

Figure S10. UCF101 Distribution.

model is effectively encoding the frames independently using SDv2. Therefore, we set  $\alpha = 1.0$  for all layers in order to use the same backbone to extract SDv2 features.

### A.2.2 V-JEPA and VideoMAEv2 Extraction

We extract VideoMAEv2 features using the official repository and pretrained encoder provided by the authors. Specifically, we compute spatiotemporal representations for each video segment and then average across the spatial dimension to obtain a single feature vector.

### A.3. Action Localization

We can also analyze which parts of a video most contribute to the prediction of a specific action. For a given video, we can crop it into fixed patches and provide as input to the model the same patch over time. Then we can evaluate the prediction of an action over all patches and produce a heatmap like in Figure S11, which shows which parts of a video are most predicted as a prey being eaten, and we can localize the action in space. Resolution is not an issue for our model since we average pool the video features.

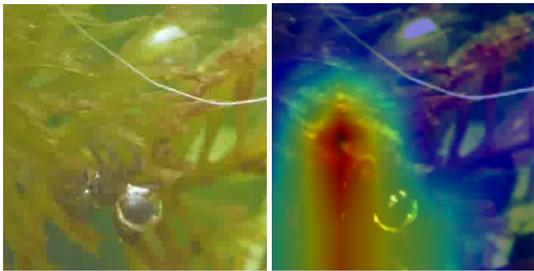


Figure S11. Action localization

## A.4. Cross-Species Generalization

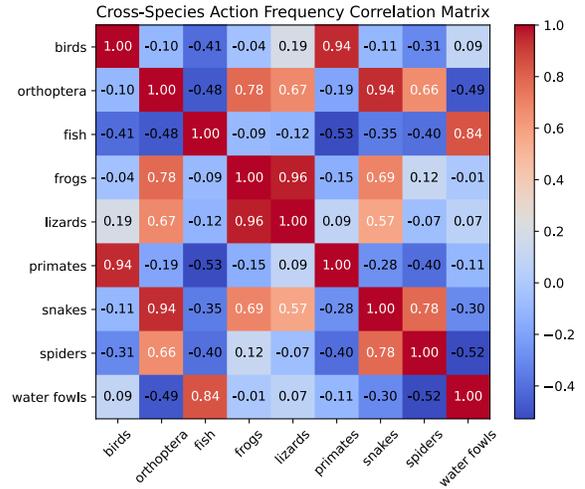


Figure S12. Action frequencies correlations across species.

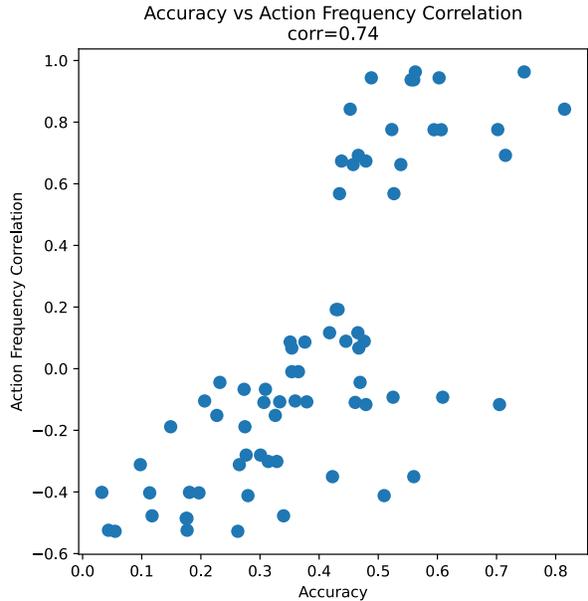


Figure S13. Accuracy vs action frequency correlation across species.

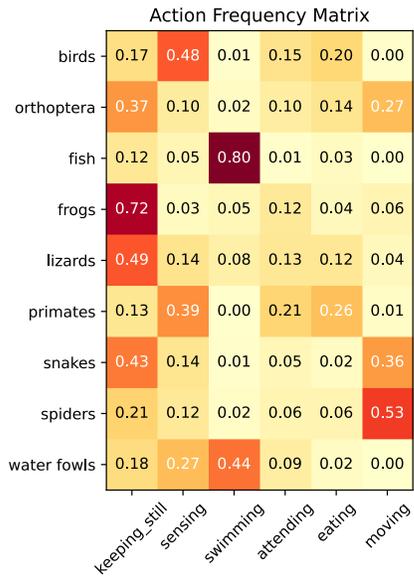


Figure S14. Action frequencies across species.

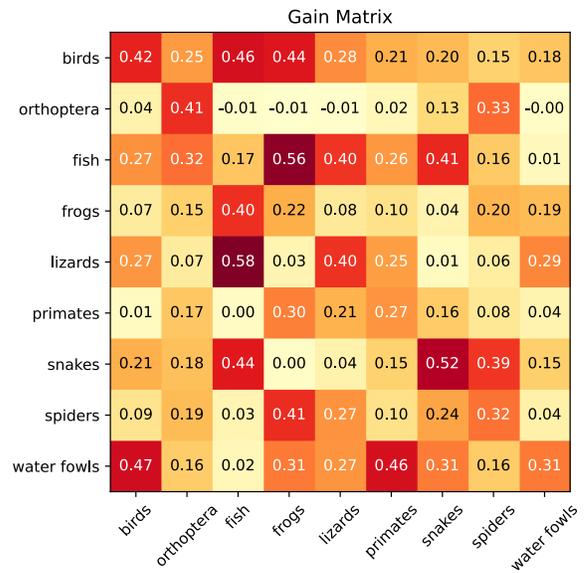


Figure S16. ActionDiff gains over predicting the most frequent class from the training species in the test species.

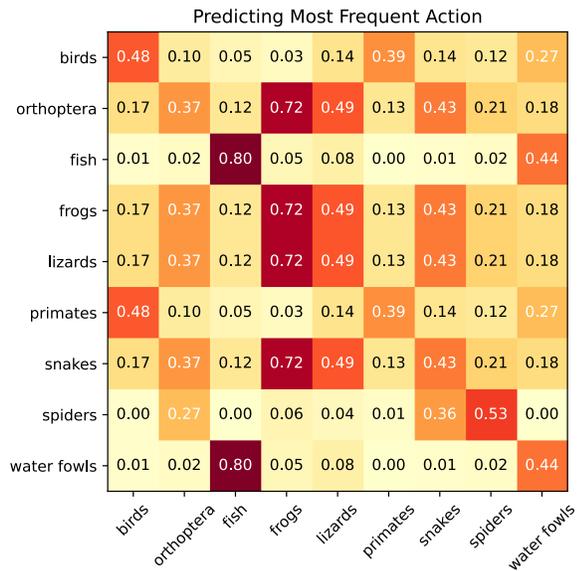


Figure S15. Results of predicting the most frequent action from the training species in the test species.

### A.5. Choice of UNet Layer

Layer	Animal Kingdom		Charades-Ego		UCF101-HMDB51	
	Full Dataset (mAP)	Unseen Species (acc)	$1^{st} \rightarrow 1^{st}$ mAP	$3^{rd} \rightarrow 1^{st}$ mAP	U $\rightarrow$ H Acc	H $\rightarrow$ U Acc
<i>ActionDiff (ours)</i>						
Layer 3	55.16	35.51	16.98	15.1	17.17	17.57
Layer 6	71.12	40.90	21.07	16.6	33.92	39.67
<b>Layer 9</b>	<b>80.79</b>	<b>51.49</b>	<b>36.5</b>	<b>30.2</b>	<b>75.6</b>	<b>81.5</b>

Table S5. **Effect of UNet layer on all datasets.** Results improve with deeper decoder layers across all metrics. **Bold** denotes best result.