

420 **A. Appendix**421 **A.1. Pseudo Code for SVD-Cache Framework****Algorithm 1: SVD-Cache Framework**

Input: Feature F_t , cached basis (V_C, σ_C) , previous low-rank feature $\hat{F}_{k,t-\Delta}$, EMA coefficient β , rank k

Output: Predicted feature $\hat{F}_{t+\Delta}$

// Stage 1: One-Time SVD on Reference Prompt (offline)

if *reference prompt (first time)* **then**

 Compute SVD: $F_{\text{ref}} = U\Sigma V^\top$;

 Cache basis: $V_C \leftarrow V, \sigma_C \leftarrow \text{diag}(\Sigma)$;

// Stage 2: Low-Rank Subspace Reconstruction

Compute $U = F_t V_C \text{diag}(\sigma_C)^{-1}$;

$(U_k, \sigma_{C,k}, V_{C,k}) \leftarrow$

 top- k components of (U, σ_C, V_C) ;

$F_{k,t} = U_k \text{diag}(\sigma_{C,k}) V_{C,k}^\top$;

$R_t = F_t - F_{k,t}$;

// Stage 3: Hybrid Caching Strategy

// EMA prediction for low-rank component

$\hat{F}_{k,t+\Delta} = \beta \hat{F}_{k,t-\Delta} + (1 - \beta) F_{k,t}$;

// Direct reuse for residual component

$\hat{R}_{t+\Delta} = R_t$;

// Stage 4: Feature Reconstruction

$\hat{F}_{t+\Delta} = \hat{F}_{k,t+\Delta} + \hat{R}_{t+\Delta}$;

return $\hat{F}_{t+\Delta}$

422 **A.2. Mathematical Derivation of Left Singular Matrix Reconstruction**

424 In this section, we provide the mathematical justification for
425 reconstructing the left singular matrix U using the cached
426 basis (V_C, σ_C) from a reference prompt.

427 **Theoretical Foundation.** Given the SVD decomposition of
428 the reference feature matrix:

$$429 \quad F_{\text{ref}} = U_{\text{ref}} \Sigma_{\text{ref}} V_{\text{ref}}^\top, \quad (12)$$

430 we cache the right singular vectors $V_C = V_{\text{ref}}$ and singular
431 values $\sigma_C = \text{diag}(\Sigma_{\text{ref}})$.

432 For a new feature matrix F_t at timestep t , if we assume
433 that F_t can be well-approximated in the subspace spanned
434 by V_C , we can express:

$$435 \quad F_t \approx U_t \text{diag}(\sigma_C) V_C^\top, \quad (13)$$

where $U_t \in \mathbb{R}^{N \times r}$ represents the coordinates of F_t in the
cached subspace.

Derivation of Reconstruction Formula. To solve for U_t ,
we multiply both sides by V_C from the right:

$$F_t V_C = U_t \text{diag}(\sigma_C) V_C^\top V_C. \quad (14)$$

Since V_C is orthonormal ($V_C^\top V_C = I$), this simplifies to:

$$F_t V_C = U_t \text{diag}(\sigma_C). \quad (15)$$

Multiplying both sides by $\text{diag}(\sigma_C)^{-1}$ yields:

$$U_t = F_t V_C \text{diag}(\sigma_C)^{-1}. \quad (16)$$

Geometric Interpretation. This reconstruction can be understood geometrically: $F_t V_C$ projects the feature matrix onto the cached principal directions, and $\text{diag}(\sigma_C)^{-1}$ normalizes these projections by the corresponding singular values. The resulting U_t represents the coefficients needed to express F_t as a linear combination of the cached basis weighted by the reference singular values.

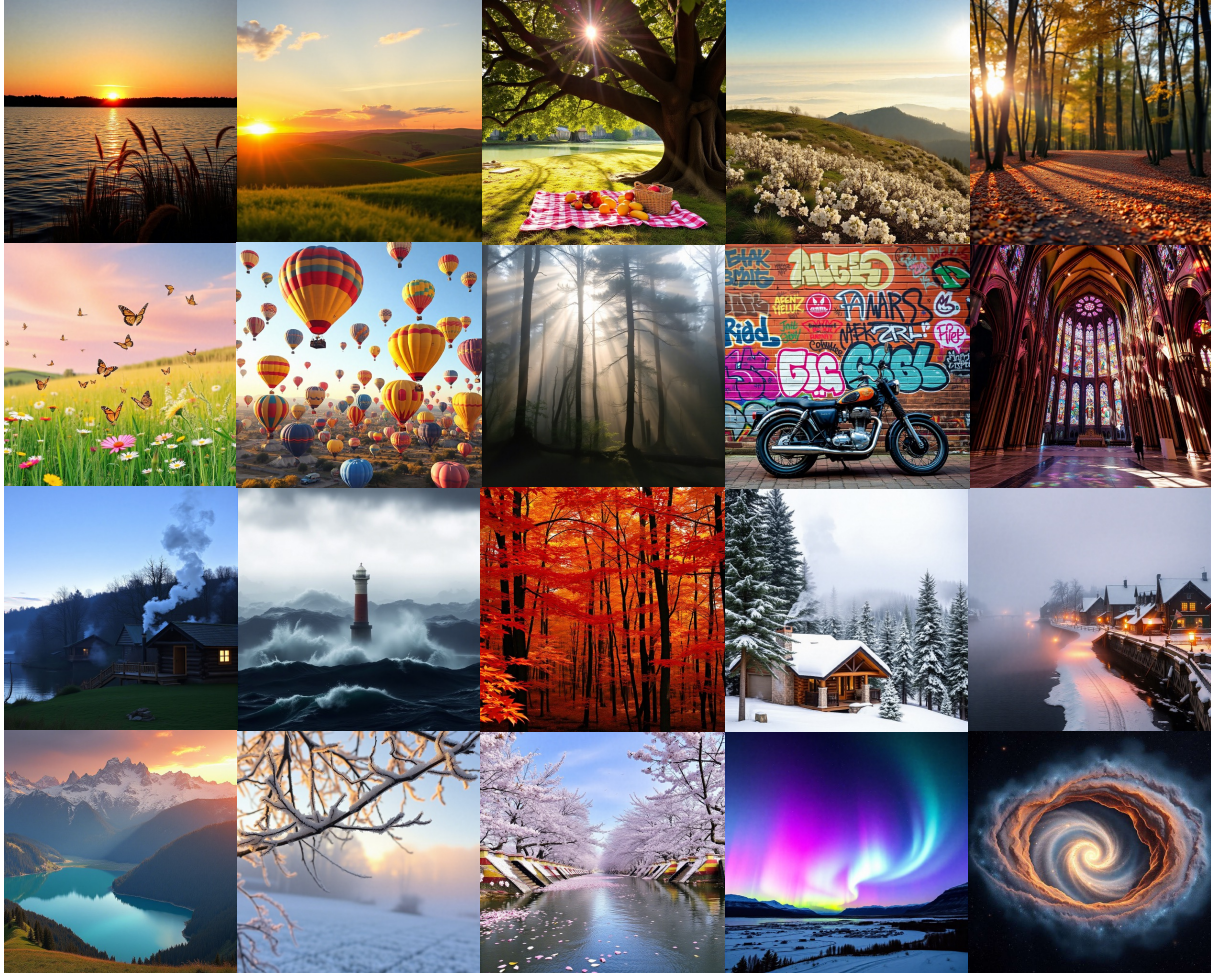
Key Advantage. This approach avoids recomputing the full SVD for each new prompt, reducing computational complexity from $\mathcal{O}(ND \min(N, D))$ for full SVD to $\mathcal{O}(NDr)$ for basis projection, where $r \ll \min(N, D)$ is the rank of the cached basis. The validity of this approximation is empirically supported by our observation in Fig. 1(b) that V and σ remain stable across different prompts.

A.3. Detailed Experiment Settings

Model Configurations and Evaluation Protocols. Our experimental evaluation encompasses six distinct generative models across two primary tasks: text-to-image generation and text-to-video generation. For text-to-image generation, we evaluate four different models: FLUX.1-dev, FLUX.1-Schnell, FLUX.1-dev-int8 and FLUX with sparse attention. Each model is assessed using standardized benchmarks and protocols to ensure fair comparisons. For text-to-video generation, we utilize HunyuanVideo, which is evaluated using VBench framework to capture various aspects of video quality and coherence. Detailed configurations for each model, including resolution settings, prompt selections, and evaluation metrics, can be found in the subsequent sections.

A.3.1. Text-to-Image Generation

FLUX.1-dev. We evaluate FLUX.1-dev using the standard DrawBench protocol, which provides a diverse set of 200 prompts spanning multiple categories including animals, colors, conflicting, and fine-grained details. All images are generated at a resolution of 1024×1024 pixels, maintaining consistency with the model’s optimal operating parameters. We assess generation quality using two metrics: ImageReward (IR) for evaluating photorealism and overall image



SVD-Cache on FLUX.1-schnell with $15.22 \times$ acceleration

Figure 7. Images of SVD-Cache on FLUX.1-Schnell. SVD-Cache achieves $15.22 \times$ acceleration while preserving image fidelity.

482 quality and CLIP Score for measuring text-image semantic
483 alignment.

484 **FLUX.1-schnell.** FLUX.1-schnell, a step-distillation model
485 from FLUX.1-dev, is designed for rapid image generation.
486 It is evaluated using the standard DrawBench protocol with
487 200 prompts at 1024×1024 resolution. Despite its focus on
488 speed optimization, we maintain the same evaluation stan-
489 dards, assessing both generation quality through ImageRe-
490 ward and CLIP Score, as well as fidelity through PSNR,
491 SSIM, and LPIPS metrics. This allows us to quantify the
492 trade-offs between generation speed and output quality.

493 **FLUX.1-dev-int8** FLUX.1-dev-int8, a quantized version of

494 FLUX.1-dev, is designed for efficient inference on resource-
495 constrained devices. It is evaluated using the standard Draw-
496 Bench protocol with 200 prompts at 1024×1024 resolution.
497 We assess generation quality using ImageReward, CLIP
498 Score, PSNR, SSIM and LPIPS.

499 **FLUX with sparse attention** We utilize SpargeAttention, a
500 universal training-free sparse attention accelerating language,
501 image, and video models. Prompts, resolution and evaluation
502 metrics are the same as FLUX.1-dev.



SVD-Cache on FLUX.1[dev]-int8, 7.61 ×



SVD-Cache with SparseAttention, 10.73 ×

Figure 8. **Images of SVD-Cache on FLUX.1[dev]-int8 and Sparse Attention.** SVD-Cache achieves 7.61 × and 10.73 × acceleration respectively while generating high-quality images.

503 A.3.2. Text-to-Video Generation

504 **HunyuanVideo.** For text-to-video generation, we evaluate
 505 HunyuanVideo using the comprehensive VBench framework,
 506 which provides multi-dimensional human-aligned assess-
 507 ments across 946 diverse prompts sourced from the VBench-
 508 full-info.json file. Videos are generated at 480×640 reso-
 509 lution with 65 frames each, providing substantial temporal
 510 content for thorough evaluation. VBench evaluates 18 criti-

cal aspects including but not limited to motion quality, visual
 511 appearance consistency, temporal coherence, semantic align-
 512 ment with text prompts, and overall video quality through
 513 human-correlated metrics. 514