

Layered Controllable Video Generation

Supplementary Material

This supplementary material is structured as follows. In Section A, we present our architectural details of our models, and the implementation details of our training/testing procedures. In Section B, we show visual results of the frames generated by our model in under different conditions. In Section C, we discuss some potential social implications of our method. We also include a [readme.html](#) to summarize our work and to show some video results.

A. Implementation Details

Architecture details. Our model consists of four main modules: the *Mask Network* \mathcal{M} , the *Encoder* \mathcal{E} , the learnable discrete code book \mathcal{Z} , and the *Decoder* \mathcal{D} . We summarize the details in Table 3. In total, our model has 47.8 M parameters, and takes about 191.2 MB to store on disk.

Training Details. In all our experiments, we use the Adam Optimizer with a fixed learning rate of 4.5×10^{-6} , which we found empirically.

- *Initial training for mask-based generation:* For this stage, we train all models for 10,000 iterations. For experiments on the *BAIR* dataset, we initialize the ratio between $\lambda_{\text{bg}} : \lambda_{\text{fg}}$ to be 80 : 1, and reduce it by a factor of 2 for every 1,000 iterations until this ratio becomes 5 : 1. For experiments on the *Tennis* dataset, this ratio is initialized to 10 : 1 and reduced by half after 5,000 iterations.
 - *Fine-tuning for controllability:* In this stage of training, we drop all losses related to the mask regularization as their gradient becomes zero (see Section 3.3), except for \mathcal{L}_{bg} . On the *BAIR* dataset, λ_{bg} is set to 5, and on the *Tennis* dataset, it's set to 0.5.
 - *Finding the “ground-truth” control signal:* In Equation. 5, we use the fully differentiable transformation operation $\mathcal{T}(\theta^t)$ to approximate the mask of the ground truth next frame using the mask of the current frame. For each frame, We optimize the affine transformation matrix that's been used to transform \mathbf{m}^t into \mathbf{m}_c^t , the optimization is performed by the ADAM optimizer for 1,000 iterations with the learning rate of 0.1.
- On the *BAIR* dataset, our model requires 80GB GPU memory and 40 hours to train on 4× RTX-6000 (24GB V-RAM per GPU) GPUs, and on the *Tennis* daaset, the model takes 22GB GPU memory and 23 hours to train on a single RTX-3090 (24GB V-RAM).

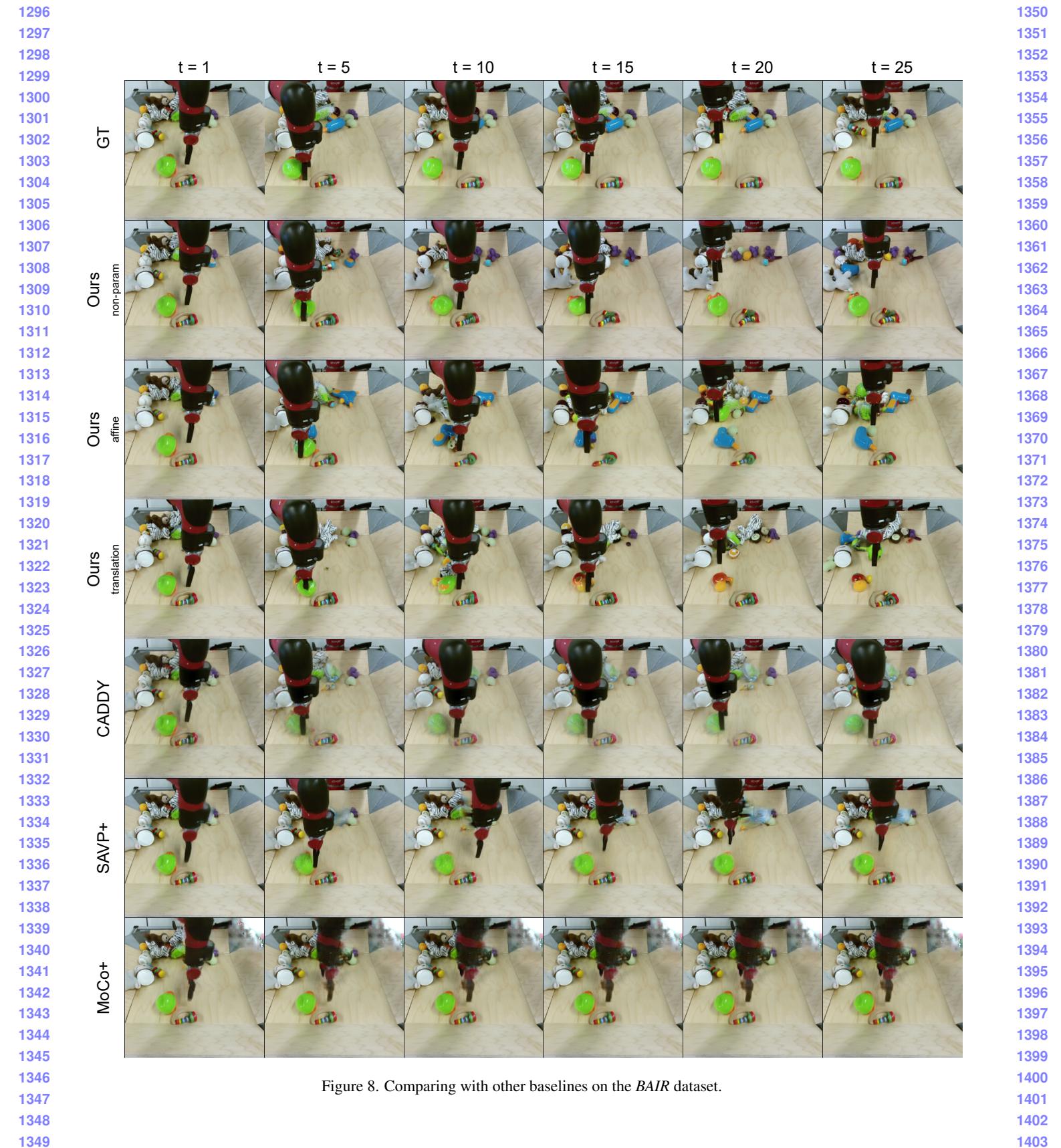
1080	Mask Network \mathcal{M}	1134
1081	$x \in \mathbb{R}^{H \times W \times 3}$	1135
1082	$\text{Conv2D} \rightarrow \mathbb{R}^{H \times W \times 64}$	1136
1083	$2 \times \{\text{Residual Blocks} + \text{Downsample}\} \rightarrow \mathbb{R}^{h \times w \times 256}$	1137
1084	$9 \times \text{Residual Blocks} \rightarrow \mathbb{R}^{h \times w \times 256}$	1138
1085	$2 \times \{\text{Upsample} + \text{Residual Blocks}\} \rightarrow \mathbb{R}^{h \times w \times 64}$	1139
1086	$\text{Conv2D} + \text{Sigmoid} \rightarrow \mathbb{R}^{H \times W \times 1}$	1140
1087		1141
1088		1142
1089		1143
1090		1144
1091		1145
1092		1146
1093		1147
1094		1148
1095		1149
1096		1150
1097		1151
1098		1152
1099		1153
1100		1154
1101		1155
1102		1156
1103		1157
1104		1158
1105		1159
1106		1160
1107		1161
1108		1162
1109		1163
1110		1164
1111		1165
1112		1166
1113		1167
1114		1168
1115		1169
1116		1170
1117		1171
1118		1172
1119		1173
1120		1174
1121		1175
1122		1176
1123		1177
1124		1178
1125		1179
1126		1180
1127		1181
1128		1182
1129		1183
1130		1184
1131		1185
1132		1186
1133		1187

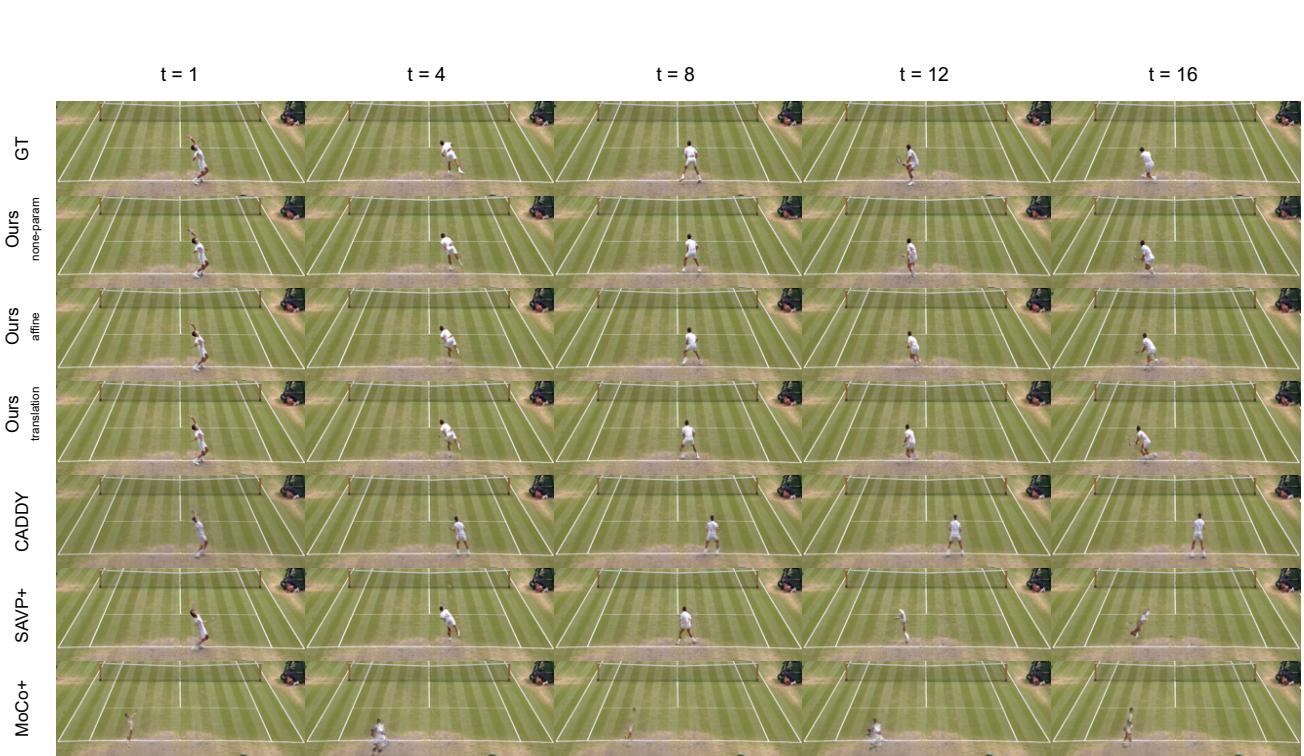
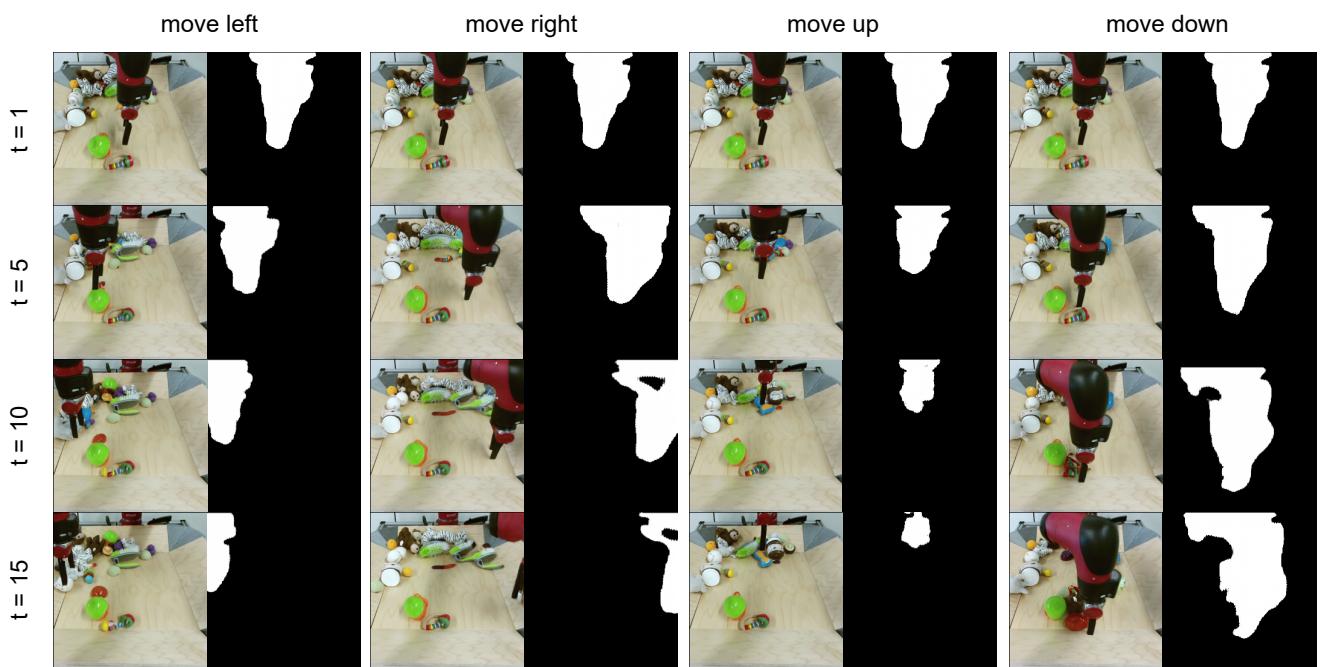
Table 3. Architecture details of the main modules of our model. The high-level design follows the architecture presented in [14]. K denotes the number of entries in the code book, n_z is the dimension of each entry, $h, w = (H, W)/4$ and $h', w' = (H, W)/16$. For the discriminator \mathcal{C} , we use the vanilla PatchGAN discriminator as described in [25].

B. Additional qualitative results

Comparing with other baselines. In Figure 8 and Figure 9, we show the comparison of the three variants of our model against other baselines that we mentioned in the main text, on the *BAIR* and *Tennis* datasets respectively. Noticeably on both datsets, our model achieves the most precise control over the generated video (see how the generated motion corresponds to the ground truth sequences), meanwhile generating one of the best frame-quality videos. For example,

1188	CADDY fails to control the robot arm, SAVP+ controls it	1242
1189	up to some degree, but as shown in frames $t=10$ and $t=15$,	1243
1190	their control is not as precise, whereas all of our variants	1244
1191	provide highly accurate control.	1245
1192		1246
1193	Generating videos conditioned on user inputs. In Fig-	1247
1194	ure 10 and Figure 11, we show how the video generated	1248
1195	by our model responds to the user inputs, on the <i>BAIR</i>	1249
1196	and <i>Tennis</i> datasets respectively. Here, for each sam-	1250
1197	ple, we keep applying the same control signal (moving	1251
1198	left/right/up/down) to the model so that it can generate	1252
1199	video sequences with single consistent motions.	1253
1200		1254
1201	Action mimicking. In Figure 12, we show that our method	1255
1202	can be used to extract the motion from a driving sequence	1256
1203	and then apply onto different appearances (staring frames).	1257
1204		1258
1205	Video results. For more video results, please refer to the at-	1259
1206	tached readme.html , where we demonstrate all above men-	1260
1207	tioned results in video format, as well as other applications	1261
1208	of our method, such as generating videos of multi players,	1262
1209	real-time user controllable generation demo, and animating	1263
1210	a single frame with different motions.	1264
1211		1265
1212	<h2>C. Broader Impact</h2>	1266
1213	Our work enables generating videos conditioned on the	1267
1214	on-the-fly user inputs at the frame-level.	1268
1215	Since it is also an unsupervised method, unlike other	1269
1216	works on conditioned video generation, it does not require	1270
1217	any kind of data annotations while training. Thus our model	1271
1218	can be applied to a large variety of raw video datasets. In	1272
1219	other words, our method has the potential to widen the ap-	1273
1220	plication area of conditional video generation, which may	1274
1221	increase the impact this line of research has.	1275
1222	Like other generative models, our method shares similar	1276
1223	benefits and risks, which have been discussed extensively	1277
1224	in [7, 50].	1278
1225	We acknowledge that our method can be used to generate	1279
1226	videos that involve humans, as shown in our Tennis exam-	1280
1227	ple, that may have potential for misuse. While it is unlikely	1281
1228	that our method can be used to create videos of completely	1282
1229	novel scenes and objects, such as in the case of deep fakes,	1283
1230	it can be used to generate a video of action combinations	1284
1231	that did not happen in reality. At the same time, we see the	1285
1232	potential benefits of having controllable generative models,	1286
1233	for example, allowing simulations without the need of an	1287
1234	explicit model, or to be used as a data augmentation tool	1288
1235	that provides both frame data and action annotations. As	1289
1236	concluding remark, we would like to emphasize the impor-	1290
1237	tance of research on detecting videos and images generated	1291
1238	from deep models [20, 59] to further mitigate any potential	1292
1239	misuse.	1293
1240		1294
1241		1295

Figure 8. Comparing with other baselines on the *BAIR* dataset.

Figure 9. Comparing with other baselines on the *Tennis* dataset.Figure 10. Generating videos conditioned on user inputs, *BAIR* dataset.

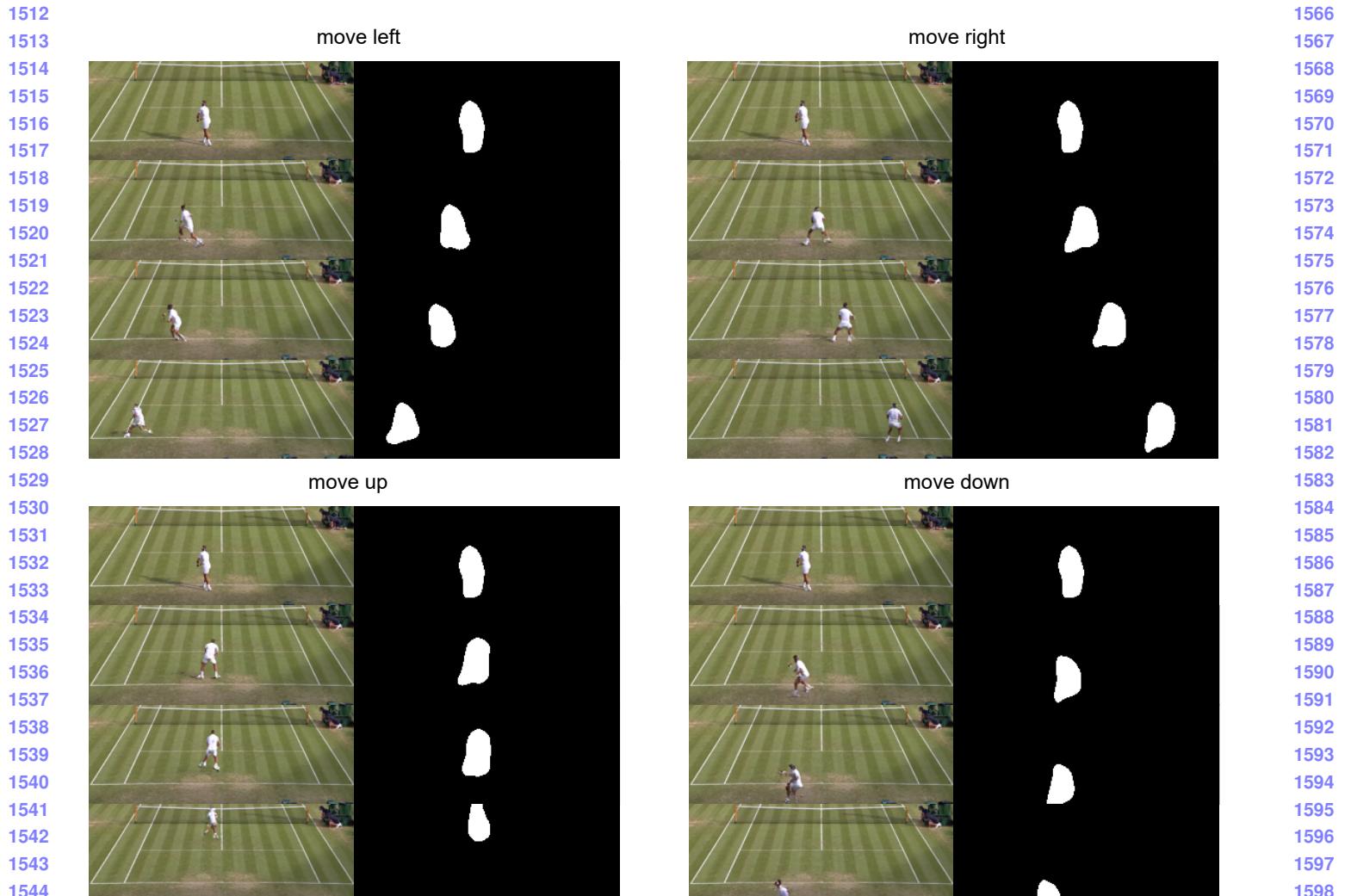
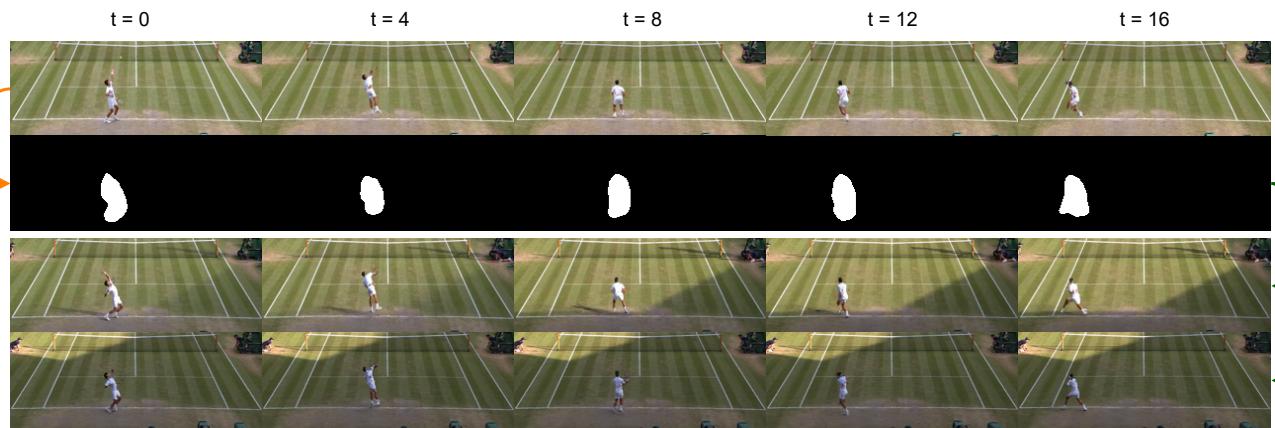
Figure 11. Generating videos conditioned on user inputs, *Tennis* dataset.

Figure 12. “Mimicking”. Our model can be used to extract the motion from a driving sequence and then apply onto different appearances (staring frames). The mask sequence is extracted from the driving video (first row), then applied onto two other staring frames (third, fourth row).