# MODEL ALREADY KNOWS THE BEST NOISE: BAYESIAN ACTIVE NOISE SELECTION VIA ATTENTION IN VIDEO DIFFUSION MODEL

#### **Anonymous authors**

Paper under double-blind review



"A joyful Corgi with a fluffy coat and perky ears frolics in a sunlit park, the golden hues of sunset casting a warm glow on the scene. The camera zooms in on the Corgi's expressive face, capturing its bright eyes and wide, happy grin...."

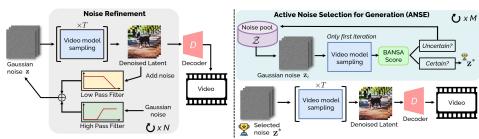


"A majestic brown bear, with its thick fur glistening in the dappled sunlight, begins its ascent up a towering pine tree in a dense forest. The bear's powerful claws grip the rough bark as it climbs higher, its muscles rippling with each movement..."

Figure 1: **Random Seed vs. Ours.** We propose ANSE, a noise selection framework, and the BANSA Score, an uncertainty-based metric. By selecting initial noise seeds with lower BANSA scores, which indicate more certain noise samples, ANSE improves video generation performance.

#### ABSTRACT

The choice of initial noise strongly affects quality and prompt alignment in video diffusion; different seeds for the same prompt can yield drastically different results. While recent methods use externally designed priors (e.g., frequency filtering or inter-frame smoothing), they often overlook internal model signals that indicate inherently preferable seeds. To address this, we propose ANSE (Active Noise Selection for Generation), a model-aware framework that selects high-quality seeds by quantifying attention-based uncertainty. At its core is BANSA (Bayesian Active Noise Selection via Attention), an acquisition function that measures entropy disagreement across multiple stochastic attention samples to estimate model confidence and consistency. For efficient inference-time deployment, we introduce a Bernoulli-masked approximation of BANSA that estimates scores from a single diffusion step and a subset of informative attention layers. Experiments across diverse text-to-video backbones demonstrate improved video quality and temporal coherence with marginal inference overhead, providing a principled and generalizable approach to noise selection in video diffusion. See our project page: https://anse-anonymous.vercel.app/.



(a) Conventional Noise Initialization approaches

(b) Our proposed framework, ANSE

Figure 2: **Conceptual comparison of noise initialization.** (a) Prior methods (Wu et al., 2024; Yuan et al., 2025) refine noise via frequency priors and full diffusion sampling, leading to high cost. (b) Our method instead selects noise seeds by estimating attention-based uncertainty at the first denoising step, enabling efficient, model-aware selection.

# 1 Introduction

Diffusion models have rapidly established themselves as a powerful class of generative models, demonstrating state-of-the-art performance across images and videos (Rombach et al., 2022; Esser et al., 2024; Xie et al., 2024; Chen et al., 2024b;a; Wang et al., 2025; Yang et al., 2024; Zheng et al., 2024; Team, 2025). In particular, Text-to-Video (T2V) diffusion models have received increasing attention for their ability to generate temporally coherent and visually rich video sequences. To achieve this, most T2V model architectures extend Text-to-Image (T2I) diffusion backbones by incorporating temporal modules or motion-aware attention layers (Blattmann et al., 2023; He et al., 2022; Wang et al., 2023; Chen et al., 2023; 2024a; Guo et al., 2024b; Wang et al., 2025). Furthermore, other works explore video generative structures, such as causal autoencoders or video autoencoderbased models, which aim to generate full video volumes rather than a sequence of independent frames (Hong et al., 2022; Yang et al., 2024; HaCohen et al., 2024; Kong et al., 2024; Zheng et al., 2024; Team, 2025).

Beyond architectural design, another promising direction lies in improving noise initialization at inference time for T2I and T2V generation (Guo et al., 2024a; Eyring et al., 2024; Chen et al., 2024c; Xu et al., 2025). This aligns with the growing trend of inference-time scaling, observed not only in Large Language Models (Brown et al., 2024; Snell et al., 2024) but also in diffusion-based generation systems (Ma et al., 2025). Due to the iterative nature of the diffusion process, the choice of initial noise profoundly influences video quality, temporal consistency, and prompt alignment (Ge et al., 2023; Wu et al., 2024; Qiu et al.; Yuan et al., 2025). As illustrated in Figure 1, the same prompt can lead to drastically different videos depending solely on the noise seed, motivating the need for intelligent noise selection.

Recent approaches tackle this by introducing external noise priors. PYoCo (Ge et al., 2023) enforces inter-frame dependent patterns for coherence but requires heavy fine-tuning. FreeNoise (Qiu et al.) reschedules noise across time with a fusion strategy, FreeInit (Wu et al., 2024) preserves low-frequency components via frequency filtering, and FreqPrior (Yuan et al., 2025) extends this with Gaussian-shaped priors and partial sampling. While effective, these methods rely on external priors and repeated full diffusion passes, while ignoring internal model signals that identify inherently preferable seeds.

To address this limitation, we propose a model-aware noise selection framework, ANSE (Active Noise Selection for Generation), grounded in Bayesian uncertainty. At the core of ANSE is BANSA (Bayesian Active Noise Selection via Attention), an acquisition function that identifies noise seeds inducing confident and consistent attention behaviors under stochastic perturbations. A conceptual comparison between our method and prior frequency-based approaches is illustrated in Figure 2, highlighting the difference between external priors and model-informed uncertainty estimates.

Unlike BALD (Gal et al., 2017), which estimates uncertainty from classification logits, BANSA measures entropy in attention maps, arguably the most informative signals in diffusion. It compares the mean of per-pass entropies with the entropy of the mean map, capturing both uncertainty and cross-pass disagreement. A lower BANSA score indicates more confident, consistent attention and empirically correlates with more coherent video generation (Fig. 1). To make BANSA inference-

friendly, we approximate it with Bernoulli-masked attention, yielding multiple stochastic attention samples from a single forward pass. We further cut cost by evaluating early denoising steps and only a subset of informative layers selected via correlation analysis. Our contributions are threefold:

- We present **ANSE**, the first active noise selection framework for video diffusion, grounded in a Bayesian formulation of attention-based uncertainty.
- We introduce **BANSA**, an acquisition function that measures attention consistency under stochastic perturbations, enabling model-aware noise selection without retraining.
- Our method improves video quality and temporal consistency across diverse text-to-video backbones, with only marginal inference overhead.

## 2 Preliminary

**Video Diffusion Models** Diffusion models (Ho et al., 2020; Song et al., 2021b) have achieved strong results across generative tasks. In T2V, pixel-space diffusion is costly, so most video diffusion models adopt latent diffusion and operate in a compressed latent space. A video autoencoder with encoder  $\mathcal{E}$  and decoder  $\mathcal{D}$  reconstructs  $\mathbf{x}$  as  $\mathbf{x} = \mathcal{D}(\mathcal{E}(\mathbf{x}))$ . Let  $\mathbf{z}_0 = \mathcal{E}(\mathbf{x})$ . The forward diffusion then progressively adds noise over time:

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \, \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}), \quad t = 1, \dots, T,$$

where  $\bar{\alpha}_t$  is a pre-defined variance schedule. To learn the reverse process, a denoising network  $\epsilon_{\theta}$  is trained using the denoising score matching loss (Vincent, 2011):

$$\mathcal{L}_{\theta} = \mathbb{E}_{\mathbf{z}_{t}, \boldsymbol{\epsilon}, t} \left[ \left\| \boldsymbol{\epsilon}_{\theta}(\mathbf{z}_{t}, \mathbf{c}, t) - \boldsymbol{\epsilon} \right\|^{2} \right],$$

where c is the conditioning text. Sampling starts from Gaussian noise  $\mathbf{z}_T \sim \mathcal{N}(0, \mathbf{I})$  and uses a deterministic DDIM solver (Song et al., 2021a). Each step updates as:

$$\mathbf{z}_{t-1} = \sqrt{\bar{\alpha}_t} \, \hat{\mathbf{z}}_0(t) + \sqrt{1 - \bar{\alpha}_{t-1}} \, \boldsymbol{\epsilon}_{\theta}(\mathbf{z}_t, \mathbf{c}, t),$$

where the denoised latent estimate  $\hat{\mathbf{z}}_0(t) := \frac{\mathbf{z}_t - \sqrt{1 - \bar{\alpha}_t} \, \epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}, t)}{\sqrt{\bar{\alpha}_t}}$  is obtained using Tweedie's formula (Efron, 2011; Kim & Ye, 2021). This iterative process continues until t = 1, yielding the final denoised latent  $\mathbf{z}_0$ , which is decoded into a video via  $\mathcal{D}$ .

Bayesian Active Learning by Disagreement (BALD) Active learning selects the most informative samples from an unlabeled pool to improve training. Acquisition functions are commonly uncertainty-based (Houlsby et al., 2011; Gal et al., 2017; Kirsch et al., 2019; Yoo & Kweon, 2019) or distribution-based (Mac Aodha et al., 2014; Yang et al., 2015; Sener & Savarese, 2017; Sinha et al., 2019), and some use external modules such as auxiliary predictors (Yoo & Kweon, 2019; Tran et al., 2019; Kim et al., 2021). While most prior work targets image classification, we adapt uncertainty-based acquisition to text-to-video generation without additional models.

Predictive entropy is widely used, but it reflects data noise and does not isolate parameter uncertainty. BALD instead measures *epistemic* uncertainty via the mutual information between predictions y and parameters  $\theta$ :

$$BALD(\mathbf{x}) = \mathcal{H}[p(\mathbf{y}|\mathbf{x})] - \mathbb{E}_{p(\theta|\mathcal{D}_U)} \left[ \mathcal{H}[p(\mathbf{y}|\mathbf{x},\theta)] \right], \tag{1}$$

where  $\mathcal{H}[p] = -\sum_y p(y) \log p(y)$  is Shannon entropy (Shannon, 1948). The first term is the entropy of the mean prediction, and the second is the average entropy across stochastic forward passes. A high BALD score means predictions are confident yet disagree, indicating high epistemic uncertainty. Since the posterior over  $\theta$  is intractable, BALD is approximated using K stochastic forward passes (e.g., Monte Carlo dropout):

$$\widehat{\text{BALD}}(\mathbf{x}) = \mathcal{H}\left[\frac{1}{K} \sum_{k=1}^{K} p^{(k)}(\mathbf{y}|\mathbf{x})\right] - \frac{1}{K} \sum_{k=1}^{K} \mathcal{H}\left[p^{(k)}(\mathbf{y}|\mathbf{x})\right].$$
(2)

We reinterpret BALD for inference-time generative modeling. Rather than selecting samples for labeling, we apply BALD to rank noise seeds by their epistemic uncertainty. Selecting seeds with lower BALD scores results in more stable model behavior and leads to higher-quality generations.

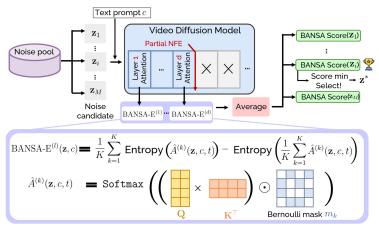


Figure 3: Overview of our BANSA-based noise selection process. Given a text prompt c, we compute BANSA scores for multiple noise seeds  $\{\mathbf{z}_1, \dots, \mathbf{z}_M\}$  using Bernoulli-masked attention maps from selected layers at an early diffusion step. The seed with the lowest score, indicating confident and consistent attention, is selected for generation.

# 3 METHODS

We propose **ANSE**, a framework for selecting high-quality noise seeds in T2V diffusion based on model uncertainty (Fig. 2). ANSE centers on **BANSA**, an acquisition function that transfers uncertainty criteria from classification to the attention space of generative diffusion models (Sec. 3.1). For efficient inference, we approximate BANSA via Bernoulli-masked attention sampling (Sec. 3.2). To avoid redundancy, we select a representative attention layer using correlation-based linear probing (Sec. 3.3). The full pipeline appears in Fig. 3.

#### 3.1 BANSA: BAYESIAN ACTIVE NOISE SELECTION VIA ATTENTION

We introduce **BANSA**, an acquisition function for selecting optimal noise seeds in T2V diffusion. Unlike classifiers with explicit predictive distributions, diffusion models do not expose such outputs, so we estimate uncertainty in the attention space where text and visual tokens align during generation. We treat attention maps as stochastic predictions conditioned on the seed z, prompt c, and timestep t. BANSA measures both disagreement and confidence across multiple attention samples, providing a BALD-style uncertainty criterion tailored to the generative setting.

**Definition 1** (BANSA Score). Let  $\mathbf{z}$  be a noise seed, c a text prompt, and t a diffusion timestep. Let  $\mathbf{Q}(\mathbf{z}, c, t), \mathbf{K}(\mathbf{z}, c, t) \in \mathbb{R}^{N \times d}$  denote the query and key matrices from a denoising network  $\epsilon_{\theta}$ . The attention map is computed as:

$$A(\mathbf{z}, c, t) := \operatorname{Softmax} \left( \mathbf{Q}(\mathbf{z}, c, t) \mathbf{K}(\mathbf{z}, c, t)^{\top} \right) \in \mathbb{R}^{N \times N}. \tag{3}$$

Let  $A(\mathbf{z}, c, t) = \{A^{(1)}, \dots, A^{(K)}\}$  denote a set of K stochastic attention maps obtained via forward passes with random perturbations (e.g., Bernoulli masking). The **BANSA** score is defined as:

$$\mathit{BANSA}(\mathbf{z}, c, t) := \mathcal{H}\left(\frac{1}{K} \sum_{k=1}^{K} A^{(k)}\right) - \frac{1}{K} \sum_{k=1}^{K} \mathcal{H}(A^{(k)}), \tag{4}$$

where  $\mathcal{H}(A) := \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{N} -A_{ij} \log A_{ij}$ . This formulation captures both the sharpness (confidence) and the consistency (agreement) of attention behavior. BANSA can be applied to various attention types (e.g., cross-, self-, or temporal) and allows layer-wise interpretability.

Given a noise pool  $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_M\}$ , we select the optimal noise seed that minimizes score:

$$\mathbf{z}^* := \arg\min_{\mathbf{z} \in \mathcal{Z}} \text{BANSA}(\mathbf{z}, c, t).$$
 (5)

A desirable property of BANSA is that its score becomes zero when all attention samples are identical, reflecting complete agreement and certainty. We formalize this as follows:

## Algorithm 1: Active Noise Selection with BANSA Score for Video Generation

**Input:** Text prompt c, noise pool  $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_M\}$ , timestep t, cutoff layer  $d^*$ 

1 foreach  $\mathbf{z}_i \in \mathcal{Z}$  do

- Compute BANSA score: BANSA- $E_{< d^*}(\mathbf{z}_i, c, t)$  via Eq. (7);
- 3 Select optimal noise:  $\mathbf{z}^* = \arg\min_{\mathbf{z}_i} \widehat{BANSA} \cdot \widehat{E}_{\leq d^*}(\mathbf{z}_i, c, t);$
- **return** Generate video:  $\hat{v} = SampleVideo(\mathbf{z}^*, c, t)$ ;

**Proposition 1** (BANSA Zero Condition). Let  $A(\mathbf{z}, c, t) = \{A^{(1)}, \dots, A^{(K)}\}$  be a set of row-stochastic attention maps. Then:

$$BANSA(\mathbf{z}, c, t) = 0 \quad \Leftrightarrow \quad A^{(1)} = \dots = A^{(K)}.$$

The proof is deferred to the Appendix B. This result implies that minimizing the BANSA score encourages attention that is both confident and consistent under stochastic perturbations. Empirically, low BANSA correlates with better prompt alignment, temporal coherence, and visual fidelity. Thus, BANSA provides a principled, model-aware criterion for noise selection in T2V.

#### 3.2 STOCHASTIC APPROXIMATION OF BANSA VIA BERNOULLI-MASKED ATTENTION

While BANSA is principled for noise selection, computing it with K independent passes per seed  $\mathbf{z}$  is costly. We instead draw K stochastic attention samples from a single pass via Bernoulli-masked attention. Rather than running K passes, we inject stochasticity directly into the attention scores. For each  $k=1,\ldots,K$ , sample a binary mask  $m_k \in \{0,1\}^{N\times N}$  i.i.d. from Bernoulli(p) and compute masked attention map,  $\hat{A}^{(k)}(\mathbf{z},c,t):=A(\mathbf{z},c,t)\odot m_k$ ,  $\odot$  denotes element-wise multiplication and rows are re-normalized. Using these K such samples, we define the approximate BANSA:

$$\text{BANSA-E}(\mathbf{z}, c, t) := \mathcal{H}\left(\frac{1}{K} \sum_{k=1}^{K} \hat{A}^{(k)}(\mathbf{z}, c, t)\right) - \frac{1}{K} \sum_{k=1}^{K} \mathcal{H}(\hat{A}^{(k)}(\mathbf{z}, c, t)). \tag{6}$$

By concavity of entropy, BANSA-E  $\geq 0$ . Although approximate, it efficiently captures attention-level epistemic uncertainty and serves as a practical surrogate for model-aware noise selection. <sup>1</sup> As shown in Table 4, experimental validation confirms that this method is sufficient for selecting optimal noisy samples from the model's perspective.

## 3.3 LAYER SELECTION VIA CUMULATIVE BANSA CORRELATION

BANSA can be computed at any layer, but behavior differs across depth. Using all layers is costly, so we truncate at the smallest depth  $d^*$  where the average BANSA up to d remains highly correlated with the full-layer score. Given a noise seed  $\mathbf{z}_i \in \mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_M\}$  and L attention layers, we compute per-layer scores BANSA- $\mathbf{E}^{(l)}(\mathbf{z}_i, c, t)$  and define the cumulative average up to layer d as:

$$\widehat{\text{BANSA-E}}_{\leq d}(\mathbf{z}_i, c, t) := \frac{1}{d} \sum_{l=1}^{d} \text{BANSA-E}^{(l)}(\mathbf{z}_i, c, t). \tag{7}$$

To find  $d^*$ , we compute the Pearson correlation (Pearson, 1895) between the partial average  $\widehat{BANSA}$ - $E_{\leq d}$  and the full-layer average  $\widehat{BANSA}$ - $E_{\leq L}$ , and choose the smallest d such that  $\widehat{Corr}\left(\widehat{BANSA}$ - $E_{\leq d}$ ,  $\widehat{BANSA}$ - $E_{\leq L}\right) \geq \tau$ . We validate this procedure using 100 prompts and 10 noise seeds across all evaluated models as detailed in Appendix D. As shown in Figure 6, the correlation quickly stabilizes at a moderate depth, allowing us to select  $d^*$  without using all layers. We then define the BANSA score as  $\widehat{BANSA}$ - $E_{\leq d^*}$  to guide noise selection, as summarized in Algorithm 1. With  $d^*$  fixed per model, it adds no runtime cost, and  $\widehat{BANSA}$ - $E_{\leq d^*}$  (truncated layer) consistently matches the full-layer score and robustly preserves generation quality across models, confirming its reliability as a practical surrogate as shown Appendix D and E (see detail).

<sup>&</sup>lt;sup>1</sup>We use "Bayesian" in the sense of epistemic uncertainty estimation, following BALD (Houlsby et al., 2011), rather than a full posterior derivation.

271 272

273

274

275

276 277

278

279

280

281

282

283

284 285

286

287

288

289

290

291

292

293

294 295

296

297

298

299

300

301

302

303

304

305

306 307

308

309

310

311

312

313 314

315 316

317

318

319

320

321

322

323

# ANALYSIS OF BANSA-SELECTED NOISE

In this section, we aim to investigate the property of noise seeds selected by BANSA, and consist of three experiments which declare our framework and interpret the noise seeds. Our framework does not seek a universal "golden" seed, but selects the most suitable one for each prompt by minimizing uncertainty. The effectiveness of a noise seed is therefore *prompt-dependent*: a seed that improves one prompt may degrade another.

**Experiment 1: Cross-Prompt Behavior.** To Table 1: Cross-prompt evaluation of a highdemonstrate this prompt-dependent behavior, We BANSA seed from A. select a high-BANSA (worst-performing) seed from prompt A and applied it to prompts B and C (Table 1) and measures Subject and Background Consistency. The seed degraded quality for prompt C but improved performance for prompt B, con-

From	То	Subject.	BackGround
A	В	$0.8491 \rightarrow 0.8703$	$0.9157 \rightarrow 0.9538$
A	C	$0.7683 { ightarrow} 0.7213$	$0.9368 { ightarrow} 0.9169$

firming that noise effectiveness depends on prompt context rather than intrinsic seed quality.

While Experiment 1 shows how noise effectiveness varies across prompts, the following two experiments examine whether similar patterns also emerge within a *fixed* prompt.

**tency.** To examine whether BANSA uncertainty remaps. lates to attention stability, we generate five samples per prompt (total 100 prompts), selected the highestand lowest-scoring seeds, and computed pairwise Euclidean distances between their attention maps (Ta-

Experiment 2: Intra-Prompt Attention Consis- Table 2: Pairwise distances within attention

Type	High (intra)	Low (intra)	Cross
Euclid.	1.635	1.567	1.735

ble 2). Low-BANSA seeds showed smaller intra-group distances, indicating more stable and consistent patterns, while high-BANSA seeds exhibited higher variability.

**Experiment 3: Latent Trajectory and Expressiveness.** Table 3: Latent stability and visual dy-Beyond attention maps, we further investigate latent-space namics. dynamics. Across 100 prompts with five repetitions each, we compare high- and low-BANSA seeds over all denoising steps (Table 3). We measured two metrics: (1) Latent Trajectory Variation, obtained by applying a Butterworth low-pass filter to isolate low-frequency structural compo-

Group	Traj. Var. ↓	Var. ↑
High BANSA	52.34	0.041
Low BANSA	51.07	0.053

nents (Wu et al., 2024; Yuan et al., 2025) and computing their temporal variation, where lower values indicate smoother and more stable trajectories; and (2) Intra-frame Variance, the average spatial variance within frames, where higher values indicate richer dynamic expressiveness. Low-BANSA seeds showed both lower trajectory variation and higher intra-frame variance, combining temporal stability with more expressive video generation.

**Takeaway and Relation to Prior Work.** The three findings indicate that BANSA identifies seeds with lower uncertainty, leading to (1) prompt-specific improvements, (2) more stable attention, and (3) smoother yet more expressive latent dynamics. This is consistent with FreeInit Wu et al. (2024) and FreeQPrior (Yuan et al., 2025), which emphasize stable latent initialization and low-frequency structure. BANSA complements these insights by offering an uncertainty-driven selector that explains why such seeds generalize within a prompt while remaining prompt-dependent across prompts.

# **EXPERIMENTS**

Experimental Setting. We evaluate ANSE on diverse T2V diffusion backbones—AnimateDiff (Guo et al., 2024b), CogVideoX-2B/5B (Yang et al., 2024), Wan2.1 (Team, 2025), and Hunyuan-Video (Kong et al., 2024). To rigorously evaluate our method, we follow each model's official sampling protocol; for resolution, Wan2.1 is evaluated at 480p and HunyuanVideo at 360p due to resource limits. Results for noise-prior refinement baselines (FreqPrior (Yuan et al., 2025)) are reported only on AnimateDiff, as they lack official support on the others and incur  $\sim 3 \times$  inference time. ANSE is orthogonal and can be combined with these priors. Unless noted, we use a noise pool of M=10 with K=10 stochastic passes per seed, and Bernoulli-masked attention with p=0.2. All experiments run on NVIDIA H100 GPUs. Further details are in the Appendix C.

324 325

Table 4: Quantitative results on VBench using AnimateDiff, CogVideoX-2B and -5B.

330 331 332

341 342 343

344

345

346

352

353

358

359 360

361 362 364

366

367

368

369

370

371

372

373

374 375

376

377

**Backbone Model** Method **Quality Score** Semantic Score **Total Score Inference Time** 28.23s Vanilla 80.22 69.03 77.98 AnimateDiff 81.66 71.09 79.33 + Ours 31.33s(+10.98%) (Guo et al., 2024b) 81.22 70.45 79.07 Freqprior 58.01(+105.36%) 82.23 73.23 80.43 + Ours 61.12s<sub>(+5.36%)</sub> CogVideoX-2B Vanilla 82.08 81.03 247.8s 76.83 (Yang et al., 2024) 269.3s<sub>(+8.67%)</sub> + Ours 82.56 78.06 81.66 CogVideoX-5B Vanilla 82.53 77.50 81.52 667.3s (Yang et al., 2024) + Ours 82.70 78.10 81.71 754.1s<sub>(+13.1%)</sub>

Table 5: Quantitative results on quality score metrics for HunyuanVideo and Wan2.1

Backbone Model	Method	Subject Consistency ↑	Background Consistency ↑	Motion Smoothness ↑	Aesthetic Quality ↑	Imaging Quality ↑	Temporal Flickering ↑	Inference Time ↓
HunyuanVideo	Vanilla	0.9562	0.9656	0.9850	0.6276	0.6137	0.9921	52.3s
(Kong et al., 2024)	+ Ours	0.9612	0.9661	0.9858	0.6268	0.6151	0.9938	59.8s <sub>(+14.34%)</sub>
Wan2.1	Vanilla	0.9562	0.9656	0.9850	0.6276	0.6137	0.9943	43.8s
(Team, 2025)	+ Ours	0.9612	0.9661	0.9858	0.6268	0.6151	0.9956	50.7s <sub>(+15.75%)</sub>

**Evaluation Metric.** We evaluate with VBench (Huang et al., 2024), which reports a quality score and a semantic score, combined into a total score on a 0-100 scale. For AnimateDiff and CogVideoX-2B/5B, we report both quality and semantic scores. For HunyuanVideo and Wan2.1, we restrict evaluation to the six quality dimensions due to computational constraints, while noting that semantic evaluation can also be applied in principle. All reported results are obtained by repeating each evaluation dimension five times and averaging the scores.

Quantitative Comparison. As shown in Table 4, ANSE consistently improves performance across diverse T2V backbones. On AnimateDiff, ANSE outperforms the vanilla baseline and further demonstrates clear advantages over noise-initialization methods such as FreqPrior (Yuan et al., 2025). Notably, ANSE is also fully compatible with FreqPrior, achieving even higher scores when combined. Table 4 also reports results on CogVideoX-2B and -5B, which adopt the more advanced MMDiT architecture. ANSE improves quality, semantic, and total VBench scores on both scales, demonstrating robust generalization across architectures and model sizes.

Table 5 presents additional results on recent state-of-the-art models, HunyuanVideo and Wan2.1. ANSE achieves consistent improvements across quality metrics, underscoring its plug-and-play applicability to a wide spectrum of video diffusion models—from U-Net to MMDiT, and from mid-scale to the latest state-of-the-art systems.

Qualitative Comparison. Figure 4 presents representative qualitative comparisons on CogVideoX-2B, CogVideoX-5B, and Wan2.1 with and without ANSE. Our approach consistently enhances semantic fidelity, motion realism, and visual clarity across diverse prompts. For example, in separate prompts such as "exploding" and "descends gracefully", ANSE captures critical semantic transitions—generating visible explosions in the former and preserving smooth temporal continuity in the latter. In other prompts like "koala playing the piano" and "tasting beer", it produces anatomically coherent bodies with natural, expressive motion. These examples highlight ANSE's ability to improve spatial-temporal fidelity while generalizing effectively to large-scale video diffusion models. Additional qualitative results, including other backbones, are provided in Appendix H.

Computational Cost. As shown in Table 4 and Table 5, ANSE introduces only a minimal increase in inference time while delivering consistent quality gains. On AnimateDiff, inference time increases by just +10.98%, compared to more than +105% when combined with FreqPrior and over +200% for FreeInit. On CogVideoX models, the overhead is similarly small: +8% on the 2B variant and +13% on the 5B variant. For more recent architectures, ANSE adds only +14% on HunyuanVideo and +15% on Wan2.1.

The overhead comes only from seed evaluation and leaves the sampling process and memory usage unchanged. In contrast, FreeInit and FreqPrior require multiple full passes, causing large slowdowns. ANSE improves quality across diverse backbones while keeping inference overhead below +15%, making it a practical plug-and-play solution for video diffusion.

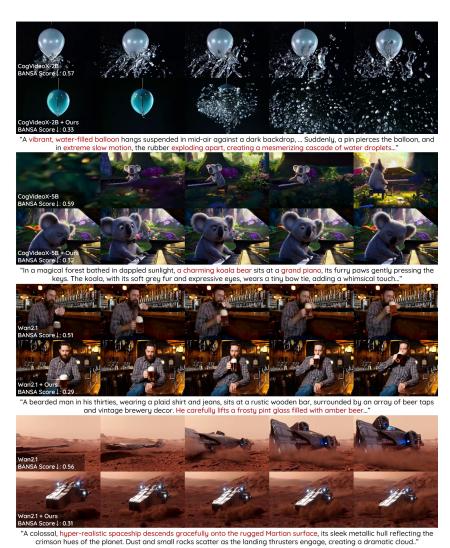


Figure 4: **Qualitative comparison with and without ANSE.** Results from CogvideoX-2B, 5B and Wan2.1. With ANSE, videos exhibit improved visual quality, better text alignment, and smoother motion transitions compared to the baseline.

# 6 ABLATION STUDY

Comparison of Acquisition Functions. Using CogVideoX-2B, we compare BANSA with random sampling, entropy-based selection, and two variants: BANSA (B) with Bernoulli masking and BANSA (D) with dropout-based stochasticity (Table 6). All methods improve over the baseline, but BANSA (B) consistently achieves the best scores across quality, semantic, and total metrics. This suggests that Bernoulli masking better captures attention-level uncertainty than dropout, underscoring the importance of modeling uncertainty in line with the model's structure.

**Effect of Ensemble Size** K. We investigate how the number of stochastic forward passes K influences subject and background consistency (Table 7). Both metrics improve as K increases from 1 to 10, indicating that larger ensembles provide more reliable noise evaluation. Performance saturates at K=10, which we set as the default in all experiments.

**Effect of Noise Pool Size** M. We analyze the effect of noise pool size M, which controls the diversity of candidate seeds in BANSA. Larger M improves the chance of finding high-quality seeds but increases inference cost. As shown in Figure 7 (Appendix), performance saturates around M=10 for CogVideoX-2B and -5B, which we adopt as the default for each model.

Table 6: Comparison of different acquisition functions for noise selection.

Method	Quality Score	Semantic Score	Total Score
Random	82.08	76.83	81.03
Entropy	82.23	76.73	81.13
BANSA (D)	82.43	76.91	81.33
BANSA (B)	82.56	78.06	81.66

Table 7: Effect of varying the number of K.

K	Subject Consistency	Background Consistnecy
1	0.9618	0.9788
3	0.9623	0.9793
5	0.9632	0.9798
7	0.9638	0.9802
10	0.9641	0.9811

Table 8: Quantitative comparison of reversed BANSA scoring on CogVideoX-2B. This presents results when selecting samples using the highest BANSA scores, compared to the default selection.

Method	Subject Consistency	Background Consistency	Temporal Flickering	Motion Smoothness	Aesthetic Quality	Imaging Quality	Dynamic Degree	Quality Score
Vanilla	0.9616	0.9788	0.9715	0.9743	0.6195	0.6267	0.6380	82.08
+ Ours (reverse)	0.9626	0.9785	0.9700	0.9741	0.6181	0.6253	0.6328	81.93
+ Ours	0.9641	0.9811	0.9775	0.9746	0.6202	0.6276	0.6511	82.56



A majestic chestnut horse with a flowing mane stands at the edge of a crystal-clear river, ... The sunlight filters through the trees, casting a golden glow on the scene. The horse gracefully bends its neck, its reflection shimmering in the gentle ripples of the water.

Figure 5: **Failure case and limitation of our method.** Although the BANSA score indicates low uncertainty, the resulting video still contains unnatural content. This represent a limitation of ours: we select optimal seeds but do not alter the generation process itself.

**Reversing the BANSA Criterion.** To further validate BANSA, we conduct a control experiment where the noise seed with the *highest* BANSA score is selected—i.e., choosing the seed associated with the greatest model uncertainty. As shown in Table 8, this reversal results in degradation of quality-related metrics, confirming that lower BANSA scores are predictive of perceptually stronger generations and supporting the validity of our selection strategy.

#### 7 DISCUSSION AND LIMITATIONS

Our approach focuses on noise seed selection via model uncertainty estimation, but has certain limitations. As shown in Figure 5, even low-BANSA seeds—reflecting high model confidence—can still yield unnatural generations. This indicates that while ANSE effectively identifies promising seeds, it does not directly influence the generation process. Another limitation is the gap between estimated uncertainty and perceptual quality. BANSA captures attention-level uncertainty, but may not fully reflect semantic or aesthetic aspects. Ideally, one could generate multiple candidates per seed and select based on strong quality metrics, but this is computationally expensive. We therefore view BANSA as a practical surrogate for such strategies. Future work could explore integrating BANSA with information-theoretic refinement or active learning to enhance performance.

## 8 CONCLUSION

We present ANSE, a framework for active noise selection in video diffusion models. At its core is BANSA, an acquisition function that leverages attention-derived uncertainty to identify noise seeds yielding confident, consistent attention and thus higher-quality generations. BANSA adapts the BALD principle to the generative setting by operating in attention space, with efficient deployment enabled through Bernoulli-masked attention and lightweight layer selection. Experiments across multiple T2V backbones show that ANSE improves video quality and prompt alignment with minimal inference overhead. This introduces an inference-time scaling paradigm, enhancing generation not by altering the model or sampling steps, but through informed seed selection.

## REFERENCES

- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22563–22575, 2023.
- Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. arXiv preprint arXiv:2407.21787, 2024.
- Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter1: Open diffusion models for high-quality video generation, 2023.
- Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7310–7320, 2024a.
- Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-σ: Weak-to-strong training of diffusion transformer for 4k text-to-image generation. In *European Conference on Computer Vision*, pp. 74–91. Springer, 2024b.
- Sherry X Chen, Yaron Vaxman, Elad Ben Baruch, David Asulin, Aviad Moreshet, Kuo-Chin Lien, Misha Sra, and Pradeep Sen. Tino-edit: Timestep and noise optimization for robust diffusion-based image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6337–6346, 2024c.
- Bradley Efron. Tweedie's formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno: Enhancing one-step text-to-image models through reward-based noise optimization. *Advances in Neural Information Processing Systems*, 37:125487–125519, 2024.
- Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International conference on machine learning*, pp. 1183–1192. PMLR, 2017.
- Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs, Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve your own correlation: A noise prior for video diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 22930–22941, 2023.
- Xiefan Guo, Jinlin Liu, Miaomiao Cui, Jiankai Li, Hongyu Yang, and Di Huang. Initno: Boosting text-to-image diffusion models via initial noise optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9380–9389, 2024a.
- Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *International Conference on Learning Representations*, 2024b.
- Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, et al. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024.
- Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity long video generation. 2022.

- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
  - Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.
  - Neil Houlsby, Ferenc Huszár, Zoubin Ghahramani, and Máté Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
  - Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. VBench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
  - Kwanyoung Kim and Jong Chul Ye. Noise2score: Tweedie's approach to self-supervised image denoising without clean images. *Advances in Neural Information Processing Systems*, 34:864–874, 2021.
  - Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational adversarial active learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8166–8175, 2021.
  - Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *Advances in neural information processing systems*, 32, 2019.
  - Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024.
  - Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang, Yandong Li, Tommi Jaakkola, Xuhui Jia, et al. Inference-time scaling for diffusion models beyond scaling denoising steps. *arXiv preprint arXiv:2501.09732*, 2025.
  - Oisin Mac Aodha, Neill DF Campbell, Jan Kautz, and Gabriel J Brostow. Hierarchical subquery evaluation for active learning on a graph. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 564–571, 2014.
  - Karl Pearson. Notes on regression and inheritance in the case of two parents. *Proceedings of the Royal Society of London*, pp. 337–344, 1895.
  - Haonan Qiu, Menghan Xia, Yong Zhang, Yingqing He, Xintao Wang, Ying Shan, and Ziwei Liu. Freenoise: Tuning-free longer video diffusion via noise rescheduling. In *The Twelfth International Conference on Learning Representations*.
  - Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
  - Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*, 2017.
  - Claude E Shannon. A mathematical theory of communication. *Bell system technical journal*, 27(3): 379–423, 1948.
  - Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 5972–5981, 2019.
  - Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv* preprint arXiv:2408.03314, 2024.
  - Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021a.

- Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021b.
- Wan Team. Wan: Open and advanced large-scale video generative models. 2025.
  - Toan Tran, Thanh-Toan Do, Ian Reid, and Gustavo Carneiro. Bayesian generative active deep learning. In *International conference on machine learning*, pp. 6295–6304. PMLR, 2019.
  - Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
  - Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
  - Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *International Journal of Computer Vision*, 133(5):3059–3078, 2025.
  - Tianxing Wu, Chenyang Si, Yuming Jiang, Ziqi Huang, and Ziwei Liu. Freeinit: Bridging initialization gap in video diffusion models. In *European Conference on Computer Vision*, pp. 378–394. Springer, 2024.
  - Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li, Ligeng Zhu, Yao Lu, et al. Sana: Efficient high-resolution image synthesis with linear diffusion transformers. *arXiv preprint arXiv:2410.10629*, 2024.
  - Katherine Xu, Lingzhi Zhang, and Jianbo Shi. Good seed makes a good crop: Discovering secret seeds in text-to-image diffusion models. In 2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pp. 3024–3034. IEEE, 2025.
  - Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and Alexander G Hauptmann. Multi-class active learning by uncertainty sampling with diversity maximization. *International Journal of Computer Vision*, 113:113–127, 2015.
  - Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
  - Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 93–102, 2019.
  - Yunlong Yuan, Yuanfan Guo, Chunwei Wang, Wei Zhang, Hang Xu, and Li Zhang. Freqprior: Improving video diffusion models with frequency filtering gaussian noise. *arXiv preprint arXiv:2502.03496*, 2025.
  - Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. *arXiv* preprint arXiv:2412.20404, 2024.

# A SUPPLEMENTARY SECTION

 In this supplementary document, we present the following:

- Proof of Proposition 1 from the main paper regarding the BANSA Zero condition in Section B.
- Implementation details of ANSE and evaluation metrics in Section C.
- Further explanation of layer selection through cumulative BANSA correlation in Section D.
- Additional ablation studies, including full-layer BANSA score analysis (Section E), temporal scope effects (Section F), and attention-space effects (Section G).
- Additional qualitative results demonstrating the impact of BANSA Score in Section H.

# B PROOF OF PROPOSITION 1

**Proposition 1** (BANSA Zero Condition). Let  $\mathcal{A}(\mathbf{z}, c, t) = \{A^{(1)}, \dots, A^{(K)}\}$  be a set of row-stochastic attention maps. Then:

$$BANSA(\mathbf{z}, c, t) = 0 \quad \Leftrightarrow \quad A^{(1)} = \dots = A^{(K)}.$$

*Proof.* BANSA is defined as the difference between the average entropy and the entropy of the average:

$$\mathrm{BANSA}(\mathbf{z}, c, t) = \mathcal{H}\left(\frac{1}{K}\sum_{k=1}^K A^{(k)}(\mathbf{z}, c, t)\right) - \frac{1}{K}\sum_{k=1}^K \mathcal{H}(A^{(k)}(\mathbf{z}, c, t)).$$

Since the Shannon entropy  $\mathcal{H}(\cdot)$  is strictly concave over the probability simplex. Therefore, by Jensen's inequality:

$$\mathcal{H}\left(\frac{1}{K}\sum_{k=1}^{K}A^{(k)}\right) \ge \frac{1}{K}\sum_{k=1}^{K}\mathcal{H}(A^{(k)}),$$

with equality if and only if  $A^{(1)} = \cdots = A^{(K)}$ . Thus, BANSA( $\mathbf{z}, c, t$ ) = 0 if and only if all attention maps are identical.

**Remark 1.** (Interpretation) This result confirms that the BANSA score quantifies disagreement among sampled attention maps. A BANSA score of zero occurs only when all stochastic attention realizations collapse to a single deterministic map—i.e., the model exhibits **no epistemic uncertainty** in its attention distribution. Higher BANSA values indicate greater variation across samples, and thus, higher uncertainty. In this sense, BANSA acts as a Jensen—Shannon-type divergence over attention maps, capturing their dispersion under stochastic masking.

## C FURTHER DETAILS ON EVALUATION METRICS AND IMPLEMENTATION

**Evaluation Metrics** To evaluate performance on Vbench, we use the Vbench-long version, where prompts are augmented using GPT-40 across all evaluation dimensions. This version is specifically designed for assessing videos longer than 4 seconds.

We rigorously evaluate our generated videos following the official evaluation protocol. The Quality Score is a weighted average of the following aspects: subject consistency, background consistency, temporal flickering, motion smoothness, aesthetic quality, imaging quality, and dynamic degree.

The Semantic Score is a weighted average of the following semantic dimensions: object class, multiple objects, human action, color, spatial relationship, scene, appearance style, temporal style, and overall consistency.

The Total Score is then computed as a weighted combination of the Quality Score and Semantic Score:

$$\text{Total Score} = \frac{w_1}{w_1 + w_2} \cdot \text{Quality Score} + \frac{w_2}{w_1 + w_2} \cdot \text{Semantic Score}$$

where  $w_1 = 4$  and  $w_2 = 1$ , following the default setting in the official implementation.

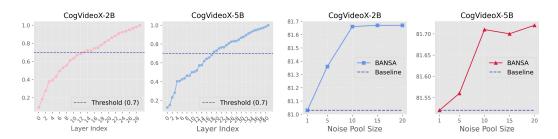


Figure 6: Correlation analysis between cu- Figure 7: Ablation study on noise pool size M. mulative BANSA score and full-layer scores. We evaluate total scores across three text-to-video The 0.7 threshold is reached around layer 14 for diffusion models with varying M, and select suit-CogVideoX-2B and layer 19 for CogVideoX-5B. able values based on computational cost.

**Implementation** We compute our BANSA score using the BALD-style formulation, which yields non-negative values. For clearer visualization, we normalize the BANSA scores from their original range (minimum: 0.45, maximum: 0.60) to the [0, 1] interval. This normalization is used solely for visual clarity in figures and plots, and does not affect the noise selection process, which operates on the raw BANSA scores.

## D FURTHER DETAIL OF BANSA LAYER-WISE CORRELATION ANALYSIS

**Prompt construction.** We evenly sampled 100 prompts from the four official VBench categories: *Subject Consistency, Overall Consistency, Temporal Flickering*, and *Scene*. Each category contains 25 prompts, selected to ensure diversity in motion, structure, and semantics.

Below are representative examples:

- Subject Consistency (e.g., "A young man with long, flowing hair sits on a rustic wooden stool in a cozy room, strumming an acoustic guitar...")
- Overall Consistency (e.g., "A mesmerizing splash of turquoise water erupts in extreme slow motion, each droplet suspended in mid-air...")
- *Temporal Flickering* (e.g., "A cozy restaurant with flickering candles and soft music. Patrons dine peacefully as snow falls outside...")
- *Scene* (e.g., "A university campus transitions from lively student life to a golden sunset behind the clock tower...")

Prompt sampling was stratified to ensure coverage of diverse visual and temporal patterns. The full list of prompts will be made publicly available upon code release.

**BANSA score computation and correlation analysis.** For each prompt, we generated 10 videos using different random noise seeds and computed BANSA scores at each attention layer. This yielded one full-layer BANSA score and a set of layer-wise scores per seed. To obtain stable estimates, we averaged the per-layer and full-layer BANSA scores across the 10 seeds, reducing noise-specific variance and capturing consistent uncertainty patterns.

We then computed Pearson correlations between the cumulative BANSA scores (summed from layer 1 to d) and the official quality scores. The optimal depth  $d^*$  was defined as the smallest d at which the correlation exceeded 0.7, a widely used threshold indicating a strong positive relationship. As shown in Figure 6, this point was reached at  $d^* = 14$  for CogVideoX-2B and  $d^* = 19$  for CogVideoX-5B. For AnimateDiff, the correlation crossed the threshold earlier at  $d^* = 10$ , while for HunyuanVideo and Wanx 2.1 it consistently appeared around  $d^* = 20$ . These model-specific depths were applied throughout all experiments. Such consistency across architectures empirically supports our choice of the 0.7 threshold and ensures that the correlation analysis reflects generalizable, noise-agnostic patterns of attention-based uncertainty.

Table 9: Comparison between full-layer and truncated BANSA score.

Backbone Model	Method	Subject Consistency ↑	Background Consistency ↑	Temporal Flickering ↑	Motion Smoothness ↑	Aesthetic Quality ↑	Imaging Quality ↑	Dynamic Degree ↑	Quality Score ↑	Inference Time ↓
CogvideoX-2B	Full-layer Truncated	0.9639 <b>0.9641</b>	0.9810 <b>0.9811</b>	<b>0.9801</b> 0.9775	0.9743 <b>0.9746</b>	0.6198 <b>0.6202</b>	0.6244 <b>0.6276</b>	<b>0.6516</b> 0.6511	<b>82.58</b> 82.56	303.7 <b>269.3</b>
CogvideoX-5B	Full-layer Truncated	<b>0.9660</b> 0.9658	0.9630 <u>0</u> .9639	<b>0.9863</b> 0.9861	0.9708 <b>0.9711</b>	0.6168 <b>0.6179</b>	0.6290 <b>0.6290</b>	<b>0.6979</b> 0.6918	<b>82.71</b> 82.70	810.3 <b>754.3</b>

# E EFFECTIVENESS OF TRUNCATED BANSA SCORE

To reduce the computational overhead of BANSA evaluation, we adopt a truncated score that aggregates attention uncertainty only up to a fixed depth  $d^*$ , rather than summing over all layers. To evaluate the effectiveness of this approximation, we compared the final generation quality when selecting noise seeds using either the full-layer or truncated BANSA scores.

As shown in Table 9, both approaches yield highly similar results across all seven dimensions of the VBench evaluation protocol (subject consistency, background consistency, aesthetic quality, imaging quality, motion smoothness, dynamic degree, and temporal flickering). Importantly, the overall quality scores are preserved despite the substantial reduction in attention layers used.

This demonstrates that truncated BANSA is sufficient to capture the key uncertainty signals for reliable noise selection while reducing inference time. The strong alignment in quality stems from the fact that our method relies on relative ranking rather than absolute values, allowing for efficient yet robust selection with significantly lower computational cost. We attribute this effectiveness to the fact that most informative attention behaviors emerge early in the denoising process, allowing accurate uncertainty estimation without full-layer computation.

Table 10: Effect of temporal scope in BANSA score on generation quality.

BANSA Scope	Subject Consistency ↑	Temporal Flickering ↑	Motion Smoothness ↑	Aesthetic Quality ↑	Imaging Quality ↑	Dynamic Degree ↑	Inference Time ↓
1-step	0.9639	0.9801	0.9743	0.6198	0.6244	0.6516	× 1
25-step avg	0.9651	0.9798	0.9746	0.6202	0.6271	0.6511	$\times$ 25
50-step avg	0.9652	0.9799	0.9751	0.6203	0.6276	0.6514	$\times$ 50

## F EFFECT OF TEMPORAL SCOPE IN BANSA SCORE

While our method computes the BANSA score only at the first denoising step to minimize cost, it is natural to ask whether incorporating more timesteps improves its predictive power for noise selection. To investigate this, we compute the average BANSA score across the first 1, 25, and 50 denoising steps and compare their effectiveness in predicting video quality.

Table 10 reports the VBench scores for subject consistency, aesthetic quality, imaging quality, motion smoothness, dynamic degree, and temporal flickering when using BANSA computed over different temporal scopes. Although using more timesteps results in slightly better quality, the gains are marginal. This indicates that most of the predictive signal for noise quality is embedded early in the generation trajectory.

More importantly, since BANSA is used solely to assess the uncertainty of the initial noise seed—not to track full-step generation behavior—our 1-step computation is sufficient to capture the core uncertainty signal. In contrast, computing BANSA over all steps requires running multiple attention forward passes across the full trajectory, resulting in substantial computational overhead that limits its practicality for real-world applications.

## G EFFECT OF ATTENTION SPACE IN BANSA SCORE

To further examine whether BANSA scores capture quality- or semantic-related uncertainty, we conducted additional experiments with the AnimateDiff backbone, which separates self- and cross-

Table 11: Quantitative results on VBench using AnimateDiff, CogVideoX-2B and -5B.

Backbone Model	Method	Quality Score	Semantic Score	Total Score
AnimateDiff	Vanilla	80.22	69.03	77.98
(Guo et al., 2024b)	+ Ours (Self-Attention)	81.68	70.72	79.19
	+ Ours (Cross-Attention)	81.66	71.09	79.33

attention layers. Unlike CogVideo and Wanx, which adopt an MMDiT architecture where multimodal attention entangles both components, AnimateDiff provides a clearer lens for disentangling semantic effects. As shown in Table 11, applying BANSA to self-attention primarily improved perceptual quality scores, whereas applying it to cross-attention enhanced semantic alignment. This indicates that uncertainty over noise manifests differently depending on the attention type, supporting the interpretation that BANSA reflects both quality- and semantics-related variations. In contrast, MMDiT-based models such as CogVideo and Wanx integrate these dimensions within their unified attention structure, making them more suitable for overall noise selection in the main paper's experiments.

# H ADDITIONAL QUALITATIVE COMPARISON

More qualitative results. Figures 8, 9, 10 11, and 12 present additional examples generated using our noise selection framework including diverse backbone such as, AnimateDiff, CogVideoX2B and 5B, HunyangVideo, and Wan2.1. Across diverse prompts, the selected seeds yield improved spatial detail, aesthetic quality, and semantic alignment, further validating the robustness of our approach. These examples complement our quantitative findings by illustrating the visual impact of BANSA-based noise selection.

**Effect of BANSA score on generation quality.** Figures 13 provide a qualitative comparison of outputs generated using three types of noise seeds: a randomly sampled seed, the seed with the highest BANSA score (lowest quality), and the seed with the lowest BANSA score (highest quality). All videos were generated using 50 denoising steps with the CogVideoX-5B backbone. The lowest-BANSA seed consistently produces sharper, more coherent, and semantically faithful videos, whereas the highest-BANSA seed often leads to structural artifacts or temporal instability. These results highlight the practical value of BANSA-guided noise selection.

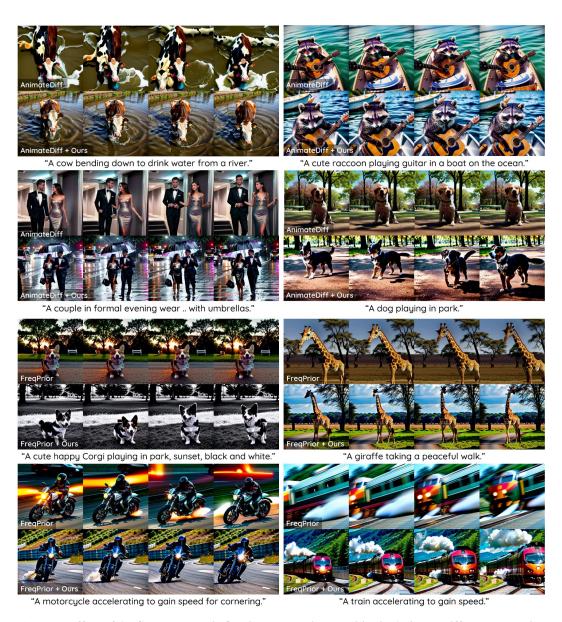


Figure 8: Effect of ANSE on semantic fidelity and motion stability in AnimateDiff and FreqPrior. Each block compares baseline generations with those using ANSE-selected noise. In addition, while FreqPrior serves as a noise-refinement baseline, our method is fully compatible and achieves further improvements when combined.



"In a charming Parisian café, an animated panda sits at a quaint wooden table, sipping coffee from a delicate porcelain cup.

The panda, wearing a stylish beret and a striped scarf, gazes out the window at the bustling Paris streets,...."



"A cool cat, sporting sleek black sunglasses and a red lifeguard vest, sits confidently on a high lifeguard chair overlooking a sparkling blue pool. The feline's fur is a mix of orange and white, and its tail flicks with authority."



"A young woman with flawless skin and a serene expression sits at a vanity, bathed in soft morning light. She begins by applying a light moisturizer, her fingers moving gently across her face. Next, she uses a foundation brush to blend a sheer.."



"An astronaut in a pristine white spacesuit, floats effortlessly against the vast, star-studded expanse of space. As the camera zooms out, the intricate details of the suit, including the life-support backpack and tether,"

Figure 9: Effect of ANSE on semantic fidelity and motion stability in CogVideoX outputs. Each block compares baseline generations with those using ANSE-selected noise. Across both CogVideoX-2B and 5B, ANSE improves semantic alignment to the prompt and reduces artifacts such as temporal flickering and object distortion.



"A sophisticated couple, dressed in elegant evening attire, walks down a dimly lit street, their formal. The man, in a tailored black tuxedo, and the woman, in a flowing red gown, share a black umbrella as the rain captures their synchronized steps.."



"A drone captures a breathtaking aerial view of a festive celebration in a snow-covered town square, centered around a towering, brilliantly lit Christmas tree adorned with twinkling lights and ornaments. The scene is alive with vibrant fireworks.."



"In a whimsical forest clearing, a raccoon with a mischievous glint in its eye stands on a tree stump, holding an electric guitar.

The raccoon, wearing a tiny leather jacket, strums the guitar with surprising skill, its tiny paws moving deftly over the strings..."



"A majestic elephant strolls gracefully through a lush, verdant forest, its massive feet gently pressing into the soft earth. The sunlight filters through the dense canopy, casting dappled shadows on its wrinkled, grey skin....."

Figure 10: Additional qualitative comparison of CogVideoX variants with and without ANSE. Results from CogVideoX-2B are shown in the first two rows; the rest show CogVideoX-5B. With ANSE, videos exhibit improved visual quality, better text alignment, and smoother motion transitions compared to the baseline.

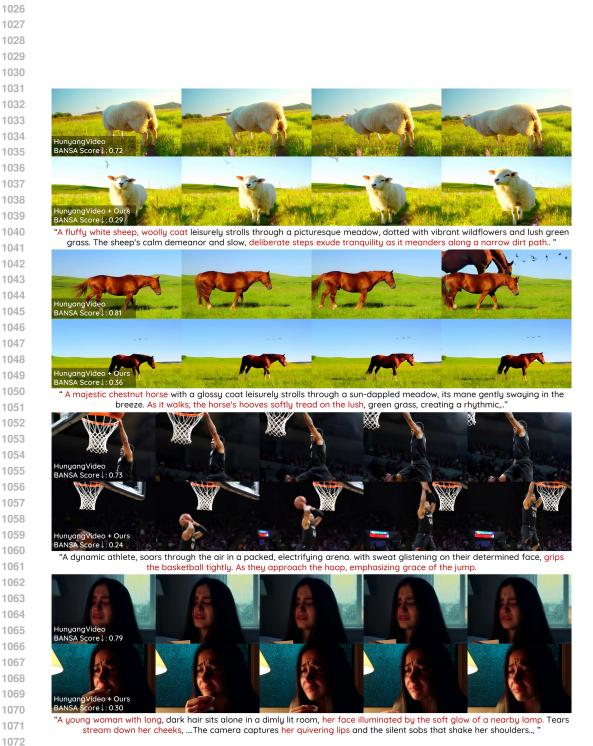


Figure 11: Additional qualitative comparison of HunyangVideo with and without ANSE. With ANSE, the generated videos achieve sharper visual fidelity, stronger alignment between text and content, and smoother motion progression than the baseline.

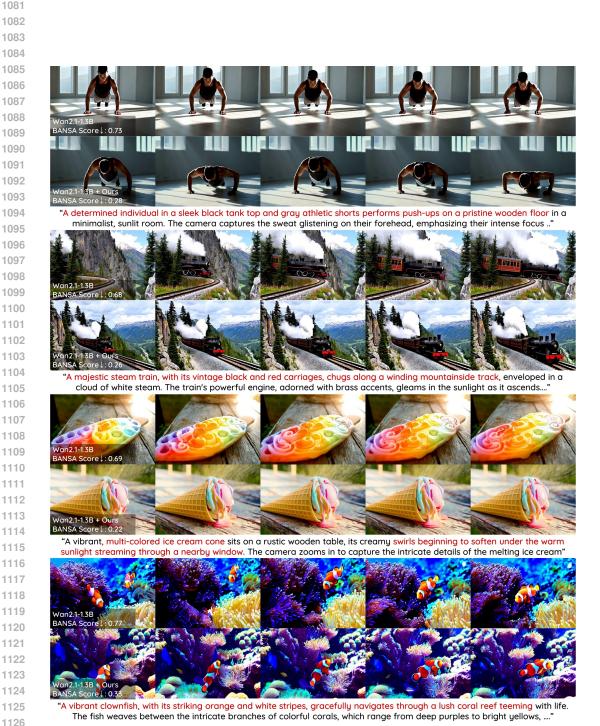


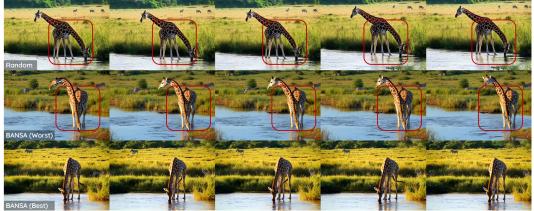
Figure 12: Additional qualitative comparison of Wan2.1 with and without ANSE. ANSE enables videos to deliver higher visual quality, more accurate adherence to the given text, and more seamless motion dynamics compared to the baseline.

Random

BANSA (Worst)

BANSA (Best)

"A lone bicycle, with its sleek frame and black tires, glides effortlessly through a vast, snow-covered field under a pale winter sky. The rider, bundled in a red parka, black gloves, and a woolen hat, pedals steadily, leaving a delicate trail in the pristine snow. The scene captures the quiet serenity of the landscape, with snowflakes gently falling and the distant silhouette of bare trees lining the horizon. As the rider continues, the sun



"A majestic giraffe, its long neck gracefully arching, bends down to drink from a serene river, surrounded by lush greenery and tall grasses. The sun casts a golden glow, highlighting the giraffe's patterned coat and the gentle ripples in the water. Nearby, a family of zebras grazes peacefully, adding to the tranquil scene. Birds flutter above, their reflections dancing on the water's surface. The giraffe's delicate movements create a sense of harmony with nature, as the river flows gently, reflecting the vibrant colors of the surrounding landscape..."

Figure 13: Qualitative comparison of generations from different noise seeds. We compare outputs generated from a randomly sampled seed (top), the seed with the highest BANSA score (middle), and the seed with the lowest score (bottom), using the same prompt and model. BANSA-selected seeds produce more coherent structure, stable motion, and stronger semantic alignment than both random and high-uncertainty seeds.