

Causality in Video Diffusers is Separable from Denoising (Supplementary Material)

Anonymous CVPR submission

Paper ID ****

001 A. Additional Experimental Setup

002 A.1. Datasets

003 **TECO-Minecraft[19].** We adopt the long-context,
004 action-conditioned prediction setup popularized by
005 TECO [19]. Each video contains 300 frames, with
006 128×128 resolution, with per-frame action annotations.
007 We evaluate on 256 video clips; for long-horizon quality
008 (FVD), each clip supplies 36 ground-truth context frames
009 followed by 264 generated frames; for frame-wise metrics
010 (LPIPS, SSIM, PSNR), each clip supplies 144 observed
011 frames followed by 156 generated frames. This set-up
012 exactly aligns with TECO.

013 **UCF0101[14].** We use the UCF-101 dataset to demon-
014 strate the models' capability in real-world motion. This
015 dataset comprises $\sim 13K$ unconstrained, unconditional ac-
016 tion videos. Following MCVD [15]/ExtDM [22] and
017 FAR [6], we randomly sample 256 videos and, for each,
018 draw 100 stochastic trajectories. Pixel metrics (LPIPS/S-
019 SIM/PSNR) are computed best-of-100 per video, and FVD
020 is averaged over all 100 trajectories.

021 **RealEstate10K[1].** We also perform unconditional
022 generation experiments on an auxiliary benchmark,
023 RealEstate10K, a ddataset consisting of real-world indoor
024 scenes. While this dataset is predominantly used in 3D
025 tasks, we use it because it is a relatively small real-world
026 dataset, where pretraining is feasible for our experiments.
027 We use a resolution of 256. Since our model is orthogonal
028 to camera-pose conditioning techniques, we simply per-
029 form unconditional prediction tasks with 16 context frames
030 and 48 generated frames. Results are shown in Tab 1.

031 **Tokenizer.** Following the common practice, we compress
032 video frames into video latent, and apply diffusion mod-
033 els in the corresponding latent space. To compress a video,
034 we adopt a series of VAE and DCAE models [3]. For
035 Minecraft and UCF, we use the DCAE trained in FAR [6];

Table 1. **Unconditional generation on RealEstate10K (128^2 , $16 \rightarrow 48$).** Numbers are measured at 400k training steps with 50 denoising steps.

Model	Sec/F	16 → 48			
		LPIPS↓	SSIM↑	PSNR↑	FVD↓
Causal DiT-B	1.07	0.172	0.594	19.35	101.64
CSD-B	0.44	0.142	0.616	19.67	102.83
CSD-B^E	0.45	<u>0.139</u>	<u>0.622</u>	<u>19.95</u>	<u>101.61</u>
CSD-B^D	1.03	0.135	0.623	20.01	85.12

for RealEstate10K, we adopt the E2E-VAE tokenizer from [8] finetuned from VA-VAE [21].

036
037

038 B. Model and Training Details

039 B.1. Design Details of Separable Causal Diffusion

Let \mathcal{E}_ϕ denote the causal reasoning encoder and \mathcal{D}_θ the frame-wise diffusion decoder. The encoder runs once per video frame outside the denoising loop to produce a latent c_t , summarizing the temporal context. Then, the decoder denoises each frame with multiple diffusion steps, conditioned on c_t :

$$c_t = \mathcal{E}_\phi(x_{<t}, a_{\leq t}) \quad , \quad \hat{v}_t^{(s)} = \mathcal{D}_\theta(x_t^{(s)}, s, c_t).$$

046

Here $x_t^{(s)}$ are noisy frame latents at step s , and $\hat{v}_t^{(s)}$ is the predicted velocity/score used by the sampler. \mathcal{E}_ϕ uses frame-wise causal attention and KV caches; \mathcal{D}_θ is perform intra-frame bidirectional attention. The amortized per-frame cost is therefore $O(\mathcal{E}_\phi) + S \cdot O(\mathcal{D}_\theta)$.

047

048

049

050

051

In our experiments, we choose $O(\mathcal{E}_\phi) \gg O(\mathcal{D}_\theta)$ because of two reasons: 1) Empirically, we observe that a large portion of layer features are sharable across the denoising process (Fig.2). 2) the encoder's cost is amortized across multiple denoising steps, so it can be made larger without significantly sacrificing efficiency.

052

053

054

055

056

057

Table 2. Model variants and depth split. Depth is split into ℓ causal blocks and m diffusion blocks. BP/frame = $\ell + S \cdot m$ with $S=50$.

Model	#Blocks	Hidden	#Heads	Params	BP / frame
DiT-B	12	768	12	131M	600
SCD-B	8+4	768	12	132M	208
SCD-B^E	12+4	768	12	174M	212
SCD-B^D	8+12	768	12	217M	608
FAR-M	12	1024	16	230	600
SCD-M	8+4	1024	16	230M	208
SCD-M^E	12+4	1024	16	306M	212
SCD-M^D	8+12	1024	16	383M	608

Algorithm 1 SCD Training

Require: Videos $\mathbf{x}_{1:N}$ with controls $\mathbf{a}_{1:N}$, where N is the number of frames; temporal reasoning module \mathcal{E}_ϕ ; frame diffusion module \mathcal{D}_θ ; diffusion loss \mathcal{L} ; noisy multi-batch size K .

```

1: repeat
2:   Choose target frame  $i \in \{1, \dots, N\}$ 
3:    $c_i \leftarrow \mathcal{E}_\phi(\mathbf{x}_{<i}, \mathbf{a}_{\leq i})$ 
4:    $\mathcal{L}_{\text{step}} \leftarrow 0$ 
5:   for  $k = 1$  to  $K$  do
6:     Sample  $t \sim \mathcal{U}[0, 1]$ ,  $\epsilon \sim \mathcal{N}(0, I)$ 
7:      $x_i^t \leftarrow (1 - t) x_i + t \epsilon$ 
8:      $\hat{u} \leftarrow \mathcal{D}_\theta(x_i^t, t, c_i)$ 
9:      $\mathcal{L}_{\text{step}} \leftarrow \mathcal{L}_{\text{step}} + \mathcal{L}(\hat{u}, x_i, \epsilon, t)$ 
10:    Take a gradient step on  $\nabla_{\theta, \phi} (\mathcal{L}_{\text{step}}/K)$ 
11: until converged

```

058

B.2. Architectures and Model Variants

We follow the Diffusion Transformer (DiT) structure [12] to implement the SCD neural network. We follow DiT's width/head configuration for hidden size and MLP. To compare with FAR [6], the SOTA model on Minecraft, we also adopt its FAR-M parametrization. Table 2 enumerates the variants used in our experiments and reports BP/frame under $S=50$.

066

B.3. Algorithmic Pipeline

We summarize the end-to-end training and inference procedures of Separable Causal Diffusion (SCD) in Algorithms 1 and 2.

070

C. Ablation Studies

071

C.1. Encoder-Decoder Interface

072

We compare two ways of providing the context latent c_t to the frame-wise diffusion decoder \mathcal{D}_θ : (i) *Channel Concatenation* we concatenates c_t with the noisy frame tokens along

Algorithm 2 SCD Generation by Roll-out over Frames

Require: Controls $\mathbf{a}_{1:N}$, where N is the number of frames to be generated; temporal reasoning module \mathcal{E}_ϕ ; frame diffusion module \mathcal{D}_θ ; sampler with T denoising steps and schedule $\{t_1, \dots, t_T\}$.

```

1:  $\hat{\mathbf{x}} \leftarrow []$  % generated frames buffer
2: for  $i = 1, \dots, N$  do
3:    $c_i \leftarrow \mathcal{E}_\phi(\hat{\mathbf{x}}_{<i}, \mathbf{a}_{\leq i})$  % AR context: previously
   generated frames
4:   Initialize  $z^T \sim \mathcal{N}(0, I)$ 
5:   for  $t = T, T-1, \dots, 1$  do
6:      $\hat{u} \leftarrow \mathcal{D}_\theta(z^t, t, c_i)$ 
7:      $z^{t-1} \leftarrow \text{SAMPLER}(z^t, \hat{u}, t)$ 
8:      $\hat{x}_i \leftarrow z^0$ ; append  $\hat{x}_i$  to  $\hat{\mathbf{x}}$ 
9: return  $\hat{\mathbf{x}}_{1:N}$ 

```

Table 3. Ablations on the encoder-decoder interface. Sequence fusion outperforms channel fusion; temporal RoPE is slightly better than identical (non-causal) RoPE. For simplicity of ablation, metrics are reported at 400 training steps, with the unified "144 context frames, 156 generated frames" evaluation setup.

Encoder-Decoder Interface	FVD \downarrow	LPIPS \downarrow
Channel dim.	25.4	0.231
Frame dim. with temporal RoPE	24.8	0.219
Frame dim. with identical RoPE	25.1	0.223

the channel dimension, projects back to the standard channel dimension with a linear layer, and feeds it to the decoder as input; (ii) *Frame Concatenation* We prepends c_t as a prefix frame for the noisy frame, and then feeds them into the decoder. This effectively positions c_t as the context frame of the current frame, performing self-attention together as a whole sequence. We also ablate the positional embedding applied to the context frame, either embedding it as "the last temporal frame" or "the current frame". As results in Tab. 3 demonstrates, sequence fusion strategy with P.E. embedded as historical frame outperforms other alternative.

075

076

077

078

079

080

081

082

083

084

085

C.2. Training Noisy Batches: Amortized Multi-Sample Decoding

Because \mathcal{E}_ϕ consumes only clean history, it is noise-agnostic. We therefore perform one efficiency trick in training: for each batch of clean videos, we encode once per frame to obtain c_t , then draw K i.i.d. noise/timestep pairs for the current frame and run \mathcal{D}_θ K times, saving the amortized cost for learning each noisy batch. As Tab. 4 demonstrates, throughput of noisy batch increases with K , which also improves the training speed.

088

089

090

091

092

093

094

095

Table 4. Amortized multi-sample decoding in Training. Encoder depth 8, decoder depth 4, as in our SCD-B and SCD-M models. For simplicity of ablation, metrics are reported at 400 training steps, with the unified "144 context frames, 156 generated frames" evaluation setup.

K	BP/clean batch	BP/noisy batch	noisy batch/s	FVD
1	$8 + 4 \times 1 = 12$	$8/1 + 4 = 12$	22.0	23.9
2	$8 + 4 \times 2 = 16$	$8/2 + 4 = 8$	39.2	23.3
4	$8 + 4 \times 4 = 24$	$8/4 + 4 = 6$	63.0	23.1

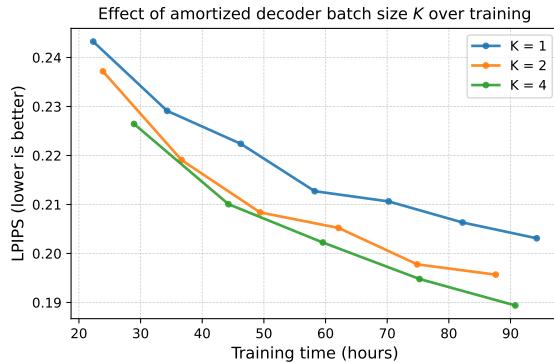


Figure 1. Effect of amortized decoder batch size K at matched training time. LPIPS on the validation set versus wall-clock training time (hours) for $K \in \{1, 2, 4\}$. Even when comparing at equal training time rather than equal optimization steps, larger K achieves lower LPIPS, indicating genuine gains from amortizing multiple noisy decoder samples per encoder pass.

096 **Training-time comparison.** The main-table ablation in
097 Sec. C.2 compares different K at matched optimization
098 steps, which slightly favors larger K because each step pro-
099 cesses more independently noised targets. To control for
100 this, Fig. 1 re-plots LPIPS against *wall-clock training time*.
101 Despite the heavier per-step compute, curves with $K=2$ and
102 especially $K=4$ reach lower LPIPS than $K=1$ at the same
103 elapsed time. This shows that amortized multi-sample de-
104 coding does more than just see more noise per step: reusing
105 the once-per-frame encoder output across multiple noisy de-
106 coder calls yields a better optimization trajectory even under
107 a fixed compute budget.

C.3. Noisy Context Latent: Context Corruption and CFG

We perturb only the context c_t by adding Gaussian noise $\tilde{c}_t = c_t + \eta_t \epsilon$, $\epsilon \sim \mathcal{N}(0, I)$, where c_t is normalized to unit variance and η_t^2 is the noise variance. Table 5 reports two complementary ablations on TECO–Minecraft at 400k steps. On the *left*, we vary the training-time corruption strength $\eta_t \in \{0, 0.05, 0.10, 0.20, 0.50\}$ and evaluate without any inference-time classifier-free guidance (CFG) on c_t : moderate noise around $\eta_t = 0.05$ improves FVD while

Table 5. Training-time causal corruption and test-time CFG with corruption (Minecraft, 400k steps). Left: vary η_t at training and evaluate with no CFG on c_t ($\eta_{cfg}=0$). Right: fix $\eta_t=0.05$ at training and sweep inference-time CFG scale η_{cfg} on a corrupted causal prior $\tilde{c}_t = c_t + 0.05 \epsilon$.

Noise η_t	Training corruption		Inference CFG with corruption		
	FVD \downarrow	LPIPS \downarrow	η_{cfg}	FVD \downarrow	LPIPS \downarrow
0.00	24.8	0.199	0.0	24.2	0.254
0.05	23.8	0.195	1.0	23.1	0.223
0.10	24.5	0.195	1.5	22.3	0.219
0.20	25.1	0.191	2.0	22.7	0.221
0.50	27.6	0.199	$(\eta_t=0.05)$		

stronger corruption eventually hurts. On the *right*, we fix training-time corruption at $\eta_t=0.05$ and ablate CFG using a negative branch with the same perturbation $\tilde{c}_t = c_t + 0.05 \epsilon$ and guidance scale $\eta_{cfg} \in \{0.0, 1.0, 1.5, 2.0\}$; $\eta_{cfg}=1.5$ gives the best trade-off. Overall, modest training-time corruption plus a small CFG weight on the corrupted causal prior ($\eta_t \approx 0.05$, $\eta_{cfg} \approx 1.5$) yields the strongest long-horizon visual quality. Because these perturbations act only on the causal interface, they do not require re-caching any context frames.

D. Additional Analysis of Section 4

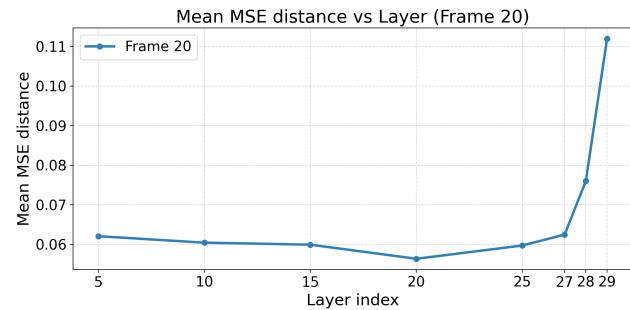
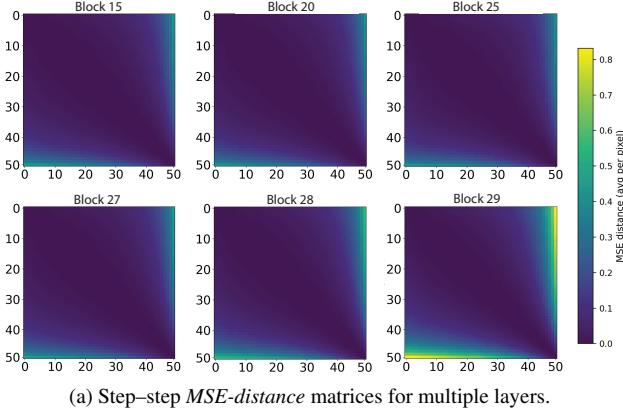
D.1. Redundancy across Denoising Steps

Feature similarity across denoising steps and depth. We extend the analysis in Fig. 2 to multiple depths of an autoregressively fine-tuned WAN-2.1 T2V-1.3B model. For a fixed set of prompts/seeds, we roll out the model for 50 denoising steps and, at every transformer block ℓ and step s , record the hidden features. We then compute a step–step *mean-squared-error (MSE) distance* matrix $\mathbf{S}_\ell \in \mathbb{R}^{T \times T}$ with entries $[\mathbf{S}_\ell]_{s,s'} = \|f_{\ell,s} - f_{\ell,s'}\|_2^2$, where $f_{\ell,s}$ denotes the layer- ℓ features at step s . Representative matrices for $\ell \in \{10, 15, 20, 25, 28, 29\}$ are shown in Fig. 2a and Fig. 2b

Layers 10–25 exhibit pronounced step-wise invariance: their similarity maps contain broad, near-uniform high-value bands (also visible at the 15th layer in the main paper), indicating that middle/early denoiser blocks repeatedly recompute almost the same features across the diffusion trajectory. In contrast, the last two layers ($\ell = 28, 29$) display markedly lower and more step-dependent similarity, consistent with these blocks performing step-specific, intra-frame rendering.

To further verify the redundancy observation, we fine-tune the baseline with a skip-layer design, in which the majority of denoising steps skip the middle-layers computation.

To further verify the redundancy finding, we fine-tune the causal baseline with a *skip-layer* schedule in which most



(b) Mean MSE distance versus block index. Each value here represents the average value across the entire corresponding matrix on the left.

Figure 2. **Redundant computation across denoising steps.** (a) Step-step feature *MSE-distance* matrices of a fine-tuned AR WAN-2.1 T2V-1.3B model at several layer depths (layers 5, 10, 15, 20, 25, 28, 29). Middle layers (10–25) show broad, *low-distance* bands across all 50 denoising steps, indicating that their features are mostly invariant along the diffusion trajectory, whereas the last few layers exhibit more step-dependent distances. (b) A complementary per-layer summary: the average MSE distance across all pairs of denoising steps remains small in the early and middle layers, but increase dramatically in late layers.. Together, these views support our claim that early/middle blocks perform largely redundant computation across denoising steps, while the deepest blocks remain step-specific for intra-frame rendering, motivating our design that amortizes the first 25 layers once per frame and reuses them across all denoising steps.

denoising steps bypass the middle of the network. Concretely, only the first five denoising steps run the full denoiser; all subsequent steps traverse just a short *prefix* of 5 early layers and a *suffix* of 10 late layers, skipping the middle chunk. We retain the prefix because, as shown in Fig. 7 of the main text, early layers are particularly important during fine-tuning. The outcome (Fig. 3 in the main text) is that the skipped model produces high-quality videos and recovers the visual fidelity of the fully causal baseline, validating that most cross-step computation in early/middle layers is redundant and can be shared. This experiment is designed to test the observation in the simplest setting. In later SCD fine-tuning, we adopt a more aggressive recipe that also works: the first 25 layers run *once per frame* (amortized across steps), and only the final 5+5 layers participate in per-step denoising under our SCD design.

D.2. Evidence on Other Models

Beyond the WAN-2.1 (1.3B) model, we also evaluate an open-source, open-weight **Self-Forcing** (SF) 1.3B model [7]—a few-step, *block-autoregressive* student distilled from a bidirectional teacher that predicts three latent frames per block. We chose this model deliberately for two reasons. **(i) Block-autoregression** is a widely used generative pattern in contemporary video systems, and even in emerging language diffusion models [2], so validating our analysis on a block-AR student makes the conclusions relevant to a broad class of architectures [9, 13, 20]. **(ii) Self-forcing** is the prevailing post-training recipe for large autoregressive video models, bridging the train-test gap and distilling many-step teachers into efficient few-step

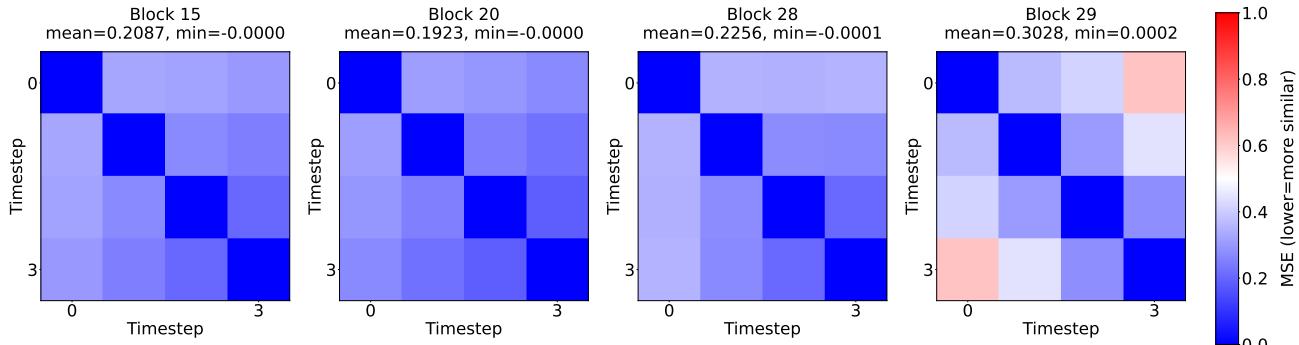
samplers [4, 5, 7, 10]. By observing our claims on both a many-step WAN and a few-step, block-autoregressive, self-forced student, we cover complementary ends of the design space; the aligned observations across these regimes would substantially strengthen the generality of our claims.

Observations. We see the same patterns as in Section 4. Early and middle *layers* produce very similar features across denoising steps. PCA views change little with the step index and already capture global structure in the *first block*. Deeper layers are temporally sparse and focus on intra-frame rendering. These results indicate that the separability trends hold for both a 50-step WAN and a distilled 4-step *block-autoregressive* self-forcing model.

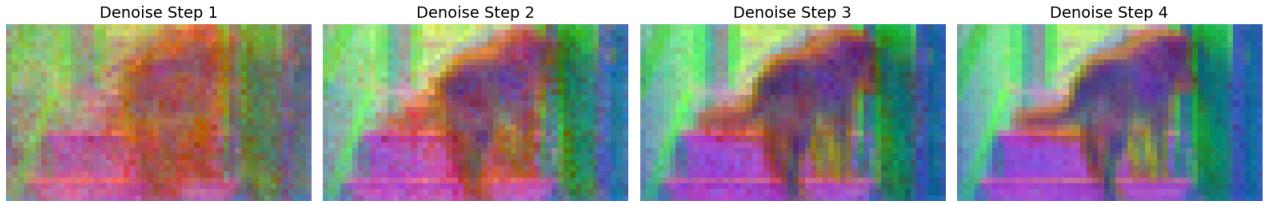
E. SCD Fine-tuning Details

This section describes fine-tuning details omitted from the main text. We first specify the teacher, data, and architecture adaptations used to convert a pretrained bidirectional video diffusion model into our Separable Causal Diffusion (SCD), which contains a causal encoder + a frame-wise diffusion decoder. We then detail the self-forcing rollout/distillation protocol used for post-training, and finally discuss capacity splits between encoder and decoder that highlight the flexibility of SCD.

Bidirectional Teacher. Unless otherwise noted, we fine-tune from a high-quality, bidirectional T2V checkpoint of WAN 2.1 T2V-1.3B [17, 18], whose weights are transplanted into our decoupled backbone (§E.1). All training

Frame 20: MSE Self-Similarity Across Blocks(a) **Self-forcing: step-step feature similarity** across multiple layers (analogous to Fig. 2(a) in Section 4). Early/middle layers show broad high-similarity (small MSE) bands over denoising steps.

PCA for block 20, frame 6 (all denoising steps)



PCA for block 20, frame 20 (all denoising steps)

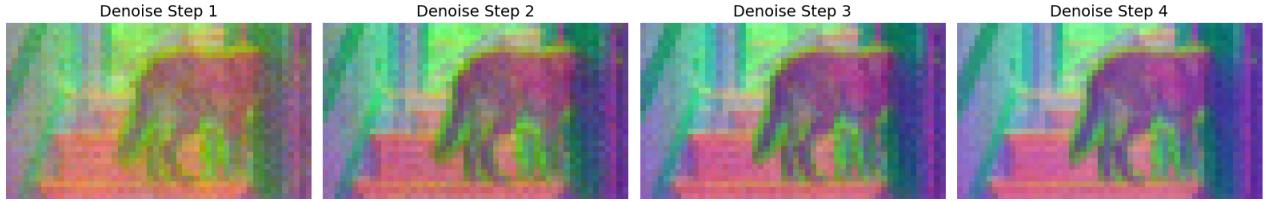
(b) **Self-forcing: PCA of activations** (combined view; analogous to Fig. 2(b) in Section 4). Principal components remain stable across denoising steps and across *blocks* in the block-autoregressive rollout.

Figure 3. Evidence on another model family (self-forcing, 4-step, block-autoregressive). We replicate the Section 4 analyses on a few-step self-forcing student. (a) Early/middle layers exhibit high step-step feature similarity. (b) PCA views confirm that principal directions stabilize early and change little across denoising steps and *blocks*.

212 lies in the latent spaces derived from the original VAE of
213 WAN 2.1.

214 **Datasets.** We use the text prompts from a 1M subset of
215 VidProM [16] following the same filtering process in [7].
216 For fine-tuning with diffusion loss with our architecture,
217 we use 70k synthetic data generated by WAN 2.1 T2V-14B
218 with the above text prompts. For self-rollout training, we
219 use the full 1M text prompts as conditions.

220 **Training Specifications.** We jointly train the encoder
221 and decoder with the conditional flow-matching objective
222 (Eq. (1) in the main text). For time step distribution, fol-

223 lowing WAN, we employ the timestep shifting $t'(k, t) =$
224 $\frac{kt}{1+(k-1)t}$ and the forward interpolation is given as $x_t =$
225 $(1 - t')x + t'\epsilon$, $\epsilon \sim \mathcal{N}(0, I)$, $t \in \mathcal{U}(0, 1)$. We use the
226 AdamW [11] optimizer for all experiments. Detailed hy-
227 perparameters can be found in the Table 6 and Table 7.

E.1. Architecture Adaptation for Decoupling

228 As described in the main paper, our decoupled backbone
229 implements once-per-frame temporal reasoning in a causal
230 encoder \mathcal{E}_ϕ and iterative rendering in a light frame-wise
231 diffusion decoder \mathcal{D}_θ , which differs from the teacher archi-
232 tecture. Therefore, to align the two architectures, in practice
233 we make two adaptations when initializing from a bidirec-
234

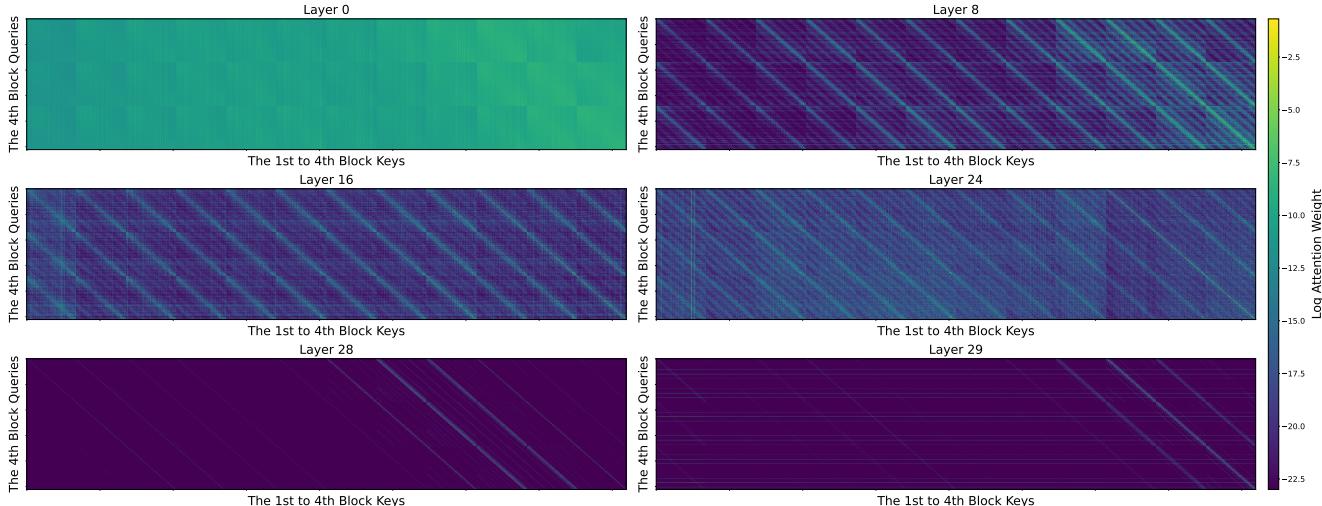


Figure 4. **Evidence on another model family (self-forcing, 4-step, block-autoregressive).** We replicate the Section 4 analyses on a few-step self-forcing student. As with the multi-step models, its deeper layers place minimal attention on past blocks, again revealing strong temporal sparsity.

235

tional teacher:

236 **(i) Input reparameterization for the encoder.** Pre-
 237 trained video diffusers consume a *noisy* current frame at
 238 each denoising step, while our encoder must operate on
 239 last generated frame instead of current frame. During
 240 fine-tuning, we therefore feed the encoder a corrupted cur-
 241 rent frame x_i^t at relatively *high* noise levels (e.g., top 20%
 242 of the diffusion/flow schedule), and at inference we replace
 243 it with pure Gaussian noise. This aligns the encoder’s input
 244 distribution with the teacher while preserving the decoupled
 245 compute pattern (once per frame for \mathcal{E}_ϕ , multi-step for \mathcal{D}_θ).

246 **(ii) Layer decomposition.** Directly treating early layers
 247 as encoder and late layers as decoder often harms genera-
 248 tion performance from finetuning. As discussed in the main
 249 text, the early layers play an important role in converting
 250 model input scale to an internal model scale. Guided by
 251 leave-one-out loss probing, we allocate the first 25 layers to
 252 \mathcal{E}_ϕ and build \mathcal{D}_θ by combining the first 5 and last 5 layers.

253 E.2. Hyperparameters

Table 6. Fine-tuning hyperparameters.

Resolution / Frames	$832 \times 480 / 81$ frames
Batch size	64
LR / WD / Optimizer	$2 \times 10^{-5} / 0.01 / \text{AdamW}(0.9, 0.99)$
EMA decay	0.99
Time sampler	$\frac{5t}{1+4t}, t \sim \mathcal{U}(0, 1)$

Table 7. Self-Forcing rollout training hyperparameters.

Resolution / Frames	$832 \times 480 / 81$ frames
Teacher	WAN 2.1 14B
Teacher CFG	3.0
Critic initialization	WAN 2.1 1.3B
Batch size	64
Student LR / WD / Optimizer	$2 \times 10^{-6} / 0.01 / \text{AdamW}(0.0, 0.99)$
Critic LR / WD / Optimizer	$4 \times 10^{-7} / 0.01 / \text{AdamW}(0.0, 0.99)$
Student EMA decay	0.99
Critic/student update ratio	5
Time sampler	$\frac{5t}{1+4t}, t \sim \mathcal{U}(0, 1)$

Throughput and latency. We report wall-clock through-
 254 put (FPS) and per-frame latency with batch size 1 on
 255 $1 \times \text{H100}$ 80 GB, charging the initial frame’s extra compute.
 256 SCD fine-tuned from a strong T2V teacher achieves ~ 11.1
 257 FPS with 0.29 s latency at 832×480 while retaining compet-
 258 itive VBench scores; the self-forcing baseline at the same
 259 scale reaches 8.9 FPS and 0.45 s latency.(Table 8).
 260

Table 8. Frame-Autoregressive Text-to-Video on VBench
 (832×480, batch 1). Reported throughput includes first-frame
 overhead.

Model	FPS ↑	Latency (s) ↓	Total ↑	Quality / Semantic ↑
Self Forcing	8.9	0.45	84.26	85.25 / 80.30
SCD (Ours)	11.1	0.29	84.03	85.14 / 79.60

E.3. Flexibility of Separable Causal Diffusion

A practical benefit of SCD is that temporal reasoning ca-
 262 pacity and per-frame rendering capacity can be traded *in-*
 263

264 dependently. Let total depth be $\ell+m$ with encoder depth ℓ
 265 (causal reasoning, amortized once per frame) and decoder
 266 depth m (frame-wise denoising, repeated S steps in infer-
 267 ence, where $S = 4$ in standard self-forcing settings). For
 268 fixed parameters:

- **Encoder-heavier variants** slightly increase reasoning cost per frame but systematically improve motion/layout adherence and long-horizon stability; it only slightly reduce the throughput since the encoder runs once per frame.
- **Decoder-heavier variants** improve per-frame detail at the cost of $S \times m$ block passes per frame; quality gains can be notable when targeting high-fidelity or very few denoising steps, but latency rises accordingly.

278 F. Additional Visualization

279 We provide an HTML viewer with additional video sam-
 280 ples in `Video_Samples/view_T2V.html`. It contains
 281 selected text-to-video generations from our 25-layer en-
 282 coder, 10-layer decoder Separable Causal Diffusion model
 283 under T2V setting; each clip is shown alongside its full text
 284 prompt so that content and conditioning can be inspected
 285 together. The folder also includes additional generations
 286 from models trained *from scratch* on TECO-Minecraft
 287 and RealEstate10K, illustrating long-horizon prediction and
 288 real-scene rendering in our decoupled architecture beyond
 289 the WAN fine-tuning regime.

290 References

- [1] Realestate10k: A large dataset of camera trajectories from videos. <https://google.github.io/realestate10k/>, 2018. 1 291
 292
 293
- [2] Marianne Arriola, Aaron Gokaslan, Justin T Chiu, Zhihan Yang, Zhixuan Qi, Jiaqi Han, Subham Sekhar Sahoo, and Volodymyr Kuleshov. Block diffusion: Interpolating between autoregressive and diffusion language models. *arXiv preprint arXiv:2503.09573*, 2025. 4 294
 295
 296
 297
 298
- [3] Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, Yao Lu, and Song Han. Deep compression autoencoder for efficient high-resolution diffusion models, 2025. 1 299
 300
 301
 302
- [4] Justin Cui, Jie Wu, Ming Li, Tao Yang, Xiaojie Li, Rui Wang, Andrew Bai, Yuanhao Ban, and Cho-Jui Hsieh. Self-forcing++: Towards minute-scale high-quality video generation. *arXiv preprint*, 2025. 4 303
 304
 305
 306
- [5] FastVideo Team. Fastvideo CausalWan2.2-I2V-A14B-Preview-Diffusers. <https://huggingface.co/FastVideo/CausalWan2.2-I2V-A14B-Preview-Diffusers>, 2025. Hugging Face model card. 4 307
 308
 309
 310
 311
- [6] Yuchao Gu, Weijia Mao, and Mike Zheng Shou. Long-context autoregressive video modeling with next-frame prediction. *arXiv preprint*, 2025. 1, 2 312
 313
 314
- [7] Xun Huang, Zhengqi Li, Guande He, Mingyuan Zhou, and Eli Shechtman. Self forcing: Bridging the train-test gap in autoregressive video diffusion. *arXiv preprint*, 2025. 4, 5 315
 316
 317
- [8] Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng. Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers. *arXiv preprint arXiv:2504.10483*, 2025. 1 318
 319
 320
 321
- [9] Shanchuan Lin, Ceyuan Yang, Hao He, Jianwen Jiang, Yuxi Ren, Xin Xia, Yang Zhao, Xuefeng Xiao, and Lu Jiang. Autoregressive adversarial post-training for real-time interactive video generation. *arXiv preprint arXiv:2506.09350*, 2025. Seaweed APT-2. 4 322
 323
 324
 325
 326
- [10] Kunhao Liu, Wenbo Hu, Jiale Xu, Ying Shan, and Shijian Lu. Rolling forcing: Autoregressive long video diffusion in real time. *arXiv preprint arXiv:2509.25161*, 2025. 4 327
 328
 329
- [11] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 5 330
 331
 332
- [12] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 2 333
 334
 335
- [13] Joonghyuk Shin, Zhengqi Li, Richard Zhang, Jun-Yan Zhu, Jaesik Park, Eli Shechtman, and Xun Huang. Motionstream: Real-time video generation with interactive motion controls. *arXiv preprint arXiv:2511.01266*, 2025. 4 336
 337
 338
 339
- [14] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human action classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 1 340
 341
 342
- [15] Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Mcvd: Masked conditional video diffusion for prediction, generation, and interpolation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 1 343
 344
 345
 346

- 347 [16] Wenhao Wang and Yi Yang. Vidprom: A million-scale real
348 prompt-gallery dataset for text-to-video diffusion models.
349 2024. 5
- 350 [17] WanTeam. Wan: Open and advanced large-scale video gen-
351 erative models. *arXiv preprint*, 2025. 4
- 352 [18] WanTeam. Wan 2.1 T2V-1.3B: Open-source text-to-video
353 model. *arXiv preprint arXiv:2503.00123*, 2025. 4
- 354 [19] Wilson Yan, Danijar Hafner, Stephen James, and Pieter
355 Abbeel. Temporally consistent transformers for video gen-
356 eration. In *Proceedings of the 40th International Conference*
357 *on Machine Learning (ICML)*, 2023. 1
- 358 [20] Shuai Yang, Wei Huang, Ruihang Chu, Yicheng Xiao,
359 Yuyang Zhao, Xianbang Wang, Muyang Li, Enze Xie, Ying-
360 cong Chen, Yao Lu, Song Han, and Yukang Chen. Longlive:
361 Real-time interactive long video generation. *arXiv preprint*,
362 2025. 4
- 363 [21] Jingfeng Yao, Bin Yang, and Xinggang Wang. Reconstruc-
364 tion vs. generation: Taming optimization dilemma in latent
365 diffusion models, 2025. 1
- 366 [22] Zhicheng Zhang, Junyao Hu, Wentao Cheng, Danda Paudel,
367 and Jufeng Yang. Extdm: Distribution extrapolation dif-
368 fusion model for video prediction. In *Proceedings of the*
369 *IEEE/CVF Conference on Computer Vision and Pattern*
370 *Recognition (CVPR)*, 2024. 1