



Self Forcing

home
papers
study
life

About me

- Self Forcing: Bridging the Train-Test Gap in Autoregressive Video Diffusion
- <https://arxiv.org/abs/2506.08009>
- NeurIPS 2025 spotlight

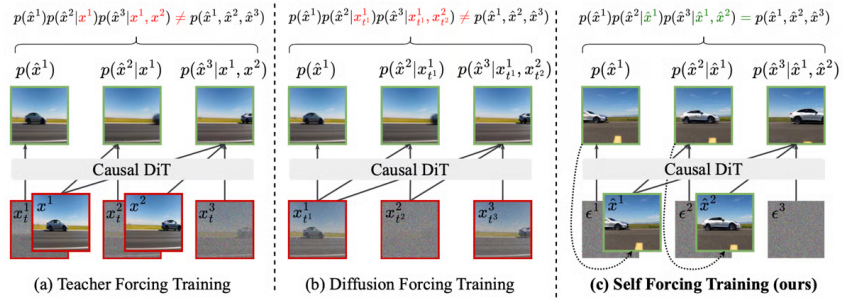


Figure 1: **Training paradigms for AR video diffusion models.** (a) In Teacher Forcing, the model is trained to denoise each frame conditioned on the preceding clean, ground-truth context frames. (b) In Diffusion Forcing, the model is trained to denoise each frame conditioned on the preceding context frames with varying noise levels. Both (a) and (b) generate outputs that do not belong to the distribution the model generates during inference. (c) Our Self Forcing approach performs autoregressive self-rollout *during training*, denoising the next frame based on previous context frames generated by itself. A distribution-matching loss (e.g., SiD, DMD, GAN) is computed on the final output video to align the distribution of generated videos with that of real videos. Our training paradigm closely mirrors the inference process, thereby bridging the train-test distribution gap.

Self Forcing
Self Forcing kv-cache
GT
kv-cache

Algorithm 1 Self Forcing Training

Require: Denoise timesteps $\{t_1, \dots, t_T\}$
Require: Number of video frames N
Require: AR diffusion model G_θ (returns KV embeddings via G_θ^{KV})

- 1: **loop**
- 2: Initialize model output $\mathbf{X}_\theta \leftarrow []$
- 3: Initialize KV cache $\mathbf{KV} \leftarrow []$
- 4: Sample $s \sim \text{Uniform}(1, 2, \dots, T)$
- 5: **for** $i = 1, \dots, N$ **do**
- 6: Initialize $x_{t_T}^i \sim \mathcal{N}(0, I)$
- 7: **for** $j = T, \dots, s$ **do**
- 8: **if** $j = s$ **then**
- 9: Enable gradient computation
- 10: Set $\hat{x}_0^i \leftarrow G_\theta(x_{t_j}^i; t_j, \mathbf{KV})$
- 11: \mathbf{X}_θ .append(\hat{x}_0^i)
- 12: Disable gradient computation
- 13: Cache $\text{kv}^i \leftarrow G_\theta^{\text{KV}}(\hat{x}_0^i; 0, \mathbf{KV})$
- 14: \mathbf{KV} .append(kv^i)
- 15: **else**
- 16: Disable gradient computation
- 17: Set $\hat{x}_0^i \leftarrow G_\theta(x_{t_j}^i; t_j, \mathbf{KV})$
- 18: Sample $\epsilon \sim \mathcal{N}(0, I)$
- 19: Set $x_{t_{j-1}}^i \leftarrow \Psi(\hat{x}_0^i, \epsilon, t_{j-1})$
- 20: **end if**
- 21: **end for**
- 22: **end for**
- 23: Update θ via distribution matching loss
- 24: **end loop**

Algorithm 2 Autoregressive Diffusion Inference with Rolling KV Cache

Require: KV cache of size L frames
Require: Denoise timesteps $\{t_1, \dots, t_T\}$
Require: Number of generated frames M
Require: AR diffusion model G_θ (returns KV embeddings via G_θ^{KV})

- 1: Initialize model output $\mathbf{X}_\theta \leftarrow []$
- 2: Initialize KV cache $\mathbf{KV} \leftarrow []$
- 3: **for** $i = 1, \dots, M$ **do**
- 4: Initialize $x_{t_T}^i \sim \mathcal{N}(0, I)$
- 5: **for** $j = T, \dots, 1$ **do**
- 6: Set $\hat{x}_0^i \leftarrow G_\theta(x_{t_j}^i; t_j, \mathbf{KV})$
- 7: **if** $j = 1$ **then**
- 8: \mathbf{X}_θ .append(\hat{x}_0^i)
- 9: Cache $\text{kv}^i \leftarrow G_\theta^{\text{KV}}(\hat{x}_0^i; 0, \mathbf{KV})$
- 10: **if** $|\mathbf{KV}| = L$ **then**
- 11: \mathbf{KV} .pop(0) \triangleright Cache eviction
- 12: **end if**
- 13: \mathbf{KV} .append(kv^i)
- 14: **else**
- 15: Sample $\epsilon \sim \mathcal{N}(0, I)$
- 16: Set $x_{t_{j-1}}^i \leftarrow \Psi(\hat{x}_0^i, \epsilon, t_{j-1})$
- 17: **end if**
- 18: **end for**
- 19: **end for**
- 20: **return** \mathbf{X}_θ

KV-Cache

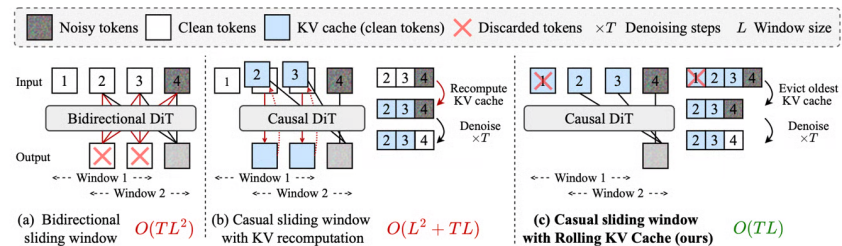


Figure 3: **Efficiency comparisons for video extrapolation.** When performing video extrapolation through sliding window inference, (a) bidirectional diffusion models trained with TF/DF [10, 73] do not support KV cache. (b) Prior causal diffusion models [69, 100] require re-computing KV when shifting the window. (c) Our method does not recompute KV and enables more efficient extrapolation.

- [xxxxxxx](#)
 - [xxx64 H100](#)
 - [xxxhttp://self-forcing.github.io](http://self-forcing.github.io)
-

[Older](#)

2026-5-6

CausVid

leicheng © 2022-2025

[Archive](#) [RSS feed](#) [GitHub](#) [Email](#) [QR Code](#)

Made with [Montaigne](#) and by [anton](#) 