

A Details and hyperparameters

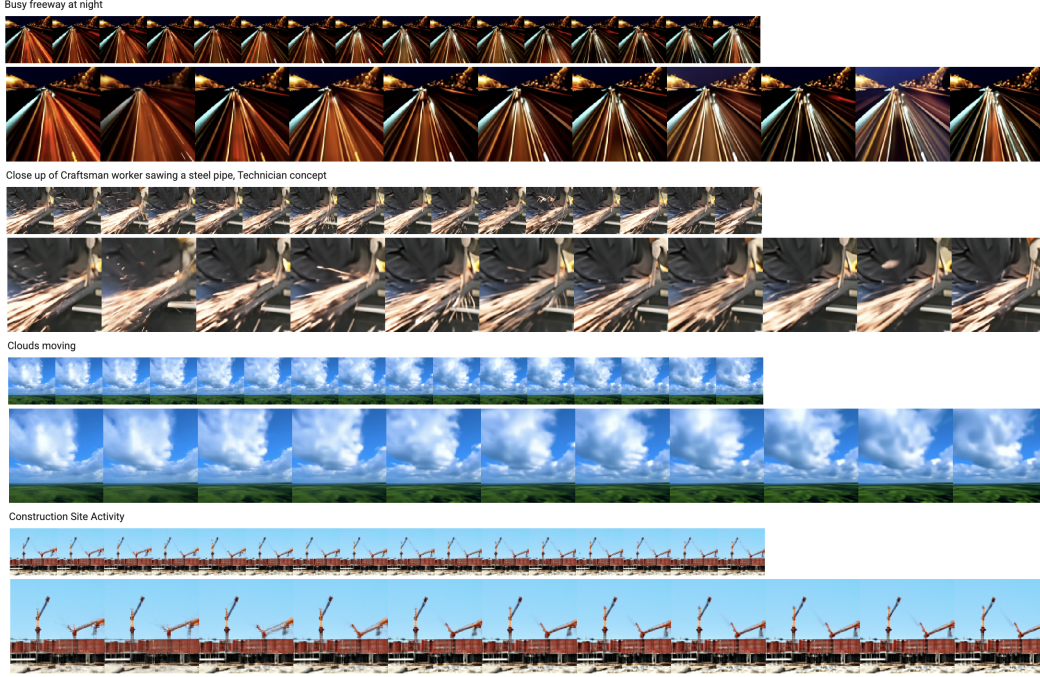


Figure 5: More samples accompanying Fig. 2.

Here, we list the hyperparameters, training details, and compute resources used for each model.

A.1 UCF101

Base channels: 256	Optimizer: Adam ($\beta_1 = 0.9, \beta_2 = 0.99$)
Channel multipliers: 1, 2, 4, 8	Learning rate: 0.0003
Blocks per resolution: 2	Batch size: 128
Attention resolutions: 8, 16, 32	EMA: 0.9999
Attention head dimension: 64	Dropout: 0.1
Conditioning embedding dimension: 1024	Training hardware: 128 TPU-v4 chips
Conditioning embedding MLP layers: 4	Training steps: 260000
Diffusion noise schedule: cosine	Joint training independent images per video: 8
Noise schedule log SNR range: $[-20, 20]$	Sampling timesteps: 256
Video resolution: 16x64x64 frameskip 1	Sampling log-variance interpolation: $\gamma = 0.1$
Weight decay: 0.0	Prediction target: ϵ

A.2 BAIR Robot Pushing

Base channels: 128	Optimizer: Adam ($\beta_1 = 0.9, \beta_2 = 0.999$)
Channel multipliers: 1, 2, 3, 4	Learning rate: 0.0002
Blocks per resolution: 3	Batch size: 128
Attention resolutions: 8, 16, 32	EMA: 0.999
Attention head dimension: 64	Dropout: 0.1
Conditioning embedding dimension: 1024	Training hardware: 128 TPU-v4 chips
Conditioning embedding MLP layers: 2	Training steps: 660000
Diffusion noise schedule: cosine	Joint training independent images per video: 8
Noise schedule log SNR range: $[-20, 20]$	Sampling timesteps: 256 (+256 Langevin cor.)
Video resolution: 16x64x64 frameskip 1	Sampling log-variance interpolation: $\gamma = 0.0$
Weight decay: 0.01	Prediction target: \mathbf{v}
Reconstruction guidance weight: 50	Data augmentation: left-right flips

A.3 Kinetics

Base channels: 256	Optimizer: Adam ($\beta_1 = 0.9, \beta_2 = 0.99$)
Channel multipliers: 1, 2, 4, 8	Learning rate: 0.0002
Blocks per resolution: 2	Batch size: 256
Attention resolutions: 8, 16, 32	EMA: 0.9999
Attention head dimension: 64	Dropout: 0.1
Conditioning embedding dimension: 1024	Training hardware: 256 TPU-v4 chips
Conditioning embedding MLP layers: 2	Training steps: 220,000
Diffusion noise schedule: cosine	Joint training independent images per video: 8
Noise schedule log SNR range: $[-20, 20]$	Sampling timesteps: 128 (+128 Langevin cor.)
Video resolution: 16x64x64 frameskip 1	Sampling log-variance interpolation: $\gamma = 0.0$
Weight decay: 0.0	Prediction target: \mathbf{v}
Reconstruction guidance weight: 9	

A.4 Text-to-video

Small 16x64x64 model

Base channels: 128	Optimizer: Adam ($\beta_1 = 0.9, \beta_2 = 0.99$)
Channel multipliers: 1, 2, 4, 8	Learning rate: 0.0003
Blocks per resolution: 2	Batch size: 128
Attention resolutions: 8, 16, 32	EMA: 0.9999
Attention head dimension: 64	Dropout: 0.0
Conditioning embedding dimension: 1024	Training hardware: 64 TPU-v4 chips
Conditioning embedding MLP layers: 4	Training steps: 200,000
Diffusion noise schedule: cosine	Joint training independent images per video: 0, 4, 8
Noise schedule log SNR range: $[-20, 20]$	Sampling timesteps: 256
Video resolution: 16x64x64 frameskip 1	Sampling log-variance interpolation: $\gamma = 0.3$
Weight decay: 0.0	Prediction target: ϵ

Large 16x64x64 model

Base channels: 256	Optimizer: Adam ($\beta_1 = 0.9, \beta_2 = 0.99$)
Channel multipliers: 1, 2, 4, 8	Learning rate: 0.0003
Blocks per resolution: 2	Batch size: 128
Attention resolutions: 8, 16, 32	EMA: 0.9999
Attention head dimension: 64	Dropout: 0.0
Conditioning embedding dimension: 1024	Training hardware: 128 TPU-v4 chips
Conditioning embedding MLP layers: 4	Training steps: 700,000
Diffusion noise schedule: cosine	Joint training independent images per video: 8
Noise schedule log SNR range: $[-20, 20]$	Sampling timesteps: 256
Video resolution: 16x64x64 frameskip 1,4	Sampling log-variance interpolation: $\gamma = 0.3$
Weight decay: 0.0	Prediction target: ϵ

Large 9x128x128 model

Base channels: 128	Optimizer: Adam ($\beta_1 = 0.9, \beta_2 = 0.99$)
Channel multipliers: 1, 2, 4, 8, 16	Learning rate: 0.0002
Blocks per resolution: 2	Batch size: 128
Attention resolutions: 8, 16, 32	EMA: 0.9999
Attention head dimension: 128	Dropout: 0.0
Conditioning embedding dimension: 1024	Training hardware: 128 TPU-v4 chips
Conditioning embedding MLP layers: 4	Training steps: 800,000
Diffusion noise schedule: cosine	Joint training independent images per video: 7
Noise schedule log SNR range: $[-20, 20]$	Sampling timesteps: 256
Video resolution: 9x128x128 frameskip 1	Sampling log-variance interpolation: $\gamma = 0.3$
Weight decay: 0.0	Prediction target: ϵ