

## 310 Appendix

### 311 A Implementation Details

312 In this section, we describe the implementation details of our algorithm for training on the training  
313 view and test time training in view generalization settings on the DMControl [30], xArm [7], and  
314 Adroit [20] environments. We utilize the official implementation of TD-MPC [11] and MoDem [8]  
315 which are available at [github.com/nicklashansen/tdmpc](https://github.com/nicklashansen/tdmpc) and [github.com/facebookresearch/modem](https://github.com/facebookresearch/modem) as  
316 the model-based reinforcement learning codebase. During training time, we use the default hyperpa-  
317 rameters in official implementation of TD-MPC and MoDem. We present relevant hyperparameters  
318 during both training and test time in Table 10 and Table 11. One seed of our experiments could be  
319 run on a single 3090 GPU with fewer than 2GB and it takes  $\sim 1$  hours for test-time training.

320 **Training time setup.** We train visual model-based policies with TD-MPC on DMControl and xArm  
321 environments, and MoDem on Adroit environments, We employ identical network architecture and  
322 hyperparameters as original TD-MPC and MoDem during training time.

323 The network architecture of the encoder in original TD-MPC is composed of a stack of 4 convolutional  
324 layers, each with 32 filters, no padding, stride of 2,  $7 \times 7$  kernels for the first one,  $5 \times 5$  kernels for  
325 the second one and  $3 \times 3$  kernels for all others, yielding a final feature map of dimension  $3 \times 3 \times 32$   
326 (inputs whose framestack is 3 have dimension  $84 \times 84 \times 9$ ). After the convolutional layers, a fully  
327 connected layer with an input size of 288 performs a linear transformation on the input and generates  
328 a 50-dimensional vector as the final output.

329 The network architecture of the encoder in original Modem is composed of a stack of 6 convolutional  
330 layers, each with 32 filters, no padding, stride of 2,  $7 \times 7$  kernels for the first one,  $5 \times 5$  kernels for  
331 the second one and  $3 \times 3$  kernels for all others, yielding a final feature map of dimension  $2 \times 2 \times 32$   
332 (inputs whose framestack is 2 have dimension  $224 \times 224 \times 6$ ). After the convolutional layers, a fully  
333 connected layer with an input size of 128 performs a linear transformation on the input and generates  
334 a 50-dimensional vector as the final output.

335 **Test time training setup.** During test time, we train spatial adaptive encoder (SAE) to adapt to view  
336 changes. We insert STN blocks before and after the first convolutional layer of the original encoders  
337 in TD-MPC and MoDem. The original encoders are augmented by inserting STN blocks, resulting in  
338 the formation of SAE. Particularly, for the STN block inserted before the first convolutional layer, the  
339 input is a single frame. This means that when the frame stack size is  $N$ ,  $N$  individual frames are fed  
340 into this STN block. This is done to apply different transformations to different frames in cases of  
341 moving and shaking view.

342 To update the SAE, we collect online data using a buffer with a size of 256. For each update, we  
343 randomly sample 32 (observation, action, next\_observation) tuples from the buffer as a batch. The  
344 optimization objective is to minimize the loss in predicting the dynamics of the latent states, as  
345 defined in Equation 2.

346 During testing on each task, we run 20 consecutive episodes, although typically only a few or even  
347 less than one episode is needed for the test-time training to converge. To make efficient use of the  
348 data collected with minimal interactions, we employ a multi-update strategy. After each interaction  
349 with the environment, the SAE is updated 32 times.

350 The following is the network architecture of the first STN block inserted into the encoder of TD-MPC.

```
351 STN_Block_0_TDMPC(  
352   (localization): Sequential(  
353     # By default, each image consists of three channels. Each frame in the  
354     ↪ observation is treated as an independent input to the STN.  
355     (0): Conv2d(in_channels=3, out_channels=8, kernel_size=7, stride=1)  
356     (1): MaxPool2d(kernel_size=4, stride=4, padding=0)  
357     (2): ReLU()  
358     (3): Conv2d(in_channels=8, out_channels=10, kernel_size=5, stride=1)
```

```

360         (4): MaxPool2d(kernel_size=4, stride=4, padding=0)
361         (5): ReLU()
362     )
363     (fc_loc): Sequential(
364         (0): Linear(in_dim=90, out_dim=32)
365         (1): ReLU()
366         (2): Linear(in_dim=32, out_dim=6)
367     )
368 )

```

369 The following is the network architecture of the second STN block inserted into the encoder of  
370 TD-MPC.

```

371 STN_Block_1_TDMPC(
372     (localization): Sequential(
373         (0): Conv2d(in_channels=32, out_channels=8, kernel_size=7, stride=1)
374         (1): MaxPool2d(kernel_size=3, stride=3, padding=0)
375         (2): ReLU()
376         (3): Conv2d(in_channels=8, out_channels=10, kernel_size=5, stride=1)
377         (4): MaxPool2d(kernel_size=2, stride=2, padding=0)
378         (5): ReLU()
379     )
380     (fc_loc): Sequential(
381         (0): Linear(in_dim=90, out_dim=32)
382         (1): ReLU()
383         (2): Linear(in_dim=32, out_dim=6)
384     )
385 )

```

386 The following is the network architecture of the first STN block inserted into the encoder of MoDem.

```

387 STN_Block_0_MoDem(
388     (localization): Sequential(
389         # By default, each image consists of three channels. Each frame in the
390         ↪ observation is treated as an independent input to the STN.
391         (0): Conv2d(in_channels=3, out_channels=5, kernel_size=7, stride=2)
392         (1): MaxPool2d(kernel_size=4, stride=4, padding=0)
393         (2): ReLU()
394         (3): Conv2d(in_channels=5, out_channels=10, kernel_size=5, stride=2)
395         (4): MaxPool2d(kernel_size=4, stride=4, padding=0)
396         (5): ReLU()
397     )
398     (fc_loc): Sequential(
399         (0): Linear(in_dim=90, out_dim=32)
400         (1): ReLU()
401         (2): Linear(in_dim=32, out_dim=6)
402     )
403 )

```

404 The following is the network architecture of the second STN block inserted into the encoder of  
405 MoDem.

```

406 STN_Block_1_MoDem(
407     (localization): Sequential(
408         (0): Conv2d(in_channels=32, out_channels=8, kernel_size=7, stride=2)
409         (1): MaxPool2d(kernel_size=3, stride=3, padding=0)
410         (2): ReLU()
411         (3): Conv2d(in_channels=8, out_channels=10, kernel_size=5, stride=2)
412         (4): MaxPool2d(kernel_size=2, stride=2, padding=0)
413         (5): ReLU()
414     )
415     (fc_loc): Sequential(
416         (0): Linear(in_dim=90, out_dim=32)
417         (1): ReLU()
418         (2): Linear(in_dim=32, out_dim=6)
419     )
420 )

```

Table 10: **Hyperparameters for training time.**

Hyperparameter	Value
Discount factor	0.99
Image size	84 × 84 (TD-MPC) 224 × 224 (MoDem)
Frame stack	3 (TD-MPC) 2 (MoDem)
Action repeat	1 (xArm) 2 (Adroit, Finger, and Walker in DMControl) 4 (otherwise)
Data augmentation	±4 pixel image shifts (TD-MPC) ±10 pixel image shifts (MoDem)
Seed steps	5000
Replay buffer size	Unlimited
Sampling technique	PER ( $\alpha = 0.6, \beta = 0.4$ )
Planning horizon	5
Latent dimension	50
Learning rate	1e-3 (TD-MPC) 3e-4 (MoDem)
Optimizer ( $\theta$ )	Adam ( $\beta_1 = 0.9, \beta_2 = 0.999$ )
Batch size	256
Number of demos	5 (MoDem only)

Table 11: **Hyperparameters for test time training.**

Hyperparameter	Value
Buffer size	256
Batch size	32
Multi-update times	32
Learning rate for encoder	1e-6 (xArm) 1e-7 (otherwise)
Learning rate for STN blocks	1e-5

## 421 B Environment Details

422 We categorize the view generalization problem into four distinct settings: novel view, moving view,  
 423 shaking view, and novel FOV. In this section, we provide descriptions of the implementation details  
 424 for each setting. The detailed camera settings can be referred to in the code of the environments that  
 425 we are committed to releasing or in the visualization available on our website [movie-rl.github.io](http://movie-rl.github.io).

426 **Novel view.** In this setting, for locomotion tasks (cheetah-run, walker-stand, walker-walk, and  
 427 walker-run), the camera always faces the moving agent, while for other tasks, the camera always  
 428 faces a fixed point in the environment. Therefore, as we change the camera position, the camera  
 429 orientation also changes accordingly.

430 **Moving view.** Similar to the previous setting, the camera also always faces the moving agent or a  
 431 fixed point in the environment. The camera position varies continuously.

432 **Shaking view.** To simulate camera shake, we applied Gaussian noise to the camera position (XYZ  
 433 coordinates in meters) at each time step. For DMControl and Adroit, the mean of the distribution  
 434 is 0, the standard deviation is 0.04, and we constrain the noise within the range of -0.07 to +0.07.  
 435 For xArm, the mean of the distribution is 0, the standard deviation is 0.4, and we constrain the noise  
 436 within the range of -0.07 to +0.07.

437 **Novel FOV.** We experiment with a larger FOV. For DMControl, we modify the FOV from 45 to 53.  
 438 For xArm, we modify the FOV from 50 to 60. For Adroit, we modify the FOV from 45 to 50. We  
 439 also experiment with a smaller FOV and results are presented in Appendix E.

## 440 C Visualization of Feature Map Transformation

441 We visualize the first layer feature map of the image encoder from TD-MPC and MoVie in Figure 7.  
 442 It is observed that the feature map from MoVie on the novel view exhibits a closer resemblance to  
 443 that on the training view.

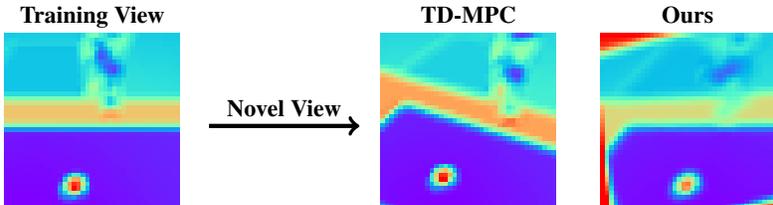


Figure 7: Visualization of the first layer feature maps of the original encoder on the training view and on the novel view, and the learned SAE on the novel view.

## 444 D Extended Description of Baselines

445 **TD-MPC.** We test the agent trained on training view without any adaptation in the view generalization  
 446 settings,

447 **DM.** This is derived from MoVie by removing STN blocks, which just adapts encoder during test  
 448 time.

449 **IDM+STN.** This is derived from MoVie by replacing the dynamics model with the inverse dynamics  
 450 model which predicts the action in between based on the latent states before and after transition. The  
 451 inverse dynamics model is finetuned together with the encoder and STN blocks during testing.

## 452 E Ablation on Different FOVs

453 In our main experiments, we consider the novel FOV as a FOV larger than the original. In Table  
 454 12, we present results for both smaller and larger FOV scenarios. Our method demonstrates the  
 455 successful handling of both cases.

Table 12: **Ablation on different FOVs.** The best method on each setting is in **bold**.

Cheetah-run	TD-MPC	DM	IDM+STN	MoVie
Small FOV	104.85±4.59	398.75±17.52	75.92±8.76	<b>530.37±12.84</b>
Large FOV	128.55±6.57	379.01±10.90	299.02±88.47	<b>532.94±19.74</b>

## 456 F Results of Original Models on Training View

457 The performance of the original agents without any adaptation under the training view is reported  
 458 in Table 13, 14, and 15 for reference. In the context of view generalization, it is evident that the  
 459 performance of agents without adaptation significantly deteriorates.

Table 13: **Training Result on DMControl.**

Task	Cheetah, run	Walker, walk	Walker, stand	Walker, run	Cup, catch	Finger, spin
Reward	658.10 $\pm$ 9.98	944.99 $\pm$ 21.71	983.53 $\pm$ 5.34	697.75 $\pm$ 11.35	980.56 $\pm$ 4.33	985.20 $\pm$ 2.25
Task	Finger, turn-easy	Finger, turn-hard	Pendulum, swingup	Reacher, easy	Reacher, hard	
Reward	756.16 $\pm$ 150.76	616.96 $\pm$ 149.44	827.26 $\pm$ 61.72	983.8 $\pm$ 0.34	937.43 $\pm$ 54.59	

Table 14: **Training Result on xArm.**

Task	Reach	Push	Peg in box	Hammer
Success rate (%)	96 $\pm$ 5	90 $\pm$ 17	80 $\pm$ 10	83 $\pm$ 20

Table 15: **Training Result on Adroit.**

Task	Door	Hammer	Pen
Success rate (%)	96 $\pm$ 3	78 $\pm$ 7	48 $\pm$ 15