

# MultiModal Action Conditioned Video Generation

Yichen Li Antonio Torralba  
MIT CSAIL

## Abstract

*Current video models fail as world model as they lack fine-grained control. General-purpose household robots require real-time fine motor control to handle delicate tasks and urgent situations. In this work, we introduce fine-grained multimodal actions to capture such precise control. We consider senses of proprioception, kinesthesia, force haptics, and muscle activation. Such multimodal senses naturally enables fine-grained interactions that are difficult to simulate with text-conditioned generative models. To effectively simulate fine-grained multisensory actions, we develop a feature learning paradigm that aligns these modalities while preserving the unique information each modality provides. We further propose a regularization scheme to enhance causality of the action trajectory features in representing intricate interaction dynamics. Experiments show that incorporating multimodal senses improves simulation accuracy and reduces temporal drift. Extensive ablation studies and downstream applications demonstrate the effectiveness and practicality of our work.<sup>†</sup>*

## 1. Introduction

For general-purpose household robots to operate dexterously and safely like humans, they need to be enabled with multi-potent sensory systems. Our interoceptive senses, including kinesthesia, proprioception, force haptics, and muscle activation, work together to enable us to dynamically engage with our surroundings. The ability to simulate such multisensory actions is crucial for developing robust embodied intelligence and guiding future directions for sensor design.

Traditionally, physics engines are used to simulate state changes of the environment [23, 36, 42, 62, 63], but creating a physics simulator with fine-grained multisensory capabilities for diverse tasks is both computationally expensive and complex in engineering. Recent works [17, 71] demonstrate the potential to use text-conditioned video models as simulators, but text struggles to capture the delicate control needed for tasks such as culinary or surgical activities. In this work, we introduce multisensory interaction signals in generative simulation to enable fine-grained control.

<sup>†</sup><https://people.csail.mit.edu/yichenl/projects/multimodalvideo/>

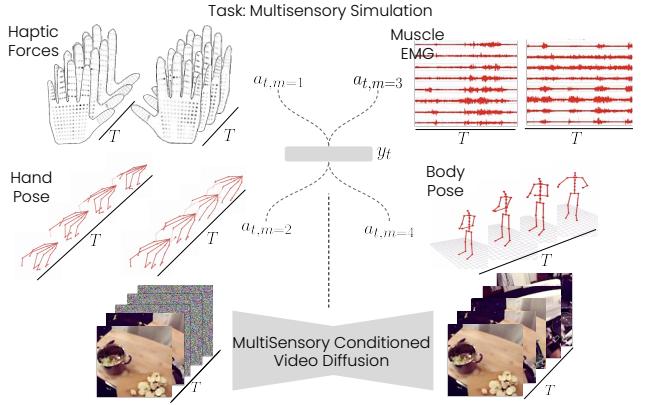


Figure 1. **Overview.** We introduce a new task for fine-grained control of video generative model using multisensory interaction signals.

We focus on learning an effective multimodal representation to control generative simulation. Prior works on multimodal feature learning [16, 19, 28, 37, 59, 76] focus the task of cross-modal retrieval. They thus emphasize multimodal alignment but overlook the unique information each modality provides. As a result, they are insufficient for conditioning generative simulators. For our task, we introduce an multimodal feature extraction paradigm that align modalities to a shared representation space while preserving the unique aspects each modality contributes. Additionally, we propose a generic feature regularization scheme to ensure the encoded action trajectories to be more context-and-consequence-aware, allowing for seamless integration with downstream video generation frameworks.

In this work, we introduce multisensory interoceptive signals of haptic forces, muscle stimulation, hand poses, and body proprioception to generative simulation for fine-grained responses. We focus on learning effective multisensory action representation to control generative video models. Our proposed multimodal feature extraction paradigm aligns different sensory signals while preserving the unique contributions from each modality. Additionally, we introduce a novel feature regularization scheme that the extracted latent representations of action trajectories to capture the intricate causality in context and consequences in interaction dynamics. Extensive comparisons to existing methods shows that our multisensory method helps increase accu-

racy by 36 percent and improve temporal consistency by 16 percent. Ablation studies and downstream applications further demonstrate the effectiveness and practicality of our proposed approach. To summarize, our contributions are:

- To the best of our knowledge, we are the first to introduce multisensory signals, including touch, pose, and muscle activity, to generative simulation for fine-grained responses.
- We devise a multimodal feature extraction paradigm that aligns modalities to a shared representation space while preserving the unique information each modality provides.
- We propose a novel feature regularization scheme to enhance encoded action trajectories to be context and consequence aware, capturing intricate interaction dynamics.
- We compare our proposed framework with prior approaches and also provide various possible downstream applications in policy optimization, planning, and more.

## 2. Simulating Multi-Sensory Interactions

We focus on two perspectives of modeling multi-sensory interactions. We first consider ways of working with **multimodal** signals, arriving at a multi-sensory action conditioning feature. We then focus on effective **interaction** modeling to capture the relationship between context and consequences in the learned representation. Finally, we cast our multisensory action feature into a generative video model to simulate accurate exteroceptive visual responses.

**Problem Statement.** Simulators, at core, are next state prediction models. They estimate the consequential state changes of the world resulted from actions. Let  $t \in [0, T]$  denote time frames, where  $t_{\text{hist}} \in [0, t - 1]$  denotes the history horizon, and  $t_{\text{future}} \in [t, T]$  are the future frames. For our task, at a snapshot of time  $t$ , we describe the state of the external world  $s_t$  as visual observations  $x_t \in \mathcal{O}$ , that are the video frames. We observe set of sensory modalities denoted as  $a_{t,m}$  of total number of  $M$  modalities,  $m \in [1, M]$ . Given past observations ( $\{a_{[0,t-1],m}\}, x_{[0,t-1]}$ ) and current action sequence  $\{a_{[t,T],m}\}$ , the goal of the simulator is to predict the consequential future states  $s_{[t,T]}$  represented as a set of frames  $x_{[t,T]}$ . We denote the encoded video frame feature as  $z_{x_t}$  that corresponds to  $x_t | t \in [1, T]$ , and we denote the encoded modality-specific features are denoted as  $z_{t,m}$ , and cross-modal feature is denoted as  $y_t$ . Under the generative simulation framework, we focus on extracting effective multimodal action representation  $y_t$  from a set of multisensory actions  $\{a_{[t,T],m}\}$  to condition a downstream generative simulator  $g_\theta$  to accurately predict future states  $x_{[t,T]}$ . We include a notation chart in Appendix Table. 4.

### 2.1. Multi-Sensory Action Representation

Multisensory actuation data are composed of temporal sequences of various sensory modalities of different granularity, dimension, and scale. How to effectively represent them, synchronize them, and combine them so they can accurately

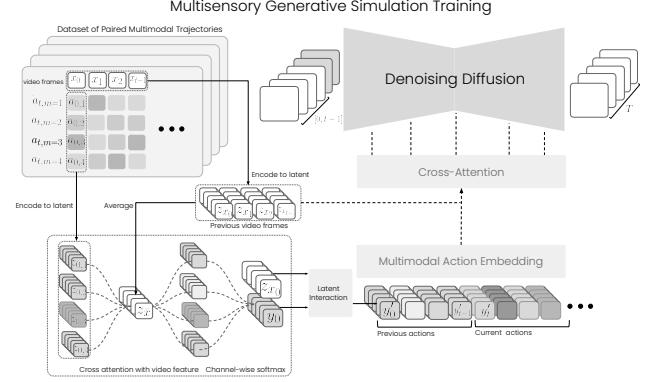


Figure 2. **Overview.** We focus on learning effective multimodal action representations and propose a generative simulation method.

control a generative simulator are the three key challenges in generative *multimodal* feature learning.

One straight-forward way to extract feature representations from various sensory modalities is through mixture-of-expert (MoE) encodings. It is a commonly employed method for encoding heterogeneous data [44, 52, 55]. Various expert encoder heads  $f_m(\cdot)$  are used to extract features  $z_{t,m} = f_m(a_{t,m})$  that represent each sensory modality  $m \in [1, M]$  at each time step  $t$ . To ensure that the encoded information in  $z_{t,m}$  is meaningful, a self-supervised reconstruction scheme is introduced through MoE decoding branches  $d_m(\cdot)$  across each sensory modality  $\hat{a}_{t,m} = d_m(f_m(a_{t,m}))$  supervised by reconstruction loss,  $\mathcal{L}_{\text{SSL}} = \|\hat{a}_{t,m} - a_{t,m}\|^2$ , which gives rise to a set of MoE features  $\{z_{t,m}\}_m^M$ , as shown in Fig. 2.

Before we combine these modality-specific features into a coherent multimodal feature, we need to synchronize them into the same representation space. Ideally, the synchronization strategy should align different MoE features to implicitly follow some shared latent structure and simultaneously preserve uniqueness of each modality, *e.g.* hand pose can inform the action direction, while forces and muscle EMG both indicate action magnitude. These information should be meaningfully packed into different dimensions of the action feature. To encourage such association, we introduce an implicit cross-modal anchoring through channel-wise cross attention. We encode context video frames into latent vectors  $z_{x_{[0,t-1]}}$  of dimension  $d$ , and obtain an anchor feature  $z_{x_{\bar{t}}}$  by averaging across frames. We then use a learnable linear layer to project MoE features  $z_{t,m}$  to anchor dimension  $d$ . Taking a channel-wise cross-attention between the anchor feature  $z_{x_{\bar{t}}}$  and action features  $\{z_{t,m}\}_m^M$  allows channels of the action latents  $\{z_{t,m}\}_m^M$  to be associated through the channels of  $z_{x_{\bar{t}}}$ . In this way, we can train the linear projection layer to implicitly encourage a shared latent structure to arise. Let  $z_{t,m,j}$  denote the  $j$ -th dimension of the action latent vector  $z_{t,m}$  of modality  $m$  and timestep  $t$ .

$$z_{t,m,j} = \sum_i^d \frac{\exp z_{x_{\bar{t}},i} \cdot z_{t,m,j}}{\sum_{l=1}^d \exp z_{x_{\bar{t}},i} \cdot z_{t,m,l}} z_{t,m,j} \quad (1)$$

We are now ready to combine this set of modality-specific features  $\{z_{t,m}\}_{m=1}^M$  into a cross-modal feature  $y_t$ . Different sensory modalities reflect different aspects of our actuation. These sensory modalities complement each other to provide comprehensive information about different actuators. This intuition suggests two properties of our multi-sensory input, over-completeness and permutation invariance. A good feature fusion function works as an information bottleneck to only select the most useful information. Moreover, unlike text sentences or image pixels, data of various sensory modalities is an unordered set. Therefore, the fusion scheme needs to be permutation-invariant regardless the modality order of the input. These properties encourage us to use symmetric functions for feature fusion. After comparing various symmetric functions (Sec. 3.3), we choose softmax weighting function to aggregate different modalities of actuation,

$$y_t = \sum_{m=1}^M w_{t,m} z_{t,m}, \quad \text{where } w_{t,m} = \frac{e^{z_{t,m}}}{\sum_{m'=1}^M e^{z_{t,m'}}}.$$

**Remark.** We avoid explicit alignment of the features through contrastive learning, as the task requires us to preserve differences between some modalities that are *complementary*. The channel-wise softmax function helps us obtain a final vector allowing *substitutional modalities* to work together on the same dimensions. We observe that hand forces and the muscle EMG are highly correlated. In this way, these latent dimensions are implicitly attributed to reflect similar action property, *e.g.* strength for muscle and haptic forces, and thus increase robustness to missing modalities at test-time.

## 2.2. Context-Aware Latent Interaction

Previous steps have taken us to learn features that represent actions. Interaction is a special subset of action that bears the notion of contexts and consequences. We take one step further to investigate ways to represent **interaction**. An effective interaction feature should not only summarize the action property itself but engage with its contexts and hint at potential consequences.

**Latent Projection Interaction.** Under our task setting, interaction describes a way to take the observed context  $x_{[0,t-1]}$  to the consequential states  $x_{[t,T]}$ . In the latent space, vectors that represent interactions are analogous to flow vectors that can be applied to various context states  $z_{x_{[0,t-1]}}$  to the consequential changes states  $z_{x_{[t,T]}}$ . We wish to capture such effects in the latent vector itself. Intuitively, the direction of latent interaction vectors  $\{y'_t\}$  should consistently introduce similar effects relative to any context frames where they are applied. In other words, a good interaction vector should be locally constrained to its context frame, at the same time when applied to different contexts, the interaction vector should introduce similar behavior relative to the new context. These observations encourage us to constrain the behavior of

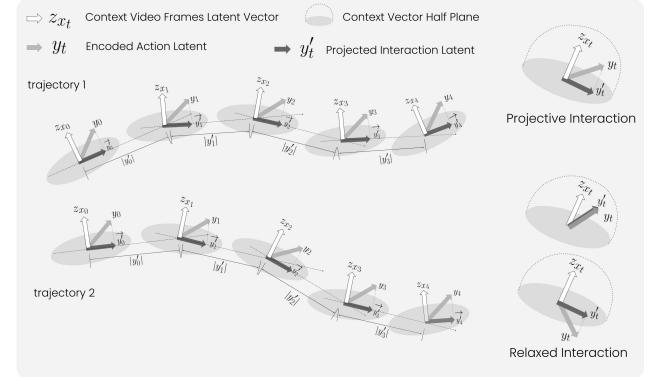


Figure 3. Latent Interaction

action vectors through projective regularization. By removing the projected components on the context vector from the action vector, we extract the orthogonal component of the actions that reflects the dominant direction of change that an action can impose onto its context

$$y'_t = y_t - \left\langle y_t, \frac{z_{x_{t-1}}}{|z_{x_{t-1}}|} \right\rangle \frac{z_{x_{t-1}}}{|z_{x_{t-1}}|}. \quad (2)$$

In addition to direction constraint, we further capture the rate of such changes through an additional supervision signal, by matching the norm of the interaction vector  $y'_t$  with the magnitude of frame-wise differences,  $\mathcal{L}_{\text{NORM}} = ||y'_t|| - |z_{x_t} - z_{x_{t-1}}|||^2$ . As shown in Fig. 3, these constraints help introduce the desired behavior in latent space. The two latent trajectories are formed by imposing the same interaction vector  $y'_t$  to two different context frames  $z_{x_0}$  and  $z_{x'_0}$ . Because the direction of change follows the orthogonal direction locally to the specific context frames and by the same magnitude, the two trajectories are similar.

**Relaxed Hyperplane Interaction.** A geometric interpretation of the latent interaction  $y'_t$  reveals that the relative angle between context  $x_{t-1}$  and interaction  $y'_t$  depicts two spaces partitioned by a hyperplane defined by the normal vector  $z_{x_{t-1}}$ . This observation encourages us to rethink latent interaction modeling. The previous projection perspective forms a hard constraint where the interaction must follow the orthogonal direction of the context. In reality, behaviors of interactions might be slightly different when context changes. Hence, we relax the hard orthogonal projection constraint. Through a geometric lens, the context vector  $z_{x_{t-1}}$  can be viewed as a normal vector that defines a partitioning hyperplane, where interaction  $y'_t$  with significant consequence to  $x_{t-1}$  lies in the positive hemisphere, and negligible interaction resides below the hyperplane is clipped and projected.

$$y'_t = i(y_t, z_{x_{t-1}}) = \begin{cases} y_t & \text{if } \langle y_t, z_{x_{t-1}} \rangle \geq 0 \\ y_t - \left\langle y_t, \frac{z_{x_{t-1}}}{|z_{x_{t-1}}|} \right\rangle \frac{z_{x_{t-1}}}{|z_{x_{t-1}}|} & \text{otherwise} \end{cases}$$

We use this formulation to regularize interaction feature vectors  $y'$  and adopt the frame-wise difference magnitude

constraint. The learned interaction feature  $y_t'$  is used to condition diffusion network to simulate future video frames.

### 2.3. Conditioning Generative Visual Simulator

Inspired by [33, 71], our simulator employs a video diffusion model to solve for future observations. Denoising diffusion [26], in the forward process, predicts noise  $\epsilon \sim \mathcal{N}(0, I)$  applied to video frames  $x_{[t,T]}$  according to a noise schedule  $\bar{\alpha}^n \in \mathbb{R}$  over several steps  $n \in [1, N]$ , where  $\bar{\alpha}^n = \Pi_{s=1}^n \alpha^s$ . The optimization objective to train model  $g_\theta$  is,

$$\mathcal{L}_{\text{VDM}} = \left\| \epsilon - g_\theta \left( \sqrt{\bar{\alpha}^n} x_{[t,T]} + \sqrt{1 - \bar{\alpha}^n} \epsilon, n \mid x_{t-1}, a \right) \right\|^2$$

For the task of future observation prediction, we use the learned model  $g_\theta$  and reverse the process by iteratively denoising an initial noise sample  $x_{[t,T]}^{n=N} \doteq \epsilon \sim \mathcal{N}(0, I)$  to recover video frames  $x_{[t,T]}^{n-1}$  at denoising step  $n-1$ . When  $n=0$ , we obtain the estimated future video frames  $\hat{x}_{[t,T]}$ .

$$x_{[t,T]}^{n-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_{[t,T]}^n - \frac{1 - \alpha^n}{\sqrt{1 - \bar{\alpha}^n}} g_\theta \left( x_{[t,T]}^n, n \mid x_{t-1}, a \right) \right) + \sigma, \sigma \sim \mathcal{N}(0, \frac{1 - \bar{\alpha}^{n-1}}{1 - \bar{\alpha}^n} (1 - \alpha) I)$$

We use I2VGen [74] as our diffusion backbone. It uses a 3D UNet [64] with dual condition architecture that generates future video frames  $x_{[t,T]}$  based on text prompt  $a$  and context image  $x_{t-1}$ . We modify I2VGen [74] replacing the single context frame with a history horizon of  $h$  context frames by concatenating in the channel dimension. We also replace the text conditioning with our learned multimodal action feature  $y_t$ , where the cross attention is applied between noise frame samples and our conditioning feature  $y_t$ . Different from text-prompted simulation [71, 74], where a single text prompt  $a$  is repeatedly used for all frames, our action condition is temporal, allowing our temporal attention to be frame-specific. (moved from end of sec. 2.2) We train the model end-to-end using a weighted sum of the aforementioned loss functions. The final supervision signal is given by  $\mathcal{L} = \lambda_1 \mathcal{L}_{\text{VDM}} + \lambda_2 \mathcal{L}_{\text{SSL}} + \lambda_3 \mathcal{L}_{\text{NORM}}$ , where  $\lambda_1 = 10.0$ ,  $\lambda_2 = 1.0$ ,  $\lambda_3 = 0.1$ . The relative weighting between different loss components  $\{\lambda\}$  are chosen to align the magnitude of each component to the same level. We provide the details of our network architecture in Appendix Sec. 7.4 and in Fig. 11.

## 3. Experiments

We design experiments to answer the following questions:

- Do we need multisensory action data to achieve fine-grained control over simulated videos?
- How do our multimodal feature extraction compare with existing ones when used for conditioning?
- Is our method robust to missing modalities at test time and how they influence prediction?

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
UniSim verb	0.131	14.1	0.332	337.9
UniSim phrase	0.118	14.6	0.321	275.9
UniSim sentence	0.117	14.6	0.317	251.7
Body-pose only	0.127	14.4	0.345	295.9
Hand-pose only	0.122	14.5	0.349	307.6
Muscle-EMG only	0.134	13.8	0.364	348.2
Hand-force only	0.120	14.5	0.334	278.9
Ours multisensory	<b>0.110</b>	<b>16.0</b>	0.276	203.5
Ours w/ phrase	0.113	<b>16.0</b>	<b>0.274</b>	<b>200.4</b>
Ours w/ sentence	0.111	<b>16.0</b>	<b>0.274</b>	201.7

Table 1. Quantitative comparison

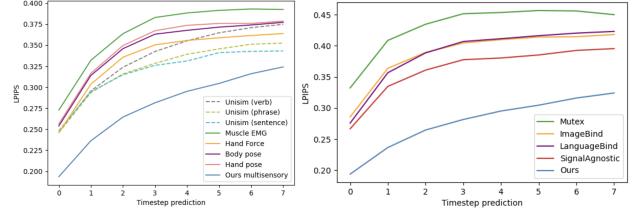


Figure 4. Temporal drift. LPIPS per frame.

**Experimental Setup.** We use the ActionSense [13] dataset for our experiments. It includes five different interoceptions, including hand haptic forces, EMG muscle activities, hand pose, body pose, and gaze tracking. To the best of our knowledge, it is the first and only multi-sensory dataset with paired actuation monitoring and video sequences. While we focus our efforts on multimodal representation, to show the generality of our proposed method, we provide additional qualitative results on other unimodal handpose datasets, H2O [34] and HoloAssist [65] in Sec. 7.8.12 in Appendix. For our main experiments, we use ActionSense dataset and subject five as our test set, and the remaining four subjects as training and validation set. We parse the dataset into paired sequences of 12 frames, and use first 4 frame as the context frames and predict the following 8 frames. All experiments and methods use the same diffusion backbone, modified I2VGen [74] (Sec. 2.3), which is a dual condition video network that predicts frames  $x_{[t,T]}$  based on conditioning prompt  $a$  and context image(S)  $x_{[0,t-1]}$ . We vary the conditioning type  $a$  for all experiments. All methods are trained from scratch on the same data with the same hardware and software setup. Due to computational constraints, our experiments and comparisons are conducted with videos of  $64 \times 64$  resolution. We provide higher resolution results of our model of  $128 \times 128$  and  $192 \times 192$  (Sec. 7.8.11). Experiments on out-of-domain generalization is shown in Sec. 7.8.8.

**Evaluation Metric.** We are interested in how various types of data and method used for conditioning can have different effects when simulating videos. We evaluate on a withheld test set from ActionSense [13], and use three different metrics to evaluate the quality of predicted video trajectories and the ground truth video trajectories, following [71]. We use MSE, PSNR, LPIPS, and FVD scores as evaluation metrics to quantify the quality and accuracy of predicted video

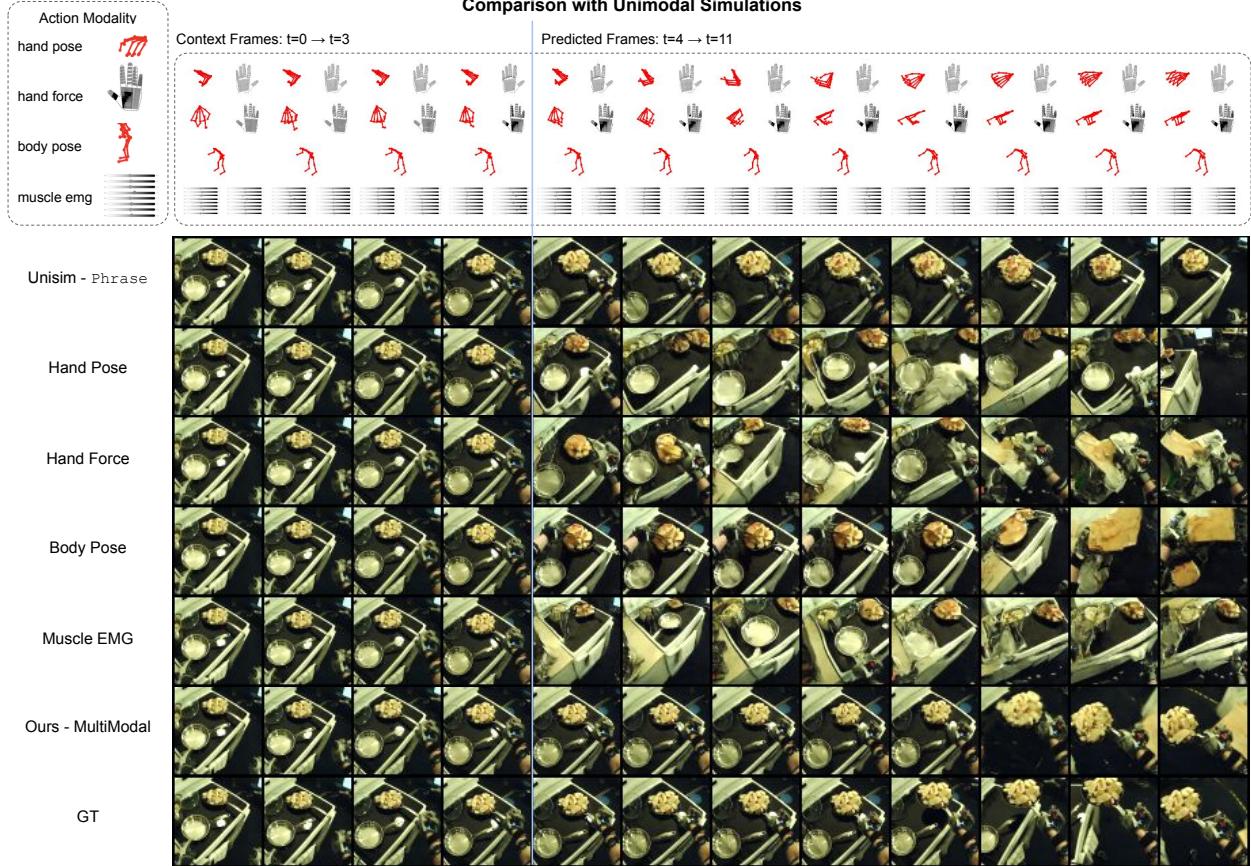


Figure 5. **Comparison to Unimodal Simulation.** We compare our proposed multisensory conditioning to unimodal conditioning, including text and each action sensory modality. The first four frames are the context frames, and the last eight frames are predictions by each method.

frames. In all tables,  $\downarrow$  means lower is better for the metric, and  $\uparrow$  indicates higher is better.

### 3.1. Conditioning Action Modalities

We are interested in understanding whether we need multisensory action data to achieve fine-grained control over simulated videos. To answer this question, we investigate the benefit of different action signal modalities, including text description, unimodal action, and multisensory action as input. For fairness of comparison, we use the same video generation model while varying the condition type.

**Comparison with Text-conditioned Simulation.** We first compare our proposed method and the state-of-the-art text-based video-diffusion simulator, UniSim [71]. We vary the input condition with increasing details in description, using verb, phrase, sentence. Phrase are composed of verbs and subjects, *e.g.* cut potato. We add more detailed descriptions to form sentences, *e.g.* person cut potato in a very fast manner, while holding it with left hand. As shown in Table. 1, our proposed method can achieve more accurate future frame prediction, because it takes temporally fine-grained action trajectories with subtle differences as inputs to control the video prediction to match the action signals for each time step, whereas the subtle differences in the action

Method	MSE $\downarrow$	PSNR $\uparrow$	LPIPS $\downarrow$	FVD $\downarrow$
Mutex	0.164	12.4	0.431	410.1
Imagebind	0.134	13.9	0.390	315.6
Languagebind	0.143	13.7	0.387	332.0
SignalAgnostic	0.127	14.3	0.361	267.5
Ours	<b>0.110</b>	<b>16.0</b>	<b>0.276</b>	<b>203.5</b>

Table 2. Quantitative comparison on multimodal feature extraction. trajectory are difficult to be accurately captured by text descriptions. Fig. 7 further demonstrates that our method can be used to generate more diverse video trajectories from the same context frames, whereas text-conditioned video simulation is more prone to mode collapse, converging to similar future frame predictions from similar context frames. These new video trajectories generated with our method can be used for data augmentation to compensate the scarcity of paired action video data. As shown in Table. 1 and Fig. 7, adding `text phrase` as an additional modality to our method can help reduce model confusion. Additional discussion is included in Appendix Sec. 7.8.1.

**Comparison with Unimodal Action Simulation.** We extend our experiments to test the necessity of **multimodal** interaction by comparing to each action modality alone. As there lacks direct baseline method that utilizes these action modalities for simulation, we use our own method for en-

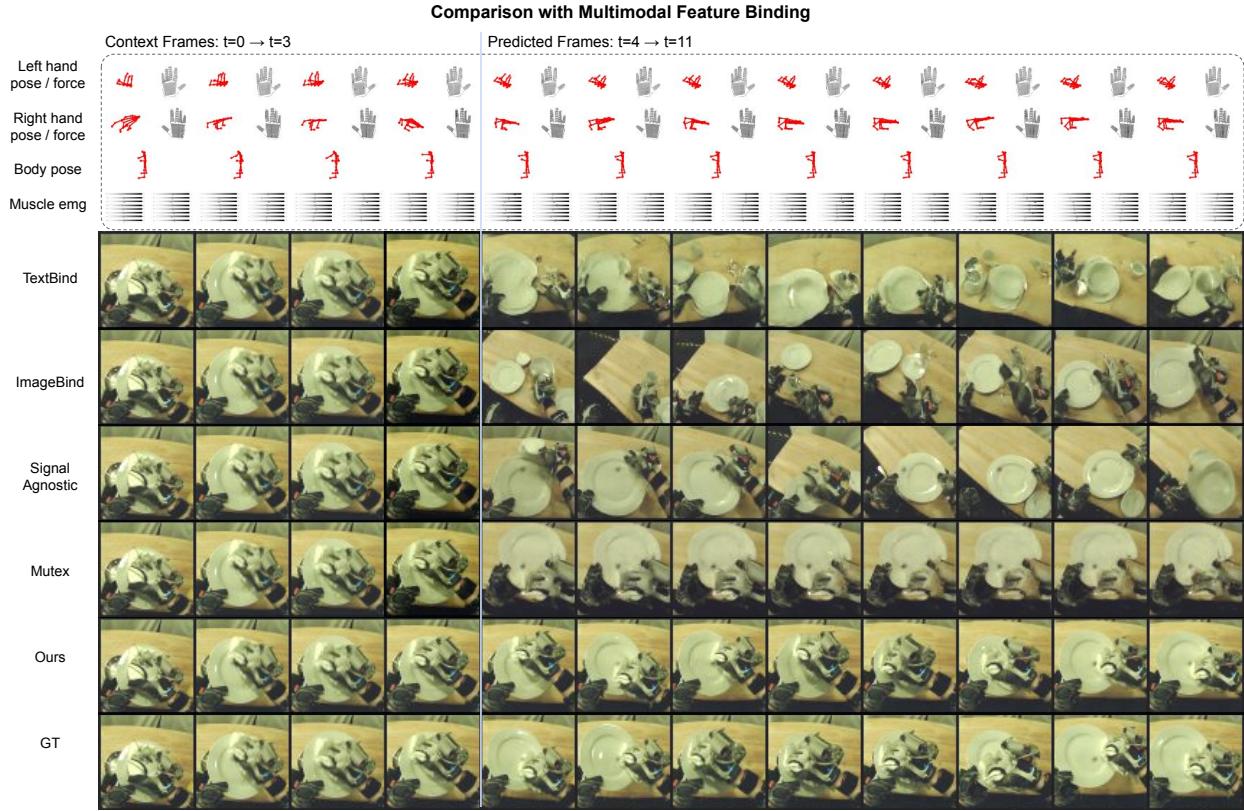


Figure 6. **Comparison with multimodal feature extraction baselines.** We compare with various multimodal feature extraction methods for conditioning the video simulator. Similarly, the first four frames are context frames and the last eight frames are predictions.

coding these modalities and conditioning video models. The closest work to one of our unimodal baseline setting is [30], which uses a two stage finetuning of stable diffusion to generate full-body videos from pixel-level dense poses assuming static camera. The assumptions of dense poses, static camera, and full-body video make it difficult and unfair for this method to tackle our task setting with egocentric videos.

The middle section in Table. 1 shows that future video frame prediction is most accurate when all modalities are combined together. This is because not all modalities are created equal, and our ability to swiftly control and operate with our surroundings is a multiplicative effect of different functions working together. As shown in Fig. 5, a simple task of removing the pan from the stove top requires us to reach to the pan (body pose), grab the pan (hand pose and force), lift the pan (muscle and body pose), and finally turn around(body pose). When training only with hand-forces, the model has no information to locate the hand, and thus generate hand holding random things in the image instead of the pan and results drift off (Fig. 5). We almost never entirely isolate one sense to interact with the world. Therefore, training with a single modality is not enough for such tasks, even when each signal is temporally fine-grained.

### 3.2. MultiModal Feature for Generative Simulation

For the task of multisensory action controlled simulation, we study how multimodal action representations impacts

explicit pixel space. We compare our method with various state-of-the-art multimodal feature extraction paradigms:

- Mutex [59] proposes to randomly mask out and project some of the input modalities and directly align and match the remaining modalities to future frames.
- LanguageBind [76] proposes to use text as a binding modality instead of using images.
- ImageBind [19] is a contrastive binding technique that leverages InfoNCE [47] contrastive loss to bind different modality of features to clip-encoded image features.
- Signal-Agnostic learning [16, 37] extracts cross-modal feature using signal-agnostic neural field.

As shown in Table 2, our multi-sensory interaction feature outperforms baseline methods for multi-modal feature extraction in controlled generative simulation. Different multimodal tasks require distinct representations. Previous approaches [19, 40, 53, 57, 76], designed mainly for cross-modal retrieval, extract shared information via contrastive learning or modality anchoring, emphasizing interchangeability between modalities. However, in generative simulation, each action modality captures unique, complementary aspects of human behavior. For example, TextBind [76] uses contrastive loss to align various modalities with text descriptions, which can erase the fine-grained temporal details of action signals, leading to compromised predictions. Similarly, ImageBind [19] and Mutex [59] align fea-

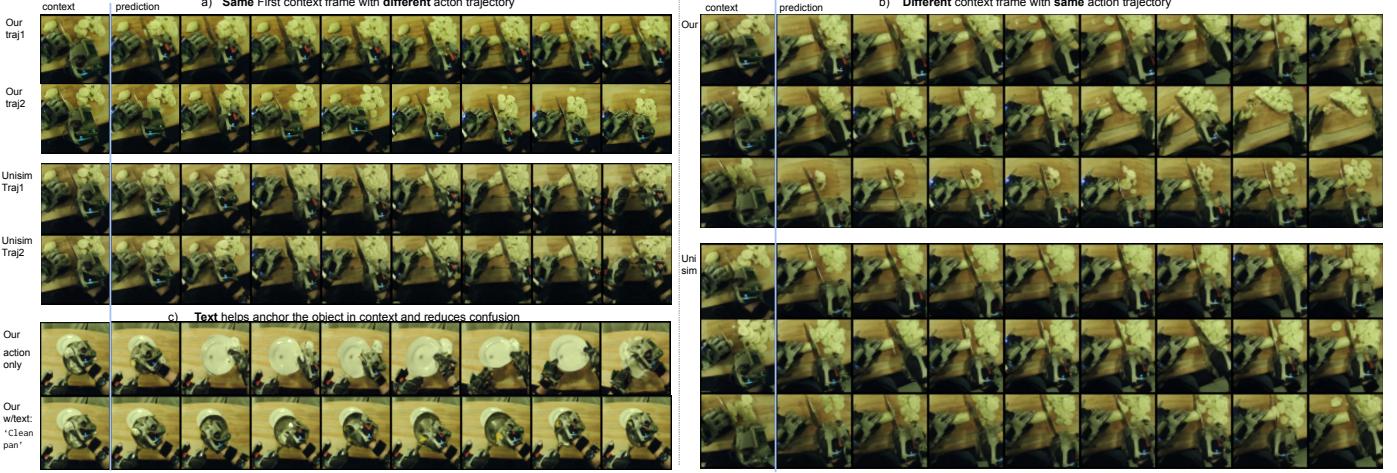


Figure 7. **Simulating new video trajectories** Comparing our multisensory method and text-based Unisim in generating diverse video trajectories from same or different context frames. We show the last context frame  $x_{t-1}$  and the predicted video frames  $x_{[t, T]}$ .

tures with visual frames, either by contrastive loss or L2 regression against pretrained CLIP features, but the inherent one-to-many mapping between similar actions and different visual contexts hampers the network’s ability to extract intrinsic motion, resulting in error accumulation and mode collapse. Signal agnostic learning [16, 37] avoids contrastive loss by letting gradients from different modalities optimize a shared latent manifold, yet its loose coupling between action and video modalities also leads to larger error. Therefore, generative simulation demands representations that preserve the complementary nature of signals. To meet these requirements, our proposed method is better suited for this task.

### 3.3. Ablation Experiments

We provide comprehensive ablation studies to show how different senses help with video prediction. We also conduct ablation studies to validate various design choices and effect of history horizon length (Appendix Sec. 7.8.5).

**Robustness to Test Time Missing Modalities.** We evaluate our model trained on all modalities with each of the modalities removed, shown in Table 3a. We can see that the prediction accuracy of our model is slightly influenced by ablated modalities during test time. From the right side of Fig. 8, we can see that our model can still make sensible predictions under missing modalities, although prediction is most accurate with all modalities included. The left side of the Fig. 8 shows a stress test evaluating our model provided with only one modality. We see when that the hand pose trajectory is more accurate compared to other ones, which hint at a task-specific critical modality. Comprehensive test-time robustness tests are included in Appendix Sec. 7.8.3. Additional results on training with ablated modalities are included in Appendix Sec. 7.8.2.

**Multimodal Feature Extraction** We investigate how different multi-sensory fusion strategies affect simulated video trajectories. To validate our softmax-ensemble approach, we compare it with common symmetric fusion functions. As shown in Table 3b, softmax outperforms mean and max

Table 3. Ablation Experiments

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
No hand pose	0.111	15.3	0.304	205.1
No hand force	0.113	15.5	0.307	205.0
No body pose	0.115	15.3	0.304	205.6
No muscle EMG	0.113	15.2	0.291	204.7
All sensory used	<b>0.110</b>	<b>16.0</b>	<b>0.276</b>	<b>203.5</b>

(a) Testing with missing modalities

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
Max	0.128	14.1	0.294	284.8
Mean	0.126	14.4	0.293	285.3
Concatenation	0.117	15.0	0.282	279.9
Without $y'$	0.142	13.7	0.327	339.0
Projection $y'$	0.116	14.5	0.288	265.5
Ours full	<b>0.110</b>	<b>16.0</b>	<b>0.276</b>	<b>203.5</b>

(b) Ablation of network components

pooling. We avoid direct feature concatenation to maintain permutation invariance and ensure robustness when some modalities are missing at test time. We also perform an ablation study on our interaction feature  $y'$  learning scheme. Table 3b shows that removing the interaction module and using the action feature  $y$  as a condition significantly degrades performance. Although the action feature contains all action information, it does not effectively modify the context frame, leading the downstream video model to focus on irrelevant details and causing mode collapse. Adding hard projection regularization on  $y'$  greatly improves video prediction accuracy, though it remains slightly inferior to our full pipeline that employs the relaxed hyperplane interaction scheme.

## 4. Downstream Applications

We show two potential downstream applications of our work in policy optimization shown below and multimodal action planning shown in Appendix Sec. 7.8.10.

**Low-level Policy Optimization** One downstream application of our proposed action-conditioned video generative simulator is to optimize a policy of low-level actuation. Inspired by [71], we set up task as goal-conditioned policy

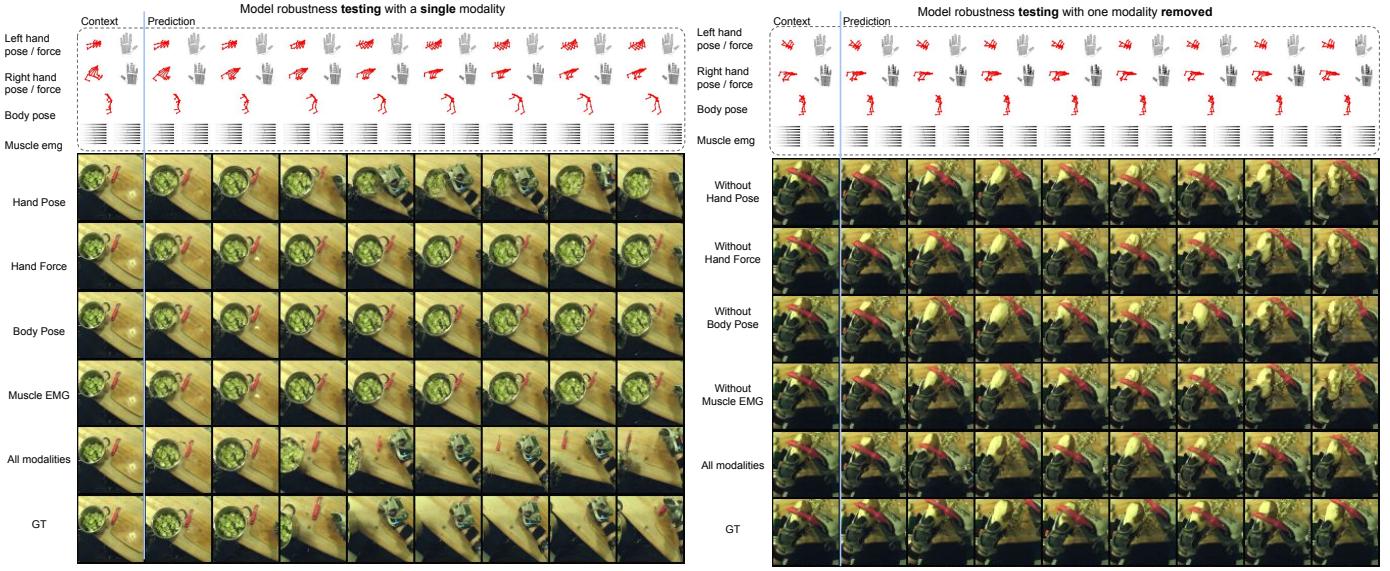


Figure 8. **Robustness to missing modalities during test time.** Left side shows stress test with evaluating with one single modality *provided*. Right side shows testing with one modality *removed*. We show the last context frame  $x_{t-1}$  and the predicted video frames  $x_{[t,T]}$ .

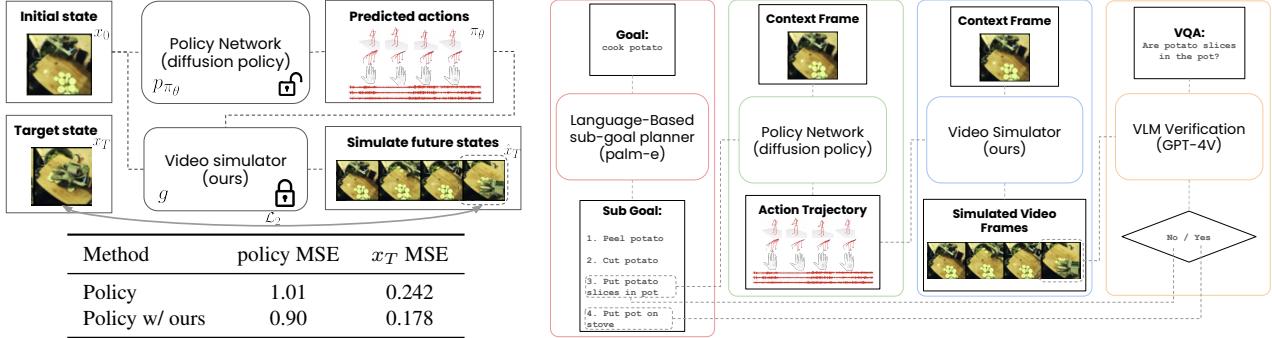


Figure 9. **Left:** Pipeline for goal-conditioned policy optimization. **Right:** Pipeline for long-term task planning.

optimization, where we optimize a policy to generate a trajectory of low-level actuation  $a_{[1,T]}$  that brings the environment from start state  $s_0$  to target  $s_T$ . States are described by images  $s_t \doteq x_t$ . We show one use case of our model in goal-conditioned policy optimization. We compare training of the same policy network  $p(\cdot)_{\pi_\theta}$  under two conditions. First, we define the baseline method using the commonly employed goal-conditioned policy training approach [10, 15, 54]. This baseline is the policy network taking the starting state and target state, depicted by two video frames  $x_0$  and  $x_T$ , and directly regress policy  $\pi_\theta$  minimizing the L2 distance between the predicted action  $\hat{a}[1,T] = \pi_\theta(x_0, x_T)$  and ground truth expert action trajectory  $a[1,T]$ . This L2 loss term is defined as  $\mathcal{L}_a = \sum_t \|\hat{a}_t - a_t\|_2 = \|p(x_0, x_T)_{\pi_\theta} - a[1,T]\|_2$ . The second condition is to train the same policy  $\pi_\theta$  in conjunction with our pretrained simulator. We feed the action trajectory predicted by policy network  $\hat{a}_{[1,T]} = \pi_\theta(x_0, x_T)$  into our pretrained simulator model  $g(\cdot)$  to predict the video frames from this action trajectory  $\hat{x}_T = g(p(x_0, x_T)_{\pi_\theta})_T$ . This additional loss term is defined as  $\mathcal{L}_{sim} = \|\hat{x}_T - x_T\|_2 = \|g(p(x_0, x_T)_{\pi_\theta})_T - x_T\|_2$ . The total loss term for the second condition is  $\mathcal{L}_{simpolicy} = \mathcal{L}_a + \mathcal{L}_{sim}$ . We evaluate the effectiveness of by using L2 distance between the predicted

action  $\hat{a}_{[1,T]}$  and ground truth action  $a_{[1,T]}$ , which is defined  $\|\hat{a}_{[1,T]} - a_{[1,T]}\|_2$ .  $x_T$  MSE is a supporting metric that compares target state and the simulated end state using our simulator. Unfortunately, no other multisensory action simulator exist to use for further validation. We see from Fig. 9 that adding our additional supervision signal helps to improve policy optimization. Directly regressing multi-sensory actions with a policy network is difficult because the action space in our task setting is quite large, 2292 dimensional. More results are shown in Fig. 15 in Appendix Sec. 7.8.9.

## 5. Conclusion

In this work, we introduce the concept of multisensory interaction for fine-grained generative simulation. We focus on learning effective multisensory feature representations to effectively control a downstream video generative simulator. Our proposed multimodal feature extraction paradigm along with regularization scheme to extract action feature vectors capable of accurately controlling video prediction and robust to missing modalities at test time. We conduct extensive comparisons, ablations, and downstream applications to showcase the merits of our method. We hope our work brings insights and inspirations to the research community.

## 6. Acknowledgement

We thank Tianwei Yin, Congyue Deng, and Shivam Dugal for the useful feedback and discussion. The work is supported by Research Grant from Delta Electronics.

## References

- [1] open-clip-torch: OpenAI CLIP Implementation in PyTorch. <https://pypi.org/project/open-clip-torch/>. Accessed: 2023-06-10. 13
- [2] Alessandro Achille and Stefano Soatto. A separation principle for control in the age of deep learning. *Annual Review of Control, Robotics, and Autonomous Systems*, 1:287–307, 2018. 12
- [3] Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*, 2023. 12
- [4] Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *Proceedings of the IEEE international conference on computer vision*, pages 609–617, 2017. 12
- [5] Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 1995. 12
- [6] Pablo Samuel Castro. Scalable methods for computing state similarity in deterministic markov decision processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 10069–10076, 2020. 12
- [7] Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learning with transformer world models. *arXiv preprint arXiv:2202.09481*, 2022. 12
- [8] Ziyang Chen, Shengyi Qian, and Andrew Owens. Sound localization from motion: Jointly learning sound direction and camera rotation. *arXiv*, 2023. 12
- [9] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023. 14, 17
- [10] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023. 8
- [11] Yung-Yu Chuang, Dan B Goldman, Ke Colin Zheng, Brian Curless, David H Salesin, and Richard Szeliski. Animating pictures with stochastic motion textures. In *ACM SIGGRAPH 2005 Papers*, pages 853–860. 2005. 12
- [12] Virginia R de Sa. Learning classification with unlabeled data. *Advances in neural information processing systems*, pages 112–112, 1994. 12
- [13] Joseph DelPreto, Chao Liu, Yiyue Luo, Michael Foshey, Yunzhu Li, Antonio Torralba, Wojciech Matusik, and Daniela Rus. Actionsense: A multimodal dataset and recording framework for human activities using wearable sensors in a kitchen environment. In *Advances in Neural Information Processing Systems*, pages 13800–13813. Curran Associates, Inc., 2022. 4, 14, 17
- [14] Karan Desai and Justin Johnson. Virtex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11162–11173, 2021. 12
- [15] Yiming Ding, Carlos Florensa, Pieter Abbeel, and Mariano Phielipp. Goal-conditioned imitation learning. *Advances in neural information processing systems*, 32, 2019. 8
- [16] Yilun Du, M. Katherine Collins, , B. Joshua Tenenbaum, and Vincent Sitzmann. Learing signal-agnostic manifolds of neural fields. In *Advances in Neural Information Processing Systems*, 2021. 1, 6, 7
- [17] Yilun Du, Mengjiao Yang, Pete Florence, Fei Xia, Ayzaan Wahid, Brian Ichter, Pierre Sermanet, Tianhe Yu, Pieter Abbeel, Joshua B Tenenbaum, et al. Video language planning. *arXiv preprint arXiv:2310.10625*, 2023. 1, 17
- [18] Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for finite markov decision processes. In *UAI*, pages 162–169, 2004. 12
- [19] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *CVPR*, 2023. 1, 6, 12, 13
- [20] Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, Marc Sun, and Benjamin Bossan. Accelerate: Training and inference at scale made simple, efficient and adaptable. <https://github.com/huggingface/accelerate>, 2022. 14
- [21] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020. 12
- [22] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023. 12
- [23] Philippe Hansen-Estruch, Amy Zhang, Ashvin Nair, Patrick Yin, and Sergey Levine. Bisimulation makes analogies in goal-conditioned reinforcement learning. In *International Conference on Machine Learning*, pages 8407–8426. PMLR, 2022. 1
- [24] Zekun Hao, Xun Huang, and Serge Belongie. Controllable video generation with sparse trajectories. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7854–7863, 2018. 12
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 13
- [26] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020. 4
- [27] Xixi Hu, Ziyang Chen, and Andrew Owens. Mix and localize: Localizing sound sources in mixtures. *Computer Vision and Pattern Recognition (CVPR)*, 2022. 12
- [28] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal

- Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. If you use this software, please cite it as below. 1, 14
- [29] Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from large weakly supervised data. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*, pages 67–84. Springer, 2016. 12
- [30] Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dreampose: Fashion image-to-video synthesis via stable diffusion. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22623–22633. IEEE, 2023. 6
- [31] Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models. *ArXiv*, abs/2312.02696, 2023. 14
- [32] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. 14
- [33] Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and Joshua B Tenenbaum. Learning to Act from Actionless Videos through Dense Correspondences. *ICLR*, 2024. 4
- [34] Taein Kwon, Bugra Tekin, Jan Stühmer, Federica Bogo, and Marc Pollefeys. H2o: Two hands manipulating objects for first person interaction recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10138–10148, 2021. 4, 13, 17
- [35] Timothée Lesort, Natalia Díaz-Rodríguez, Jean-François Goudou, and David Filliat. State representation learning for control: An overview. *Neural Networks*, 108:379–392, 2018. 12
- [36] Mengxi Li, Rika Antonova, Dorsa Sadigh, and Jeannette Bohg. Learning Tool Morphology for Contact-Rich Manipulation Tasks with Differentiable Simulation. In *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023. 1
- [37] Yichen Li, Yilun Du, Chao Liu, Francis Williams, Michael Foshey, Benjamin Eckart, Jan Kautz, Joshua B Tenenbaum, Antonio Torralba, and Wojciech Matusik. Learning to jointly understand visual and tactile signals. In *The Twelfth International Conference on Learning Representations*, 2023. 1, 6, 7, 12
- [38] Zhengqi Li, Richard Tucker, Noah Snavely, and Aleksander Holynski. Generative image dynamics, 2023. 12
- [39] Lennart Ljung and Torkel Glad. *Modeling of dynamic systems*. Prentice-Hall, Inc., 1994. 12
- [40] Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*, 2023. 6
- [41] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European conference on computer vision (ECCV)*, pages 181–196, 2018. 12
- [42] Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, and Deepak Pathak. Discovering and achieving goals via world models. In *NeurIPS*, 2021. 1
- [43] Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample efficient world models. *arXiv preprint arXiv:2209.00588*, 2022. 12
- [44] Basil Mustafa, Carlos Riquelme, Joan Puigcerver, Rodolphe Jenatton, and Neil Houlsby. Multimodal contrastive learning with limoe: the language-image mixture of experts. *Advances in Neural Information Processing Systems*, 35:9564–9576, 2022. 2
- [45] Medhini Narasimhan, Shiry Ginosar, Andrew Owens, Alexei A Efros, and Trevor Darrell. Strumming to the beat: Audio-conditioned contrastive video textures. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3761–3770, 2022. 12
- [46] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y Ng. Multimodal deep learning. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 689–696, 2011. 12
- [47] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 6
- [48] OpenAI and Josh Achiam et al. Gpt-4 technical report, 2024. 17
- [49] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H Adelson, and William T Freeman. Visually indicated sounds. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2405–2413, 2016. 12
- [50] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*. 2019. 13
- [51] Wentian Qu, Zhaopeng Cui, Yinda Zhang, Chenyu Meng, Cuixia Ma, Xiaoming Deng, and Hongan Wang. Novel-view synthesis and pose estimation for hand-object interaction from sparse views. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15100–15111, 2023. 13
- [52] Gorjan Radevski, Dusan Gruijicic, Matthew Blaschko, Marie-Francine Moens, and Tinne Tuytelaars. Multimodal distillation for egocentric action recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5213–5224, 2023. 2
- [53] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 6, 12

- [54] Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal conditioned imitation learning using score-based diffusion policies. In *Robotics: Science and Systems*, 2023. 8
- [55] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021. 2
- [56] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 13
- [57] Ludan Ruan, Yiyang Ma, Huan Yang, Huigu He, Bei Liu, Jianlong Fu, Nicholas Jing Yuan, Qin Jin, and Baining Guo. Mm-diffusion: Learning multi-modal diffusion models for joint audio and video generation. In *CVPR*, 2023. 6, 13
- [58] Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement learning with action-free pre-training from videos. In *International Conference on Machine Learning*, pages 19561–19579. PMLR, 2022. 12
- [59] Rutav Shah, Roberto Martín-Martín, and Yuke Zhu. Mutex: Learning unified policies from multimodal task specifications. In *7th Annual Conference on Robot Learning*, 2023. 1, 6
- [60] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. Animating arbitrary objects via deep motion transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2377–2386, 2019. 12
- [61] Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991. 12
- [62] Bingjie Tang, Michael A Lin, Iretiayo Akinola, Ankur Handa, Gaurav S Sukhatme, Fabio Ramos, Dieter Fox, and Yashraj Narang. Industreal: Transferring contact-rich assembly tasks from simulation to reality. In *Robotics: Science and Systems*, 2023. 1
- [63] Yunsheng Tian, Jie Xu, Yichen Li, Jieliang Luo, Shinjiro Sueda, Hui Li, Karl D.D. Willis, and Wojciech Matusik. Assemble them all: Physics-based planning for generalizable assembly by disassembly. *ACM Trans. Graph.*, 41(6), 2022. 1
- [64] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023. 4
- [65] Xin Wang, Taein Kwon, Mahdi Rad, Bowen Pan, Ishani Chakraborty, Sean Andrist, Dan Bohus, Ashley Feniello, Bugra Tekin, Felipe Vieira Frujeri, Neel Joshi, and Marc Pollefeys. Holoassist: an egocentric human interaction dataset for interactive ai assistants in the real world. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 20270–20281, 2023. 4
- [66] Yaohui Wang, Di Yang, Francois Bremond, and Antitza Dantcheva. Latent image animator: Learning to animate images via latent space navigation. *arXiv preprint arXiv:2203.09043*, 2022. 12
- [67] Chung-Yi Weng, Brian Curless, and Ira Kemelmacher-Shlizerman. Photo wake-up: 3d character animation from a single photo. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5908–5917, 2019. 12
- [68] Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. Slotformer: Unsupervised visual dynamics simulation with object-centric models. *arXiv preprint arXiv:2210.05861*, 2022. 12
- [69] Tianfan Xue, Jiajun Wu, Katherine L Bouman, and William T Freeman. Visual dynamics: Stochastic future generation via layered cross convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2236–2250, 2018. 12
- [70] Fengyu Yang, Chenyang Ma, Jiacheng Zhang, Jing Zhu, Wen-zhen Yuan, and Andrew Owens. Touch and go: Learning from human-collected vision and touch. *Neural Information Processing Systems (NeurIPS) - Datasets and Benchmarks Track*, 2022. 12
- [71] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. *ICLR*, 2023. 1, 4, 5, 7, 12, 13, 16, 17
- [72] Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in projected latent space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18456–18466, 2023. 12
- [73] Lingzhi Zhang, Shenghao Zhou, Simon Stent, and Jianbo Shi. Fine-grained egocentric hand-object segmentation: Dataset, model, and applications. In *European Conference on Computer Vision*, pages 127–145. Springer, 2022. 13
- [74] Shiwei\* Zhang, Jiayu\* Wang, Yingya\* Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qing, Xiang Wang, Deli Zhao, and Jingen Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. 2023. 4, 12, 13
- [75] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022. 12
- [76] Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, Wang HongFa, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Cai Wan Zhang, Zhifeng Li, Wei Liu, and Li Yuan. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment, 2023. 1, 6, 12
- [77] Karl J. Åström and Björn Wittenmark. Adaptive control of linear time-invariant systems. *Automatica*, 9(6):551–564, 1973. 12

## 7. Appendix

1. Disclaimer
2. Notation
3. Related Work
4. Implementation Details
  - (a) Network Architecture
  - (b) Hardware, Software, Training Setup
  - (c) Experimental Setup
5. Additional Pipeline Figure
6. Model Size
7. Discussion of Limitations and Future Work
8. Additional Experiments and Discussions
  - (a) Text as addition to multisensory action
  - (b) Additional ablation experiment: comprehensive test-time robustness
  - (c) Additional comparison between training with limited modalities and testing with missing modalities
  - (d) Additional ablation experiment: effect of history horizon length
  - (e) Comprehensive Cross-subject testing: using other subject as test set and train on the rest.
  - (f) Examples of fine-grained control
  - (g) Generalization to out-of-domain OOD data through model finetuning
  - (h) Additional qualitative results on other dataset
  - (i) Downstream application 2: multisensory action planning
  - (j) Additional discussion and results on and downstream application
9. Higher Resolution Results
10. Additional Qualitative Results
11. Discussion of failure cases

### 7.1. Disclaimer

This is a research work where the primary focus is introducing a new task and a method to learn effective multimodal representation for generative simulation. We devise our multimodal feature extraction as generic to be combined when stronger video generation backbone is invented. High-resolution videos are **not** the main focus of this work. We **provide higher resolution** results of our model in Sec. 7.8.11, and we conduct all experiments shown in the paper using the same video resolution, including our model and all baseline methods trained. We hope our work can inspire future research works and industrial efforts to build foundational digital twin of our world with fine-grained control. We hope that our work can be used to scale with more abundant resources.

### 7.2. Notation Chart

We summarize the notation used in our paper in Table. 4.

### 7.3. Related Work

**Learning Multi-Modal Representations.** Learning shared representations across various modalities has been instrumental in a variety of research areas. Early research by De Sa et al. [12] pioneered the exploration of correlations between vision and audio. Since then, many deep learning techniques have been proposed to learn shared multi-modal representations, including vision-language [14, 29, 41, 53], audio-text [3], vision-audio [4, 27, 45, 46, 49], vision-touch [37, 70], and sound with Inertial Measurement Unit (IMU) [8]. Recently, ImageBind [19] and LanguageBind [76] demonstrate that images and text could successfully bind multiple modalities, including audio, depth, thermal, and IMU, into a shared representation. However, these previous efforts take bind-all fuse-all perspective, which takes away many of the inherent differences brought by various sensory modalities. Our work takes a different perspective. By differentiating between the active and passive senses, we allow a bilateral model to arise and capture the interaction between the two. The prior fuse-all strategy also overshadows an inherent need in multi-modal representation learning, which is interaction. We propose a representation learning scheme to capture the nature of multi-modal interactions.

**Learning World Models.** Learning accurate dynamics models to predict environmental changes from control inputs has long challenged system identification [39], model-based reinforcement learning [61], and optimal control [5, 77]. Most approaches learn separate lower-dimensional state space models per system instead of directly modeling the high-dimensional pixel space [2, 6, 18, 35]. While simplifying modeling, this limits cross-system knowledge sharing. Recent large transformer architectures enable learning image-based world models, but mostly in visually simplistic, data-abundant simulated games/environments [7, 21, 22, 43, 58, 68]. Prior generative video modeling works leverage text prompts [72, 75], driving motions [60, 66], 3D geometries [67, 69], physical simulations [11], frequency data [38], and user annotations [24] to introduce video movements. Recently, Yang et al. [71] proposes Unisim, which uses text conditioned video diffusion model as an interactive visual world simulator. However, these prior works focus on using text as condition to control video generation, which limits their ability to precisely control the generated video output, as many fine-grained interactions and subtle variations in control are difficult to be accurately described only using text. We propose to use complementary multi-sensory data to achieve more fine-grained temporal control over video generation through multi-sensory action conditioning.

### 7.4. Implementation Details

**Network Architecture Detail** We use the open-source I2VGen [74] video diffusion network as our backbone. We

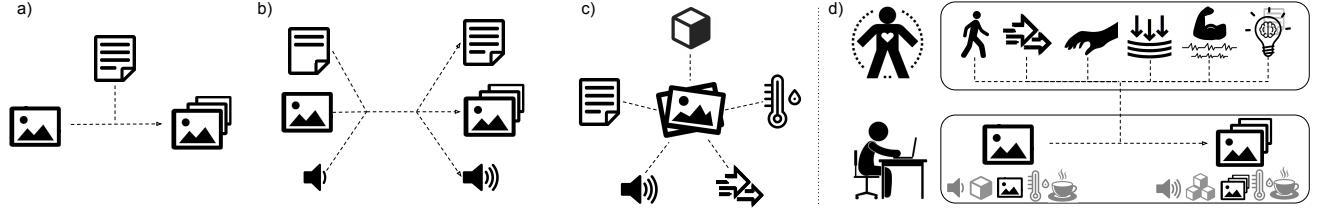


Figure 10. Existing multimodal learning tasks focus on vision-language binding, cross-modal retrieval, and modality anchoring focuses on mining the similarity between different modalities of data (a, b, c) [19, 57, 71]. On the other hand, the task of multisensory action conditioned generative simulation (d) need to understand the unique aspect of each interoceptive action modalities (top) and combine the synchronously to change the exteroception of the external world (bottom).

modify original I2VGen to take pixel space data by changing the input channel to 3 (originally set to 4) and change input image size to  $64 \times 64$ . We keep all other parameters unmodified, and vary the input condition type. We note that single condition models that only use image or text such as Stable Diffusion [56] and etc. are not sufficient for our purpose.

All text input are encoded using CLIP text encoder from the open-source OpenClip [1] library. Images are encoded also using OpenClip Image encoder. Specifically, we use the *ViT-H-14* version with *laion2b\_s32b\_b79k* weights. Please refer to the original papers [1, 74] their architecture details. We describe the architecture of the remaining modules of our model.

Signal specific encoder heads for hand pose, body pose, emg uses the same MLP architecture with different input dimension. The input dimension for hand pose is  $24 \times 3 \times 8$ , body pose is  $28 \times 3 \times 8$ , emg is  $8 \times$ , hand force is  $32 \times 32 \times 8$ . MLP is composed of four layers, with GeLU activation. We set the hidden and output dimension of 128. We apply a dropout with  $p=0.1$ , with batchnorm applied in the first two layers. All encoded signals then goes through a three-layer MLP projection head to project the encoded feature to the same space  $\mathbb{R}^{1024}$  as the clip image feature. The projection MLP also uses GeLU activation with dimensions of [input\_dim, 512, 768, 1024]. We apply batchnorm after the first layer. The set of features are then aggregated across the sensory modalities and masked by a softmax in the modality dimension.

For the latent interaction layers, we use each context frame vector and the action vector for the correponding timestep  $t$  for the context frame feature regularization, we use the aggregated average context frame feature  $z_{x_t}$  to form the context vector for the current action features.

For the experiments comparing to unimodal action sensories, we use our own method for encoding these modalities and conditioning video model. For the sensory modalities of muscle EMG and hand forces, there lacks research works concerning the senses of muscle activation and haptic forces. For hand poses, most works concerning hand poses tackle the task of detection of hand regions from videos [34, 51, 73]. Therefore they also cannot be directly adapted to compare with our work. For this reason, we use our own method for encoding these modalities and conditioning video model.

For experiments on down stream application, we follow the original diffusion policy implementation. The image prompted DP (Sec. 4) uses ResNet [25]-18 image encoder, and the text prompted DP (Sec. 7.8.10) uses OpenClip [1] text-encoder. We modify the original 1D UNet to be four layers with hidden dimensions set to [128, 256, 512, 1024]. The dimension of action space comes to 2292, with two hand poses  $24 \times 3 \times 2$ , one body pose  $28 \times 3 \times 1$ , two arm muscle emg  $8 \times 2$ , two hand forces is  $32 \times 32 \times 2$ .

**Hardware, Software, Training Setup** We use a server with 8 NVIDIA H100 GPU, 127 core CPU, and 1T RAM to train our models for 15 days. We implement all models using the Pytorch [50] library of version 2.2.1 with CUDA

time frame	$t$
history horizon	$[0, t - 1]$
future frames	$[t - 1, T]$
video frame	$x_t$
encoded video frame	$z_{x_t}$
action modality	$m$
action modality signal	$a_{t,m}$
encoded action modality $m$ signal at time step $t$	$z_{t,m}$
j-th dimension of encoded action modality $m$ signal at time step $t$	$z_{t,m,j}$
cross-modal feature	$y_t$
regularized cross-modal feature	$y'_t$

Table 4. Notation Chart

12.1, and accelerator [20] and EMA [31]. We train our models with batch size of 18 per GPU. We use the Adam [32] optimizer with learning rate of  $1e - 4$  and betas (0.9, 0.99), ema decay at 0.995 every 10 iterations.

**Experimental Setup** The ActionSense [13] dataset does not contain the detailed text description used in Sec. 3.1. We generate these text descriptions by using several metrics. We augment the original dataset by resampling video frames, three-ways, every frame, every other frame, and every three frames. We add description of *slow* in speed to the first chunk of data, and *fast* in speed to the third chunk of data. Additionally we also calculate the average hand force magnitude for every task. If the hand force sequence contains frames that are significantly larger than the average frame we add *holding tightly* and add *holding gently* to the lowest force data sequences.

## 7.5. Additional Pipeline Figure

We provide additional pipeline Fig. 11.

## 7.6. Model Size

We report the modules of our model in Table. 5. We can see that the multimodal action signal module is fairly small compared to the video module. Each signal average to around 18044828 parameters which is only 5 percent of the total model weights. The lightweight action signal heads highlights the advantage of our method for low computational cost added for each action signal modality

module	parameter count	percentage of total
signal expert encoder	43780932	0.13
signal projection	11537408	0.03
signal decoder	28398382	0.08
signal Total	83716722	0.25
video model	252380168	0.75
total model	336096890	1.00

Table 5. Parameter Count on  $64 \times 64$  model.

Additionally, people are frequently concerned the real-time execution and edge device computing. We would like to highlight that our work proposes a multisensory conditioned video simulator. When employed in robotics applications, simulator are used in to train policy networks. Normally, only the trained policy network, rather than the simulator itself, needs to be deployed on edge devices / robots. In general, simulators, including ours, do not require to be executed on edge devices or robots for real-time deployment.

We show such application in Sec. 4 Downstream application. Similar to UniSim or any other robotic simulators, we train a goal-conditioned policy network using our pretrained video model. We directly adopt diffusion policy [9] as our policy network, which is lightweight (shown below) and can

be executed on Jetsons as shown below, the parameter count for the policy network trained in Table. 6.

## 7.7. Discussion of Limitations and Future Work

Our experiments are conducted on datasets of human actuation and activities. Ideally, it would be interesting to see the deployment of planned and optimized policies on real humanoid robots with similar multi-sensory capabilities. Because we currently do not have such hardware setup that enables dense force readings on human-hand-like robotic hands or various other fine-grained interoceptive modalities on humanoid robots. We leave this direction for a future research.

There are other passive exteroceptive senses that can be combined with vision, such as depth, 3D and audio etc. One can directly leverage a multi-branch visual-audio or visual-depth UNet diffusion model as the backbone to achieve such multi-modal experception responses. However, due to limited availability of such data, we leave this direction as future work.

Additionally, because of limited computational resources, we limit our video diffusion model to be very low resolution. However, one can employ upsampling approaches to map low-resolution video predictions to higher resolution. Our work is less concerned with the specifics of image quality but more with the application of using multi-sensory interoception data. Therefore, we leave the study of low-cost video upsampling or better video diffusion backbone as future work.

## 7.8. Additional Experiments and Discussion

### 7.8.1. Text as addition to multisensory actions

We are also interested in learning whether multi-sensory action can entirely replace text as condition. We integrate an additional text-encoder head to the MoE feature encoding branches to incorporate simple text phrases, *e.g.* *cut potato*. The encoded text features are aggregated with other multi-sensory action features in the same manner as described in Sec. 2.1. We use the pretrained OpenClip [28] text encoder to encode text in all baselines and our model.

As depicted in the bottom half of Figure. 7, when multiple objects (pan and plate) appear in context image and when the action trajectory can be applied to both objects, the network is uncertain about which object to apply the action. It cleans the plate instead of the pan. When we add text description *clean pan* as an extra piece of information, ambiguity is removed and accurate video can be generated. We also observe that when the context frame is not ambiguous, multi-sensory action provides enough information to generate accurate video trajectories. Adding additional text feature induces a temporal smoothing effect generating similar images across frames.

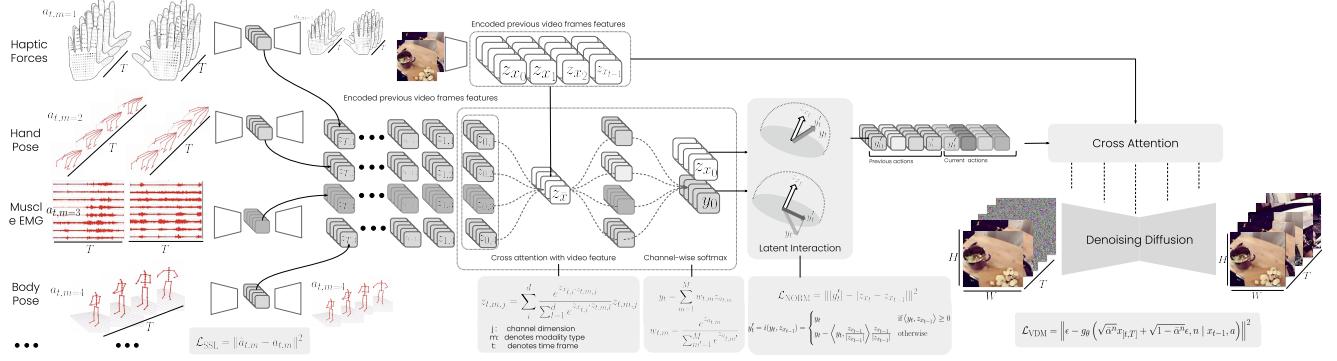


Figure 11. Additional pipeline figure.

module	parameter count	float16 in MB	float32 in MB
policy network (to be deployed on edge devices)	120690484	241MB	482 MB

Table 6. Parameter Count for the policy network model used in Downstream application section.

Hardware Type	NVIDIA Jetson Nano	Jetson Xavier	Jetson Orin NX	Jetson AGX Orin	RTX 4090	H100
Throughput (FPS)	166 ~ 111	415 ~ 290	1,725 ~ 1,293	2,555 ~ 1,916	26,528 ~ 19,896	315,141
Latency (ms)	6.6 ~ 9.2	2.4 ~ 3.44	0.57 ~ 0.77	0.39 ~ 0.52	0.037 ~ 0.050	0.00317
Energy Cost(J)	0.06 ~ 0.09	0.036 ~ 0.051	0.0114 ~ 0.0154	0.0195 ~ 0.026	0.01665 ~ 0.02250	0.02219

Table 7. Table shows that the trained policy can be deployed onto Edge devices.

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
No hand pose	0.138	14.1	0.314	264.0
No hand force	0.129	14.5	0.317	256.3
No body pose	0.137	14.5	0.322	273.1
No muscle EMG	0.121	15.2	0.311	217.1
All sensory used	0.110	16.0	0.276	203.5

Table 8. Training with ablated modalities

Table 9. Testing with single modality available

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
Hand pose	0.121	14.6	0.309	210.2
Hand force	0.117	14.7	0.307	208.0
Body pose	0.123	14.6	0.310	210.5
Muscle EMG	0.132	13.9	0.312	214.8
All sensory used	0.110	16.0	0.276	203.5

### 7.8.2. Additional results on training with missing modalities

We first ablate different sensory signal input, when training our video simulator. We observe that body pose is crucial for larger motions that involve moving in space such as turning or walking. For more delicate manipulations such as cutting or peeling, hand poses and haptic forces get us most of the way. Results in Table 8 suggests that contribution of muscle EMG is minimal. A closer look into the dataset reveals that muscle EMG is highly correlated with hand force magnitude, but it provides extra information in scenarios where hands are fully engaged.

### 7.8.3. Additional results on test-time robustness

As we see from the Table. 9 that when one modality is provided, our model can still produce higher prediction accuracy

compared to text-based models or single-model models. Comparing this result with Table. 1 shows that our proposed multisensory action training strategy induces higher quality action feature compared to training with a single modality. This comparison indicates that through implicit association between different modalities, both feature alignment and information preservation is achieved. That is, the complementary information is preserved in the feature representation such that when only one action modality is provided, the model might have access to commonly co-activated feature dimensions and thus produce better result than training with single modality.

To provide a comprehensive set of ablation studies on testing with missing modalities, we show Table 10 that includes all possible pairs of modalities used during testing. The results in Table. 10 along with Table. 9 and Table. 3a makes a

comprehensive study cross all possible ablated experiments. We can from Table 10, that the model achieves better performance when different aspect of information is provided.

Table 10. Testing with paired modality available

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
Hand Pose and Hand Force	0.115	14.9	0.304	206.4
Body Pose and Muscle EMG	0.122	14.6	0.309	210.1
Hand Force and Muscle EMG	0.117	14.7	0.307	207.6
Hand Pose and Body Pose	0.113	15.0	0.297	206.2
All sensory used	0.110	16.0	0.276	203.5

#### 7.8.4. Comparison between Training and Testing with Ablated Modalities

The critical difference between the above two experiments, training with ablated modalities (Table. 8) and testing with missing modalities (Table. 3a) is the modalities used during training. The latter ablation experiment, testing with missing modalities, employs a model trained with all modalities, whereas the former is trained only on a subset of modalities. Comparing the performance decrease in Table. 8 and Table. 3a, we can see that the latter experiment, testing with missing modalities, induces very minimal drop in prediction accuracy. This comparison confirms the advantage of training on multimodal action signals. We believe that this test-time robustness is induced by channel-wise attention and channel-wise softmax module, as these design choices allows the model to leverage substitutional information in the given modalities to bridge different modalities to allow for robustness during inference.

#### 7.8.5. History Horizon.

Finally, we study the effect of history horizon length on our model with comparison to text-conditioned simulation. We follow prior works [71] to compare context frame length  $h(x)=4$  and  $h(x)=1$ , shown in Table 11. We can see that increased history frame length reduces prediction error for all methods. Additionally, our proposed multisensory action condition is temporally fine-grained, which allows the cross attention between action and observation history  $h(x, a) = 4$  to help further increase simulation accuracy.

#### 7.8.6. Cross Subject Testing

We report the cross subject testing, where we use three other different subjects for testing and training with the rest using the ActionSense dataset, result can be found in Table. 12.

#### 7.8.7. Examples of fine-grained control

We can see from Fig. 12 where hand force together with hand pose helps accurately controls the timing of the hand grabbing the pan.

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
Unisim $h(x) = 1$	0.177	12.7	0.408	674.9
Unisim $h(x) = 4$	0.118	14.6	0.321	275.9
Ours $h(x) = 1$	0.142	12.9	0.362	535.1
Ours $h(x, a) = 1$	0.138	12.7	0.356	529.1
Ours $h(x) = 4$	0.114	15.4	0.306	256.3
Ours $h(x, a_h) = 4$	<b>0.110</b>	<b>16.0</b>	<b>0.276</b>	<b>203.5</b>

Table 11. Effects of history horizon length

Table 12. Cross Subject Testing

Method	MSE ↓	PSNR ↑	LPIPS ↓	FVD ↓
subject 2	0.115	15.8	0.301	206.7
subject 4	0.112	16.0	0.282	204.6
subject 5	0.110	16.0	0.276	203.5

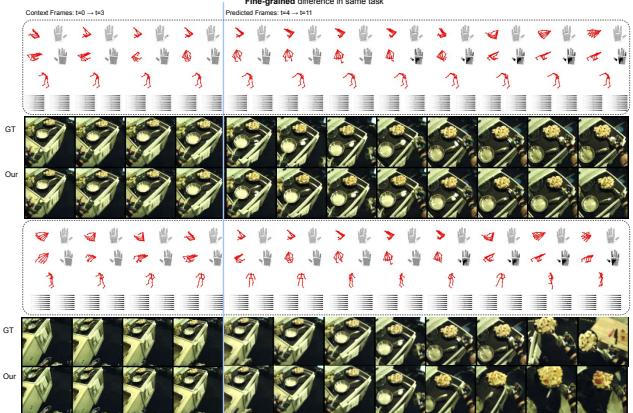


Figure 12. Temporally fine-grained control

#### 7.8.8. Additional Experiment on Generalization to OOD data through Finetuning

We present a second experiment to demonstrate that our method can handle specific out-of-distribution (OOD) scenarios through fine-tuning. For this experiment, we modified the original ActionSense dataset to create OOD data. Using LangSAM, we extracted segmentation masks for "potatoes" and recolored them to appear as "tomatoes." Since the video model had not encountered red vegetables or fruits during training, we fine-tuned our pretrained model on a small dataset of approximately 600 frames (30 seconds) and evaluated it on the test split of this "tomato" data. The data creation procedure is shown in Fig. 13 and results on this experiment can also be found in 14. The results show that the model achieves reasonable performance after fine-tuning. While we acknowledge that robust in-the-wild generalization requires training on larger-scale datasets with diverse domain coverage, this experiment illustrates a practical use

case for addressing OOD data. Specifically, it demonstrates that by collecting a small, specialized dataset, our pretrained model can be effectively fine-tuned to adapt to new domains.

### 7.8.9. Additional discussion and results on downstream application

Sample results visualization can be found in Fig. 15. We also observe from the figure that the policy optimized by our proposed approach can be different from the ground truth action trajectory, yet the simulated visual observations still closely resemble the ground truth state observations. We believe that the softmax aggregation learns to pick out information deemed useful by the simulator, leaving freedom in irrelevant dimensions in the action space.

### 7.8.10. Downstream Application2: Multi-Sensory Action Planning

Another potential downstream application is long-term planning. Inspired by [17], we use text to describe high-level goals to generate a set of executable next-step actions. Our video model takes an image observation and the generated actions to simulate future image sequences, which can be further evaluated for next-step execution planning. As shown in Fig. 15, our model can potentially be used for low-level actuation planning through iterative action roll outs. We adapt diffusion policy (DP) [9] to take in both first frame image feature  $x_0$  and high-level goal  $\gamma$  described by a text feature  $f_\gamma$  as the context conditions to generate multi-sensory trajectories of fine-grained actions  $a_{[1,T]} = p(x_0, f_\gamma)$ . The action steps are then fed into our action-conditioned video generative model  $g(\cdot)$  to generate sequences of future video frames  $\hat{x}_{[1,t]} = g(x_0, a_{[1,t]})$ . To decide whether the subtask  $\tau$  has been achieved, we use a vision language model  $f_v(\cdot)$  as a heuristic function [48], which can be prompted with the end state of the current roll out  $\hat{x}_t$  to evaluate whether subgoal  $\tau$  has been achieved. If more steps are needed, we can further iterate the process  $a_{[t,it]} = p(\hat{x}_t, \gamma)$ ,  $x_{[t,it]} = g(\hat{x}_t, a_{[t,it]})$ . A sample result from text-promted diffusion policy is shown in Figure. 15. We observe long iterations result in accumulative error, as shown in the bottom row of Fig. 19 in Appendix Sec. 7.8). A larger-scale dataset can further boost performance for this task. This downstream application hints at fully automated low-level motion planning and dexterous manipulation, enabling realization of household robots.

### 7.8.11. Higher Resolution Results

We include some sample results for higher resolution model of video size  $128 \times 128 \times 12$  and  $192 \times 192 \times 12$ , matching the video resolution of existing generative video simulation paper, such as Unisim [71]. The results are shown in Fig. 18

### 7.8.12. Additioanl qualitative results on other dataset

To show that our proposed method is generic is not designed for the ActionSense [13] dataset, we conducted an experi-

ment by directly applying our proposed approach on another dataset, H2O dataset [34]. H2O [34] dataset is a unimodal action-video dataset that includes paired video and hand pose sequences. We would love to expand our our training on larger and more diverse dataset, However, to the best of our knowledge, ActionSense [13] is the only dataset that includes paired multisensory action signal monitoring sequences alongside video sequences. We show experiment on H2O [34] in Figure 16. We provide additional sample test on the holoAssist dataset 17, which is also a hand-pose video dataset in Fig. 16. These results demonstrate that our system is generic, not dataset specific, and can achieve reasonable performance. These results indicate that our model is capable of training and testing on unimodal action datasets, highlighting its generalizability beyond the ActionSense dataset. This demonstrates that our method is not specifically tailored to ActionSense and can adapt to various scenarios. We believe our proposed method offers a generalizable framework that can serve as a reference and can be applied more broadly as additional datasets of this nature become available.

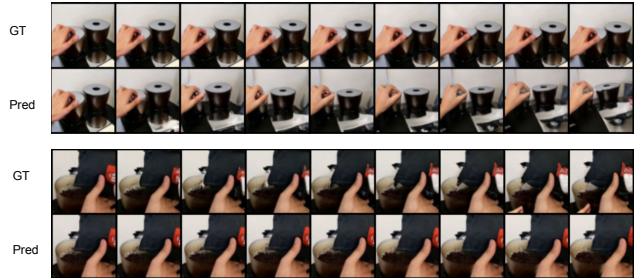


Figure 17. Test on HoloAssist dataset

## 7.9. Additional Qualitative Results

Additional Qualitative Results are shown in Fig. 19, Fig. 20, and Fig. 21. Fig. 19 and Fig. 20 show additional qualitative results of context frames and predicted video frames from our proposed multisensory action signals. Fig. 21 shows demonstrations of failure cases, policy optimization, and long-trajectory planning. We show one most recent context frame and the eight prediction frames. Fig. 21 shows results paired in two rows, where the top row shows ground truth trajectory the bottom row shows predicted trajectory.

### 7.9.1. Failure Cases

We show the failure cases on the top right section. Common failure cases include false hallucination of environment with large motion. Failure to identify object with similar apperance to background. The wooden chopboard gradually disappear into the wooden table background and fails to pick it up in simulation. Failure in identify object to act on (also hallucates pan handle on plate and cleaning the plate). The last five rows in Fig. 19 show additional results on down stream tasks of policy planning, shown in the middle rows, and long-trajectory simulation, show in the bottom row.

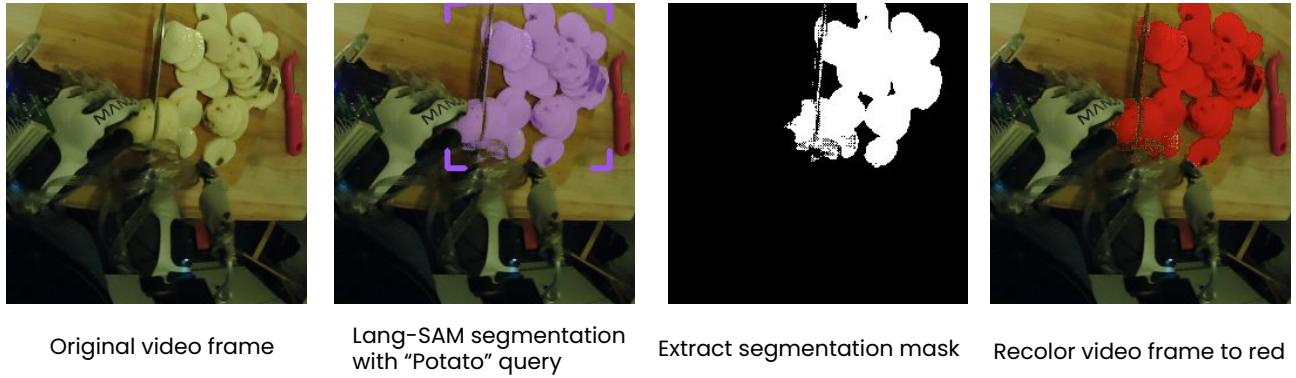


Figure 13. Experimental set up on OOD testing.

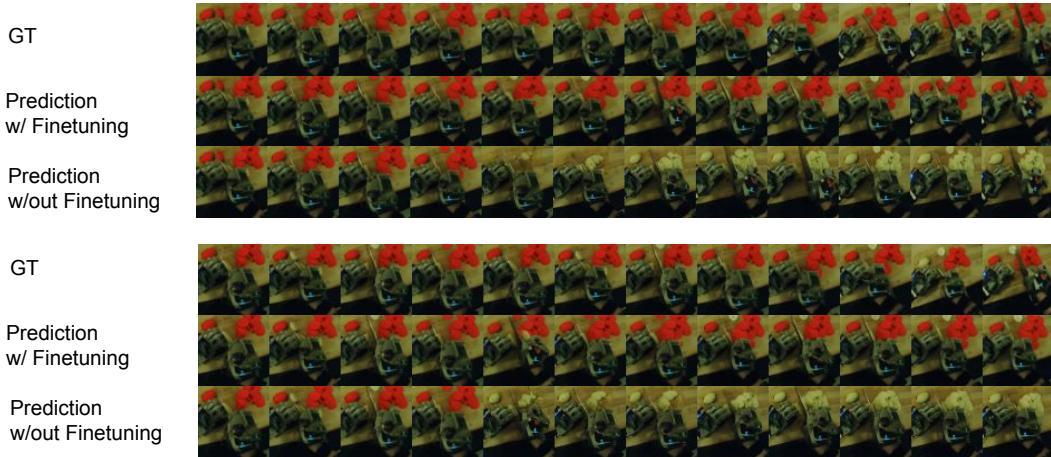


Figure 14. Experimental results on OOD testing.

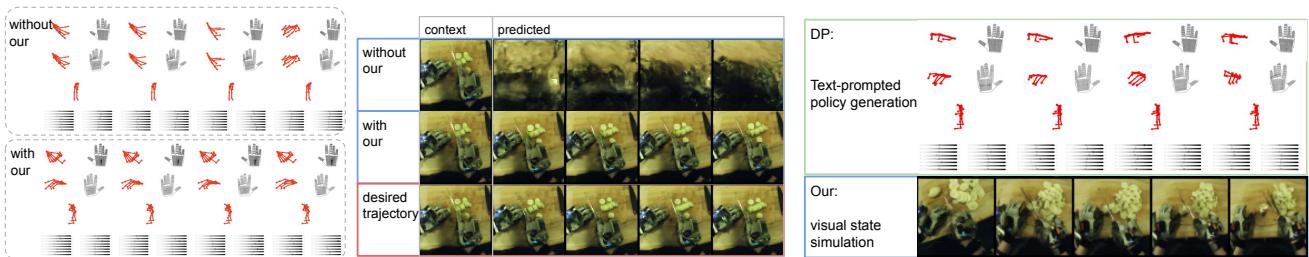


Figure 15. **Left:** Results on goal-conditioned policy optimization. **Right:** Results on long-term task planning.

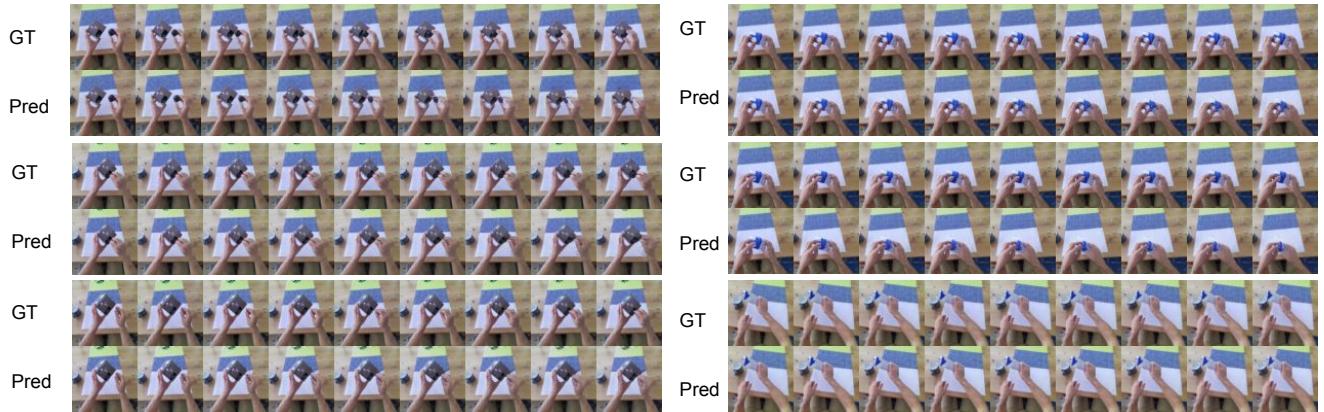


Figure 16. Test on H2O dataset

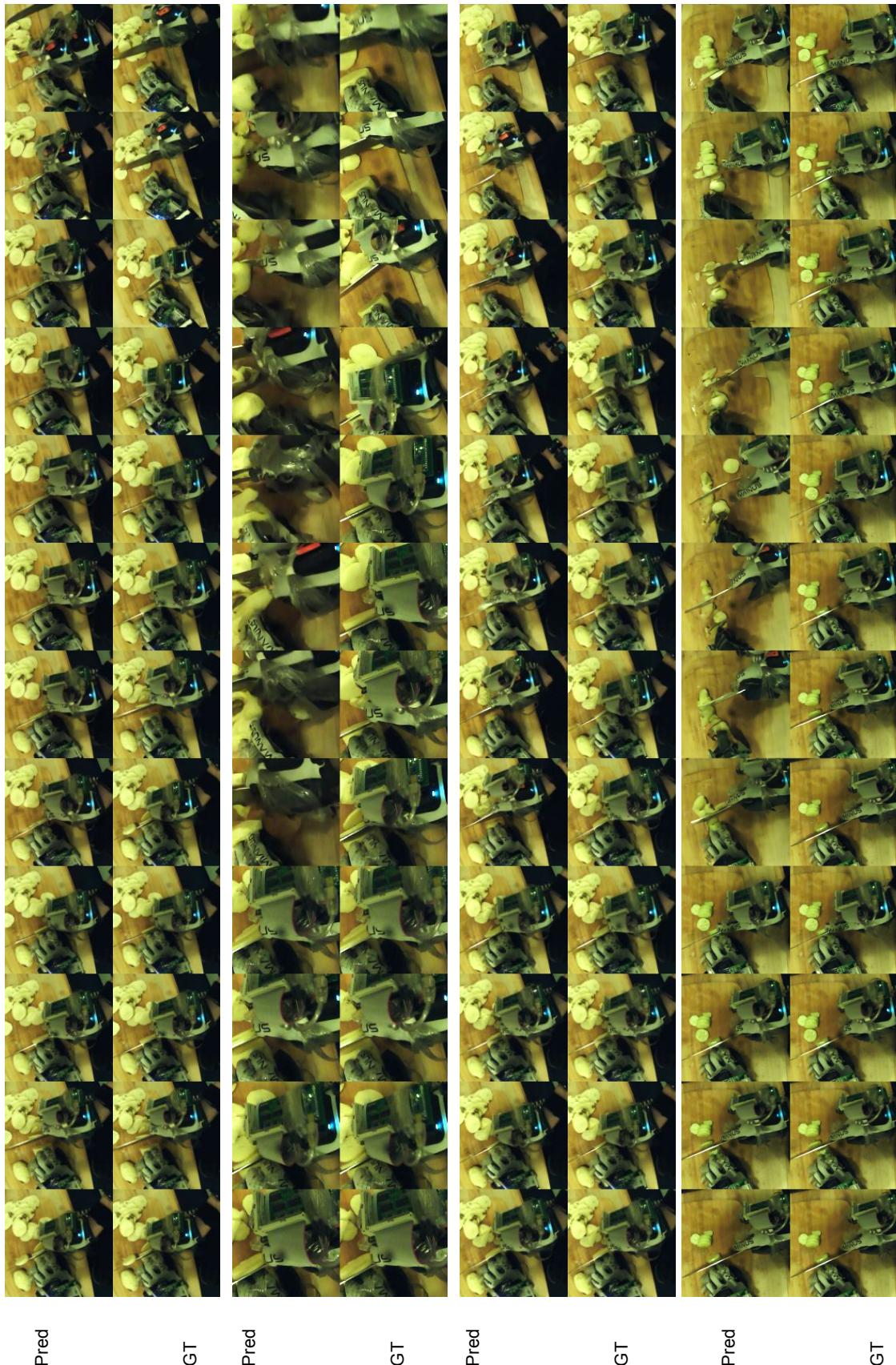


Figure 18. The left three are of resolution 128 and the last one is of resolution 192

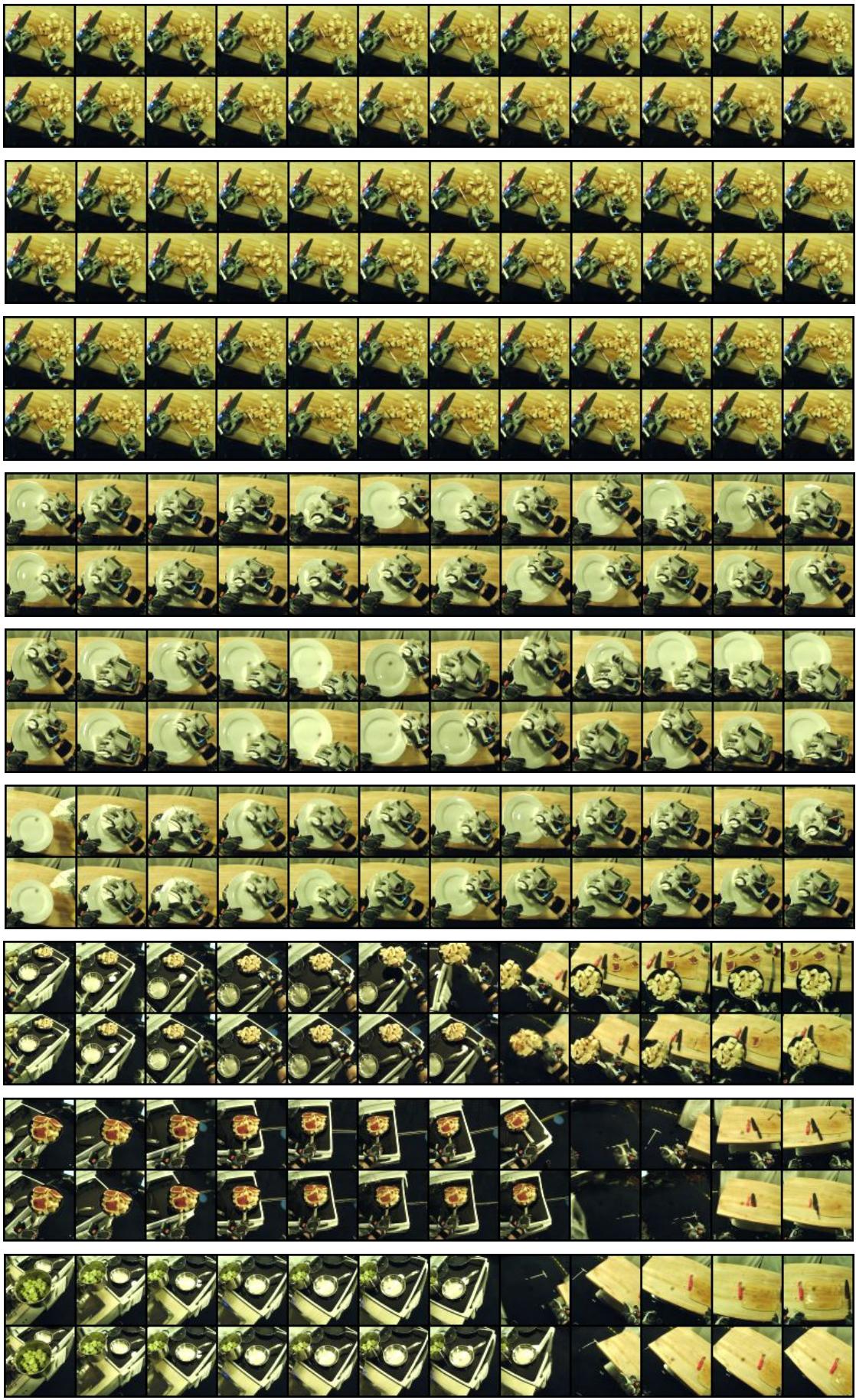


Figure 19. Additional qualitative results

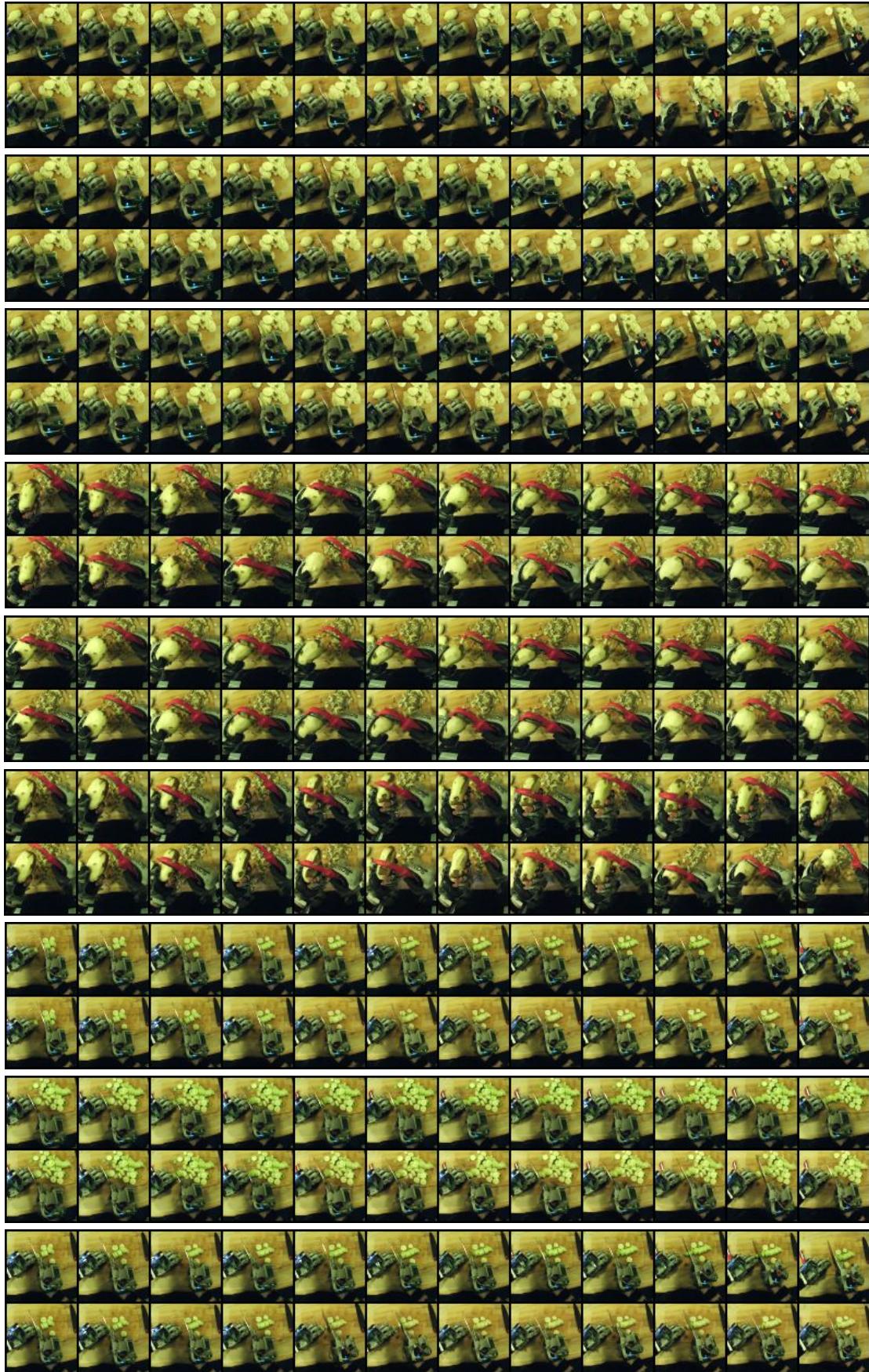


Figure 20. Additional qualitative results

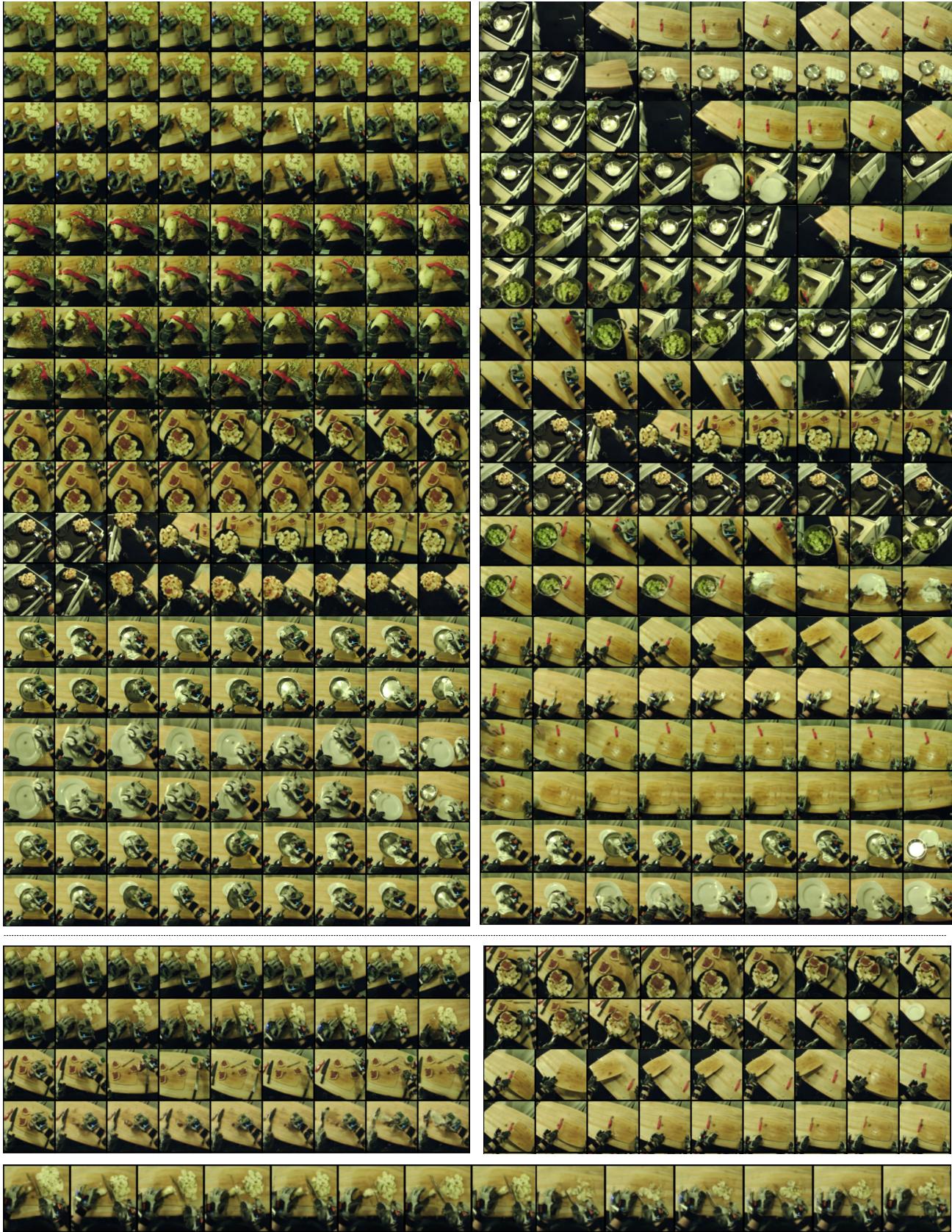


Figure 21. **Top left:** Additional qualitative results. **Top right:** Failuare cases. **Middle left and right:** Additional results on policy optimization. **Bottom:** long-trajectory policy planning.