreamer [76]. [85, 9] predict the future visual states using goal image in  
 116 robot data. GR-1 [97] conducts language-conditioned video prediction for policy model pre-training  
 117 in a frozen visual representation space. [96] proposed dynamics-aware representation learning,  
 118 and [82, 72] employed forward dynamics for self-supervised pre-training. Some works explored to  
 119 train video prediction models and utilize visual foresight [32], inverse dynamics models [18],  
 120 goal-conditioned policy learning [23], and geometry estimation [51] methods for motor control,  
 121 respectively. [92] fine-tuned pre-trained representations into dynamic and functional distance  
 122 modules for manipulation tasks. Unlike these works, we utilize the public large-scale egocentric  
 123 video data and employ masked spatiotemporal predictive learning as a *self-supervised proxy task*  
 124 (*without any language or action annotations*) for robotic motor control representation learning,  
 125 instead of designing elaborate architectures or methods for specific predictive tasks [28, 37].

126 **Vision-based Robot Learning.** Vision-based robot learning plays a crucial role in robotics com-  
 127 munity. Recently some related works focus on studying model architectures [44, 12, 47], observa-  
 128 tion spaces [107], downstream policy learning methods [41], sim-to-real transfer [79], designing  
 129 adapters [78, 56], learning-from-scratch baseline [36], and affordance model [6, 105, 45, 60], in  
 130 visuo-motor representation learning. Other related works [70, 5, 91, 101, 48] attempt to learn ma-  
 131 nipulation skills from small-scale and in-domain human videos. In addition, language-conditioned  
 132 vision robot learning has received significant attention. Some works scale multimodal robotic  
 133 data [42, 11, 34, 90, 24, 68, 84] or introduce Internet data and knowledge [81, 103, 10, 54, 43, 94, 64]  
 134 for end-to-end robot learning. In our study, we pre-train a off-the-shelf visual representation from  
 135 large-scale egocentric video datasets for robotic motor control tasks. Our method is more simple and  
 136 general for different downstream tasks of motor control.

### 137 3 Method

138 In this section, we describe our method in details. First, we give an overview of our spatiotemporal  
 139 predictive pretraining (STP) framework. Then, we give a technical description on our core components  
 140 during pre-training: the masked image encoder and dual decoders scheme. Finally, we describe how  
 141 to adapt our pre-trained encoder to downstream robotic motor control tasks.

#### 142 3.1 Overview of STP

143 As illustrated in Figure 1, our STP aims to pre-train an image encoder for robotic motor control from  
 144 video datasets. This pre-trained image encoder is subsequently frozen and directly transferred to solve  
 145 motor control tasks. Specifically, given a video dataset  $\mathcal{D}$ , our goal is to learn an image encoder  $\Phi_{enc}$ ,  
 146 that maps images to the visual representations. During pre-training and post-pre-training,  $\mathcal{D}$  represents  
 147 large-scale out-of-domain videos and task-specific demonstration videos, respectively. After pre-  
 148 training, we reuse  $\Phi_{enc}$  for downstream motor control policy learning. Specifically, the downstream

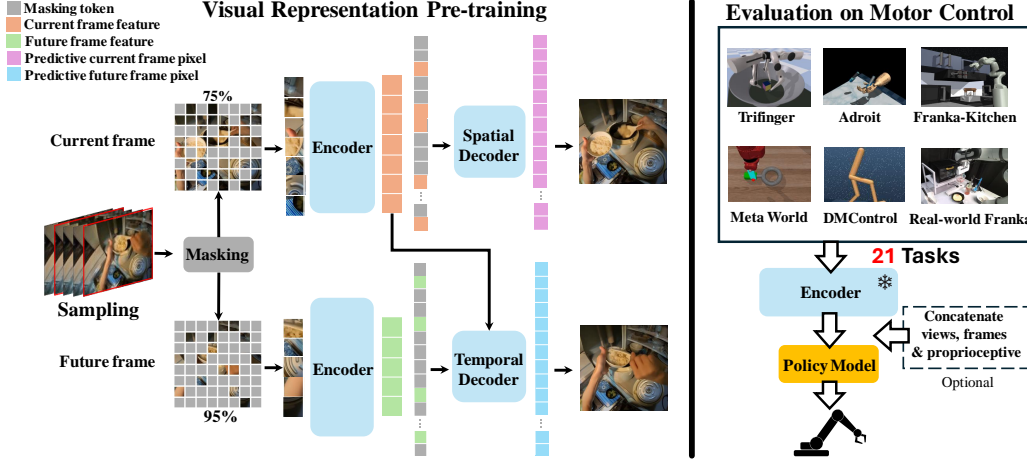


Figure 1: **STP framework.** **Left:** During pre-training, we sample the current frame and the future frame from the video clip, and carry out spatiotemporal predictive pre-training. **Right:** During motor control tasks evaluation, we freeze the pre-trained encoder to extract visual state representations and discard the decoders.

149 task will require an agent to make sequential action decisions based on visual observations  $\mathcal{O}$ . Instead  
 150 of using the raw observation images as direct input like end-to-end policy learning from pixel, the  
 151 agent will employ the pre-trained  $\Phi_{enc}$  to extract its visual state representation  $\Phi_{enc}(\mathcal{O})$  for the  
 152 subsequent policy learning module.

### 153 3.2 Masked Image Encoder

154 We first introduce the pipeline of our image encoder. Our image encoder processes image frames  
 155 using a vanilla vision transformer [22]. Given an image  $\mathbf{I} \in \mathbb{R}^{C \times H \times W}$ , we initially process it by the  
 156 patch embedding layer to obtain its token sequences  $\mathbf{T}$ , where  $\mathbf{T} = \{P_i\}_{i=1}^N$  and  $N$  is the total  
 157 token number, (e.g.,  $N = 196$  for a  $224 \times 224$  image with a patch size of  $16 \times 16$ ). Then we add the  
 158 fixed 2D sine-cosine positional embeddings for all tokens. Following this, we mask and remove a  
 159 part of tokens, according to a randomly generated masking map  $\mathbb{M}(\rho)$ , where  $\rho$  is the masking ratio.  
 160 The encoder applies several transformer blocks (consisting of a global self-attention layer and a FFN  
 161 layer) on all unmasked tokens:  $\mathbf{Z} = \Phi_{enc}(\mathbf{T}^u)$ , where  $\mathbf{T}^u = \{T_i\}_{i \in (1 - \mathbb{M}(\rho))}$ . During this process, a  
 162 [CLS] token is added at the beginning.

163 Then we describe our encoding process during pre-training. We randomly sample two frames from a  
 164 video clip based on an interval: the current frame  $\mathbf{I}_c$  and the future frame  $\mathbf{I}_f$ . Following the above  
 165 pipeline, we randomly generate two asymmetric masking maps for the current frame and the future  
 166 frame, denoted as  $\mathbb{M}_c = \mathcal{M}_c(\rho^c)$  and  $\mathbb{M}_f = \mathcal{M}_f(\rho^f)$ , respectively. Each of these maps has a  
 167 different masking ratio. We then use these maps to separately process the two frames and obtain their  
 168 features,  $\mathbf{Z}_c$  and  $\mathbf{Z}_f$ . As analyzed above, our STP aims to jointly learn content and motion features  
 169 by spatiotemporal predictive learning. For content feature learning, we follow MAE [38], masking a  
 170 portion of the current frame based on  $\mathbb{M}_c$ , with  $\rho^c = 75\%$ , and predict the masked parts during the  
 171 decoding process. This encourages the model to learn spatial and geometric structure priors from the  
 172 current frame data through spatial reasoning. For motion feature learning, we establish an objective to  
 173 predict the future frame based on the masked current frame. However, predicting the future frame  
 174 without any conditions could be meaningless and extremely challenging. Therefore, we use the future  
 175 frame with an extremely high masking ratio as a condition, specifically  $\rho^f = 95\%$ , which reveals  
 176 some behavior and dynamic priors. In the experiments section, we will further discuss different  
 177 condition schemes, including language narration and the combination between them. In summary,  
 178 our encoding process during pre-training can be formally described as follows:

$$\begin{cases} \mathbf{Z}_c = \Phi_{enc}(\mathbf{I}_c, \mathbb{M}_c), \\ \mathbf{Z}_f = \Phi_{enc}(\mathbf{I}_f, \mathbb{M}_f). \end{cases} \quad (1)$$

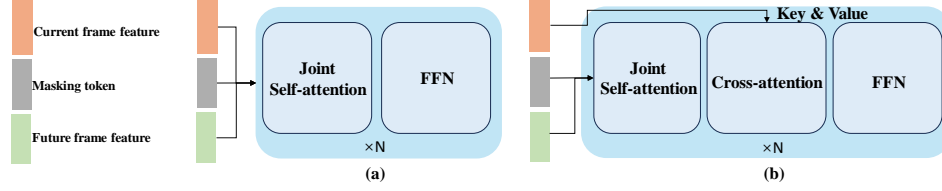


Figure 2: Temporal decoder design. (a) Standard joint-self architecture. (b) Our self-cross architecture.

### 179 3.3 Dual Decoders

180 To jointly capture static content and object motion features for better spatiotemporal understanding,  
 181 our STP present a dual decoders scheme to predict both the pixel of current and future frame  
 182 simultaneously in a multi-task learning manner. As shown in Figure 1, our dual decoder scheme  
 183 includes a spatial decoder  $\Phi_{dec-s}$  for spatial prediction and a temporal decoder  $\Phi_{dec-t}$  for temporal  
 184 prediction. We firstly give a technical description on them, respectively. Then we describe how we  
 185 combine them into our final method.

186 **Spatial Decoder.** To capture static content features, our spatial decoder is solely utilized for pro-  
 187 cessing the current frame visual feature. Specifically, after obtaining the masked current frame  
 188 visual feature  $\mathbf{Z}_c$ , we concatenate it with some learnable masking tokens, leading to the formation  
 189 of  $\mathbf{Z}_c^d = \mathbf{Z}_c \cup \{\mathbf{M}_i\}_{i \in \mathbb{M}_c}$ , where  $\mathbb{M}_c$  is the current frame masking map. Then, each of these tokens  
 190 further adds a corresponding positional embedding. Subsequently,  $\mathbf{Z}_c^d$  undergoes decoding in the  
 191 decoder and is continuously updated. The architecture of the spatial decoder block aligns with the  
 192 standard transformer encoder block, comprised of a global self-attention layer and a FFN layer.  
 193 Finally, with the decoded token sequence  $\mathbf{Z}_c^d$ , our spatial decoder predicts the invisible tokens of the  
 194 current frame  $\hat{\mathbf{I}}_c^d$ , operating under the current frame masking map  $\mathbb{M}_c$ .

195 **Temporal Decoder.** To capture motion features, our temporal decoder jointly processes the current  
 196 frame and the future frame which serves as the temporal prediction condition. To elaborate, we  
 197 firstly obtain the masked current frame feature  $\mathbf{Z}_c$  and the masked future frame feature  $\mathbf{Z}_f$ . We then  
 198 concatenate  $\mathbf{Z}_f$  with the masking tokens that have the positional embedding added, resulting in  $\mathbf{Z}_f^d$ .  
 199 Following this,  $\mathbf{Z}_f^d$  and  $\mathbf{Z}_c$  interact within the temporal decoder for decoding. The architecture of our  
 200 temporal decoder block is in alignment with the standard transformer decoder block [89], consisting  
 201 of a self-attention layer, a cross-attention layer, and a FFN layer, as shown in Figure 2 (b). During  
 202 decoding, the self-attention layer and FFN are solely used to process  $\mathbf{Z}_f^d$ . For the cross-attention  
 203 layer,  $\mathbf{Z}_f^d$  is continuously updated as the query, while  $\mathbf{Z}_c$ , acting as the key and value, is kept constant.  
 204 Compared to standard architecture, it ensures that the past frame representation space will not be  
 205 updated in the temporal decoder and are specifically used for temporal correlation and prediction.  
 206 This asymmetric interact architecture not only achieves more efficient training but also produces better  
 207 results. Finally, with the decoded token sequence  $\mathbf{Z}_f^d$ , our temporal decoder predicts the invisible  
 208 tokens of the future frame  $\hat{\mathbf{I}}_f^d$ , operating under the future frame masking map  $\mathbb{M}_f$ .

209 **Multi-task Predictive Learning.** As mentioned above, our STP jointly conducts spatiotemporal  
 210 prediction by asymmetric masking ratio and dual decoders scheme, the whole decoding pipeline can  
 211 be formally described as follows:

$$\begin{cases} \hat{\mathbf{I}}_c^d = \Phi_{dec-s}(\mathbf{Z}_c^d), \\ \hat{\mathbf{I}}_f^d = \Phi_{dec-t}(\mathbf{Z}_c, \mathbf{Z}_f^d). \end{cases} \quad (2)$$

212 Our loss function is the mean squared error (MSE) loss between the normalized masked pixels and  
 213 the predicted pixels. So our loss function  $\ell$  is as follows:

$$\ell = \text{MSE}(\hat{\mathbf{I}}_c, \mathbf{I}_c) + \text{MSE}(\hat{\mathbf{I}}_f, \mathbf{I}_f). \quad (3)$$

### 214 3.4 Downstream Policy Learning

215 To enable data and computation efficiency during the policy learning process, we adopt the paradigm  
 216 of few-shot behavior cloning by learning from demonstrations (Lfd), and we keep the image encoder

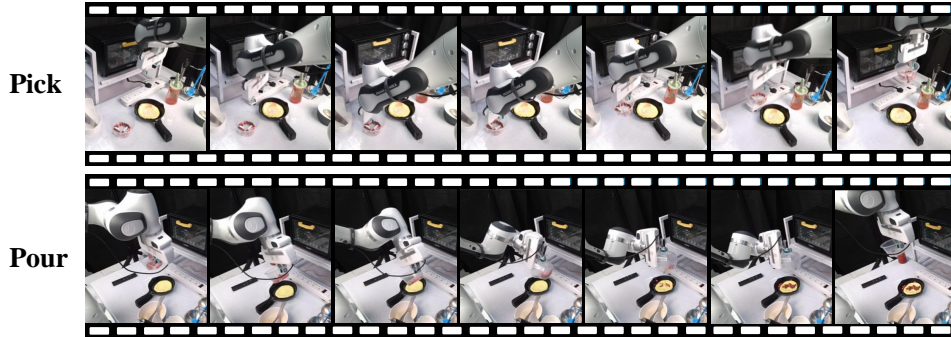


Figure 3: **The evaluation demonstrations of our real-world tasks.** For picking, the robot arm needs to pick up the bowl on the desktop. For pouring, the robot arm needs to pour the ingredients from the bowl into the pot.

217 frozen. Concretely, for each task, we are given offline expert demonstrations  $\mathcal{S} = \{\tau_1, \dots, \tau_n\}$ , where  
 218 each  $\tau_i$  is a trajectory of robot observations and actions, denoted as  $\tau_i = [(o_0, a_0), \dots, (o_T, a_T)]$ .  
 219 Based on the  $\mathcal{S}$ , we train a policy model,  $\pi_\theta(a|\mathcal{C}(\Phi_{enc}(o)))$ , parameterized by  $\theta$ , which maps from  
 220 robot’s state representations to actions. Here,  $\mathcal{C}$  represents an optional concatenation operation that  
 221 effectively fuses multi-view and multi-frame visual features, along with the robot’s proprioceptive  
 222 state in the channel dimension. We optimize the  $\pi_\theta$  through a standard behavior cloning MSE loss:

$$\min_{\theta} \sum_{(o,a) \sim \mathcal{S}} \text{MSE}(a, \pi_\theta(\mathcal{C}(\Phi_{enc}(o))))). \quad (4)$$

## 223 4 Experiments

### 224 4.1 Implementation on Pre-training

225 We execute pre-training with data from EgoVLP [55] for comprehensive ablation and fair comparison.  
 226 It processes untrimmed videos of Ego4D and filters out that miss language narrations and belong  
 227 to validation or test sets, resulting in a total of 3.8 million clips, called as Egoclip. In pre-training,  
 228 we sample a frame pair from each clip for training. As for all experiments, we employ ViT [22] as  
 229 backbone. Additionally, we maintain consistency with prior works [73, 59], directly using the [CLS]  
 230 token as the global representation. The pre-training hyperparameters can be found in section A.3.

### 231 4.2 Implementation on Downstream Policy

232 **Evaluation Scheme.** Following popular settings on PVRs for robotic motor control [65, 46, 59], for  
 233 each task, we learn a single policy  $\pi$  which is structured as a MLPs network. The policy models  
 234 utilize both the history of visual observation embeddings and optional robot proprioceptive as inputs,  
 235 subsequently generating executable actions as outputs.

236 **Simulation Tasks.** We select the union of manipulation and locomotion tasks from prior  
 237 works [65, 59] for evaluation, encompassing 19 tasks across 5 simulated environments. These include  
 238 Meta-World [104] (Assembly, Bin-Picking, Button-Press, Drawer-Open, and Hammer), Franka-  
 239 Kitchen [31] (Sliding Door, Turning Light On, Opening Door, Turning Knob, and Opening Mi-  
 240 crowave), Adroit [74] (Relocate and Reorient-Pen), DMControl [83] (Finger-Spin, Reacher-Hard,  
 241 Cheetah-Run, Walker-Stand, and Walker-Walk), and Trifinger [98] (Reach-Cube and Push-Cube).  
 242 More detailed simulation evaluation details can be found in section A.4.

243 **Real-World Tasks.** In our real-world experiments, we evaluate contact-rich picking and pouring  
 244 tasks using a Franka Emika Research 3 robot arm in a tabletop environment, ensuring no duplication  
 245 with simulation Franka-Kitchen [31]. For each task, we collect 100 noise demonstrations for training,  
 246 and we conduct 20 trials per task during evaluation phase. The robotic arm and objects have different  
 247 initial pose between training and testing. The evaluation demonstrations of our real-world tasks is  
 248 shown in Figure 3. Please see section A.5 for more real-world setup details.

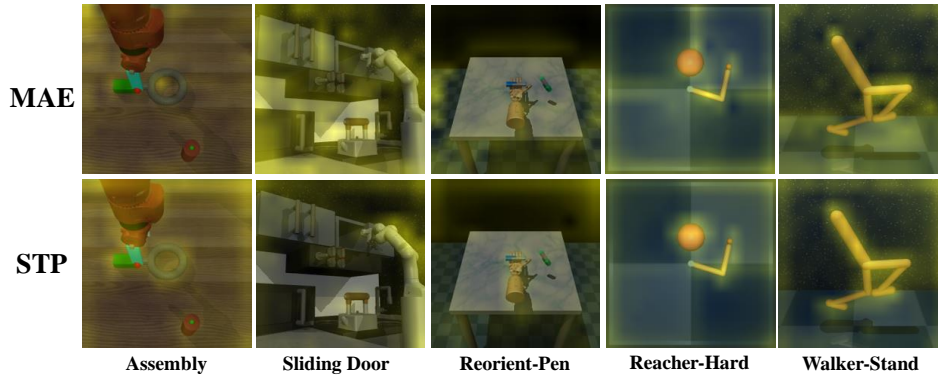


Figure 4: **Attention Visualization.** We use the [CLS] token as query, average the attention of all heads at the last layer of the frozen ViT encoder, and perform min-max normalization. We then upsample the attention map and overlay it on the original image, where the size of the attention value is directly proportional to the intensity of the yellow light. **Top:** MAE pre-training. **Bottom:** STP pre-training.

### 249 4.3 Performance on Downstream Simulation Tasks

250 In this section, we mainly analyze the performance of some pre-trained image representations on  
 251 reproducible simulation tasks. Specifically, we first evaluate the following models: (1) public  
 252 DINOv2 [67] that combines masked image modeling with self-distillation on large-scale image  
 253 datasets; (2) public CLIP [71] that conducts contrastive learning on large-scale image-text pairs;  
 254 (3) R3M trained based on Egoclip [55]; (4) public VC-1 [59]; (5) MAE trained based on Egoclip;  
 255 (6) STP trained based on Egoclip. (7) STP that conducts hybrid pre-training with initialization  
 256 using ImageNet-MAE [59]. Among them, (1) and (2) achieve excellent performance on core visual  
 257 understanding tasks using zero-shot or linear probing evaluation settings. (3) and (4) utilize egocentric  
 258 videos for robotic motor control. (5), (6) and (7) are used for fair comparison and exploring the  
 259 potential benefits of STP from more diverse image data, respectively. The experimental results are  
 260 presented in Table 1. Consistent with prior findings [41, 59], there is not a universal foundation model  
 261 that performs optimally across all benchmarks. However, on the whole, the MAE method is superior  
 262 due to its effective modeling of low-level geometry and spatial structure, especially for the MetaWorld  
 263 tasks that demand fine-grained control. Another intriguing observation is that MAE underperforms  
 264 in the Franka-Kitchen and Adroit tasks. We believe that this could be due to its relatively weaker  
 265 semantic representation. Under a fair comparison, our STP outperforms MAE by 4.1 (59.6  $\rightarrow$  63.7),  
 266 and additionally benefits from a more diverse image data, improving by 0.5 (63.7  $\rightarrow$  64.2). This is  
 267 attributed to that our STP not only captures static content features but also effectively models motion  
 268 information by extracting temporal clues from videos of interactions and manipulations with the  
 269 environment and objects. Additionally, we provide the visualization of the attention maps (model (5)  
 270 and (6)) of several specific tasks in Figure 4. The results indicate that, on top of effectively capturing  
 271 spatial information, our method further encourages the model to focus on motion areas or objects,  
 272 thereby providing a more *sparse and compact* representation for downstream low-data BC paradigm.

273 Next, we also evaluate and compare the adaptation results of our representations to downstream motor  
 274 control tasks. Specifically, we evaluate following settings: (a) The MAE pre-trained representation  
 275 undergoes further MAE post-pre-training with task-specific data, and is frozen during policy training;  
 276 (b) The STP pre-trained representation undergoes further STP post-pre-training with task-specific  
 277 data, and is frozen during policy training; (c) The STP pre-trained representation undergoes end-to-  
 278 end fine-tuning with task-specific data; (d) STP pre-training is performed directly using task-specific  
 279 data and the resulting representation is frozen during policy training. The results show that end-to-end  
 280 fine-tuning fails to yield the best results, suggesting that naively fine-tuning ViT-base could still lead  
 281 to overfitting under few-shot behavior cloning scheme. Conversely, (a) and (b) achieve competitive  
 282 results, with our STP achieving a 3.9 (72.5  $\rightarrow$  76.4) improvement on the weight average success rate  
 283 than MAE, further demonstrating the effectiveness and data efficiency of our STP for in-domain data.  
 284 In addition, the comparison between (a) and (d) also proves the effectiveness of pre-training with  
 285 out-of-domain data. Finally, we also scale up both MAE and our STP to ViT-L/16, and the results  
 286 still demonstrate the superiority of STP. Among them, compared to ViT-B/16, ViT-L/16 brings a  
 287 smaller performance improvement, which may be due to the task’s performance saturation. However,  
 288 the ViT-L/16 of STP does not show improvement in Meta-World and Trifinger, indicating that simply

Table 1: Performance comparisons of visual representations on simulation benchmarks. We report the average score across all tasks for each simulation environment. DINOv2 uses **ViT-B/14**, CLIP uses **ViT-B/32**, and unless otherwise specified, others use **ViT-B/16**. Mt-Wd, Fr-Ki, DMC, Adro, Tr-fi, and WA respectively represent MetaWorld, Franka-Kitchen, DMControl, Adroit, Trifinger, and weight average. \* denotes that public VC-1 samples image frames from full Ego4D dataset.

	Pre-training Data	Mt-Wd	Fr-Ki	DMC	Adro	Tr-fi	WA
DINOv2 [67]	LVD-142M	77.9	41.2	59.4	50.7	69.0	59.6
CLIP [71]	Image-text pairs	75.5	39.8	52.2	<b>51.3</b>	57.7	55.6
R3M [65]	Ego	81.3	30.6	52.2	46.7	64.7	54.9
VC-1 [59]	Ego*+MNI	88.8	38.4	60.9	46.0	70.5	61.8
MAE [38]	Ego	85.1	36.7	59.2	43.4	<b>70.6</b>	59.6
STP	Ego	92.0	40.9	<b>62.1</b>	48.0	69.3	63.7
STP	Ego+I	<b>94.1</b>	<b>42.5</b>	61.6	47.3	66.7	<b>64.2</b>
MAE (Post PT)	Ego+Demo	93.6	46.9	81.1	58.0	76.8	72.5
STP (Post PT)	Ego+Demo	<b>97.3</b>	<b>53.6</b>	<b>82.8</b>	<b>63.3</b>	<b>78.0</b>	<b>76.4</b>
STP (E2E FT)	Ego	87.2	52.4	55.2	40.0	70.4	62.9
STP	Demo	70.3	30.4	52.5	38.0	70.8	51.8
MAE-L/16 (Post PT)	Ego+Demo	95.7	54.7	83.5	66.0	<b>77.6</b>	76.7
STP-L/16 (Post PT)	Ego+Demo	<b>97.3</b>	<b>57.4</b>	<b>85.0</b>	<b>70.0</b>	75.4	<b>78.4</b>

Table 2: The ablation experiment results. Me, Fra, DMC, Adr, Tri, and WA respectively represent MetaWorld, Franka-Kitchen, DMControl, Adroit, Trifinger, and weight average. All models use **ViT-B/16**.

(a) Current Frame Masking and Spatial Prediction.

$\rho^c$	Predict	Me	Fra	DMC	Adr	Tri	WA
75%	✓	<b>92.0</b>	<b>40.9</b>	<b>62.1</b>	<b>48.0</b>	<b>69.3</b>	<b>63.7</b>
75%		84.5	34.7	55.4	43.3	65.3	57.4
50%	✓	82.1	36.0	60.3	<b>48.0</b>	66.8	59.0
0%		79.2	39.7	54.8	44.0	63.1	57.0

(c) Temporal Decoder Architecture Design.

Decoder	Me	Fra	DMC	Adr	Tri	WA
8 joint-self	87.7	36.9	55.7	46.0	71.3	59.8
12 joint-self	88.5	35.0	55.7	46.0	67.0	59.1
8 self-cross	<b>92.0</b>	<b>40.9</b>	<b>62.1</b>	<b>48.0</b>	69.3	<b>63.7</b>

(b) Temporal Prediction Condition Design.

Condition	Me	Fra	DMC	Adr	Tri	WA
L-E	82.1	30.7	55.5	42.0	63.8	55.4
95%	<b>92.0</b>	40.9	62.1	<b>48.0</b>	<b>69.3</b>	<b>63.7</b>
90%	91.2	<b>42.5</b>	62.8	44.7	65.9	63.4
L-E + 95%	91.0	37.7	<b>64.1</b>	46.7	<b>70.8</b>	63.1
L-D + 95%	88.0	34.3	62.6	46.7	69.3	60.9

(d) Frame Sampling Strategy.

Frame interval	Me	Fra	DMC	Adr	Tri	WA
8	89.6	39.9	58.4	46.0	67.0	61.3
16	92.0	40.9	<b>62.1</b>	<b>48.0</b>	<b>69.3</b>	<b>63.7</b>
24	89.1	<b>41.1</b>	61.5	46.0	68.1	62.5
8, 24	<b>92.3</b>	37.1	57.3	42.0	68.4	60.8

289 scaling up model capacity does not necessarily lead to performance gains. In the few-shot BC setting,  
 290 there is a risk of overfitting in both policy and backbone training.

#### 291 4.4 Ablation on Downstream Simulation Tasks

292 In this section, we perform extensive ablation studies to further demonstrate the effectiveness of our  
 293 joint spatial and temporal prediction, as well as temporal prediction condition design. In addition, we  
 294 also study the influence of temporal decoder architecture design and future frame sampling strategy.

295 **Current frame masking.** The design of the current frame masking is crucial. On one hand, similar  
 296 to MAE [38], masking some patches and predicting the missing parts can effectively promote the  
 297 learning of image content features. On the other hand, the visible patches of the current frame need  
 298 to interact with the condition to predict the future frame. Specifically, we mask the current frame at  
 299 masking rates of 75%, 50%, and 0%, respectively, and optionally predict the missing parts through  
 300 the spatial decoder. The results are shown in Table 2 (a). From results, we see that the masking ratio  
 301 of 75% and performing spatial prediction still lead to the best performance. This demonstrates the  
 302 importance of retaining MAE [38] for content features learning, especially for low-level manipulation  
 303 in Meta-World, while a current frame with a high masking ratio (75%) is sufficient to interact with  
 304 other conditions to predict the future frame.

305 **Temporal prediction condition design.** Subsequently, we discuss the influence of temporal predic-  
 306 tion condition design. We implicitly model motion in actionless video data by predicting the pixels of  
 307 the future frame. A direct and simple idea is to use language narration as a condition. The text tokens  
 308 can be flexibly utilized as inputs to ViT [22], forming a multimodal encoder. Language narration



309 provides a high-level behavior description, but lacks low-level visual dynamic priors for pixel-level  
 310 prediction. However, leaking part of the future frame can effectively provide these priors. In order  
 311 to explore how to construct a more meaningful temporal prediction proxy task, we compare the  
 312 following schemes: (1) only language narration, (2) masking 95% of the future frame, (3) masking  
 313 90% of the future frame, (4) masking 95% of the future frame and language narration, and (5) masking  
 314 95% of the future frame and language narration, but the language is added in the temporal decoder,  
 315 instead of being fused with the visible image patches in the multimodal encoder. We tokenize all  
 316 language narration by pre-trained DistilBERT [75]. The results are shown in Table 2 (b). From  
 317 results, we see that using only language as a prediction condition leads to a significant decline in  
 318 performance, while leaking a small amount of future frame (masking 95%) in the temporal decoder  
 319 can achieve competitive results. As for joint conditions of language and future frame with 95%  
 320 masking ratio, adding language in the encoder is better than in the decoder. Additionally, adding  
 321 language performs better on DMControl (64.1 vs. 62.1) and Trifinger (70.8 vs. 69.3), while not  
 322 adding language performs better on Meta-World (92.0 vs. 91.0), Franka-Kitchen (40.9 vs. 37.7) and  
 323 Adroit (48.0 vs. 46.7). We speculate the reasons for language hurts performance are as follows: (i)  
 324 The input gap (multi-modal and single-modal) between upstream and downstream; (ii) Extra language  
 325 in ViT may result in the loss of some fine-grained information capture. Furthermore, the latter does  
 326 not require language supervision, and can provide a more scalable self-supervised solution.

327 **Temporal decoder design.** We also investigate the impact of the temporal decoder design. Specifi-  
 328 cally, we consider two types of decoder blocks. One is the joint-self architecture, as shown in Figure 2  
 329 (a), and similar joint architecture are adopted in [26, 102]. The other is the self-cross architecture, as  
 330 shown in Figure 2 (b), and similar cross architecture are adopted in [3, 33]. We consider the following  
 331 settings: (1) 8 joint-self decoder blocks, (2) 12 joint-self decoder blocks, (3) 8 self-cross decoder  
 332 blocks. Among them, setting (2) and (3) have similar amounts of parameters for a fairer comparison.  
 333 The results are shown in Table 2 (c). The results demonstrate the importance of maintaining a fixed  
 334 representation space of the past frame during temporal prediction.

335 **Frame sampling strategy.** Finally, we investigate the impact of the sampling strategy between the  
 336 current frame and future frame. The difficulty of temporal prediction is directly proportional to the  
 337 frame interval values. We establish four settings where we fix the sampling intervals at 8, 16, and 24  
 338 respectively, and for the fourth setting, we randomly select an interval within the range of [8, 24].  
 339 The results are shown in Table 2 (d). The results show that an interval of 16 achieves the best balance  
 340 for building temporal prediction proxy task.

#### 341 4.5 Performance on Downstream Real-world Tasks

342 In this section, we report our experiment results on  
 343 real-world picking and pouring tasks. We report the  
 344 average success rate for each task. Specifically, we  
 345 compare STP with the baseline MAE, both of which  
 346 are trained on out-of-domain videos and kept frozen  
 347 during policy training. The results are shown in Ta-  
 348 ble 3. From the results, it can be seen that STP has  
 349 achieved significant advantages in the pouring task.  
 350 It can more accurately align with the moving bowl  
 351 and the pot. In addition, although MAE and STP have a same success rate in picking tasks, STP tends  
 352 to execute grasping in a better position. This indicates that the trend and conclusion of our STP are  
 353 consistent in both simulation and the real-world, which also aligns with the findings of [79].

Table 3: Performance comparisons on real-world tasks.

Method	Picking	Pouring	Average
MAE	65.0	45.0	55.0
STP	65.0	65.0	<b>65.0</b>

## 354 5 Conclusion

355 In this work, we have proposed the STP, a simple, efficient and effective self-supervised visual repre-  
 356 sentation pre-training framework for robotic motor control. Our STP jointly performs spatiotemporal  
 357 predictive learning on large-scale videos within a multi-task learning manner. Our STP captures  
 358 content features by predicting the invisible areas within the masked current frame, and simultaneously  
 359 captures motion features by using a future frame with an extremely high masking ratio as a condition  
 360 to predict the invisible areas within that future frame. We carry out the largest-scale BC evaluation of  
 361 PVRs for robotic motor control to date to demonstrate the effectiveness of STP. Furthermore, as for  
 362 pre-training data, we also prove that extending STP to hybrid pre-training and post-pre-training could  
 363 further unleash its generality and data efficiency.

364 **References**

- 365 [1] OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob  
366 McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al.  
367 Learning dexterous in-hand manipulation. *IJRR*, 39(1):3–20, 2020.
- 368 [2] Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael  
369 Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a  
370 joint-embedding predictive architecture. In *CVPR*, pages 15619–15629, 2023.
- 371 [3] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimaes: Multi-modal  
372 multi-task masked autoencoders. In *ECCV*, pages 348–367. Springer, 2022.
- 373 [4] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli.  
374 Data2vec: A general framework for self-supervised learning in speech, vision and language.  
375 In *ICML*, pages 1298–1312. PMLR, 2022.
- 376 [5] Shikhar Bahl, Abhinav Gupta, and Deepak Pathak. Human-to-robot imitation in the wild. In  
377 Kris Hauser, Dylan A. Shell, and Shoudong Huang, editors, *RSS*, 2022.
- 378 [6] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances  
379 from human videos as a versatile representation for robotics. In *CVPR*, pages 13778–13790,  
380 2023.
- 381 [7] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image  
382 transformers. In *ICLR*, 2021.
- 383 [8] Adrien Bardes, Jean Ponce, and Yann LeCun. Mc-jepa: A joint-embedding predictive  
384 architecture for self-supervised learning of motion and content features. *arXiv preprint*  
385 *arXiv:2307.12698*, 2023.
- 386 [9] Konstantinos Bousmalis, Giulia Vezzani, Dushyant Rao, Coline Devin, Alex X Lee, Maria  
387 Bauza, Todor Davchev, Yuxiang Zhou, Agrim Gupta, Akhil Raju, et al. Robocat: A self-  
388 improving foundation agent for robotic manipulation. *arXiv preprint arXiv:2306.11706*,  
389 2023.
- 390 [10] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof  
391 Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-  
392 language-action models transfer web knowledge to robotic control. In *CoRL*, 2023.
- 393 [11] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea  
394 Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian  
395 Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J. Joshi, Ryan Julian,  
396 Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu,  
397 Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn  
398 Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael S. Ryoo, Grecia  
399 Salazar, Pannag R. Sanketi, Kevin Sayed, Jaspier Singh, Sumedh Sontakke, Austin Stone,  
400 Clayton Tan, Huong T. Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted  
401 Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: robotics transformer for  
402 real-world control at scale. In Kostas E. Bekris, Kris Hauser, Sylvia L. Herbert, and Jingjin  
403 Yu, editors, *RSS*, 2023.
- 404 [12] Shaofei Cai, Zihao Wang, Xiaojian Ma, Anji Liu, and Yitao Liang. Open-world multi-task  
405 control through goal-aware representation learning and adaptive horizon prediction. In *CVPR*,  
406 pages 13734–13744, 2023.
- 407 [13] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,  
408 and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*,  
409 pages 9650–9660, 2021.
- 410 [14] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
411 for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020.

- 412 [15] X Chen, S Xie, and K He. An empirical study of training self-supervised vision transformers.  
413 in 2021 *ieee*. In *ICCV*, pages 9620–9629.
- 414 [16] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and  
415 Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. In Kostas E.  
416 Bekris, Kris Hauser, Sylvia L. Herbert, and Jingjin Yu, editors, *RSS*, 2023.
- 417 [17] Yutao Cui, Cheng Jiang, Gangshan Wu, and Limin Wang. Mixformer: End-to-end tracking  
418 with iterative mixed attention. *arXiv preprint arXiv:2302.02814*, 2023.
- 419 [18] Yilun Dai, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuur-  
420 mans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *arXiv*  
421 *preprint arXiv:2302.00111*, 2023.
- 422 [19] Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Abhinav Gupta. An unbiased look at  
423 datasets for visuo-motor pre-training. In *CoRL*, 2023.
- 424 [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
425 hierarchical image database. In *CVPR*, pages 248–255. *ieee*, 2009.
- 426 [21] AM Derrington and P Lennie. Spatial and temporal contrast sensitivities of neurones in lateral  
427 geniculate nucleus of macaque. *The Journal of physiology*, 357(1):219–240, 1984.
- 428 [22] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
429 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
430 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
431 recognition at scale. In *ICLR*.
- 432 [23] Yilun Du, Mengjiao Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet,  
433 Tianhe Yu, Pieter Abbeel, Joshua B Tenenbaum, et al. Video language planning. *arXiv preprint*  
434 *arXiv:2310.10625*, 2023.
- 435 [24] Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, Chenxi Wang, Junbo Wang, Haoyi  
436 Zhu, and Cewu Lu. Rh20t: A comprehensive robotic dataset for learning diverse skills in  
437 one-shot. In *Towards Generalist Robots: Learning Paradigms for Scalable Skill Acquisition@*  
438 *CoRL2023*, 2023.
- 439 [25] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang,  
440 Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation  
441 learning at scale. In *CVPR*, pages 19358–19369, 2023.
- 442 [26] Xinyang Geng, Hao Liu, Lisa Lee, Dale Schuurams, Sergey Levine, and Pieter Abbeel.  
443 Multimodal masked autoencoders learn transferable representations. *arXiv preprint*  
444 *arXiv:2205.14204*, 2022.
- 445 [27] Abraham George, Alison Bartsch, and Amir Barati Farimani. Openvr: Teleoperation for  
446 manipulation. *arXiv preprint arXiv:2305.09765*, 2023.
- 447 [28] Rohit Girdhar and Kristen Grauman. Anticipative video transformer. In *ICCV*, pages 13505–  
448 13515, 2021.
- 449 [29] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit  
450 Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the  
451 world in 3,000 hours of egocentric video. In *CVPR*, pages 18995–19012, 2022.
- 452 [30] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Tri-  
453 antafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al.  
454 Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives.  
455 *arXiv preprint arXiv:2311.18259*, 2023.
- 456 [31] Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay  
457 policy learning: Solving long-horizon tasks via imitation and reinforcement learning. In *CoRL*,  
458 pages 1025–1037. PMLR, 2020.

- 459 [32] Agrim Gupta, Stephen Tian, Yunzhi Zhang, Jiajun Wu, Roberto Martín-Martín, and Li Fei-Fei.  
460 Maskvit: Masked visual pre-training for video prediction. In *ICLR*, 2022.
- 461 [33] Agrim Gupta, Jiajun Wu, Jia Deng, and Li Fei-Fei. Siamese masked autoencoders. 2023.
- 462 [34] Huy Ha, Pete Florence, and Shuran Song. Scaling up and distilling down: Language-guided  
463 robot skill acquisition. In *Conference on Robot Learning*, pages 3766–3777. PMLR, 2023.
- 464 [35] Tengda Han, Weidi Xie, and Andrew Zisserman. Video representation learning by dense  
465 predictive coding. In *ICCV Workshops*, pages 0–0, 2019.
- 466 [36] Nicklas Hansen, Zhecheng Yuan, Yanjie Ze, Tongzhou Mu, Aravind Rajeswaran, Hao Su,  
467 Huazhe Xu, and Xiaolong Wang. On pre-training for visuo-motor control: Revisiting a  
468 learning-from-scratch baseline. *arXiv preprint arXiv:2212.05749*, 2022.
- 469 [37] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood.  
470 Flexible diffusion modeling of long videos. *NeurIPS*, 35:27953–27965, 2022.
- 471 [38] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
472 autoencoders are scalable vision learners. In *CVPR*, pages 16000–16009, 2022.
- 473 [39] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
474 unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020.
- 475 [40] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
476 recognition. In *CVPR*, pages 770–778, 2016.
- 477 [41] Yingdong Hu, Renhao Wang, Li Erran Li, and Yang Gao. For pre-trained vision models in mo-  
478 tor control, not all policy learning methods are created equal. *arXiv preprint arXiv:2304.04591*,  
479 2023.
- 480 [42] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey  
481 Levine, and Chelsea Finn. BC-Z: zero-shot task generalization with robotic imitation learning.  
482 In Aleksandra Faust, David Hsu, and Gerhard Neumann, editors, *CoRL*, volume 164 of  
483 *Proceedings of Machine Learning Research*, pages 991–1002. PMLR, 2021.
- 484 [43] Chuhao Jin, Wenhui Tan, Jiange Yang, Bei Liu, Ruihua Song, Limin Wang, and Jianlong Fu.  
485 Alphablock: Embodied finetuning for vision-language reasoning in robot manipulation. *arXiv*  
486 *preprint arXiv:2305.18898*, 2023.
- 487 [44] Ya Jing, Xuelin Zhu, Xingbin Liu, Qie Sima, Taozheng Yang, Yunhai Feng, and Tao Kong.  
488 Exploring visual pre-training for robot manipulation: Datasets, models and methods. *arXiv*  
489 *preprint arXiv:2308.03620*, 2023.
- 490 [45] Yuanchen Ju, Kaizhe Hu, Guowei Zhang, Gu Zhang, Mingrun Jiang, and Huazhe Xu. Robo-  
491 abc: Affordance generalization beyond categories via semantic correspondence for robot  
492 manipulation. *arXiv preprint arXiv:2401.07487*, 2024.
- 493 [46] Siddharth Karamcheti, Suraj Nair, Annie S. Chen, Thomas Kollar, Chelsea Finn, Dorsa Sadigh,  
494 and Percy Liang. Language-driven representation learning for robotics. In Kostas E. Bekris,  
495 Kris Hauser, Sylvia L. Herbert, and Jingjin Yu, editors, *RSS*, 2023.
- 496 [47] Heecheol Kim, Yoshiyuki Ohmura, and Yasuo Kuniyoshi. Multi-task robot data for dual-arm  
497 fine manipulation. *arXiv preprint arXiv:2401.07603*, 2024.
- 498 [48] Moo Jin Kim, Jiajun Wu, and Chelsea Finn. Giving robots a hand: Learning generalizable  
499 manipulation with eye-in-hand human video demonstrations. *arXiv preprint arXiv:2307.05959*,  
500 2023.
- 501 [49] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*  
502 *preprint arXiv:1412.6980*, 2014.
- 503 [50] Andreas Kleinschmidt, Kai V Thilo, Christian Büchel, Michael A Gresty, Adolfo M Bronstein,  
504 and Richard SJ Frackowiak. Neural correlates of visual-motion perception as object-or self-  
505 motion. *Neuroimage*, 16(4):873–882, 2002.

- 506 [51] Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and Joshua B Tenenbaum. Learning to  
507 act from actionless videos through dense correspondences. *arXiv preprint arXiv:2310.08576*,  
508 2023.
- 509 [52] Xiangwen Kong and Xiangyu Zhang. Understanding masked image modeling via learning  
510 occlusion invariant feature. In *CVPR*, pages 6241–6251, 2023.
- 511 [53] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep  
512 visuomotor policies. *JMLR*, 17(1):1334–1373, 2016.
- 513 [54] Xinghang Li, Minghuan Liu, Hanbo Zhang, Cunjun Yu, Jie Xu, Hongtao Wu, Chilam Cheang,  
514 Ya Jing, Weinan Zhang, Huaping Liu, et al. Vision-language foundation models as effective  
515 robot imitators. *arXiv preprint arXiv:2311.01378*, 2023.
- 516 [55] Kevin Qinghong Lin, Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Z XU,  
517 Difei Gao, Rong-Cheng Tu, Wenzhe Zhao, Weijie Kong, et al. Egocentric video-language  
518 pretraining. *NeurIPS*, 35:7575–7586, 2022.
- 519 [56] Xingyu Lin, John So, Sashwat Mahalingam, Fangchen Liu, and Pieter Abbeel. Spawn-  
520 net: Learning generalizable visuomotor skills from pre-trained networks. *arXiv preprint*  
521 *arXiv:2307.03567*, 2023.
- 522 [57] Yecheng Jason Ma, William Liang, Vaidehi Som, Vikash Kumar, Amy Zhang, Osbert Bastani,  
523 and Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control.  
524 *arXiv preprint arXiv:2306.00958*, 2023.
- 525 [58] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and  
526 Amy Zhang. Vip: Towards universal visual reward and representation via value-implicit  
527 pre-training. In *ICLR*, 2023.
- 528 [59] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Yecheng Jason Ma, Claire Chen, Sneha  
529 Silwal, Aryan Jain, Vincent-Pierre Berges, Pieter Abbeel, Jitendra Malik, et al. Where are we  
530 in the search for an artificial visual cortex for embodied intelligence? 2023.
- 531 [60] Pierre Marza, Laetitia Matignon, Olivier Simonin, and Christian Wolf. Task-conditioned  
532 adaptation of visual features in multi-task policy learning. *arXiv preprint arXiv:2402.07739*,  
533 2024.
- 534 [61] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction  
535 beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015.
- 536 [62] Russell Mendonca, Shikhar Bahl, and Deepak Pathak. Structured world models from human  
537 videos. In Kostas E. Bekris, Kris Hauser, Sylvia L. Herbert, and Jingjin Yu, editors, *RSS*, 2023.
- 538 [63] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and  
539 Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million  
540 narrated video clips. In *CVPR*, pages 2630–2640, 2019.
- 541 [64] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang,  
542 Jifeng Dai, Yu Qiao, and Ping Luo. EmbodiedGPT: Vision-language pre-training via embodied  
543 chain of thought. In *NeurIPS*, 2023.
- 544 [65] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A  
545 universal visual representation for robot manipulation. In *CoRL*, 2022.
- 546 [66] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive  
547 predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 548 [67] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
549 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2:  
550 Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 551 [68] Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan,  
552 Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment:  
553 Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.

- 554 [69] Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, and Abhinav Gupta. The unsur-  
555 prising effectiveness of pre-trained vision models for control. In *ICML*, pages 17359–17371.  
556 PMLR, 2022.
- 557 [70] Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong  
558 Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *ECCV*,  
559 pages 570–587. Springer, 2022.
- 560 [71] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
561 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
562 models from natural language supervision. In *ICLR*, pages 8748–8763. PMLR, 2021.
- 563 [72] Ilija Radosavovic, Baifeng Shi, Letian Fu, Ken Goldberg, Trevor Darrell, and Jitendra Malik.  
564 Robot learning with sensorimotor pre-training. 2023.
- 565 [73] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell.  
566 Real-world robot learning with masked visual pre-training. In *CoRL*, pages 416–426. PMLR,  
567 2023.
- 568 [74] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman,  
569 Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep  
570 reinforcement learning and demonstrations. *RSS*, 2018.
- 571 [75] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled  
572 version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 573 [76] Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter  
574 Abbeel. Masked world models for visual control. In *CoRL*, pages 1332–1344. PMLR, 2023.
- 575 [77] Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with  
576 action-free pre-training from videos. In *ICML*, pages 19561–19579. PMLR, 2022.
- 577 [78] Mohit Sharma, Claudio Fantacci, Yuxiang Zhou, Skanda Koppula, Nicolas Heess, Jon Scholz,  
578 and Yusuf Aytar. Lossless adaptation of pretrained vision models for robotic manipulation. In  
579 *ICLR*, 2023.
- 580 [79] Sneha Silwal, Karmesh Yadav, Tingfan Wu, Jay Vakil, Arjun Majumdar, Sergio Arnaud, Claire  
581 Chen, Vincent-Pierre Berges, Dhruv Batra, Aravind Rajeswaran, et al. What do we learn from  
582 a large-scale study of pre-trained visual representations in sim and real environments? *arXiv*  
583 *preprint arXiv:2310.02219*, 2023.
- 584 [80] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of  
585 video representations using lstms. In *ICML*, pages 843–852. PMLR, 2015.
- 586 [81] Austin Stone, Ted Xiao, Yao Lu, Keerthana Gopalakrishnan, Kuang-Huei Lee, Quan Vuong,  
587 Paul Wohlhart, Brianna Zitkovich, Fei Xia, Chelsea Finn, et al. Open-world object manipula-  
588 tion using pre-trained vision-language models. In *CoRL*, 2023.
- 589 [82] Yanchao Sun, Shuang Ma, Ratnesh Madaan, Rogerio Bonatti, Furong Huang, and Ashish  
590 Kapoor. Smart: Self-supervised multi-task pretraining with control transformers. In *ICLR*,  
591 2023.
- 592 [83] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David  
593 Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite.  
594 *arXiv preprint arXiv:1801.00690*, 2018.
- 595 [84] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep  
596 Dasari, Joey Hejna, Charles Xu, Jianlan Luo, et al. Octo: An open-source generalist robot  
597 policy, 2023.
- 598 [85] Garrett Thomas, Ching-An Cheng, Ricky Loynd, Felipe Vieira Frujeri, Vibhav Vineet, Mihai  
599 Jalobeanu, and Andrey Kolobov. Plex: Making the most of the available data for robotic  
600 manipulation pretraining. In *CoRL*, pages 2624–2641. PMLR, 2023.

- 601 [86] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are  
602 data-efficient learners for self-supervised video pre-training. *NeurIPS*, 35:10078–10093, 2022.
- 603 [87] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and  
604 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In  
605 *ICLR*, pages 10347–10357. PMLR, 2021.
- 606 [88] David C Van Essen and Jack L Gallant. Neural mechanisms of form and motion processing in  
607 the primate visual system. *Neuron*, 13(1):1–10, 1994.
- 608 [89] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
609 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017.
- 610 [90] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan Vuong, Chongyi Zheng, Philippe  
611 Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al. Bridgedata v2:  
612 A dataset for robot learning at scale. In *Conference on Robot Learning*, pages 1723–1736.  
613 PMLR, 2023.
- 614 [91] Chen Wang, Linxi Fan, Jiankai Sun, Ruohan Zhang, Li Fei-Fei, Danfei Xu, Yuke Zhu, and  
615 Anima Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play.  
616 *arXiv preprint arXiv:2302.12422*, 2023.
- 617 [92] Jianren Wang, Sudeep Dasari, Mohan Kumar Srirama, Shubham Tulsiani, and Abhinav Gupta.  
618 Manipulate by seeing: Creating manipulation controllers from pre-trained representations. In  
619 *ICCV*, pages 3859–3868, 2023.
- 620 [93] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yanan He, Yi Wang, Yali Wang, and  
621 Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *CVPR*,  
622 pages 14549–14560, 2023.
- 623 [94] Lirui Wang, Jialiang Zhao, Yilun Du, Edward H Adelson, and Russ Tedrake. Poco: Policy  
624 composition from and for heterogeneous robot learning. *arXiv preprint arXiv:2402.02511*,  
625 2024.
- 626 [95] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer.  
627 Masked feature prediction for self-supervised visual pre-training. In *CVPR*, pages 14668–  
628 14678, 2022.
- 629 [96] William F Whitney, Rajat Agarwal, Kyunghyun Cho, and Abhinav Gupta. Dynamics-aware  
630 embeddings. In *ICLR*, 2020.
- 631 [97] Hongtao Wu, Ya Jing, Chilam Cheang, Guangzeng Chen, Jiafeng Xu, Xinghang Li, Minghuan  
632 Liu, Hang Li, and Tao Kong. Unleashing large-scale video generative pre-training for visual  
633 robot manipulation. *arXiv preprint arXiv:2312.13139*, 2023.
- 634 [98] Manuel Wuthrich, Felix Widmaier, Felix Grimminger, Shruti Joshi, Vaibhav Agrawal, Bi-  
635 lal Hammoud, Majid Khadiv, Miroslav Bogdanovic, Vincent Berenz, Julian Viereck, et al.  
636 Trifinger: An open-source robot for learning dexterity. In *CoRL*, pages 1871–1882. PMLR,  
637 2021.
- 638 [99] Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked visual pre-training  
639 for motor control. *arXiv preprint arXiv:2203.06173*, 2022.
- 640 [100] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han  
641 Hu. Simmim: A simple framework for masked image modeling. In *CVPR*, pages 9653–9663,  
642 2022.
- 643 [101] Mengda Xu, Zhenjia Xu, Cheng Chi, Manuela Veloso, and Shuran Song. Xskill: Cross  
644 embodiment skill discovery. In *CoRL*, 2023.
- 645 [102] Jiange Yang, Sheng Guo, Gangshan Wu, and Limin Wang. Comae: Single model hybrid  
646 pre-training on small-scale rgb-d datasets. *arXiv preprint arXiv:2302.06148*, 2023.

- 647 [103] Jiange Yang, Wenhui Tan, Chuhao Jin, Bei Liu, Jianlong Fu, Ruihua Song, and Limin Wang.  
648 Pave the way to grasp anything: Transferring foundation models for universal pick-place  
649 robots. *arXiv preprint arXiv:2306.05716*, 2023.
- 650 [104] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn,  
651 and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta  
652 reinforcement learning. In *CoRL*, pages 1094–1100. PMLR, 2020.
- 653 [105] Chengbo Yuan, Chuan Wen, Tong Zhang, and Yang Gao. General flow as foundation affordance  
654 for scalable robot learning. *arXiv preprint arXiv:2401.11439*, 2024.
- 655 [106] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong.  
656 ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021.
- 657 [107] Haoyi Zhu, Yating Wang, Di Huang, Weicai Ye, Wanli Ouyang, and Tong He. Point cloud  
658 matters: Rethinking the impact of different observation spaces on robot learning. *arXiv  
659 preprint arXiv:2402.02500*, 2024.

## 660 A Appendix

### 661 A.1 Limitations and Discussion

662 Although STP has demonstrated superior performance in extensive experiments, there remain some  
663 challenges and future works. From the perspective of pre-training data, Ego4D provides numerous  
664 human-object interaction scenes and good motion clues. Building larger-scale and more diverse  
665 potential datasets such as [63, 30] to scale up STP is worth exploring. Regarding pre-training methods,  
666 exploring predictive targets outside of pixel space and more effective sampling and masking strategies  
667 present intriguing research directions. From an evaluation standpoint, we utilize a frozen ViT to  
668 extract agent state representations and adopt the paradigm of few-shot behavior cloning, other policy  
669 learning methods (reinforcement learning, visual reward function, visual task specification), have not  
670 been explored. In conclusion, as the first method of performing temporal prediction on large-scale  
671 videos for self-supervised visual representation learning intended for robotic motor control tasks, we  
672 hope STP can be taken as a strong baseline and facilitate further research along this direction.

### 673 A.2 The influence of the loss weight ratio between temporal prediction and spatial prediction

674 In this section, we further explore the influence of the loss weight ratio between temporal prediction  
675 and spatial prediction. Specifically, taking five tasks from Franka-Kitchen as examples, we load the  
676 pre-trained STP and perform post-pre-training with three different loss weight ratios (temporal to  
677 spatial). The results, as shown in Figure 5, are 54.7, 55.2, and 57.4 for the average results of the ratios  
678 3:1, 1:3, and 1:1, respectively. The results indicate that due to the different attributes of the tasks, the  
679 trends are not consistent. However, overall, the 1:1 ratio achieves the best balance and results. We  
680 chose it as a universal setting.

### 681 A.3 Pre-training Details

682 In this section, we describe the details of our STP pre-training. Specifically, we list some key training  
683 and architectural hyperparameters of STP in Table 4. In addition, as for our MAE [38] baseline, we  
684 mainly follow the publicly available code of MAE<sup>1</sup>. Additionally, we train MAE and STP using the  
685 same data and number of epochs to ensure that the comparison between them is **completely fair**.  
686 Finally, we also provide some STP prediction results in Figure 6.

### 687 A.4 Simulation Environments Details

688 In this section, we first present further details of the STP post-pre-training on downstream simula-  
689 tion environments. Subsequently, we delineate the specific hyperparameters used in the behavior  
690 cloning policy training within these simulation environments. Finally, we provide the comprehensive  
691 evaluation scheme for each simulation environment.

<sup>1</sup><https://github.com/facebookresearch/mae>



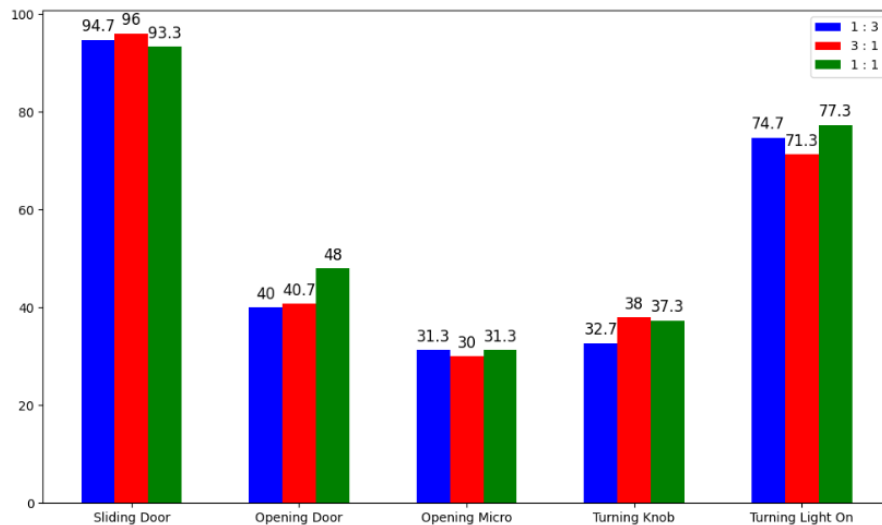


Figure 5: The results of different loss weight ratios between temporal prediction and spatial prediction.

Table 4: Training and architectural hyperparameters for STP pre-training.

Hyperparameter	Value
<i>STP Pre-training</i>	
optimizer	AdamW [49]
base learning rate	0.00015
weight decay	0.05
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.95$
effective batch size	4096
learning rate schedule	cosine decay
total epochs	50
warmup epochs	5
augmentation	RandomResizedCrop (0.8, 1)
<i>Encoder ViT-base Architecture</i>	
patch size	16
#layers	12
#MHSA heads	12
hidden dim	768
positional embedding	sin-cos initialization and fix
<i>Dual Decoder ViT-base Architecture</i>	
#layers	8
#MHSA heads	16
hidden dim	512
positional embedding	sin-cos initialization and fix

692 In regards to the STP post-pre-training, we utilize data that aligns with the policy training, and the  
 693 specific architecture hyperparameters correspond to those listed in Table 4. Depending on the specific  
 694 demonstration data, we adjust the values of total epochs, warmup epochs, effective batch size, and  
 695 the frame interval, as shown in Table 5.

696 As for policy training and evaluation schemes, we primarily refer to the publicly available code<sup>2</sup>  
 697 and training data of VC-1 [59] for Metaworld [104], DMControl [83], Adroit [74] and Trifinger [98].  
 698 Similarly, for Franka-Kitchen [31], we follow the public code<sup>3</sup> and training data of R3M [65].  
 699 Specifically, the policy training hyperparameters and evaluation schemes are shown in Table 6 and  
 700 Table 7, respectively. About policy training, we completely follow the setting of prior works [65, 59]  
 701 when freezing the encoder; when performing end-to-end fine-tuning, we make appropriate adjustments  
 702 to the batch size and learning rate. About evaluation details, similar to prior works[65, 59], we  
 703 establish all evaluation details such as the number of expert demonstrations and test trajectories,  
 704 environmental viewpoints, optimization hyperparameters, base seeds, history windows size, and  
 705 the use of robot proprioceptive. In Table 7, the term ‘prop.’ stands for whether proprioceptive  
 706 information is used or not, and ‘history window size’ signifies the number of frames received by  
 707 the policy model at each step, with features between frames being fused through concatenation.  
 708 ‘Number of trajectories’ represents the quantity of trajectories evaluated. For tasks in Meta-World,  
 709 Franka-Kitchen, Adroit, and Trifinger, we report the maximum success rate, whereas for tasks in  
 710 DMControl, we report the maximum reward score, rescaling to be in the range of [0, 100] by dividing  
 711 by 10. We report the average metric across tasks for each environment. In addition, it is worth noting  
 712 that the metrics we report are the **average value across all base seeds and camera viewpoints**.  
 713 Finally, we also report the results of our post-pre-training STP (ViT-B/16) on each task in Table 8.

714 In addition, we emphasize that different random seeds primarily affect the rendering of the initial  
 715 frame in the sampled trajectories, as shown in Figure 7. During evaluation, the seed value we provide  
 716 serves as the base seed, and the trajectory sampling process is depicted in Algorithm 1. **Therefore,**  
 717 **the actual number of trajectories we evaluate is the number of trajectories multiplied by the**  
 718 **number of base seeds.** For instance, for MetaWorld, we evaluate  $25 \times 3 = 75$  trajectories, with  
 719 random seeds for rendering being 100-124, 200-224, and 300-324.

720 Finally, for Franka-Kitchen, we utilize MuJoCo210, while all other simulation environments are  
 721 based on MuJoCo200. Our policy training and evaluation environments are conducted on Cuda 11.3,  
 722 NVIDIA TITAN Xp GPUs, and OpenGL 3.1.

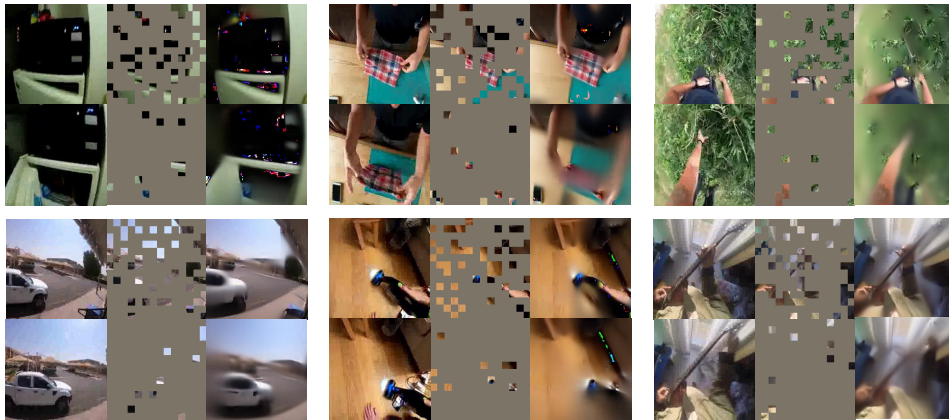


Figure 6: Some examples of our STP prediction result on Ego4D videos. For each six tuple, we show the ground-truth (left), masked frames (middle), STP prediction results (right), current frames (top), and future frames (bottom). We simply overlay the output with the visible patches to improve visual quality.

<sup>2</sup><https://github.com/facebookresearch/eai-vc/tree/main/cortexbench>

<sup>3</sup><https://github.com/facebookresearch/r3m/tree/eval/evaluation>

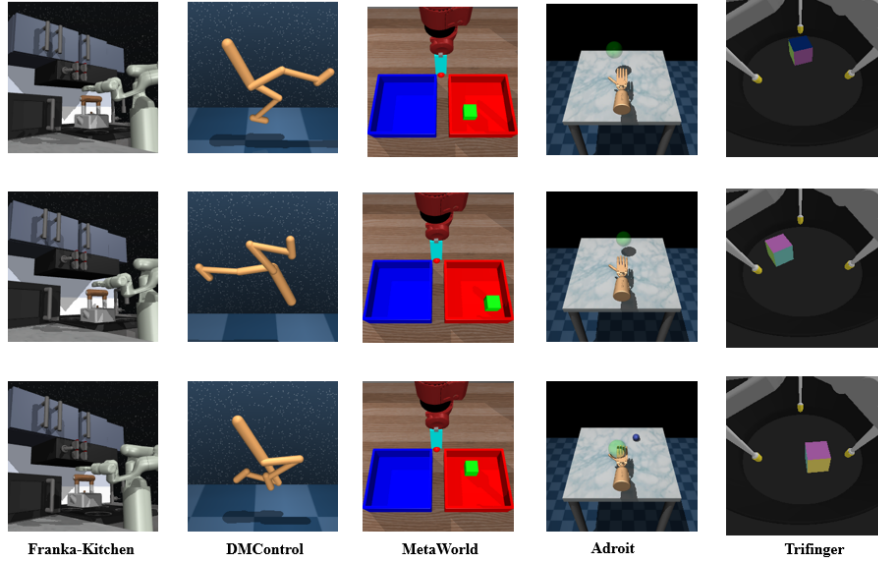


Figure 7: The visualization of initial frame rendering under different random seeds.

---

**Algorithm 1** Trajectories Sampling Pseudocode

---

```

# num_traj: the number of evaluation trajectories
# base_seed: base seed for rollouts

# rollout to sample trajectories
for ep in range(num_traj):
    seed = base_seed + ep
    env.set_seed(seed)
    o = env.reset()

```

---

723 **A.5 Real-World Environments Details**

724 In this section, we outline the details of our real-world setup and evaluation scheme. As depicted  
725 in Figure 8, our real-world scenario includes four camera viewpoints: top, left, right, and wrist. It  
726 includes two Kinect DK and two RealSense cameras. An example of four views is shown in Figure 9.  
727 Specifically, we utilize four different camera views and resize their resolution uniformly to  $224 \times 224$ .  
728 To effectively model the complex and multimodal action distribution in our real-world tasks, we  
729 select diffusion policy [16] as our policy model. In accordance with this approach, we concatenate  
730 the visual embeddings of all views from two sequential frames. Following the approach in [27], we  
731 collect robot data using a VR tele-operation setup. In this way, we collect 100 continuous trajectories  
732 for each task. It is worth noting that the quality of these demonstrations leaves room for improvement  
733 and contains a lot of noise. During the evaluation phase, we primarily evaluate two contact-rich  
734 tasks that have not appeared in Franka-Kitchen [31] benchmark: (1) Picking. It requires the robot  
735 arm to pick up the transparent bowl off the table; (2) Pouring. It requires the robot arm to pour  
736 at least three-quarter of the ingredients from the transparent bowl into the black pot. For each task,  
737 we change the initial pose of the robot arm and objects within a certain range as well as conduct 20  
738 trials. In addition, there are different distractors on the desktop during training and testing, which  
739 also evaluates the robustness of the model to distractors. Throughout the process, we use ROS and  
740 MoveIt for hardware communication and motion planning.



Figure 8: Our real-world scene with four cameras and a Franka Emika robot arm.

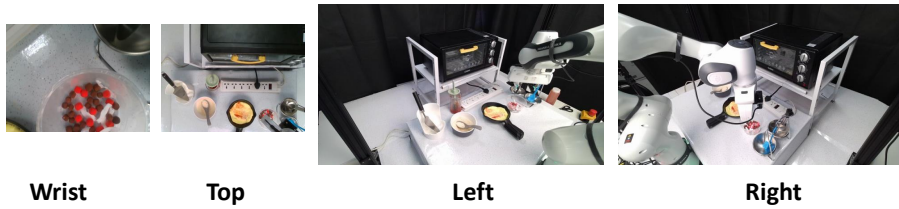


Figure 9: An example of four views.

Table 5: STP post-pre-training hyperparameters on simulation environments.

	MetaWorld	Franka-Kitchen	DMControl	Adroit	Trifinger
total epochs	50	100	50	50	50
warmup epochs	5	5	5	5	5
effective batch size	1024	128	2048	1024	1024
number of demonstrations	25	25	100	100	100
frame interval	4	4	4	4	16

Table 6: Policy training hyperparameters on simulation environments.

		MetaWorld	Franka-Kitchen	DMControl	Adroit	Trifinger
epochs		100	480	100	100	100 / 1000
batch size	frozen	256	32	256	256	32
	fine-tuning	64	32	64	64	16
learning rate	frozen	0.001	0.001	0.001	0.001	0.0001
	fine-tuning	0.00005	0.0001	0.00005	0.00005	0.0001

Table 7: Evaluation schemes on simulation environments.

Benchmark	Observation Space	History Window Size	Camera ViewPoints	Base Seeds	Number of Trajectories
Metaworld	RGB + prop.	3	top_cap2	100, 200, 300	25
Franka-Kitchen	RGB + prop.	1	left, right	123, 124, 125	50
DMControl	RGB	3	0	100, 200, 300	25
Adroit	RGB + prop.	1	vil_camera	100, 200, 300	25
Trifinger	RGB + prop.	1	default	10	25

Table 8: The success rate for each task on simulation bechmarks.

Assembly 94.7	Bin-Picking 97.3	Button-Press 94.7	Drawer-Open 100.0	Hammer 100.0
Sliding Door 96.0	Turning Light on 72.7	Opening Door 39.0	Turning Knob 31.3	Opening Microwave 29.0
Relocate 49.3	Reorient-Pen 77.3	Finger-Spin 69.6	Cheetah-Run 71.9	Reacher-Hard 87.7
Walker-Stand 95.9	Walker-Walk 89.0	Reach-Cube 85.3	Push-Cube 70.6	

## 741 **NeurIPS Paper Checklist**

742 The checklist is designed to encourage best practices for responsible machine learning research,  
743 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
744 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
745 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
746 towards the page limit.

747 Please read the checklist guidelines carefully for information on how to answer these questions. For  
748 each question in the checklist:

- 749 • You should answer [Yes], [No], or [NA].
- 750 • [NA] means either that the question is Not Applicable for that particular paper or the  
751 relevant information is Not Available.
- 752 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

753 **The checklist answers are an integral part of your paper submission.** They are visible to the  
754 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
755 (after eventual revisions) with the final version of your paper, and its final version will be published  
756 with the paper.

757 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
758 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
759 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
760 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
761 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
762 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
763 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
764 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
765 please point to the section(s) where related material for the question can be found.

766 **IMPORTANT, please:**

- 767 • **Delete this instruction block, but keep the section heading "NeurIPS paper checklist",**
- 768 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 769 • **Do not modify the questions and only use the provided macros for your answers.**

### 770 **1. Claims**

771 Question: Do the main claims made in the abstract and introduction accurately reflect the  
772 paper's contributions and scope?

773 Answer: [Yes]

774 Justification: Yes, the main claims made in the abstract and introduction accurately reflect  
775 the paper's contributions and scope.

776 Guidelines:

- 777 • The answer NA means that the abstract and introduction do not include the claims  
778 made in the paper.
- 779 • The abstract and/or introduction should clearly state the claims made, including the  
780 contributions made in the paper and important assumptions and limitations. A No or  
781 NA answer to this question will not be perceived well by the reviewers.
- 782 • The claims made should match theoretical and experimental results, and reflect how  
783 much the results can be expected to generalize to other settings.
- 784 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
785 are not attained by the paper.

### 786 **2. Limitations**

787 Question: Does the paper discuss the limitations of the work performed by the authors?

788 Answer: [Yes]

789 Justification: Yes, the paper discusses the limitations of the work performed by the authors.

790 Guidelines:

- 791 • The answer NA means that the paper has no limitation while the answer No means that
- 792 the paper has limitations, but those are not discussed in the paper.
- 793 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 794 • The paper should point out any strong assumptions and how robust the results are to
- 795 violations of these assumptions (e.g., independence assumptions, noiseless settings,
- 796 model well-specification, asymptotic approximations only holding locally). The authors
- 797 should reflect on how these assumptions might be violated in practice and what the
- 798 implications would be.
- 799 • The authors should reflect on the scope of the claims made, e.g., if the approach was
- 800 only tested on a few datasets or with a few runs. In general, empirical results often
- 801 depend on implicit assumptions, which should be articulated.
- 802 • The authors should reflect on the factors that influence the performance of the approach.
- 803 For example, a facial recognition algorithm may perform poorly when image resolution
- 804 is low or images are taken in low lighting. Or a speech-to-text system might not be
- 805 used reliably to provide closed captions for online lectures because it fails to handle
- 806 technical jargon.
- 807 • The authors should discuss the computational efficiency of the proposed algorithms
- 808 and how they scale with dataset size.
- 809 • If applicable, the authors should discuss possible limitations of their approach to
- 810 address problems of privacy and fairness.
- 811 • While the authors might fear that complete honesty about limitations might be used by
- 812 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
- 813 limitations that aren't acknowledged in the paper. The authors should use their best
- 814 judgment and recognize that individual actions in favor of transparency play an impor-
- 815 tant role in developing norms that preserve the integrity of the community. Reviewers
- 816 will be specifically instructed to not penalize honesty concerning limitations.

### 817 3. Theory Assumptions and Proofs

818 Question: For each theoretical result, does the paper provide the full set of assumptions and

819 a complete (and correct) proof?

820 Answer: [\[Yes\]](#)

821 Justification: The paper does not include theoretical results.

822 Guidelines:

- 823 • The answer NA means that the paper does not include theoretical results.
- 824 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 825 referenced.
- 826 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 827 • The proofs can either appear in the main paper or the supplemental material, but if
- 828 they appear in the supplemental material, the authors are encouraged to provide a short
- 829 proof sketch to provide intuition.
- 830 • Inversely, any informal proof provided in the core of the paper should be complemented
- 831 by formal proofs provided in appendix or supplemental material.
- 832 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 833 4. Experimental Result Reproducibility

834 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

835 perimental results of the paper to the extent that it affects the main claims and/or conclusions

836 of the paper (regardless of whether the code and data are provided or not)?

837 Answer: [\[Yes\]](#)

838 Justification: The paper fully disclose all the information needed to reproduce results.

839 Guidelines:

- 840 • The answer NA means that the paper does not include experiments.

- 841 • If the paper includes experiments, a No answer to this question will not be perceived  
842 well by the reviewers: Making the paper reproducible is important, regardless of  
843 whether the code and data are provided or not.
- 844 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
845 to make their results reproducible or verifiable.
- 846 • Depending on the contribution, reproducibility can be accomplished in various ways.  
847 For example, if the contribution is a novel architecture, describing the architecture fully  
848 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
849 be necessary to either make it possible for others to replicate the model with the same  
850 dataset, or provide access to the model. In general, releasing code and data is often  
851 one good way to accomplish this, but reproducibility can also be provided via detailed  
852 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
853 of a large language model), releasing of a model checkpoint, or other means that are  
854 appropriate to the research performed.
- 855 • While NeurIPS does not require releasing code, the conference does require all submis-  
856 sions to provide some reasonable avenue for reproducibility, which may depend on the  
857 nature of the contribution. For example
  - 858 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
859 to reproduce that algorithm.
  - 860 (b) If the contribution is primarily a new model architecture, the paper should describe  
861 the architecture clearly and fully.
  - 862 (c) If the contribution is a new model (e.g., a large language model), then there should  
863 either be a way to access this model for reproducing the results or a way to reproduce  
864 the model (e.g., with an open-source dataset or instructions for how to construct  
865 the dataset).
  - 866 (d) We recognize that reproducibility may be tricky in some cases, in which case  
867 authors are welcome to describe the particular way they provide for reproducibility.  
868 In the case of closed-source models, it may be that access to the model is limited in  
869 some way (e.g., to registered users), but it should be possible for other researchers  
870 to have some path to reproducing or verifying the results.

## 871 5. Open access to data and code

872 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
873 tions to faithfully reproduce the main experimental results, as described in supplemental  
874 material?

875 Answer: [Yes]

876 Justification: We will release all codes and model weights on github.

877 Guidelines:

- 878 • The answer NA means that paper does not include experiments requiring code.
- 879 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
880 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 881 • While we encourage the release of code and data, we understand that this might not be  
882 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
883 including code, unless this is central to the contribution (e.g., for a new open-source  
884 benchmark).
- 885 • The instructions should contain the exact command and environment needed to run to  
886 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
887 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 888 • The authors should provide instructions on data access and preparation, including how  
889 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 890 • The authors should provide scripts to reproduce all experimental results for the new  
891 proposed method and baselines. If only a subset of experiments are reproducible, they  
892 should state which ones are omitted from the script and why.
- 893 • At submission time, to preserve anonymity, the authors should release anonymized  
894 versions (if applicable).



- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper specify all the training and test details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: In our experiments, different seeds are primarily used for rendering different initial frames. Therefore, our evaluation is comprehensive and sufficient, while our comparisons are absolutely fair.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide sufficient information on the computer resources.

Guidelines:

- The answer NA means that the paper does not include experiments.

- 945
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - 946
  - 947
  - 948
  - 949
  - 950
  - 951
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

952

953 Question: Does the research conducted in the paper conform, in every respect, with the

954 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

955 Answer: [Yes]

956 Justification: Our paper aligns with these.

957 Guidelines:

- 958 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 959 • If the authors answer No, they should explain the special circumstances that require a
- 960 deviation from the Code of Ethics.
- 961 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 962 eration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

963

964 Question: Does the paper discuss both potential positive societal impacts and negative

965 societal impacts of the work performed?

966 Answer: [NA]

967 Justification: There is no societal impact of the work performed.

968 Guidelines:

- 969 • The answer NA means that there is no societal impact of the work performed.
- 970 • If the authors answer NA or No, they should explain why their work has no societal
- 971 impact or why the paper does not address societal impact.
- 972 • Examples of negative societal impacts include potential malicious or unintended uses
- 973 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 974 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 975 groups), privacy considerations, and security considerations.
- 976 • The conference expects that many papers will be foundational research and not tied
- 977 to particular applications, let alone deployments. However, if there is a direct path to
- 978 any negative applications, the authors should point it out. For example, it is legitimate
- 979 to point out that an improvement in the quality of generative models could be used to
- 980 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 981 that a generic algorithm for optimizing neural networks could enable people to train
- 982 models that generate Deepfakes faster.
- 983 • The authors should consider possible harms that could arise when the technology is
- 984 being used as intended and functioning correctly, harms that could arise when the
- 985 technology is being used as intended but gives incorrect results, and harms following
- 986 from (intentional or unintentional) misuse of the technology.
- 987 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 988 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 989 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 990 feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

991

992 Question: Does the paper describe safeguards that have been put in place for responsible

993 release of data or models that have a high risk for misuse (e.g., pretrained language models,

994 image generators, or scraped datasets)?

995 Answer: [NA]

996 Justification: The paper poses no such risks.

997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

**15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.