IGOR: IMAGE-GOAL REPRESENTATIONS ARE THE ATOMIC CONTROL UNITS FOR FOUNDATION MODELS IN EMBODIED AI

Anonymous authors

Paper under double-blind review



Figure 1: **Image-GOal Representations (IGOR) based training framework for embodied AI.** IGOR learns a unified latent action space for humans and robots by compressing visual changes between an image and its goal state on data from both robot and human activities. By labeling latent actions, IGOR facilitates the learning of foundation policy and world models from internet-scale human video data, covering a diverse range of embodied AI tasks. With a semantically consistent latent action space, IGOR enables human-to-robot generalization. The foundation policy model acts as a high-level controller at the latent action level, which is then integrated with a low-level policy to achieve effective robot control.

Abstract

We introduce Image-GOal Representations (IGOR), aiming to learn a unified, semantically consistent action space across human and various robots. Through this unified latent action space, IGOR enables knowledge transfer among large-scale robot and human activity data. We achieve this by compressing visual changes between an initial image and its goal state into latent actions. IGOR allows us to generate latent action labels for internet-scale video data. This unified latent action space enables the training of foundation policy and world models across a wide variety of tasks performed by both robots and humans. We demonstrate that: (1) IGOR learns a semantically consistent action space for both human and robots, characterizing various possible motions of objects representing the physical interaction knowledge; (2) IGOR can "migrate" the movements of the object in the one video to other videos, even across human and robots, by jointly using the latent action model and world model; (3) IGOR can learn to align latent actions with natural language through the foundation policy model, and integrate latent actions with a low-level policy model to achieve effective robot control. We believe IGOR opens new possibilities for human-to-robot knowledge transfer and control. See video demonstrations on our anonymous webpage.

043

044

045

046

048

000

002

004

006

1 INTRODUCTION

055 056

Learning foundation models for embodied AI has been notably constrained by a lack of interaction 057 data. Unlike text or video data, which are abundantly available, interaction data is much scarcer. Research efforts have been devoted to creating large-scale interaction datasets, such as Open-X-Embodiment (Collaboration et al., 2023) and DROID (Khazatsky et al., 2024). Based on multi-task 060 interaction data, a series of generalist agents (or foundation policy models) have been proposed, 061 such as RT-1 (Brohan et al., 2022), Robocat (Bousmalis et al., 2023), RT-2 (Brohan et al., 2023), 062 Octo (Team et al., 2024), and OpenVLA (Kim et al., 2024). However, the volume of interaction data remains several orders of magnitude smaller than that of internet text or video data. Given that 063 the success of foundation models relies on scaling up datasets and extracting knowledge from such 064 large-scale datasets, it is essential to design methods for building embodied AI foundation models 065 that can effectively utilize internet-scale video data. 066

067 Internet-scale video data contains abundant sequential records of human activities and perfect 068 demonstrations of how human perform various tasks by interacting with the real world. When the human brain extracts information from videos, instead of doing it frame by frame, it modular-069 izes the differences between frames into a single word such as "move", "open", "close". We refer to these highly compressed, modularized actions as latent actions that are shared across different 071 tasks. The question to ask here is, is it possible to recover latent actions from video datasets with 072 humans and robots performing various real embodied AI tasks? While recent works such as 073 Genie (Bruce et al., 2024) and LAPO (Schmidt & Jiang, 2023) made attempts in recovering such 074 latent actions from videos, they primarily focus on 2D platformer games where each latent action 075 a_t corresponds to a specific control button. The action space is highly designed to fit a specific sce-076 nario and incomparable to the complex human and robot action space in various embodied AI tasks. 077 To take a step further, the question would be, can we learn a unified, semantically consistent latent action space, allowing the transfer of knowledge across different tasks, and embodiments including human and various robots? 079

In this paper, we propose Image-GOal Representations (IGOR), which learns a unified and seman-081 tically consistent latent action space shared across different tasks and embodiments, enabling the knowledge transfer among internet-scale video data. We propose a latent action model designed 083 to capture robot and human actions across various embodied AI tasks. IGOR compresses the vi-084 sual changes between an image and its goal state into latent actions, which are also embeddings of 085 sub-tasks defined by reaching the goal from the initial image. IGOR is trained by minimizing the reconstruction loss of the goal state, which is predicted based on the image and the latent action. 086 The core insight behind IGOR is that if compressed sufficiently, the image-goal pairs with similar 087 visual changes will have similar embeddings. 880

We argue that, besides text embeddings for human instruction understanding and im-

age/video embeddings for state understanding, image-goal representations for latent action

learning and sub-task understanding are yet another crucial building blocks, which may

089

090

091 092

002

094

With the latent action model, we can transform internet-scale human video data into interaction 096 data labeled with latent actions, which largely expands the data available to building embodied 097 AI foundation models. This unified latent action space allows us to train foundation policy and 098 world models on nearly arbitrary tasks performed by robots and humans. Specifically, we train a foundation policy model on large-scale video data with text labels. This model uses text to describe tasks and makes decisions, generating the next latent action to perform. Additionally, we train a 100 foundation world model on the same dataset, learning to simulate the outcome of executing the 101 foundation policy model. Image-Goal Representations can be viewed as atomic control units in 102 visual space. They function both as latent actions for a foundational policy model to predict in 103 visual trajectory planning and as sub-tasks for a robot-specific low-level policy to execute. 104

hold great potential for becoming the atomic control unit in embodied AI.

We train our models on human video data and robot data with actions removed, with RT-1 dataset held out for OOD evaluation. First, we evaluate the latent action model qualitatively, and find that image-goal pairs with similar latent actions have similar visual changes, corresponding to semantically consistent movements, even on OOD scenarios. Then we evaluate the world model by

2



Figure 2: We extract latent actions from Image-Goal pairs in the solid line boxes, and apply the latent actions to different initial frames, generating subsequent videos via world model as shown in the corresponding dashed boxes. The first half illustrates examples from real-world videos with diverse object categories, while the second half demonstrates generalization from human to robot arms. Full videos are available on our website.

135

136

137

138

141 extracting latent actions from a video and applying such latent action (or action sequence) to the 142 initial frames of other videos, generating the rest of frames. We find that, jointly with the latent 143 action model and world model, IGOR successfully "migrates" the movements of the object in the one video to other videos, as shown in Figure 2. We also apply different latent actions to the same 144 initial image, and find that the world model has learned various possible movements of the object in 145 the image, suggesting that it has absorbed the physical interaction knowledge. For the foundation 146 policy model, we show its ability in following diverse language instruction via iteratively rolling out 147 the foundation policy and world model using latent actions. We further integrate it with a low-level 148 policy, and show that IGOR-based policy training can improve performance on Google Robot tasks 149 in low-data regime with the SIMPLER (Li et al., 2024) simulator.

150 151 152

153

2 Methodology

154 2.1 LATENT ACTION MODEL

The primary objective of the latent action model is to label latent actions from unlabeled opendomain videos in an unsupervised manner. Given a sequence of video frames $o_{1:t+1}$, the goal is to derive the latent action a_t , which captures the key information describing only the changes that occur at time step t, removing other redundant information. In contrast to prior works (Schmidt & Jiang, 2023; Bruce et al., 2024), which primarily focus on 2D platformer games where each latent action a_t corresponds to a specific control button, we aim to develop a more generalizable model. Our model is designed to handle the significantly greater complexity of open-world scenarios, where latent actions may not correspond to any specific underlying actions. This presents several additional
 challenges.

First, rather than focusing solely on absolute position of pixel changes, the latent action model must learn to capture semantic movements that remain consistent across varying scenarios. Moreover, due to the temporal redundancy, actions are often sparse given long contexts, which can lead the model to infer o_{t+1} directly from the history, bypassing the need for a more informative latent action a_t .

To address these issues, we propose a novel model architecture. Our latent action model consists of a pair of Inverse Dynamics Model (IDM) and Forward Dynamics Model (FDM). IDM I is trained 170 to predict the latent action a_t based on the full sequence of observations $o_{1:t+1}$. Instead of using 171 the raw observations, we first apply random cropping c_1 to the inputs: $a_t = I(c_1[o_{1:t+1}])$. For 172 the architecture of I, we first extract features for each frame through Vision Transformer (ViT) 173 (Dosovitskiy et al., 2021) and then adopt a Spatio-Temporal transformer (ST-transformer) (Bruce 174 et al., 2024; Xu et al., 2021) with a temporal causal mask as the backbone. Learnable readout tokens 175 are then used to extract and compress the visual changes into N tokens. To further compress the 176 information stored in latent action, we apply vector quantization to each token, restricting them to a discrete codebook of size |C|. Finally, we derive the latent action $a_t \in \mathbf{R}^{N \times D}$ where D is the 177 dimension of each code. We refer to a_t as the latent action embedding, or sub-task embedding, as 178 179 they describe the information that takes the observation o_t to the next observation o_{t+1} .

For the FDM F, we propose using a single-frame Vision Transformer to reconstruct o_{t+1} , in contrast 181 to previous works (Schmidt & Jiang, 2023; Bruce et al., 2024), which reconstruct the next frame 182 given the entire context $o_{1:t}$. This approach mitigates the case where the model might predict the 183 next frame directly from the context, bypassing the latent action. By conditioning on a single frame, 184 it encourages more information to flow into the latent action a_t . For reconstruction, we apply another 185 random cropping c_2 , and the next frame is predicted as $\hat{o}_{t+1} = F(c_2[o_t], a_t)$. By using different croppings c_1 and c_2 , the model is encouraged to learn a more semantically invariant latent action across different trajectories. The models are trained jointly with the reconstruction loss $||c_2|o_{t+1}|$ – 187 $\hat{o}_{t+1} \parallel^2$ and the commitment loss in vector quantization. 188

189

211 212

215

190 2.2 FOUNDATION WORLD MODEL

Our foundation world model is a continuous-time Rectified Flow (Liu et al., 2023b; Esser et al., 2024) that learns to predict the future frames $o_{t+1:T}$ conditioned on the history observation frame $o_{1:t}$, and future latent actions $a_{t:T-1}$. To achieve this goal, there are two key challenges: 1) Generating the photo-realistic frame that describes the states precisely; 2) Controlling the generated frames by the latent actions.

Accordingly, we start our foundation world model with the pre-trained Open-Sora (Zheng et al., 197 2024). It consists of two components: a 3D Variational AutoEncoder (VAE) that encodes the raw observation into latent space with 8×8 times downsampling in spatial dimension and $4 \times$ times 199 downsampling in temporal dimension; a Spatial-Temporal Rectified Flow Transformer (ST-RFT) 200 that generates the latent from the text conditions. To enable the control from the observation and 201 action, we make two modifications to the original Open-Sora: 1) We replace the original text input 202 of the pre-trained model with our latent actions $a_{1:T}$ obtained from LAM. Zero-padding is applied 203 for the last action. For each frame, we map the latent actions into a single token and feed it to the ST-204 RFT via the cross-attention mechanism; 2) We also make the generation conditioned on the output of FDM $\hat{o}_{t+1:T}$, which provides a coarse-grained prediction according to the input latent action. For 205 the conditioning of $\hat{o}_{t+1:T}$, we encode it to the latent space with the same 3D VAE and directly add 206 it to the noisy input element-wise. 207

Formally, Rectified Flow (Liu et al., 2023b; Albergo & Vanden-Eijnden, 2023; Esser et al., 2024) aims at directly regressing a vector field that generates a probability path between noise distribution and data distribution. For $n \in [0, 1]$, we define the interpolation between the two distributions as:

$$\boldsymbol{x}_n = (1-n)\boldsymbol{x}_0 + n\boldsymbol{x}_1,\tag{1}$$

where x_0 is the clean data, x_1 is the sampled noise, and x_n is the noisy data. During training, we train a vector-valued neural network x_{θ} with L2 loss:

$$\mathbb{E}_{n,\boldsymbol{x}_0,\boldsymbol{x}_1} \| \boldsymbol{x}_0 - \boldsymbol{x}_{\theta}(\boldsymbol{x}_n, n, a_{t:T-1}, \hat{o}_{t+1:T}) \|^2.$$
(2)

Instead of predicting the conditional expectation directly, we follow Liu et al. (2023b) to parameterize the velocity with a neural network v_{θ} and train it on:

$$L_{\text{world}}(\theta) = \mathbb{E}_{n,\boldsymbol{x}_0,\boldsymbol{x}_1} \| (\boldsymbol{x}_1 - \boldsymbol{x}_0) - \boldsymbol{v}_{\theta}(\boldsymbol{x}_n, n, a_{t:T-1}, \hat{o}_{t+1:T}) \|^2.$$
(3)

It should be noted that, our foundation world model can be fine-tuned to accommodate the different
 action spaces of robots with various embodiments. The fine-tuning of the foundational world model
 is left as future work.

2.3 FOUNDATION POLICY MODEL AND LOW-LEVEL POLICY MODEL

226 The training of the policy model consists of two stages. In the first pretraining stage, taken as input 227 the raw observation frames $o_{1:t}$ and a textual description s for the task, the foundation policy model 228 predicts latent actions $a_t = I([o_{1:t+1}])$ labeled by the IDM in the latent action model at each step. 229 The training dataset of this stage is the same as that used for the latent action model, i.e., with 230 large-scale and diverse sources of videos. In the second finetuning stage, we add an extra prediction 231 component on the foundation policy model to predict real continuous robot actions, with taking the raw observations as well as the latent actions predicted by the first stage model as input. In this 232 stage, only the prediction component (i.e., the low-level policy model) is optimized on small-scale 233 and task-specific downstream datasets, while other components are frozen. 234

Specifically, similar to the latent action model, the backbone of foundation policy model is also a ST-transformer equipped with a ViT image encoder, with a feed-forward layer as the final prediction layer. The textual description *s* is encoded to a latent representation by a pre-trained text encoder, which is concatenated with the observation representation encoded by the ViT encoder as the joint input to the model. We use the L2 distance between the predicted hidden output and the latent action as the loss function. Given a trajectory consists of *t* observations $o_{1:t}$, the training objective can be written as:

$$L_{policy} = \|P([s; o_{1:t}]) - a_t\|^2, \tag{4}$$

where $P(\cdot)$ denotes the policy model.

219

224

225

242

254 255

256 257

258 259

260 261

262

263

244 In the second stage, we train the low-level policy model to predict the real continuous actions 245 within each latent action, where the image-goal latent actions can be seen as representations for 246 sub-tasks defined by reaching a goal from an initial image. The low-level policy model is also an ST-transformer with a prediction layer. The input consists of the textual representation s, the ob-247 servation $o_{1:t}$ and latent actions predicted by the foundation policy model $P([s; o_{1:t}])$, which are 248 concatenated together at the patch level as one part. The latent action $P([s; o_{1:t}])$ predicted by the 249 foundation policy model also serves as sub-task embedding for the low-level policy model. We de-250 note that each latent action corresponds to τ real robot actions, and the latent action a_t corresponds 251 to real robot action $u_t^{1:\tau}$. Denote the low-level policy model as $P_f(\cdot)$, we train the second stage 252 model also by L2 distance: 253

$$L_{ft} = \|P_f([s; P([s; o_{1:t}]); o_{1:t}]) - u_t^{1:\tau}\|^2,$$
(5)

where only the parameters of the low-level policy are optimized.

3 EXPERIMENTS

3.1 DATASET

In the pretraining stage, we construct a large-scale dataset comprising diverse domains, including robotic data from various embodiments and a substantial amount of human activity videos.

Data Mixture. For the robotic data, we select a subset of Open-X Embodiment dataset (Collaboration et al., 2023) with single arm end-effector control, excluding RT-1 dataset for out-of-distribution (OOD) evaluation. We follow the preprocessing and data mixture weights from Team et al. (2024);
Kim et al. (2024). In total, we utilize approximately 0.8M robot trajectories. While our dataset includes data from real robots, we discard the associated actions and proprio-states, using only image frames and text instructions during pretraining. Additionally, we incorporate large-scale open-world videos with language instructions, including human daily activities from Something-Something v2

<section-header>Laten Action I (pene the grippen)Image: A perpension of the perpension of

Figure 3: Image-goal pairs with similar latent actions in OOD RT-1 dataset. In each row, we choose the leftmost image-goal pair, and retrieve 3 nearest pairs on latent action embedding. The original task instructions of the pairs are shown under the images. We find that each row shares the similar visual changes semantically, and the latent actions generalize across different raw language tasks.

(Goyal et al., 2017), and egocentric videos such as EGTEA (Li et al., 2018), Epic Kitchen (Damen et al., 2020), and Ego4D (Grauman et al., 2022; Pramanick et al., 2023). In total, we derive approximately 2.0M human activity video clips with high quality. Overall, our dataset for pretraining comprises around 2.8M trajectories and video clips, where each trajectory contains a language instruction and a sequence of observations.

Data Preprocessing. In practice, we found that the video quality has a big impact on the model performance. We exclude low-quality videos characterized by excessive shakiness or rapid camera movement, and apply stabilization techniques to the remaining videos. To ensure proper amount of changes between frames in the latent action model, we choose the optimal frame rates for robotics dataset and human activity videos.

In the finetuning stage, we use RT-1 dataset, a large-scale dataset for real-world robotic experiences. We uniformly sample 1% number of episodes from RT-1 dataset for finetuning, where each episode comprises a language instruction, a sequence of image observations, and a sequence of low-level actions. The action space is 7-dimensional, including 3 dimensions of robot arm movement ΔPos , 3 dimensions of robot arm rotation ΔRot , and 1 dimension of robot gripper action ΔGrp . We provide more details in Appendix A.

311 3.2 TRAINING DETAILS

We first pretrain our latent action model on our pretraining dataset. Then, we use the pretrained latent action model to label latent actions on our pretraining dataset, and pretrain foundation policy model and foundation world model on the labeled dataset. Finally, we finetune our low-level policy model on top of our pretrained models on RT-1.

For latent action model, we use a codebook with N = 4 tokens, and codebook size of |C| = 32, each with an embedding size of D = 128. We use a sub-task length of $\tau = 3$ for finetuning the low-level policy model on RT-1 dataset. Please refer to Appendix B for more training details.

320

322

310

312

270

281

283 284

286

287 288

289

290

291

292 293

295

296

297

298

321 3.3 QUALITATIVE RESULTS ON LATENT ACTIONS

We present qualitative results on latent actions learned from robotics and human activity dataset. Specifically, we answer the following questions on learned latent actions:



Figure 4: Controllability of latent action among multiple objects. The last two rows show the generated image by applying 6 different latent actions to the initial frame. Effects of applying different latent actions are highlighted in dashed squares: (a,b) move the apple, (c,d) move the tennis, (e,f) move the orange. Full generated videos from the world model are available on our webpage.

- Do similar latent actions reflect similar visual changes?
- Can latent actions encode semantically consistent changes across different tasks, and embodiments including human and robots? If so, are we able to migrate movements in videos across embodiments and tasks via latent action?
- Does the policy foundation model properly follow language instructions for task solving?

VISUALIZATION OF IMAGE-GOAL PAIRS WITH SIMILAR LATENT ACTIONS 3.3.1

354 We investigate whether similar learned latent actions reflect similar visual changes on robotics ma-355 nipulation dataset. We use RT-1 dataset, which was excluded from the latent action model training 356 and serves as out-of-domain samples for evaluation. We randomly select image-goal pairs from 357 RT-1 dataset, and present the image-goal pairs with smallest euclidean distance in latent action em-358 bedding in RT-1 dataset in Figure 3. We observe that pairs with similar embeddings indeed have 359 similar visual changes, and also similar sub-tasks in semantic, for example, "open the gripper", 360 "move left", and "close the gripper". Furthermore, each sub-task appears in different raw language 361 tasks, suggesting the latent actions are reused, thereby facilitating generalization in model learning.

CONTROLLABILITY OF LATENT ACTIONS 3.3.2

364 We demonstrate that latent actions are able to control the changes in objects on different real world scenes, and the effects of latent actions generalize across tasks and embodiments. Specially, the 366 generalizability of latent actions enable IGOR to successfully migrate human movement videos into robot movements provided the initial image, despite they largely differ in embodiments.

367 368

362

324

330

341

342

343

344

345

347

348

349

350 351

352 353

369 **Object Controllability Among Multiple Objects.** We evaluate the controllability of the latent 370 actions on object movements among multiple objects on the same image. In Figure 4, we generate 371 subsequent images by applying 6 different actions to the same original image on the foundation 372 world model. We observe that the latent action model and foundation world model learn to control 373 specific object's movement among multiple objects.

374

375 **Object Controllability Across Embodiments and Tasks.** We evaluate the semantic consistency of the latent actions across different setups, including embodiments and tasks. We use pairs of 376 image-goal in the real world manipulation videos to generate latent actions, and apply the same set of 377 actions to other images in different scene setups with foundation world model to generate subsequent



Figure 5: Generated image sequence jointly by the foundation policy and world model via only latent actions, following 3 different instructions from the same initial image. Full generated videos from the world model are available on our webpage.

399

400

401

402

403 404

394

videos. The results are shown in Figure 2. Impressively, we observe that (1) latent actions are semantically consistent across different tasks with different object categories; (2) latent actions are semantically consistent across human and robots. By applying latent actions extracted from human demonstrations, we generate videos of robot arm movements. With only one demonstration, the robot arm can successfully migrate human behaviors, which opens up new possibilities for few-shot human-to-robot transfer and control.

3.3.3 COUNTERFACTUAL VIDEO GENERATION WITH DIVERSE INSTRUCTIONS 405

406 We analyze whether the foundation policy model has the ability to follow human instructions. To 407 this end, we interpret the effect of latent actions visually with the foundation world model. Starting 408 from a single initial image, the foundation policy and world model can jointly generate diverse 409 behaviors in videos that follow diverse instructions using only latent actions. We experiment with 410 initial images from RT-1 and Bridge dataset and manually written instructions, and show the image 411 clips of generated videos in Figure 5. The results show that the foundation policy model can properly 412 follow different language instructions for task solving.

413 414 415

3.4 **QUANTITATIVE RESULTS**

416 EVALUATION ON THE GOOGLE ROBOT TASKS IN SIMPLER 3.4.1

417 We evaluate our IGOR-based training framework on the Google robot tasks in the SIMPLER sim-418 ulator under a low-data regime, utilizing only 1% of the data from the large RT-1 dataset for the 419 low-level policy learning stage. 420

421

Evaluation Setups. We test different model's ability to control the Google Robot following 422 language tasks with RGB images as observations, where all robots are controlled with low-level endeffector control actions, after finetuning on the same amount of data from RT-1 dataset. We evaluate 424 the success rate on three tasks: "Pick Coke Can", "Move Near", and "Open / Close Drawer".

425 426

423

Baseline Method. We compare our method with the same low-level policy model architecture 427 with ST-Transformer, without latent action embedding concatenated on the observation feature em-428 bedding. 429

We present the success rate of different methods in Figure 6(a). From the figure, we observe that 430 IGOR achieves higher or equal success rate than the model trained from scratch, showing the gen-431 eralizability of the learned latent action to real robotics actions.



Figure 6: (a). Success rate of IGOR and the low-level policy trained from scratch methods on Google Robot tasks under SIMPLER simulator, finetuned on 1% data of RT-1. (b). Predictiveness of latent action on robot action. X-axis: log(N), where N is the number of nearest neighbours in latent action embedding. Y-axis: normalized standard deviation in action embedding with respect to movement actions (orange), rotation actions (blue), and gripper actions (green).

448
4493.4.2PREDICTIVENESS OF LATENT ACTIONS ON ROBOT ACTIONS

We analyze whether our learned latent actions are predictive of real robot actions. On RT-1 dataset, we randomly sample a number of M = 15,000 pairs of images, and compute their latent action embeddings. For each pair of image, we find N nearest neighbours of image pairs in RT-1 dataset with the closest latent action embedding, and compute the standard deviation of real robot actions among N neighbours on each action dimension, normalized by the standard deviation of robot actions over each dimension over the whole RT-1 dataset. By varying N, we assess whether closer latent actions correspond to more similar downstream actions.

The results are shown in Figure 6(b). The fact that smaller N leading to lower normalized standard deviation, and all normalized standard deviation being below 1.0, show that the latent actions are predictive of real robot actions including robot movements, rotations and gripper actions. It is also shown that the latent actions are more predictive of the robot movement than rotations and gripper actions, suggesting that the IGOR learned action space reflects more information in robot movements than robot arm rotations and gripping.

463 3.5 ABLATION STUDIES

443

444

445

446

447

468

469

We provide additional ablation studies on the pretraining dataset of latent action model, showing that using a mixture of robotics and human activity dataset benefits the generalization of latent action model. Detailed ablation studies results are provided in Appendix C.

4 RELATED WORK

Foundation Agents for Robots Open-ended task-agnostic training and high-capacity neural net-470 work architectures have been recognized as key to the success of foundation models. In this context, 471 a series of generalist agents have been proposed as the foundation policy models for robots (Brohan 472 et al., 2022; Bousmalis et al., 2023; Brohan et al., 2023; Team et al., 2024; Kim et al., 2024). RT-473 1 (Brohan et al., 2022) contributes a large-scale multi-task dataset and a robotic transformer architec-474 ture, facilitating and assessing generalization across multiple tasks. RoboCat builds on Gato (Reed 475 et al., 2022), further enabling multi-embodiment generalization. RT-2 highlights the importance of 476 leveraging vision-language models trained on internet-scale data (Brohan et al., 2023). Octo (Team 477 et al., 2024) and OpenVLA (Kim et al., 2024) can be seen as open versions of RoboCat and RT-2 respectively, with some additional technical contributions. IGOR is similar to RT-2 and OpenVLA 478 in the sense that we both leverage Internet-scale data. The difference lies in that we use video data of 479 human/robot performing embodied AI tasks, while they use text data and visual question answering 480 data for the training of vision language models. To the best of our knowledge, we present the first 481 foundation policy model that performs decision making at the sub-task (i.e. latent action) level. 482

Image-Goal Visual Changes Tracking visual changes and establishing correspondence be tween an image and its goal state is crucial for dynamic visual understanding in embodied AI.
 SiamMAE (Gupta et al., 2023) proposes to use a siamese encoder on the image and goal to learn visual correspondence. Voltron (Karamcheti et al., 2023) introduces language-guided visual represen-

486 tation learning on image-goal pairs. Lin et al. (2024) and Ko et al. (2023) leverage optical flow be-487 tween image and goal to capture visual changes and correspondence, while Video-LaVIT (Jin et al., 488 2024) utilizes motion vectors. iVideoGPT (Wu et al., 2024) proposes using image-conditioned goal 489 representations as state representations to predict within a world model. VPT (Baker et al., 2022) 490 proposes to recover latent actions in videos using an inverse dynamics model trained on interaction data to predict real actions. Perhaps the most similar approaches to our methods are LAPO (Schmidt 491 & Jiang, 2023) and Genie (Bruce et al., 2024). Both works primarily focus on 2D platformer games 492 where each latent action corresponds to a specific control button. By contrast, we aim to develop 493 a more generalizable model to handle the significantly greater complexity of open-world scenarios, 494 where latent actions may not correspond to any specific underlying actions. 495

496 Video Generation for Embodied AI Video generation is another research topic closely related 497 to embodied AI. It has been proposed that video can be seen as the new language for real-world 498 decision making (Yang et al., 2024b). Many works on world models build on video generation techniques (Bruce et al., 2024; Wu et al., 2024; Hu et al., 2023; Yang et al., 2024a; Xiang et al., 499 2024). Some text-to-video works claim to be real-world simulators, such as Sora (Brooks et al., 500 2024) and WorldDreamer (Wang et al., 2024). Unipi (Du et al., 2023) proposes to first predict 501 the next goal state, then infer real robot actions with an inverse dynamics model. By contrast, our 502 foundation policy model first predicts the latent action, which can specify the goal state, and then uses the latent action to enable sub-task level generalization. We argue that forward prediction in 504 latent action space, rather than the original image space, offers several advantages. For example, 505 we can perform sub-task understanding for image-goal representations, and the compressed latent 506 action could be easier to predict than the entire image.

Pre-trained Visual Representations Pre-trained Visual Representations target on training representations for images/videos in self-supervised learning manner (He et al., 2021; Xiao et al., 2022; Radosavovic et al., 2022; Majumdar et al., 2023; Radford et al., 2021; Nair et al., 2022; Ma et al., 2023; Oquab et al., 2023; Darcet et al., 2023; Kirillov et al., 2023; Assran et al., 2023; Bardes et al., 2024), and has been demonstrated to be very effective for state understanding in embodied AI. By contrast, IGOR learns image-goal representations for sub-task understanding, which we believe are another crucial building blocks, that may significantly enhance generalization in embodied AI.

514 515

520

521

522

523

524

525

527

528

529 530

531

532

5 CONCLUSIONS, LIMITATIONS, AND FUTURE WORK

In this paper, we propose IGOR, a novel training framework, taking the first step towards learning a unified action space for humans and robots in various embodied AI tasks.

519 Qualitatively, we demonstrate that:

- IGOR learns similar representations for image pairs with similar visual changes.
- The learned latent action has control over the next state given the current image.
- The foundation world model acquires knowledge about objects and their potential movements.
 - The foundation policy model learns to follow instructions across different states.
- Quantitatively, we show that:
 - On the RT-1 dataset, image-goal pairs with similar latent actions have similar low-level robot actions.
 - The IGOR framework improves policy learning, potentially due to its capability to predict the next sub-task by leveraging internet-scale data, thereby enabling sub-task level generalization.

The IGOR framework is limited in the following perspective: we cannot separate visual changes caused by the agents, other agents (such as dogs), or the shakiness of the camera. To address this, we mitigated camera shakiness and used only ego-centric videos without other agents in view. Just like any other representation learning methods, scaling up the dataset and model size is always most straightforward and effective. To facilitate the usage of more data, incorporating image processing methods such as object segmentation with IGOR will be part of future works. For better applications in embodied AI, the foundation world model can also be tuned to match real world scenarios, along with other improvements such as adapting the latent action model for multi-agent scenarios.

540 REFERENCES

- 542 Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic
 543 interpolants. In *International Conference on Learning Representations*, 2023.
- Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat,
 Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding
 predictive architecture, 2023. URL https://arxiv.org/abs/2301.08243.
- 547
 548
 549
 549
 550
 Bowen Baker, Ilge Akkaya, Peter Zhokhov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos, 2022. URL https://arxiv.org/abs/2206.11795.
- Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun, Mahmoud
 Assran, and Nicolas Ballas. Revisiting feature prediction for learning visual representations from
 video, 2024. URL https://arxiv.org/abs/2404.08471.
- Suneel Belkhale, Yuchen Cui, and Dorsa Sadigh. Hydra: Hybrid robot actions for imitation learning.
 arxiv, 2023.
- Homanga Bharadhwaj, Debidatta Dwibedi, Abhinav Gupta, Shubham Tulsiani, Carl Doersch, Ted Xiao, Dhruv Shah, Fei Xia, Dorsa Sadigh, and Sean Kirmani. Gen2act: Human video generation in novel scenarios enables generalizable robot manipulation. *arXiv preprint arXiv:2409.16283*, 2024a.
- Homanga Bharadhwaj, Abhinav Gupta, Vikash Kumar, and Shubham Tulsiani. Towards generaliz able zero-shot manipulation via translating human interaction plans. In 2024 IEEE International
 Conference on Robotics and Automation (ICRA), pp. 6904–6911. IEEE, 2024b.
- Homanga Bharadhwaj, Roozbeh Mottaghi, Abhinav Gupta, and Shubham Tulsiani. Track2act: Predicting point tracks from internet videos enables generalizable robot manipulation. *arXiv preprint arXiv:2405.01527*, 2024c.
- 568 Konstantinos Bousmalis, Giulia Vezzani, Dushyant Rao, Coline Devin, Alex X. Lee, Maria Bauza, Todor Davchev, Yuxiang Zhou, Agrim Gupta, Akhil Raju, Antoine Laurens, Claudio Fantacci, 569 Valentin Dalibard, Martina Zambelli, Murilo Martins, Rugile Pevceviciute, Michiel Blokzijl, 570 Misha Denil, Nathan Batchelor, Thomas Lampe, Emilio Parisotto, Konrad Żołna, Scott Reed, 571 Sergio Gómez Colmenarejo, Jon Scholz, Abbas Abdolmaleki, Oliver Groth, Jean-Baptiste Regli, 572 Oleg Sushkov, Tom Rothörl, José Enrique Chen, Yusuf Aytar, Dave Barker, Joy Ortiz, Mar-573 tin Riedmiller, Jost Tobias Springenberg, Raia Hadsell, Francesco Nori, and Nicolas Heess. 574 Robocat: A self-improving generalist agent for robotic manipulation, 2023. URL https: 575 //arxiv.org/abs/2306.11706. 576
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- 580 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choro-581 manski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, 582 Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alex Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, 583 Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Hen-584 ryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, 585 Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricut, 586 Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-588 2: Vision-language-action models transfer web knowledge to robotic control. In arXiv preprint 589 arXiv:2307.15818, 2023. 590
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video
 generation models as world simulators. 2024. URL https://openai.com/research/
 video-generation-models-as-world-simulators.

- Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
- 596 597 598

594

595

Lawrence Yunliang Chen, Simeon Adebola, and Ken Goldberg. Berkeley UR5 demonstration dataset. https://sites.google.com/view/berkeley-ur5/home.

Open X-Embodiment Collaboration, Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Ab-601 hishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, 602 Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant 603 Rai, Anchit Gupta, Andrew Wang, Andrey Kolobov, Anikait Singh, Animesh Garg, Aniruddha 604 Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan 605 Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen 607 Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher Agia, Chuer 608 Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak 609 Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Ed-610 ward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Felipe Vieira Fru-611 jeri, Freek Stulp, Gaovue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guangwen Yang, Guanzhi Wang, Hao Su, Hao-Shu 612 Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homer 613 Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad 614 Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jeannette 615 Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, 616 Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra 617 Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Sal-618 vador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, 619 Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento 620 Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty El-621 lis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle 622 Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi "Jim" Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Mem-623 mel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, 624 Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Sri-625 rama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J 626 Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier 627 Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick "Tree" Miller, Patrick Yin, 628 Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiuyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Mart'in-630 Mart'in, Rohan Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, 631 Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, 632 Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, 633 Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, 634 Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, 635 Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, 636 Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, 637 Vidhi Jain, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiangyu 638 Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yan-639 song Pang, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, 640 Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yongqiang Dou, Yoonyoung Cho, Youngwoon 641 Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, 642 Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff 643 Cui, Zichen Zhang, Zipeng Fu, and Zipeng Lin. Open X-Embodiment: Robotic learning datasets and RT-X models. https://arxiv.org/abs/2310.08864, 2023. 644

645

Zichen Jeff Cui, Yibin Wang, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. From play to policy: Conditional behavior generation from uncurated robot data. *arXiv preprint arXiv:2210.10047*, 2022.

648 649 650 651	Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. The epic-kitchens dataset: Collection, challenges and baselines. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 43(11):4125–4141, 2020.
652 653 654	Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers, 2023.
655 656 657	Shivin Dass, Jullian Yapeter, Jesse Zhang, Jiahui Zhang, Karl Pertsch, Stefanos Nikolaidis, and Joseph J. Lim. CLVR jaco play dataset, 2023. URL https://github.com/clvrai/clvr_jaco_play_dataset.
658 659 660 661 662	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL https://arxiv.org/abs/2010.11929.
663 664 665	Yilun Du, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Joshua B. Tenenbaum, Dale Schu- urmans, and Pieter Abbeel. Learning universal policies via text-guided video generation, 2023. URL https://arxiv.org/abs/2302.00111.
666 667 668	Frederik Ebert, Yanlai Yang, Karl Schmeckpeper, Bernadette Bucher, Georgios Georgakis, Kostas Daniilidis, Chelsea Finn, and Sergey Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets. <i>arXiv preprint arXiv:2109.13396</i> , 2021.
670 671 672 673	Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In <i>Forty-first International Conference on Machine Learning</i> , 2024.
674 675 676 677	Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne West- phal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The" something something" video database for learning and evaluating visual common sense. In <i>Proceedings of the IEEE international conference on computer vision</i> , pp. 5842–5850, 2017.
678 679 680 681	Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 18995–19012, 2022.
683 684	Agrim Gupta, Jiajun Wu, Jia Deng, and Li Fei-Fei. Siamese masked autoencoders, 2023. URL https://arxiv.org/abs/2305.14344.
685 686 687	Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners, 2021. URL https://arxiv.org/abs/2111.06377.
689 690 691	Minho Heo, Youngwoon Lee, Doohyun Lee, and Joseph J. Lim. Furniturebench: Reproducible real-world benchmark for long-horizon complex manipulation. In <i>Robotics: Science and Systems</i> , 2023.
692 693 694	Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shot- ton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving, 2023. URL https://arxiv.org/abs/2309.17080.
695 696 697 698	Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. Bc-z: Zero-shot task generalization with robotic imitation learning. In <i>Conference on Robot Learning</i> , pp. 991–1002. PMLR, 2022.
699 700 701	Yang Jin, Zhicheng Sun, Kun Xu, Kun Xu, Liwei Chen, Hao Jiang, Quzhe Huang, Chengru Song, Yuliang Liu, Di Zhang, Yang Song, Kun Gai, and Yadong Mu. Video-lavit: Unified video-language pre-training with decoupled visual-motional tokenization, 2024. URL https: //arxiv.org/abs/2402.03161.

- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. In *CoRL*, pp. 651–673, 2018.
- Siddharth Karamcheti, Suraj Nair, Annie S. Chen, Thomas Kollar, Chelsea Finn, Dorsa Sadigh, and
 Percy Liang. Language-driven representation learning for robotics. In *Robotics: Science and Systems (RSS)*, 2023.
- 709 Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth 710 Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, 711 Peter David Fagan, Joey Hejna, Masha Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree 712 Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, Arhan Jain, Abraham Lee, Youngwoon Lee, Marius Memmel, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Albert Zhan, Kevin 713 Black, Cheng Chi, Kyle Beltran Hatch, Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pan-714 nag R Sanketi, Archit Sharma, Cody Simpson, Quan Vuong, Homer Rich Walke, Blake Wulfe, 715 Ted Xiao, Jonathan Heewon Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia, Rohan Bai-716 jal, Mateo Guaman Castro, Daphne Chen, Qiuyu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul 717 Foster, Jensen Gao, David Antonio Herrera, Minho Heo, Kyle Hsu, Jiaheng Hu, Donovon Jack-718 son, Charlotte Le, Yunshuang Li, Kevin Lin, Roy Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mir-719 chandani, Daniel Morton, Tony Nguyen, Abigail O'Neill, Rosario Scalise, Derick Seale, Victor 720 Son, Stephen Tian, Emi Tran, Andrew E. Wang, Yilin Wu, Annie Xie, Jingyun Yang, Patrick Yin, 721 Yunchu Zhang, Osbert Bastani, Glen Berseth, Jeannette Bohg, Ken Goldberg, Abhinav Gupta, 722 Abhishek Gupta, Dinesh Jayaraman, Joseph J Lim, Jitendra Malik, Roberto Martín-Martín, Subramanian Ramamoorthy, Dorsa Sadigh, Shuran Song, Jiajun Wu, Michael C. Yip, Yuke Zhu, 723 Thomas Kollar, Sergey Levine, and Chelsea Finn. Droid: A large-scale in-the-wild robot manip-724 ulation dataset. 2024. 725
- Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.
 Segment anything, 2023. URL https://arxiv.org/abs/2304.02643.
- Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and Joshua B. Tenenbaum. Learning to act from actionless videos through dense correspondences, 2023. URL https://arxiv.org/abs/2310.08576.
- Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.
- Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, Sergey Levine, Jiajun Wu, Chelsea Finn, Hao Su, Quan Vuong, and Ted Xiao. Evaluating real-world robot manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024.
- Yin Li, Miao Liu, and James M Rehg. In the eye of beholder: Joint learning of gaze and actions in first person video. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 619–635, 2018.
- Li-Heng Lin, Yuchen Cui, Amber Xie, Tianyu Hua, and Dorsa Sadigh. Flowretrieval: Flow guided data retrieval for few-shot imitation learning, 2024. URL https://arxiv.org/abs/
 2408.16944.
- Huihan Liu, Soroush Nasiriany, Lance Zhang, Zhiyao Bao, and Yuke Zhu. Robot learning on the job: Human-in-the-loop autonomy and learning during deployment. In *Robotics: Science and Systems (RSS)*, 2023a.
- 755 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *International Conference on Learning Representations*, 2023b.

756 757 758	Jianlan Luo, Charles Xu, Xinyang Geng, Gilbert Feng, Kuan Fang, Liam Tan, Stefan Schaal, and Sergey Levine. Multi-stage cable routing through hierarchical imitation learning. <i>arXiv preprint arXiv:2307.08927</i> , 2023.
759 760 761 762	Jianlan Luo, Charles Xu, Fangchen Liu, Liam Tan, Zipeng Lin, Jeffrey Wu, Pieter Abbeel, and Sergey Levine. Fmb: a functional manipulation benchmark for generalizable robotic learning. <i>arXiv preprint arXiv:2401.08553</i> , 2024.
763 764 765 766	Corey Lynch, Ayzaan Wahid, Jonathan Tompson, Tianli Ding, James Betker, Robert Baruch, Travis Armstrong, and Pete Florence. Interactive language: Talking to robots in real time. <i>IEEE Robotics and Automation Letters</i> , 2023.
767 768 769	Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training, 2023. URL https://arxiv.org/abs/2210.00030.
770 771 772 773	Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Yecheng Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Pieter Abbeel, Jitendra Malik, Dhruv Batra, Yixin Lin, Olek-sandr Maksymets, Aravind Rajeswaran, and Franziska Meier. Where are we in the search for an artificial visual cortex for embodied intelligence? 2023.
774 775 776 777	Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emmons, Anchit Gupta, Emre Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In <i>Conference on Robot Learning</i> , pp. 879–893. PMLR, 2018.
778 779 780	Oier Mees, Jessica Borja-Diaz, and Wolfram Burgard. Grounding language with visual affordances over unstructured data. In <i>Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)</i> , London, UK, 2023.
781 782 783	Russell Mendonca, Shikhar Bahl, and Deepak Pathak. Structured world models from human videos. <i>CoRL</i> , 2023.
784 785 786	Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation, 2022. URL https://arxiv.org/abs/2203.12601.
787 788 789	Soroush Nasiriany, Tian Gao, Ajay Mandlekar, and Yuke Zhu. Learning and retrieval from prior data for skill-based imitation learning. In <i>Conference on Robot Learning (CoRL)</i> , 2022.
790 791 792 793 794 795	Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Ar- mand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023.
796 797 798 799 800	Shraman Pramanick, Yale Song, Sayan Nag, Kevin Qinghong Lin, Hardik Shah, Mike Zheng Shou, Rama Chellappa, and Pengchuan Zhang. Egovlpv2: Egocentric video-language pre-training with fusion in the backbone. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 5285–5297, 2023.
801 802 803	Gabriel Quere, Annette Hagengruber, Maged Iskandar, Samuel Bustamante, Daniel Leidner, Freek Stulp, and Joern Vogel. Shared Control Templates for Assistive Robotics. In 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 7, Paris, France, 2020.
804 805 806 807	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar- wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL https://arxiv.org/abs/2103.00020.
808	Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Real-world robot learning with masked visual pre-training. <i>CoRL</i> , 2022.

822

823

824

847

848

849

850

851

856

- Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. A generalist agent, 2022. URL https://arxiv.org/abs/2205.06175.
- Erick Rosete-Beas, Oier Mees, Gabriel Kalweit, Joschka Boedecker, and Wolfram Burgard. Latent
 plans for task agnostic offline reinforcement learning. In *Proceedings of the 6th Conference on Robot Learning (CoRL)*, 2022.
- 819 Saumya Saxena, Mohit Sharma, and Oliver Kroemer. Multi-resolution sensing for real-time control
 820 with vision-language models. In 7th Annual Conference on Robot Learning, 2023. URL https:
 821 //openreview.net/forum?id=WuBv9-IGDUA.
 - Dominik Schmidt and Minqi Jiang. Learning to act without actions. *arXiv preprint arXiv:2312.10812*, 2023.
- Nur Muhammad Mahi Shafiullah, Anant Rai, Haritheja Etukuru, Yiqian Liu, Ishan Misra, Soumith
 Chintala, and Lerrel Pinto. On bringing robots home, 2023.
- Rutav Shah, Roberto Martín-Martín, and Yuke Zhu. MUTEX: Learning unified policies from multimodal task specifications. In *7th Annual Conference on Robot Learning*, 2023. URL https://openreview.net/forum?id=PwqiqaaEzJ.
- Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep
 Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot
 policy. *arXiv preprint arXiv:2405.12213*, 2024.
- Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony Zhao,
 Philippe Hansen-Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang, Chelsea Finn, and
 Sergey Levine. Bridgedata v2: A dataset for robot learning at scale, 2023.
- Xiaofeng Wang, Zheng Zhu, Guan Huang, Boyuan Wang, Xinze Chen, and Jiwen Lu. Worlddreamer: Towards general world models for video generation via predicting masked tokens, 2024. URL https://arxiv.org/abs/2401.09985.
- Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point trajectory modeling for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- Jialong Wu, Shaofeng Yin, Ningya Feng, Xu He, Dong Li, Jianye Hao, and Mingsheng Long.
 ivideogpt: Interactive videogpts are scalable world models. *arXiv preprint arXiv:2405.15223*, 2024.
 - Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning, Yuheng Zha, Zeyu Feng, Tianhua Tao, Shibo Hao, Yemin Shi, Zhengzhong Liu, Eric P. Xing, and Zhiting Hu. Pandora: Towards general world model with natural language actions and video states, 2024. URL https://arxiv.org/abs/2406.09455.
- Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked visual pre-training for
 motor control. *arXiv:2203.06173*, 2022.
- Mengda Xu, Zhenjia Xu, Yinghao Xu, Cheng Chi, Gordon Wetzstein, Manuela Veloso, and Shuran
 Song. Flow as the cross-domain manipulation interface. *arXiv preprint arXiv:2407.15208*, 2024.
- Mingxing Xu, Wenrui Dai, Chunmiao Liu, Xing Gao, Weiyao Lin, Guo-Jun Qi, and Hongkai
 Xiong. Spatial-temporal transformer networks for traffic flow forecasting, 2021. URL https: //arxiv.org/abs/2001.02908.
- Ge Yan, Kris Wu, and Xiaolong Wang. ucsd kitchens Dataset. August 2023.
- Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Leslie Kaelbling, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators, 2024a. URL https://arxiv.org/abs/2310.06114.

864 865	Sherry Yang, Jacob Walker, Jack Parker-Holder, Yilun Du, Jake Bruce, Andre Barreto, Pieter Abbeel, and Dale Schuurmans. Video as the new language for real-world decision making, 2024b.
867	Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all
869	March 2024. URL https://github.com/hpcaitech/Open-Sora.
870	Georgia Zhou, Viatoria Daen, Mahan Kumar Srirama, Aravind Baiagwaran, Justhich Dari, Kula
871	Hatch Arvan Jain Tianhe Yu Pieter Abbeel Lerrel Pinto, Chelsea Finn, and Abbinay Gunta
872	Train offline, test online: A real robot learning benchmark, 2023.
073 974	Xinghao Zhu, Ran Tian, Chenfeng Xu, Mingyu Ding, Wei Zhan, and Masayoshi Tomizuka. Fanuc
875	manipulation: A dataset for learning-based manipulation with fanuc mate 200id robot. 2023a.
876	Yifeng Zhu, Peter Stone, and Yuke Zhu. Bottom-up skill discovery from unsegmented demonstra-
877 878	tions for long-horizon robot manipulation. <i>IEEE Robotics and Automation Letters</i> , 7(2):4126–4133, 2022
879	1155, 2022.
880	Yifeng Zhu, Abhishek Joshi, Peter Stone, and Yuke Zhu. Viola: Imitation learning for vision-based
881	manipulation with object proposal priors, 2023b.
882	
883	
884	
885	
886	
887	
888	
889	
890	
891	
892	
893	
894	
895	
896	
897	
898	
899	
900	
901	
903	
904	
905	
906	
907	
908	
909	
910	
911	
912	
913	
914	
915	
916	
917	

918 A DATASET

954

955 956

We present the datasets used for pre-training in Table 1. In total, these datasets comprise approximately 0.8 million robot trajectories and 2.0 million filtered human activity video clips. The robot data ratios are from (Team et al., 2024).

924	Robot Dataset	Mix Ratio (%)
925	Kuka (Kalashrikay at al. 2018)	7.72
926	Ruka (KalasiiiiKov et al., 2016) Bridge (Welke et al., 2022; Ebert et al., 2021)	2.12
927	Taao Diay (Desete Dese et al., 2023, Ebert et al., 2021)	0.00
928	Taco Play (Rosele-Deas et al., 2022; Mees et al., 2025)	1.62
929	Jaco Flay (Dass et al., 2023) Barkeley Cable Pouting (Luc et al., 2023)	0.24
930	Poboturk (Mandlekar et al. 2018)	1.40
031	Viola (Zhu et al. 2023b)	0.55
000	Berkely Autolab UB5 (Chen et al.)	0.55
932	Toto (Zhou et al. 2023)	1.21
933	Language Table (Lynch et al. 2023)	2.67
934	Stanford Hydra Dataset (Belkhale et al. 2023)	2.67
935	Austin Buds Dataset (Zhu et al. 2022)	0.12
936	NYLL Franka Play Dataset (Cui et al. 2022)	0.12
937	Furniture Bench Dataset (Heo et al. 2023)	1 46
938	UCSD Kitchen Dataset (Yan et al., 2023)	0.06
939	Austin Sailor Dataset (Nasiriany et al., 2022)	1.34
940	Austin Sirius Dataset (Liu et al., 2023a)	1.03
941	DLR EDAN Shared Control (Ouere et al., 2020)	0.06
942	IAMLab CMU Pickup Insert (Saxena et al., 2023)	0.55
943	UTAustin Mutex (Shah et al., 2023)	1.34
944	Berkeley Fanuc Manipulation (Zhu et al., 2023a)	0.43
945	CMU Stretch (Mendonca et al., 2023)	0.12
0/6	BC-Z (Jang et al., 2022)	4.56
047	FMB Dataset (Luo et al., 2024)	4.31
947	DobbE (Shafiullah et al., 2023)	0.85
948	DROID (Khazatsky et al., 2024)	6.07
949	Ego4D (Grauman et al., 2022)	32.1
950	Something-Something V2 (Goyal et al., 2017)	9.5
951	EPIC-KITCHENS (Damen et al., 2020)	8.0
952	EGTEA Gaze+ (Li et al., 2018)	0.4
953		

Table 1: Dataset, mixture weights, and number of training examples after filtering in the pre-training stage in IGOR.

Data Filtering We observed that video quality significantly affects the action model, particularly for human activities video. Excessive shakiness in videos can introduce visual changes between consecutive frames that are unrelated to the agent's actions.

We calculate the camera motion over the videos, and filter approximately 40% percent of open-world video data. For the remaining data, we further stabilize the videos. Although we retain only 60% percent of open-world video data, we find that the action model improves dramatically.

964 Frame Interval A noticeable amount of visual changes is crucial for our latent action model. If 965 we select two frames that are too close in time, the agent may barely move, resulting in visual 966 changes that are not significant enough for inferring meaningful actions. Conversely, if the frames 967 are too far apart, the changes might be too large to model accurately. To address this issue, we 968 tune the sampling interval. For robot data, we choose frames that are three intervals apart, using s_t and s_{t+3} as the image-goal pair. For real world videos, we control the sampling. For real world 969 data, we control the sample interval to be within [0.1s, 0.5s]. For the action and policy model, the 970 context frames follow the same interval, ensuring that each pair of consecutive frames maintains this 971 consistent spacing.

972 B TRAINING DETAILS

974 B.1 LATENT ACTION MODEL TRAINING 975

The latent action model uses an ST-transformer equipped with a frozen DINO-v2 pretrained ViT image encoder. The latent action model uses 258 M parameters, a patch size of 14, and a codebook with N = 4 tokens and size |C| = 32, each with an embedding size of D = 128. We train the latent action model with batch size B = 512, training iterations of 140K steps, and learning rate of 1.5e - 4 with Adam optimizer.

980 981 982

B.2 FOUNDATION WORLD MODEL TRAINING

We start on the top of the OpenSora (Zheng et al., 2024) model with newly initialized projection layers. The foundation world model with batch size B = 12, training iterations of 48K, and learning rate of 1e - 4 with Adam optimizer.

986

987 B.3 FOUNDATION POLICY MODEL AND LOW-LEVEL POLICY MODEL

The latent action model uses an ST-transformer equipped with a frozen DINO-v2 pretrained ViT image encoder, following the latent action model's image encoder. The foundation policy model consists of 12 layers of spatial and temporal attentions, each with 12 attention heads and hidden dimension as 768 and a patch size of 14. In total the policy model has 138M parameters. We use frozen CLIP features for text instructions. We pretrain the foundation policy model with batch size B = 128, training iterations of 124K, and learning rate of 1e - 4 with Adam optimizer.

The low-level policy model adds extra 118M parameters on top of the foundation policy model. We use a sub-task length of $\tau = 3$ for finetuning the low-level policy model on RT-1 dataset. We finetune the low-level policy model with batch size B = 128, training iterations of 32K, and learning rate of 1e - 4 with Adam optimizer.

999

1000 C ADDITIONAL ABLATION RESULTS

1002 C.1 DATASET ABLATION FOR LATENT ACTION MODEL

1003 We compare two different settings for the pre-training dataset: only use the robotic dataset (robot 1004 data), and use both robotic and human activity dataset (mixed data). We evaluate the validation loss 1005 on the latent action model on RT-1 dataset, which is held out from the pretraining dataset and serves for OOD evaluation. Validation loss of the latent action model assesses the extent to which the 1007 IDM and FDM can jointly generate latent actions and recover goal states from these latent actions 1008 conditioned on states on the unseen dataset. The results are shown in Table 2. We find that the OOD 1009 validation loss is greatly reduced by adding human activity dataset. This may be due to the diversity 1010 of human videos, which comprise real daily life environments with lots of diverse backgrounds and 1011 objects. These results demonstrate that it is promising to leverage human data for improving robot tasks under the IGOR framework. 1012

	Validation Loss
Robot Data	0.145
Mixed Data	0.112

1017 1018

1020 1021

Table 2: Validation loss on held-out dataset (RT-1) with different training data.

D MODEL STRUCTURE OF IGOR

1023 1024

1025 We illustrate the model architecture of latent action model, policy model and world model in IGOR in Figure 7.



ADDITIONAL DISCUSSION ON RELATED WORK

Any-Point Trajectory Modeling (ATM) (Wen et al., 2023) and Tract2Act (Bharadhwaj et al., 2024c)
generates point tracking from the action-free videos for pretraining the policy to predict future trajectories of the tracked point, and learns a robot policy conditioned on generated trajectories. HOPMan (Bharadhwaj et al., 2024b) generates future human plans as conditions for facilitating robot policy learning. Im2Flow2Act (Xu et al., 2024) conditions the robot policy on complete generated object flows, which captures movement information only for the object, improving its crossembodiment transfer capability.

Same as ATM, HOPMan, Im2Flow2Act, and Tract2Act, IGOR also pretrain a foundation policy on action-free videos that generalize across embodiments. There are two major differences: (1) IGOR uses an unsupervised way to learn and compress the visual changes, while existing work needs an additional pretrained video model for point tracking; (2) IGOR uses a compact latent representation for visual changes, while existing works uses an explicit representation for visual changes. The compact action representation enables IGOR to transfer human actions on robots directly, as shown in Figure 2 in our paper.

In the world model aspect, Gen2Act (Bharadhwaj et al., 2024a) uses the VideoPoet model (Kondratyuk et al., 2023) for text-to-video generation, and the generated human videos are used for downstream policy learning. Compared to existing world model works, IGOR can learn generalizable latent actions for fine-grained control on manipulation scenarios for both human and robots, while existing works use coarse grained text conditions for generating video predictions.

Ε