

# SPARK: Simple Post-training for Adapting pRetrained Knowledge to Robot Control

Anonymous CVPR submission

Paper ID 23

## Abstract

001 *Large scale representation learning has produced public*  
002 *foundation models that provide strong general purpose vi-*  
003 *sual features. However, their pretraining objectives are not*  
004 *designed for robot representation learning, so the resulting*  
005 *embeddings tend to emphasize global semantics rather than*  
006 *the temporally sensitive representations required for robot*  
007 *control. Training robot models from scratch is also unde-*  
008 *sirable due to the limited scale and high collection cost*  
009 *of robotics data, as well as the substantial computational*  
010 *burden of pre-training. In this paper, we propose **SPARK***  
011 *(Simple Post-training for Adapting pRetrained Knowledge),*  
012 *a post training method that adapts foundation models to*  
013 *robot control. SPARK is built on two principles: dynamics-*  
014 *aware abstraction, which encourages the encoder to de-*  
015 *rive compact features that preserve information essential*  
016 *for understanding temporal changes and action-relevant*  
017 *scene structure, and knowledge preservation, which aligns*  
018 *patch-level representations with those of the original foun-*  
019 *dition model to retain useful pretrained semantics. These*  
020 *objectives yield compact visual state representations that re-*  
021 *main semantically meaningful while preserving useful prior*  
022 *knowledge. Experiments across multiple robotics bench-*  
023 *marks show that SPARK consistently improves success rates*  
024 *and generalization over vanilla foundation models, and fur-*  
025 *ther demonstrate that these gains transfer to real-world*  
026 *robot manipulation.*

## 027 1. Introduction

028 The evolution of robot representation learning [5, 17–19,  
029 21, 23] has primarily centered on modeling dynamics in  
030 spatiotemporal visual observations. This line of research  
031 has produced representations effective for imitation learn-  
032 ing and robot control, where success depends on capturing  
033 object motion, scene transitions, and task progress from vi-  
034 sual input. In parallel, general-purpose visual representa-  
035 tion learning has increasingly emphasized visual-language

alignment for transfer. From this perspective, conventional 036  
robot models remain limited, being largely optimized for 037  
narrow control distributions rather than language-grounded 038  
semantics. 039

Learning robot representations directly from paired 040  
video-language data offers a direct route to unifying em- 041  
bodied dynamics with language-level semantics. However, 042  
collecting robot vision-action data at scale is expensive and 043  
typically limited to laboratory settings, making it difficult to 044  
cover diverse scenes, objects, and interactions. This has led 045  
prior work to learn visual representations from large-scale 046  
video data without robot actions, in order to capture phys- 047  
ical dynamics that may benefit downstream robot learning. 048  
While such approaches improve scalability and provide tem- 049  
poral supervision, they still lack the broad semantic cover- 050  
age and vision-language alignment available in foundation 051  
models trained on web-scale corpora. These limitations moti- 052  
vate a different strategy that adapts already well-aligned 053  
foundation models [20, 22, 24, 26, 28] to the requirements 054  
of robot control. 055

In this paper, we propose **SPARK** (Simple Post-training 056  
for Adapting pRetrained Knowledge), a post-training 057  
method that adapts foundation models and general-purpose 058  
vision models to robot control. The central goal of SPARK 059  
is to enable pretrained representations more suitable for con- 060  
trol without sacrificing the general visual knowledge ac- 061  
quired during large-scale pretraining. To this end, SPARK 062  
is built on two key principles: **dynamics-aware abstrac-** 063  
**tion** and **knowledge preservation**. Dynamics-aware ab- 064  
straction encourages the encoder to produce compact fea- 065  
tures that preserve information essential for understanding 066  
temporal changes by leveraging the prior knowledge ac- 067  
quired during large-scale pretraining. However, naively en- 068  
forcing such compact abstraction can degrade useful prior 069  
knowledge inherited from large-scale pretraining. Thus, we 070  
adopt a retention objective that aligns the patch-wise rep- 071  
resentations of the adapted encoder with those of the original 072  
foundation model. 073

We empirically demonstrate that SPARK yields consis- 074  
tent improvements across multiple robotics benchmarks. 075

076	Models adapted with SPARK yield higher success rates and	<b>Post-training and Adaptation.</b> Post-training has been	124
077	better generalization than their vanilla models on average	used to reduce the cost of training when adapting pre-	125
078	on robot manipulation tasks, highlighting that cost-efficient	trained encoders to new domains or tasks through post-	126
079	post-training can substantially enhance the generalizability	training. Parameter-efficient adaptation methods introduce	127
080	of the foundation model on robot control. Moreover, exper-	small learnable components such as prompts and adapter	128
081	iments on real-world robot tasks show that these gains ext-	layers, as in CoOP and CoCoOp, CLIP Adapter and visual	129
082	end beyond simulation. Consequently, SPARK presents a	prompt tuning, or apply low-rank updates like LoRA [7, 13,	130
083	promising recipe for mitigating the mismatch between found-	29, 30]. In robotics, a prior work [18] has explored post-	131
084	ation models and the visual representations required for	training using a generic reconstruction objective based on	132
085	embodied intelligence.	masked autoencoding [11] approaches to enhance the visual	133
		understanding capacity of the original models. Besides, our	134
		approach follows the spirit of post-training adaptation but	135
		targets the representation more directly. SPARK introduces	136
		a conservative summarization loss that trains a single class	137
		token to act as a compact state, together with a retention	138
		loss that aligns token-level features with those of the origi-	139
		nal foundation model.	140
086	<b>2. Related Work</b>	<b>3. Method</b>	141
		<b>3.1. Overview</b>	142
087	<b>Pretrained Vision and Foundation Models.</b> Large-scale	The goal of SPARK (Simple Post-training for Adapting	143
088	vision models or vision-language models, including CLIP	pRetrained Knowledge) is to adapt pretrained foundation	144
089	and SigLIP families, are trained on numerous image and	models to robot control without sacrificing the broad vi-	145
090	text datasets and provide strong zero-shot and few-shot per-	sual knowledge acquired during large-scale pretraining. To	146
091	formance [1–4, 10, 11, 20, 22, 26, 28]. These founda-	this end, SPARK is built on two complementary princi-	147
092	tion models capture high-level semantics and diverse vi-	ples: <b>dynamics-aware abstraction</b> and <b>knowledge preser-</b>	148
093	visual regularities and have become standard backbones for	<b>vation.</b> Dynamics-aware abstraction encourages the en-	149
094	many vision tasks. However, their pre-training objectives	coder to derive compact features that preserve information	150
095	are designed for recognition and language alignment rather	essential for understanding temporal changes and action-	151
096	than robot representation learning, resulting in embeddings	relevant scene structure. Knowledge preservation regular-	152
097	suitable for global semantics and caption-level alignment,	izes this adaptation so that the learned representation re-	153
098	which may not be effective for temporally sensitive struc-	mains aligned with the semantically meaningful features of	154
099	ture within consecutive observations, which is crucial for	the original foundation model. Together, these principles	155
100	control. This objective mismatch motivates post-training	enable pretrained representations to become more suitable	156
101	recipes that align these models with the demands of robotic	for robot control without requiring training from scratch.	157
102	perception without pre-training from scratch.	To realize dynamics-aware abstraction, SPARK trains	158
		the encoder so that a compact feature extracted from a refer-	159
		ence observation can support the prediction of a temporally	160
		offset target scene. Specifically, the encoder processes a	161
		reference scene and produces patch tokens together with a	162
		class token, where the class token serves as a compact fea-	163
		ture. A decoder then reconstructs token representations of	164
		the target scene from this compact feature, conditioned on	165
		only a small subset of visible target patches as hints. Be-	166
		cause accurate prediction requires information that explains	167
		how the scene changes over time, this objective encour-	168
		ages the compact feature to preserve temporal and action-	169
		relevant structure rather than only static global semantics.	170
		However, adapting the encoder solely through this ab-	171
		straction objective can distort useful pretrained knowledge.	172
		To mitigate this issue, SPARK introduces a retention objec-	173
		tive that aligns the patch-wise representations of the adapted	174
103	<b>Visual Representations in Robotics.</b> In robotic percep-		
104	tion, pretrained vision backbones are widely used to encode		
105	observations for downstream manipulation, navigation, and		
106	policy learning [5, 14]. Several studies have shown that		
107	even frozen features from generic encoders can already pro-		
108	vide strong representations for imitation learning and con-		
109	trol [21]. Beyond direct reuse, robotics-oriented visual pre-		
110	training has been explored to further improve transfer by		
111	exploiting egocentric viewpoints or temporal structure in		
112	robot data. Examples include R3M, which combines time		
113	contrastive objectives with language supervision on human-		
114	centric egocentric video [19], VIP, which learns value im-		
115	plicit representations that connect initial and goal obser-		
116	vations [17], and MAE-style variants such as MVP and		
117	VC-1, which curate large ego-centric manipulation corpora		
118	for masked reconstruction [5, 18, 23]. However, these		
119	works mainly introduce pre-training recipes that require		
120	huge computation costs, often with paired language annota-		
121	tions, which are expensive to collect and still much smaller		
122	than web-scale image and text data, limiting their accessi-		
123	bility and coverage.		

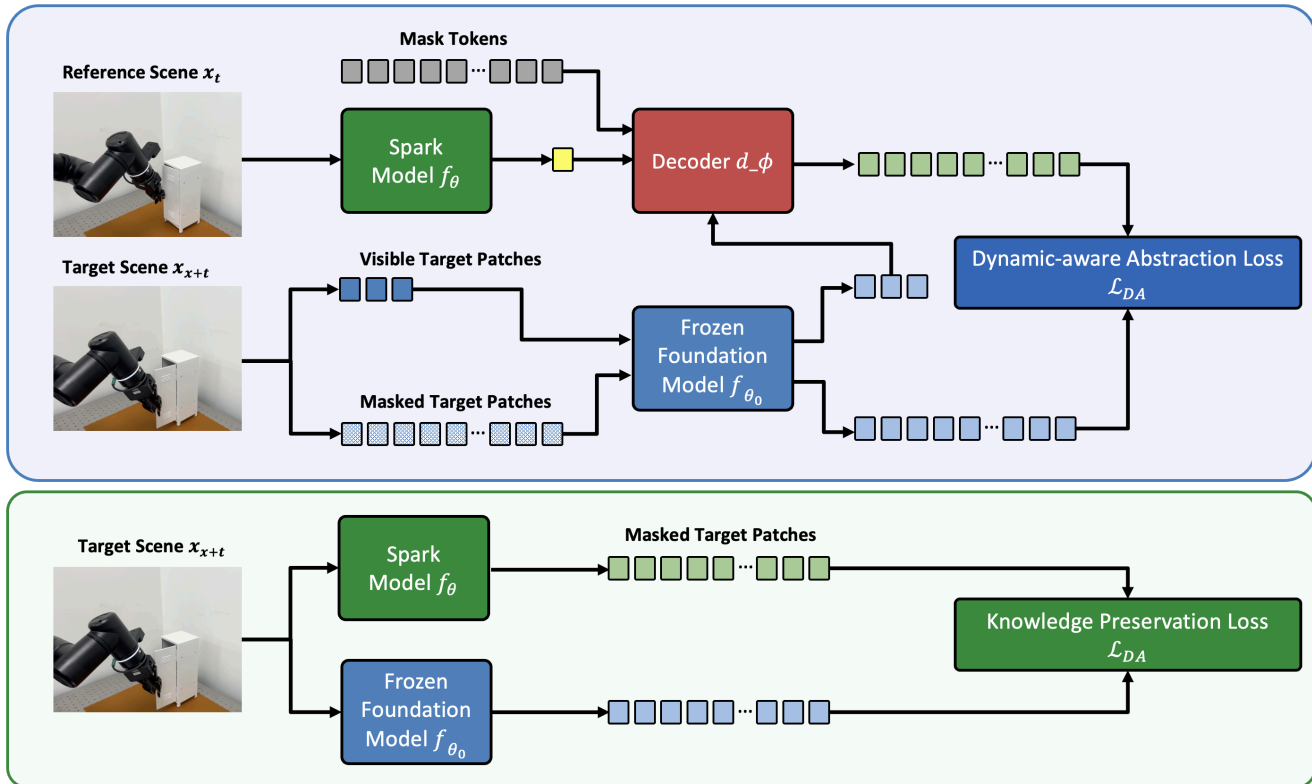


Figure 1. **Overview of SPARK (Simple Post-training for Adapting pRetrained Knowledge)**. SPARK aims to encourage conservative summarization of patch-wise well-understood semantics while preserving the prior knowledge of the foundation model obtained during pre-training. Specifically, given a reference scene  $x_t$ , the SPARK encoder  $f_\theta$  encodes the reference scene into patch tokens and a class token. Then, the decoder predicts the target token representations from the class token with hints from a few target patches. The conservative summarization loss  $\mathcal{L}_{DA}$  minimizes the discrepancy between these predicted target token representations and the corresponding target token representations produced by a frozen baseline encoder  $f_{\theta_0}$ . On the other hand, the retention loss aligns the reference token representations encoded by the SPARK encoder with the reference token representations encoded by the frozen baseline encoder  $f_{\theta_0}$  at the patch level. This encourages the class token to act as a compact visual state for control while preserving the patch-level semantic knowledge of the original foundation model.

175 encoder with those of the frozen original foundation model  
 176 on the same input. This regularization preserves semanti-  
 177 cally meaningful local features while allowing the compact  
 178 feature to specialize toward dynamic understanding for control.  
 179 The overall pipeline is illustrated in Fig. 1.

### 180 3.2. Formulation

181 **Encoding procedure.** Let  $f_\theta$  denote the SPARK encoder  
 182 and let  $f_{\theta_0}$  denote the frozen foundation model used for ini-  
 183 tialization. We sample a reference scene  $x_t \in \mathbb{R}^{3 \times H \times W}$   
 184 and a temporally offset target scene  $x_{t+k} \in \mathbb{R}^{3 \times H \times W}$ . We  
 185 patchify both scenes into  $N$  non-overlapping patches, de-  
 186 noted by  $\{x_t^i\}_{i=1}^N$  and  $\{x_{t+k}^i\}_{i=1}^N$ , respectively. The refer-  
 187 ence scene is processed by both the SPARK encoder  $f_\theta$  and  
 188 the frozen encoder  $f_{\theta_0}$ , yielding adapted token representa-  
 189 tions  $\{u_t^i\}_{i=1}^N$  and frozen token representations  $\{v_{t+k}^i\}_{i=1}^N$ . In  
 190 addition, the SPARK encoder produces a class token  $u_t^{[CLS]}$ ,

191 which serves as the compact feature associated with the refer-  
 192 ence scene. For the target scene, the frozen encoder pro-  
 193 cesses all target patches to produce  $\{v_{t+k}^i\}_{i=1}^N$ , while the  
 194 SPARK encoder processes only a subset of visible target  
 195 patches indexed by  $M \subset \{1, \dots, N\}$ , where  $|M| = \lfloor rN \rfloor$   
 196 for a visible ratio  $r \in (0, 1)$ . This yields adapted target to-  
 197 kens  $\{u_{t+k}^i\}_{i \in M}$ , which are used as hints for target-scene  
 198 prediction.

199 **Dynamics-aware abstraction loss.** The abstraction ob-  
 200 jective encourages the compact feature  $u_t^{[CLS]}$  to retain  
 201 the information required to predict the target scene. To  
 202 this end, we reconstruct target token representations from  
 203 the compact feature together with the visible target hints  
 204  $\{u_{t+k}^i\}_{i \in M}$ . Let  $M^c$  denote the masked target-patch in-  
 205 dices. We introduce learned mask tokens  $\{m_i\}_{i \in M^c}$  and  
 206 feed the concatenation of the compact feature, visible tar-

207 get tokens, and mask tokens into a decoder  $d_\phi$ . The de-  
208 coder predicts target token representations  $\{\hat{u}_{t+k}^i\}_{i \in M^c}$  for  
209 the masked regions. We then align these predicted target  
210 tokens with the corresponding frozen target representations  
211 produced by the original foundation model:

$$212 \quad \mathcal{L}_{\text{DA}} = \sum_{i \in M^c} d(\hat{u}_{t+k}^i, v_{t+k}^i), \quad (1)$$

213 where  $d(\cdot, \cdot)$  is a patch-level distance function. By requiring  
214 the compact feature to support prediction of the target scene,  
215 this loss encourages the encoder to abstract temporal and  
216 action-relevant information from the reference observation  
217 into a compact representation.

218 **Knowledge preservation loss.** While the abstraction loss  
219 specializes the representation toward temporal understand-  
220 ing, it may also distort the useful semantic structure inher-  
221 ited from the pretrained model. To prevent this drift, we  
222 introduce a retention loss that aligns the patch-wise repre-  
223 sentations of the adapted encoder with those of the frozen  
224 original encoder on the reference scene:

$$225 \quad \mathcal{L}_{\text{R}} = \sum_{i=1}^N d'(u_t^i, v_t^i), \quad (2)$$

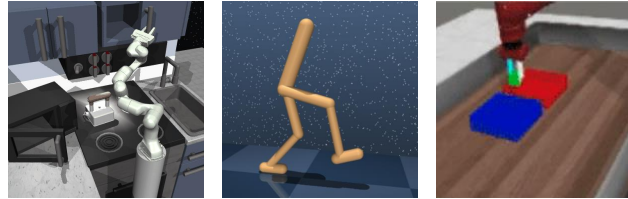
226 where  $d'(\cdot, \cdot)$  is a patch-level distance function. Since  $f_{\theta_0}$   
227 is frozen, the resulting representations provide a stable se-  
228 mantic target during post-training. This regularization pre-  
229 serves general visual knowledge while allowing the adapted  
230 encoder to specialize in robot control.

231 **Post-training objective.** The overall SPARK objective  
232 combines dynamics-aware abstraction with knowledge  
233 preservation:

$$234 \quad \mathcal{L}_{\text{SPARK}} = \mathcal{L}_{\text{DA}} + \alpha \mathcal{L}_{\text{R}}, \quad (3)$$

235 where  $\alpha > 0$  balances adaptation and preservation. Op-  
236 timizing this objective yields compact visual features that  
237 capture temporal and action-relevant information while re-  
238 maining grounded in the semantic knowledge of the original  
239 foundation model.

240 **Policy learning.** For downstream control, we train a pol-  
241 icy by behavior cloning on trajectories of image obser-  
242 vations and actions. Let  $f_\theta$  be the adapted SPARK en-  
243 coder and let  $\pi_\phi$  be a policy network. Given a trajectory  
244  $\{(I_t, o_t, a_t)\}_{t=1}^T$  with image observation  $I_t$ , optional propri-  
245 oceptive feature  $o_t$ , and action  $a_t$ , we extract a compact vi-  
246 sual feature  $z_t = f_\theta(I_t)^{[\text{CLS}]}$  and concatenate it with propri-  
247 oception to form the policy input  $s_t = [z_t; o_t]$ . The policy



(a) Franka Kitchen (b) DeepMind Control (c) Metaworld

Figure 2. **Environments for evaluation.** We visualize three vision-based robot learning benchmarks employed for the validation of our method: (a) Franka Kitchen [8], (b) DeepMind Control Suite [25], and (c) Metaworld [14].

predicts  $\hat{a}_t = \pi_\phi(s_t)$  and is trained with a standard behavior cloning objective:

$$248 \quad \mathcal{L}_{\text{BC}}(\phi) = \frac{1}{N} \sum_{i=1}^N \frac{1}{T_i} \sum_{t=1}^{T_i} \ell(\pi_\phi(s_t^{(i)}), a_t^{(i)}), \quad (4) \quad 249$$

where  $\ell$  denotes a regression loss such as mean squared error. Unless stated otherwise, the encoder  $f_\theta$  is kept fixed after post-training, and only the policy network  $\pi_\phi$  is optimized.

## 255 4. Experiments

In this section, we evaluate the impacts of the proposed method in vision-based policy learning for robotic manipulation across various environments [8, 14, 18]. Moreover, we compare our method with commonly used post-training recipes.

### 257 4.1. Evaluation suites

We evaluate on three vision-based robot learning benchmarks encompassing Franka Kitchen, DeepMind Control Suite, and Metaworld, covering a total of 15 tasks. The environments and sampled tasks for each benchmark are shown in Fig. 2.

**Franka Kitchen** [8] is a challenging robotic manipulation environment, which consists of a 9 DoF position-controlled Franka robot interacting with a kitchen scene. The environment contains various kitchenwares including an openable microwave, four turnable oven burners, an oven light switch, two hinged cabinets, and a sliding cabinet door. We consider the imitation learning evaluation setup with five imitation tasks: *Knob on*, *Light on*, *Sdoor open*, *Ldoor open*, and *Micro open*.

**DeepMind Control Suite** [25] is a set of continuous control tasks with simulated robots. We use five imitation learning cases: *Walker-stand*, *Walker-walk*, *Reacher-easy*, *Cheetah-run*, and *Finger-spin*. We report the normalized

Table 1. **Experimental results on vision-based robot policy learning on Franka Kitchen.** We report the performance of imitation learning agents on Franka Kitchen [8], which are trained upon representations of CLIP [22], DINOv2 [20], SigLIP [28], and Theia [24] post-trained on Kinetics-400 [15] dataset. The success rates (%) are reported for all the tasks. We report the gains of our method over the baselines.

Method		Knob1 on	Light on	Sdoor open	Ldoor open	Micro open	Average
CLIP	Baseline	23.0±2.6	29.5±4.7	69.5±4.4	13.5±3.4	22.0±2.8	31.5±2.4
	SPARK	31.5±1.9	38.5±6.0	74.5±11.4	18.0±2.0	29.0±6.2	39.8±5.6
	Gain	+ 8.5	+ 9.0	+ 5.0	+ 4.5	+ 7.0	+ 8.3
SigLIP	Baseline	17.0±1.2	38.5±1.9	75.5±3.4	8.5±1.0	16.5±1.9	31.2±0.6
	SPARK	27.5±5.0	43.0±5.0	87.5±1.9	20.0±2.8	33.0±5.3	42.2±2.4
	Gain	+ 10.5	+ 4.5	+ 12.0	+ 11.5	+ 16.5	+ 10.0
DINOv2	Baseline	26.0±4.3	38.5±2.5	71.0±6.8	9.5±2.5	37.0±6.2	36.4±2.7
	SPARK	21.0±4.2	51.0±8.4	81.5±2.5	17.0±2.0	33.0±3.5	40.7±2.7
	Gain	- 5.0	+ 12.5	+ 10.5	+ 7.5	- 4.0	+ 4.3
Theia	Baseline	33.5±10.5	35.5±6.6	82.5±4.7	27.5±4.7	14.0±2.0	40.1±2.3
	SPARK	35.0±2.6	56.0±6.5	75.0±5.3	34.5±2.5	21.0±2.6	44.3±2.6
	Gain	+ 1.5	+ 20.5	- 7.5	+ 7.0	+ 7.0	+ 4.2

Table 2. **Experimental results on vision-based robot policy learning on DeepMind Control Suite (DMC).** We report the performance of imitation learning agents on the DMC benchmark [25], which are trained upon representations of CLIP [22], DINOv2 [20], and SigLIP [28] post-trained on Kinetics-400 [15] dataset. The normalized scores are reported for all the tasks. We report the gains of our method over the baseline.

Method		Walker-stand	Walker-work	Reacher-easy	Cheetah-run	Finger-spin	Average
CLIP	Baseline	58.1±5.8	24.9±2.8	75.2±4.0	12.7±2.0	68.4±3.1	47.9±2.1
	SPARK	62.0±3.5	29.5±4.5	75.6±4.7	22.0±2.1	63.7±1.6	50.6±0.9
	Gain	+ 3.9	+ 4.6	+ 0.4	+ 9.3	- 4.7	+ 2.7
SigLIP	Baseline	43.9±5.9	19.9±2.8	74.6±2.5	9.0±2.5	63.6±0.4	42.1±1.5
	SPARK	45.2±6.2	17.5±2.4	75.2±3.0	8.4±2.1	66.8±1.2	42.6±1.8
	Gain	+ 1.3	-2.4	+ 0.6	-0.6	+ 3.2	+ 0.5
DINOv2	Baseline	74.2±5.6	41.5±3.4	76.8±5.9	15.9±3.7	68.3±1.6	55.3±0.4
	SPARK	78.8±3.5	45.2±5.0	74.3±6.5	22.3±2.9	70.8±1.0	58.3±2.4
	Gain	+ 4.6	+ 3.7	- 2.5	+ 6.4	+ 2.5	+ 3.0

282 scores for all tasks.

283

284 **MetaWorld** [27] is a suite of simulated robotic manipula-  
 285 tion tasks with a Sawyer robot arm. We focus on a subset  
 286 of five representative tasks: *Assembly*, *Bin-picking*, *Button-*  
 287 *press-topdown*, *Drawer-open*, and *Hammer*. We measure  
 288 the best success rates among the online evaluation trials.

## 289 4.2. Implementation details

290 **Initialization and post-training setup.** For post-training,  
 291 we initialize the SPARK encoder from public foundation  
 292 models, including CLIP [22], DINOv2 [20], SigLIP [28],  
 293 and Theia [24]. Unless otherwise stated, post-training is  
 294 conducted on Kinetics-400 [15] for 50 epochs. Following  
 295 repeated sampling [6, 12], this corresponds to an effective  
 296 training length of 200 epochs. We use AdamW [16] with

a batch size of 1536 and train on dynamic scenes at a res-  
 297 olution of  $224 \times 224$ . Frames are sampled at 30 FPS, and  
 298 the temporal offset between the reference and target scenes  
 299 is uniformly sampled between 4 and 96 frames. We apply  
 300 random resized cropping and horizontal flipping, using the  
 301 same crop for both reference and target scenes. The target  
 302 scene is masked with a ratio of 0.9. The decoder consists  
 303 of eight Vision Transformer blocks. All remaining hyperpa-  
 304 rameters follow the default pretraining settings of the corre-  
 305 sponding backbone. 306

**Policy learning setup.** For downstream control, each  
 307 agent consists of a frozen visual encoder and a policy net-  
 308 work trained with behavior cloning. The policy takes as  
 309 input the compact visual feature extracted from the encoder,  
 310 together with proprioceptive input when available. For  
 311

Table 3. **Experimental results on vision-based robot policy learning on Metaworld.** We report the performance of imitation learning agents on Metaworld [27], which are trained upon representations of CLIP [22], DINOv2 [20], and SigLIP [28] post-trained on Kinetics-400 [15] dataset. The success rates (%) are reported for all the tasks. We report the gains of our method over the baseline.

Method		assembly	bin-picking	button-press	drawer-open	hammer	Average
CLIP	Baseline	55.2±6.6	63.2±7.7	58.4±13.4	100.0±0.0	87.2±15.6	72.8±1.1
	SPARK	72.0±12.0	63.2±7.0	58.8±14.5	100.0±0.0	94.4±6.1	77.7±7.0
	Gain	+ 16.8	+ 0.0	+ 0.4	0.0	+ 7.2	+ 4.9
SigLIP	Baseline	69.6±9.2	52.0±12.3	59.2±15.6	100.0±0.0	84.0±9.4	73.0±7.6
	SPARK	74.6±10.3	60.1±9.5	58.5±12.7	100.0±0.0	91.5±5.8	76.9±6.8
	Gain	+ 5.0	+ 8.1	-0.7	0.0	+ 7.5	+ 3.9
DINOv2	Baseline	64.8±5.9	70.4±14.9	64.0±19.4	100.0±0.0	92.8±5.9	78.4±7.7
	SPARK	80.8±5.2	67.2±7.2	68.8±16.3	100.0±0.0	99.2±1.8	83.2±5.5
	Gain	+ 16.0	- 3.2	+ 4.8	0.0	+ 6.4	+ 4.8

Table 4. **Experimental results on vision-based robot policy learning on Franka Kitchen.** We report the performance of imitation learning agents on Franka Kitchen [8], which are trained upon representations from the ViT-S/16 model pre-trained on Kinetics-400 [15] dataset. The success rates (%) are reported for all the tasks. We underline the second-best performance. We report the gains of our method over the second-best baseline.

Methods	Architecture	Franka Kitchen	DeepMind Control	Metaworld	Average
SimCLR	ViT-S	38.7	39.7	78.4	52.3
MoCo v3	ViT-S	25.4	43.7	65.4	44.8
DINO	ViT-S	38.7	50.9	82.4	57.3
MAE	ViT-S	26.1	43.7	65.4	45.1
SiamMAE	ViT-S	30.4	56.0	81.1	55.8
CLIP	ViT-B	31.5	47.9	72.8	50.7
SigLIP	ViT-B	31.2	42.1	73.0	48.8
DINOv2	ViT-B	36.4	55.3	78.4	56.7
SPARK (CLIP)	ViT-B	39.8	50.6	77.7	56.0
SPARK (SigLIP)	ViT-B	42.2	42.6	76.9	53.9
SPARK (DINOv2)	ViT-B	<b>40.7</b>	<b>58.3</b>	<b>83.2</b>	<b>60.7</b>

312 Franka Kitchen, the policy network is implemented as a  
 313 two-layer MLP with batch normalization. For MetaWorld  
 314 and DMC, we follow the policy-learning protocols of [18].  
 315 Unless otherwise stated, the visual encoder is frozen after  
 316 post-training and only the policy network is optimized.

317 **Benchmark-specific settings.** For Franka Kitchen, we  
 318 evaluate five imitation learning tasks following prior pro-  
 319 tocols [19, 21]. We use either the left or right camera at a  
 320 resolution of  $224 \times 224$  without depth input. Training runs  
 321 for 20,000 steps, with online evaluation every 1,000 steps.  
 322 We report the best success rate during training, averaged  
 323 over four random seeds.

324 For MetaWorld and DMC, we follow the setup of [18].  
 325 Proprioceptive input is provided for all benchmarks except  
 326 DMC. Each agent is trained for 100 epochs and evalu-  
 327 ated online every 5 epochs. We report normalized scores  
 328 for DMC and final success rates for MetaWorld, averaged  
 329 across five random seeds.

### 4.3. Vision-based Robot Policy Learning in Real-world Environments

330 **Quantitative results.** To evaluate whether the benefits of  
 331 SPARK extend beyond simulation, we additionally conduct  
 332 experiments on real-world robot manipulation tasks. We  
 333 consider four imitation learning tasks: *Cloth Folding*, *Cup*  
 334 *Stacking*, *Cabinet Opening*, and *Drawer Closing*. For each  
 335 task, we collect 20 demonstration episodes for behavior  
 336 cloning and evaluate policies learned from visual represen-  
 337 tations extracted by the frozen backbone. The robot oper-  
 338 ates at 20 Hz with a policy’s chunk size of 20. Following  
 339 the same evaluation protocol as in simulation, we compare  
 340 the vanilla CLIP encoder with its SPARK-adapted counter-  
 341 part with 20 trials for each task. The results are reported  
 342 in Table 5. SPARK improves the success rate on three out  
 343 of four tasks, increasing performance from 65.0 to 75.0 on  
 344 *Cloth Folding*, from 45.0 to 55.0 on *Cup Stacking*, and from  
 345 50.0 to 55.0 on *Cabinet Opening*. On *Drawer Closing*, the  
 346 SPARK-adapted model achieves comparable performance  
 347  
 348

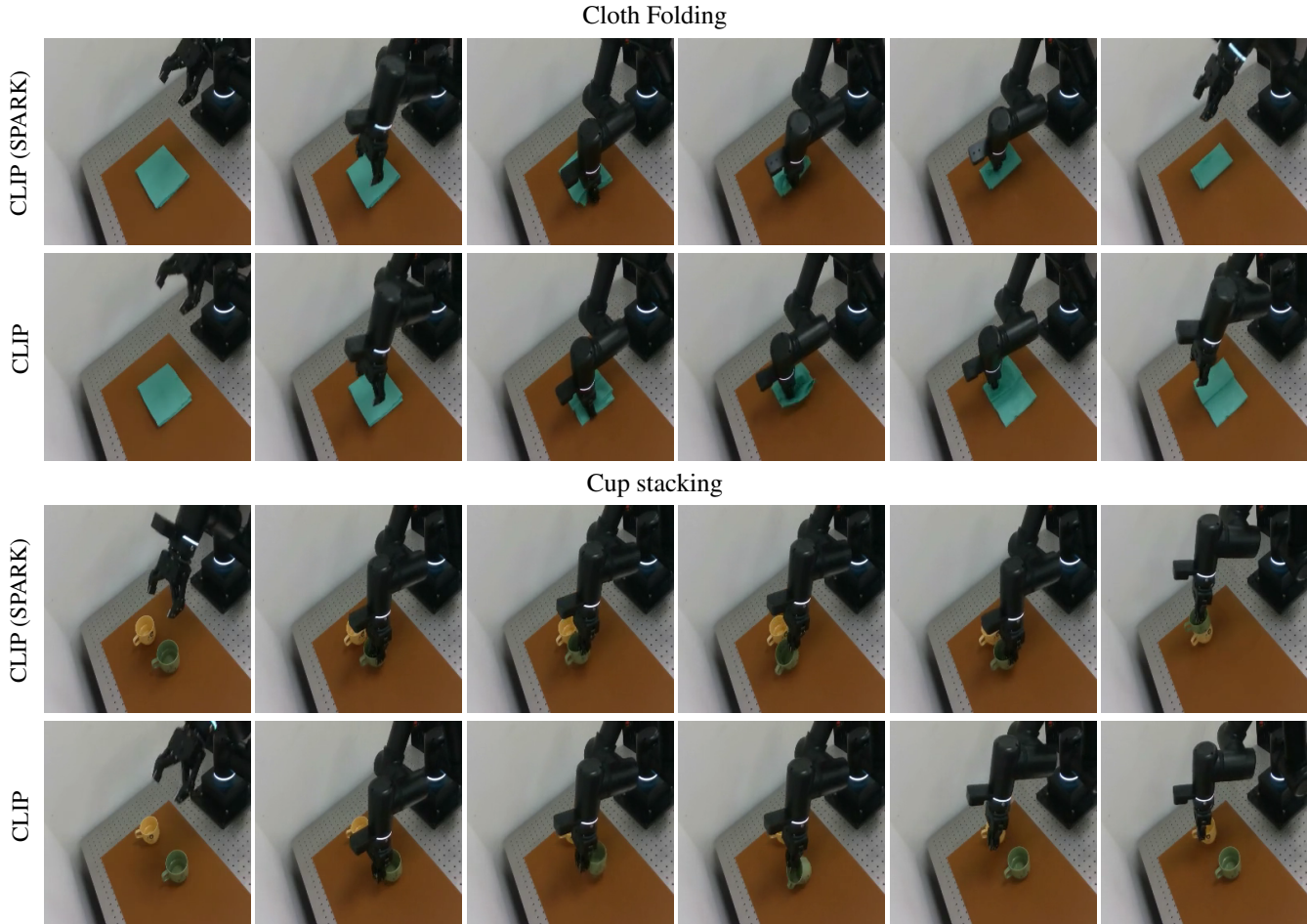


Figure 3. **Real-world demonstration trajectories.** Each row shows the temporal progression of a policy from the initial observation to intermediate interaction states and the final outcome. Across tasks, the SPARK-adapted policy exhibits stable approach, contact, and object manipulation behaviors, illustrating that the adapted visual representations capture control-relevant scene information that transfers reliably to real-world settings.

Table 5. **Performance on real-world vision-based robot policy learning.** Success rates (%) of imitation learning agents on three manipulation tasks: Cloth Folding, Cup Stacking, Cabinet Opening, and Drawer Closing. We compared CLIP [22] with our SPARK model, post-trained upon representations of CLIP on Kinetics-400 [15]. The results demonstrate the generalizability of SPARK in the real-world.

Method	Cloth Folding	Cup Stacking	Cabinet Opening	Drawer Closing
CLIP	65.0	45.0	50.0	55.0
SPARK (CLIP)	75.0	55.0	55.0	50.0

349 to the baseline. Overall, these results show that the gains obtained by SPARK are not limited to simulation, but transfer  
350 to real-world robot manipulation as well. This suggests that  
351 post-training foundation models with dynamics-aware ab-  
352

353 straction and knowledge preservation improves the useful-  
354 ness of pretrained visual representations in practical robotic  
355 settings.

356 **Qualitative results.** We further visualize representative  
357 real-world rollouts of the SPARK-adapted policy across  
358 multiple manipulation tasks in Fig. 3. Each row presents a  
359 temporal sequence from the initial scene, through interme-  
360 diate interaction stages, to the final task outcome. The visu-  
361 alizations show that the policy maintains coherent manipu-  
362 lation progress throughout the rollout, including approach,  
363 contact, and task completion. Across different tasks and  
364 scene configurations, these examples suggest that SPARK  
365 produces visual representations that preserve the informa-  
366 tion most relevant for control, enabling robust real-world  
367 execution under variations in object pose, scene layout, and  
368 appearance.

Table 6. **Ablation study varying post-training recipes** We compared our SPARK to other post-training recipes on Franka Kitchen [8] using CLIP [22] as an initial model. The success rates (%) are reported for all the tasks.

Post-training	Franka Kitchen
Baseline	31.5
Masked autoencoding	23.0
Correspondence matching	25.5
SPARK	39.8

#### 4.4. Comparison to prior visual pretraining methods

We compare SPARK against conventional self-supervised pretraining methods for visual representations pretrained on dynamic scenes (e.g., consecutive frames), including SimCLR [2], MoCo v3 [4], DINO [1], MAE [11], and a dynamics-oriented variant, SiamMAE [9]. This evaluation verifies whether explicitly learning a compact state representation via conservative summarization is preferable to generic SSL objectives trained from scratch. As shown in Table 4, applying SPARK to a strong public encoder yields the best overall performance across suites, with SPARK (DINOv2) ranking first and SPARK (CLIP) substantially narrowing the gap to the strongest SSL baselines. Overall, adapting public encoders with SPARK is more effective for downstream robotics tasks than relying on generic SSL features, including methods designed for dynamic scenes.

## 5. Analysis

**Comparison to various post-training recipes** We conduct comparative analysis with various post-training recipes on Franka Kitchen in Table 6. We compared naive post-training using masked autoencoding (MAE) [11, 18] and correspondence matching among consecutive scenes [9]. All the post-trained models are initialized by CLIP [22]. As shown in the table, standard MAE style post-training and correspondence matching style post-training, which focus on generic reconstruction and temporal correspondence respectively, perform worse than the baseline with averages of 23.0% and 25.5%. In contrast, SPARK attains an average success rate of 39.8%, substantially outperforming both the vanilla CLIP and other post-trained models. This highlights that explicitly enforcing conservative summarization and knowledge preservation is more effective for adapting foundation models to robot control than generic reconstruction based objectives.

## 6. Conclusion

In this paper, we proposed SPARK, a simple post-training method for adapting pretrained foundation models to robot

control. SPARK is based on two key principles: dynamics-aware abstraction and knowledge preservation. Dynamics-aware abstraction encourages the learned representation to capture temporal changes relevant to control, while knowledge preservation prevents the adapted model from losing useful semantic knowledge inherited from large-scale pre-training. In this way, SPARK makes pretrained visual representations more suitable for robot control without requiring training from scratch. Experiments on multiple robot learning benchmarks show that SPARK consistently improves over vanilla pretrained backbones. We further show that these gains extend beyond simulation to real-world robot manipulation tasks. These results suggest that cost-efficient post-training is a practical way to improve the usefulness of foundation models for embodied intelligence.

## References

- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2021. 2, 8
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, 2020. 8
- Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. *arXiv preprint arXiv:2104.02057*, 2021. 2, 8
- Sudeep Dasari, Mohan Kumar Srirama, Unnat Jain, and Abhinav Gupta. An unbiased look at datasets for visuo-motor pre-training. In *Conference on Robot Learning*. PMLR, 2023. 1, 2
- Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3299–3309, 2021. 5
- Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *arXiv preprint arXiv:2110.04544*, 2021. 2
- Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. In *Proceedings of the Conference on Robot Learning*, pages 1025–1037. PMLR, 2020. 4, 5, 6, 8
- Agrim Gupta, Jiajun Wu, Jia Deng, and Li Fei-Fei. Siamese masked autoencoders. In *Advances in Neural Information Processing Systems*, 2023. 8
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual repre-

- 461 sentation learning. *arXiv preprint arXiv:1911.05722*, 2019. 2
- 462
- 463 [11] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr
- 464 Dollár, and Ross Girshick. Masked autoencoders are scalable
- 465 vision learners. In *Proceedings of the IEEE/CVF conference*
- 466 *on computer vision and pattern recognition*, 2022. 2, 8
- 467 [12] Elad Hoffer, Tal Ben-Nun, Itay Hubara, Niv Giladi, Torsten
- 468 Hoeffler, and Daniel Soudry. Augment your batch: Improving
- 469 generalization through instance repetition. In *Proceedings of*
- 470 *the IEEE/CVF Conference on Computer Vision and Pattern*
- 471 *Recognition*, 2020. 5
- 472 [13] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-
- 473 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
- 474 LoRA: Low-rank adaptation of large language models. In *Inter-*
- 475 *national Conference on Learning Representations*, 2022.
- 476 2
- 477 [14] Stephen James, Zicong Ma, David Rovick Arrojo, and An-
- 478 drew J Davison. Rlbenc: The robot learning benchmark &
- 479 learning environment. *IEEE Robotics and Automation Let-*
- 480 *ters*, 5(2):3019–3026, 2020. 2, 4
- 481 [15] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang,
- 482 Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola,
- 483 Tim Green, Trevor Back, Paul Natsev, et al. The kinetics hu-
- 484 man action video dataset. *arXiv preprint arXiv:1705.06950*,
- 485 2017. 5, 6, 7
- 486 [16] Ilya Loshchilov and Frank Hutter. Decoupled weight decay
- 487 regularization. In *International Conference on Learning Rep-*
- 488 *resentations*, 2019. 5
- 489 [17] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Os-
- 490 bert Bastani, Vikash Kumar, and Amy Zhang. Vip: Towards
- 491 universal visual reward and representation via value-implicit
- 492 pre-training. *arXiv preprint arXiv:2210.00030*, 2022. 1, 2
- 493 [18] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud,
- 494 Yecheng Jason Ma, Claire Chen, Sneha Silwal, Aryan
- 495 Jain, Vincent-Pierre Berges, Pieter Abbeel, Jitendra Malik,
- 496 et al. Where are we in the search for an artificial visual
- 497 cortex for embodied intelligence? In *Advances in neural*
- 498 *information processing systems*, 2023. 2, 4, 6, 8
- 499 [19] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea
- 500 Finn, and Abhinav Gupta. R3m: A universal visual repre-
- 501 sentation for robot manipulation. In *6th Annual Conference*
- 502 *on Robot Learning*, 2022. 1, 2, 6
- 503 [20] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy
- 504 Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,
- 505 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al.
- 506 Dinov2: Learning robust visual features without supervision.
- 507 *arXiv preprint arXiv:2304.07193*, 2023. 1, 2, 5, 6
- 508 [21] Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam,
- 509 and Abhinav Gupta. The unsurprising effectiveness of pre-
- 510 trained vision models for control. In *Proceedings of the*
- 511 *39th International Conference on Machine Learning*, pages
- 512 17359–17371. PMLR, 2022. 1, 2, 6
- 513 [22] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
- 514 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
- 515 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learn-
- 516 ing transferable visual models from natural language supervi-
- 517 sion. In *International conference on machine learning*, pages
- 518 8748–8763. PmLR, 2021. 1, 2, 5, 6, 7, 8
- [23] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel,
- Jitendra Malik, and Trevor Darrell. Real-world robot learn-
- ing with masked visual pre-training. In *Proceedings of The*
- 6th Conference on Robot Learning*, pages 416–426. PMLR,
2023. 1, 2
- [24] Jinghuan Shang, Karl Schmeckpeper, Brandon B. May,
- Maria Vittoria Minniti, Tarik Kelestemur, David Watkins,
- and Laura Herlant. Theia: Distilling diverse vision founda-
- tion models for robot learning. 2024. 1, 5
- [25] Yuval Tassa, Saran Tunyasuvunakool, Alistair Muldal,
- Yotam Doron, Siqi Liu, Steven Bohez, Josh Merel, Tom
- Erez, Timothy Lillicrap, and Nicolas Heess. dm\_control:
- Software and tasks for continuous control. *arXiv preprint*
- arXiv:2006.12983*, 2020. 4, 5
- [26] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muham-
- mad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil
- Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil
- Mustafa, Olivier Hénaff, Jeremiah Harmsen, Andreas
- Steiner, and Xiaohua Zhai. Siglip 2: Multilingual vision-
- language encoders with improved semantic understand-
- ing, localization, and dense features. *arXiv preprint*
- arXiv:2502.14786*, 2025. 1, 2
- [27] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian,
- Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-
- world: A benchmark and evaluation for multi-task and meta
- reinforcement learning. In *Conference on Robot Learning*,
2020. 5, 6
- [28] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lu-
- cas Beyer. Sigmoid loss for language image pre-training. In
- Proceedings of the IEEE/CVF International Conference on*
- Computer Vision (ICCV)*, pages 11975–11986, 2023. 1, 2, 5,
- 6
- [29] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Zi-
- wei Liu. Conditional prompt learning for vision-language
- models. In *IEEE/CVF Conference on Computer Vision and*
- Pattern Recognition (CVPR)*, 2022. 2
- [30] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei
- Liu. Learning to prompt for vision-language models. *Inter-*
- national Journal of Computer Vision (IJCV)*, 2022. 2