

ENTROPY-GUIDED k -GUARD SAMPLING FOR LONG-HORIZON AUTOREGRESSIVE VIDEO GENERATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Autoregressive (AR) architectures have achieved significant successes in LLM, inspiring explorations for video generation. In LLMs, top- p /top- k sampling strategies work exceptionally well: language tokens have high semantic density and low redundancy, so a fixed size of token candidates already strike a balance between semantic accuracy and generation diversity. In contrast, video tokens have low semantic density and high spatio-temporal redundancy. This mismatch makes static top- k /top- p strategies ineffective for video decoders: they either introduce unnecessary randomness for low-uncertainty regions (static backgrounds) or stuck in early errors for high-uncertainty regions (foreground objects). Prediction errors will accumulate as more frames are generated and eventually severely degrade long-horizon quality. To address this, we propose Entropy-Guided k -Guard (ENkG) sampling, a simple yet effective strategy that adapts sampling to token-wise dispersion, quantified by the entropy of each token’s predicted distribution. ENkG uses adaptive token candidate sizes: for low-entropy regions, it employs fewer candidates to suppress redundant noise and preserve structural integrity; for high-entropy regions, it uses more candidates to mitigate error compounding. ENkG is model-agnostic, training-free, and adds negligible overhead. Experiments demonstrate consistent improvements in perceptual quality and structural stability compared to static top- k /top- p strategies.

1 INTRODUCTION

The field of video-based world models has witnessed explosive growth in recent years, with significant advancements in generating high-fidelity, temporally coherent, and physically plausible video sequences Villegas et al. (2022); Wang et al. (2024a). These models aim to build an internal representation of the world’s dynamics, enabling applications from realistic simulation for robotics to advanced content creation He et al. (2025). This progress has paved the way for models that can not only synthesize video from text but also begin to understand and simulate interactive environments Mo et al. (2025).

Among the various architectural paradigms, autoregressive (AR) models have become a cornerstone for video generation. By factorizing the joint probability distribution of video frames into a product of conditional probabilities, AR models excel at capturing temporal causality and allow for fine-grained, frame-by-frame control during generation. This sequential approach is inherently flexible, supporting variable-length video generation and compatibility with scalable transformer architectures (Dosovitskiy et al., 2020; Weissenborn et al.). However, the sequential nature of AR models also introduces significant challenges: error accumulation (Parthipan et al., 2024; Bengio et al., 2015) and exposure bias (Schmidt, 2019). Minor inaccuracies or suboptimal choices in generating a single frame can propagate and amplify over time, leading to a degradation of quality, loss of coherence, and “drifting” from the intended content in longer video sequences (Hu et al., 2023).

Several strategies have been proposed to mitigate this effect. Huang et al. (2025) simulate inference during training by feeding the model its own previous predictions, allowing it to learn to correct mistakes. Other works introduce noisy or masked contexts, encouraging the model to be robust to imperfect inputs (Ren et al., 2025a; Zhou et al., 2025). While effective, these approaches often require modifications to the model architecture or additional training complexity, which may limit their applicability to existing large-scale video generation models.

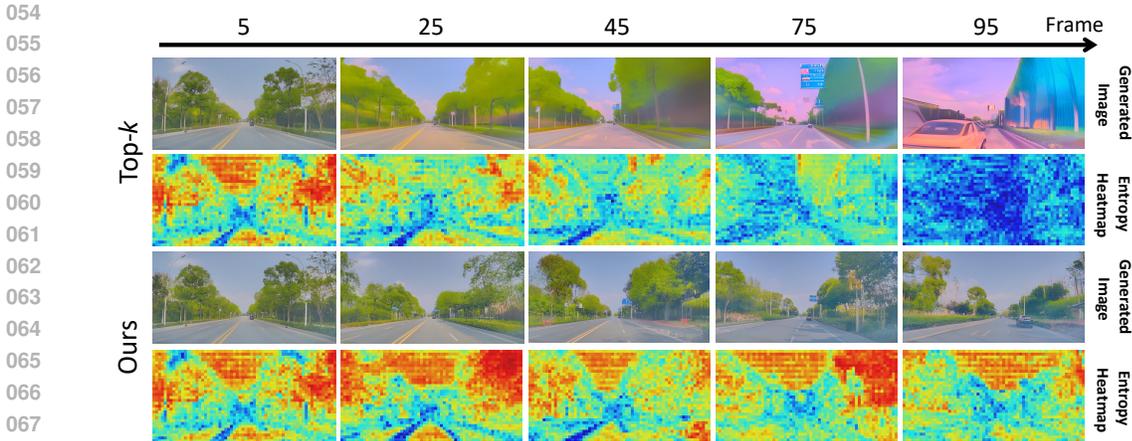


Figure 1: The results illustrate the phenomenon of entropy collapse in standard AR decoding, where blue regions indicate low entropy and red regions indicate high entropy. Our method effectively alleviates this issue.

In contrast, we focus on the often-overlooked role of the *sampling process* in autoregressive video generation. Our analysis reveals that conventional strategies such as fixed top- k or nucleus (top- p) sampling fail to account for the spatially structured uncertainty inherent in video tokens. Specifically, we observe that high-entropy regions, corresponding to complex textures like foliage or road markings, are prone to brittleness, whereas low-entropy regions representing structured geometry can suffer from overconfidence and texture wash-out. This motivates an *adaptive* sampling policy that modulates token diversity based on entropy, effectively balancing stability and richness in generated content.

Specifically, we introduce **Entropy-guided k -Guard sampling**, a model-agnostic algorithm that dynamically adjusts the size of candidates for each token according to its entropy. First, our method measure the entropy of each video token that indicates the dispersion of video token probability. For low-entropy regions, our strategy employs fewer candidates to suppress redundant noise and preserve structural integrity; for high-entropy regions, it uses more candidates to mitigate error compounding. Unlike previous solutions that modify training or rely on multiple candidate evaluations, our approach operates purely at the inference stage, making it widely applicable to existing autoregressive video models.

We validate our method on several state-of-the-art autoregressive video generation architectures, demonstrating that it significantly reduces error accumulation, preserves fine-grained textures, and improves temporal coherence over extended sequences. Quantitative metrics and qualitative results confirm that our adaptive sampling strategy enables longer, more realistic video generation without retraining or architectural changes. These findings suggest that carefully designed inference-time strategies can be a powerful tool for improving autoregressive video generation, complementing existing advances in model design and training.

In summary, our contributions are threefold: (i) we identify the limitations of fixed sampling strategies in autoregressive video generation and highlight the role of spatially structured uncertainty in error accumulation, (ii) we propose a simple yet effective entropy-guided adaptive sampling strategy with a k -guard mechanism, and (iii) we empirically demonstrate that this method improves long-sequence video quality across multiple benchmark models. This work highlights the potential of uncertainty-aware inference as a practical and generalizable solution for high-fidelity video synthesis.

2 RELATED WORK

2.1 VIDEO WORLD MODELS

Video-based world models aim to learn an internal representation of an environment, allowing the system to predict future states, simulate interactions, and support planning (Ha & Schmidhuber,

2018; Ding et al., 2024; Long et al., 2025; Zhang et al., 2025). Recent progress in large-scale video generation has enabled the creation of high-fidelity world simulators capable of producing visually realistic and physically plausible sequences (OpenAI, 2024), which are particularly valuable for applications such as autonomous driving and robotics (Wang et al., 2024b; Li et al., 2025).

A wide range of generative architectures have been explored for video synthesis. Diffusion models and Generative Adversarial Networks (GANs) have shown success in producing high-quality frame sequences (Ho & Jain, 2022; Clark & Fidler, 2019), while autoregressive (AR) models have emerged as a powerful alternative due to their inherent capacity to model temporal dependencies. By factorizing the joint distribution of video tokens or frames into a product of conditional probabilities, AR models generate coherent, long-duration videos either frame-by-frame (Gu et al., 2025) or token-by-token (Wu et al., 2024), enabling fine-grained temporal control and interactive generation.

Despite their strengths, AR models suffer from a fundamental limitation: *error accumulation*. During inference, each predicted frame or token is fed back as input for subsequent steps, so any imperfection can propagate and amplify over time (Yu & Chen, 2024; Feng & Li, 2021; Walker & Gupta, 2016). This leads to quality degradation, manifesting as flickering, unnatural motion, or drift from the intended scene (Yu & Chen, 2024; Saxena & Kumar, 2024; Benjamin & Smith, 2018; Kong et al., 2025). Error accumulation has been identified as a core challenge in theoretical analyses of autoregressive video generation, alongside issues such as memory bottlenecks (Saxena & Kumar, 2024; Goyal & Lee, 2022).

To address this limitation, our work proposes a *token-level adaptive sampling* strategy that dynamically modulates the sampling distribution according to the model’s predictive uncertainty. By integrating uncertainty into the generation process, we directly mitigate the compounding of errors, preserving both temporal coherence and visual fidelity in long-duration video sequences.

2.2 AUTOREGRESSIVE SAMPLING ALGORITHMS

Sampling strategies are central to autoregressive (AR) generation. *Greedy decoding* selects the most likely token at each step, but often yields low diversity. *Beam search* explores multiple hypotheses in parallel and improves likelihood-based metrics, yet typically reduces diversity. A widely used family of stochastic methods is *truncated probability sampling*, including top- k (Noarov et al., 2025) and nucleus (top- p) sampling (Ravfogel et al., 2023). Both restrict sampling to a subset of the distribution, balancing diversity and quality, but their hard truncation can occasionally introduce rare erroneous tokens, causing catastrophic errors in long sequences or generate duplicate and overconfident content with low threshold. *Best-of- N sampling* (Snell et al., 2024) generates multiple candidates and selects the best according to a reward model or predefined metric. While effective, it operates at a coarser granularity, is model-dependent, and incurs significant additional computation. Recently, entropy-based strategies have been explored in large language models (LLMs), where entropy guides model size switch (Simonds, 2025) or retrieval augmentation (Qiu et al., 2025). Similar entropy-guided temperature scaling methods have been proposed for LLMs (Zhang et al., 2024) and image generation (Ma et al., 2025) to balance diversity and fidelity. In contrast to these temperature-scaling approaches, our method leverages entropy to adaptively adjust the Top- p threshold with a k -guard mechanism. Crucially, we address the video-specific challenge of temporal error accumulation, rather than the static trade-offs focused on in text and image modalities.

3 MOTIVATION

We identify three critical findings in autoregressive (AR) video generation models. These observations explain why static sampling (effective for LLMs) fails for video and therefore motivate our dispersion-aware adaptive sampling strategy.

F1. Video token probability distributions are inherently flat. As shown in Figure ??, generated video tokens usually have quite small probability values, where dozens of tokens could achieve total probability. Tiny gaps between top candidates mean small logit perturbations (model noise, temporal drift) easily flip the argmax, breaking temporal coherence and accumulating visual artifacts. Truncated sampling eases brittleness but remains one-size-fits-all.

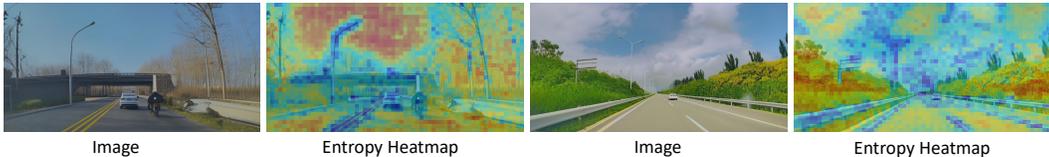


Figure 2: The visualization of entropy heatmaps. High-entropy regions form repeating textures (e.g., sky, foliage, and road), while low-entropy regions cluster in structured content and distinguishable textures (e.g., boundaries, edges between sky and trees, road markers and lines).

The root cause lies in fundamental differences between video and language tokens. Language tokens, with high semantic specificity (each mapping to a clear meaning, e.g., "car") and low redundancy (rarely repeating adjacent tokens). The language tokens have sharp probability distributions where top-1 probabilities often reach 0.7–0.8 (e.g., "street" in "I walk along the ___"). In contrast, video tokens lack direct semantic grounding and have high spatiotemporal redundancy. With no single token carrying unique meaning, their distributions are flat and diffuse, with an average top-1 probability of just 0.2.

F2. Video token probability dispersion is inherently tied to the semantic structure of the image. Given an AR model, the predicted distribution at the token i is P_i , and the entropy is as follows.

$$H_i = - \sum_j p_i(j) \log p_i(j) \quad (1)$$

where j is the vocabulary index.

Entropy is a tool to measure the model’s uncertainty Stolfo et al. (2024); Kang et al. (2025) and the token probability dispersions in predictions. Higher entropy indicates low confidence, thus showing greater dispersion in distributions, while lower entropy (high confidence) presents more concentrated distributions. Figure 1 demonstrates how the entropy heatmaps links the generation quality. High-dispersion (high-entropy) regions form repeating textures (e.g., sky, foliage, and road), where multiple tokens are equally plausible due to subtle texture variations. In contrast, low-dispersion (low-entropy) regions cluster in structured content and distinguishable textures (e.g., boundaries, edges between sky and trees, road markers and lines), where only a few tokens match the stable pixel patterns. However, existing static sampling (top-k or top-p) ignores this structure, forcing large candidate pools on low-dispersion regions (introducing redundant noise) and small pools on high-dispersion ones (discarding valid tokens), thereby exacerbating error accumulation.

F3. Entropy Collapse in Long-Horizon Autoregressive Video Generation. A third critical finding is that AR video models suffer from entropy collapse during long-horizon generation—an issue tied to evolving token dispersion patterns. As shown in Figure 1, this collapse manifests in two ways: temporally, the share of low-dispersion (low-entropy) tokens grows rapidly with each frame, driving down frame-averaged entropy; spatially, low-dispersion regions expand outward, gradually encroaching on high-dispersion areas, which erodes fine textures (e.g., foliage, road cracks) and replaces them with oversmoothed, uniform blocks (e.g., a detailed tree reduced to solid green). This collapse stems from the model overcommitting to low-dispersion token choices as generation proceeds, which reinforces structural drift and texture wash-out. Notably, entropy collapse is unique to video generation: LLMs avoid it because the high semantic density of language tokens prevents overconfidence in redundant sequences.

Insights. To address these challenges, we propose a locally adaptive, entropy-guided sampling strategy: align candidate pool size with token dispersion. For low-dispersion regions, small pools (with a minimal guard-n) suppress redundant noise and prevent entropy collapse, preserving structural stability while retaining baseline stochasticity. For high-dispersion regions, large pools include all plausible tokens to avoid brittle argmax flips and mitigate early error accumulation. The minimal guard-n is key—it avoids the extremes of greedy decoding (accelerating texture wash-out) and over-large pools (introducing noise), balancing stability and richness. Details of the entropy-to-k mapping for efficient implementation are provided in Section 4.

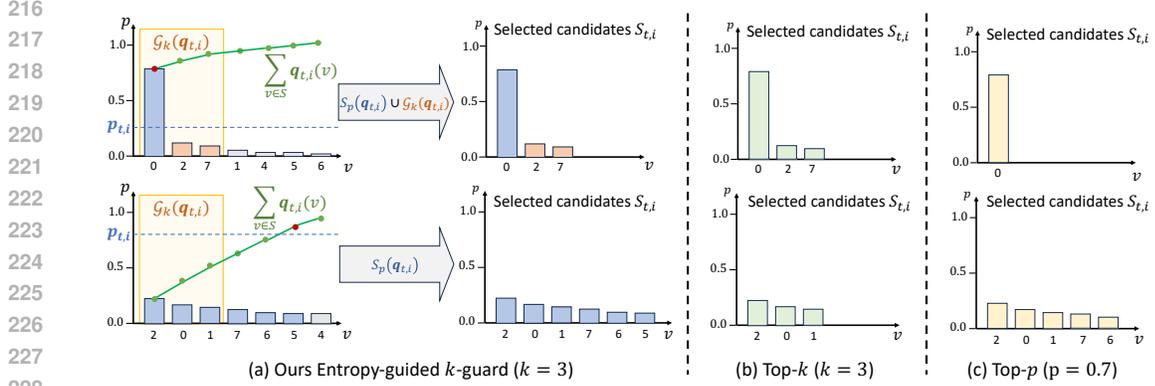


Figure 3: The overall illustration of our sampling strategy.

4 METHOD

Inspired by the findings in Section 1, we propose an Uncertainty-aware Adaptive Sampling strategy. The core idea is to leverage the model’s predictive uncertainty at each token to dynamically guide the sampling process. Specifically, for regions where the model is confident, we enforce a conservative, near-greedy sampling to preserve structure. Conversely, for ambiguous regions, we encourage more diversity to mitigate brittle decisions and enrich textures. This is implemented through a three-stage process: (1) quantifying token-wise uncertainty using entropy, (2) mapping this uncertainty to an adaptive nucleus threshold, and (3) incorporating a “ k -guard” to ensure robust exploration. We provide a pseudocode in Alg. 1.

4.1 PRELIMINARY

We begin by formalizing the **autoregressive (AR) formulation** for video generation. Let \mathcal{V} denote a discrete codebook of size V , obtained via a learned tokenizer such as VQ-VAE. Each video frame is represented as a grid of tokens $\{z_{t,i} \in \mathcal{V}\}$, where t indexes the temporal step and $i \in \{1, 2, \dots, m\}$ indexes the spatial positions within a frame (each frame contains m VQ tokens).

An AR world model factorizes the joint distribution of tokens as a product of conditional probabilities. Specifically, the probability of generating the i -th token in frame t is conditioned on all previously decoded tokens in the same frame, as well as historical context and actions:

$$p(z_{t,i} | z_{t,<i}, c_{<t}, a_{<t}) = \prod_{j=1}^i p(z_{t,j} | z_{t,j-1}, c_{<t}, a_{<t}), \quad (2)$$

where $z_{t,<i}$ denotes the previously decoded tokens within the current frame, $c_{<t}$ represents observed context such as conditioning frames, and $a_{<t}$ denotes historical actions. When $i = m$, the above product gives the joint probability of generating the entire t -th frame.

4.2 INSTABILITY IN AUTOREGRESSIVE VIDEO GENERATION

However, autoregressive token generation is inherently prone to instability due to **error accumulation**. Let $\hat{z}_{t,j}$ denote the token actually generated at step j . Then the conditional probability for the next token depends on previously generated (potentially erroneous) tokens:

$$p(z_{t,i} | \hat{z}_{t,<i}, c_{<t}, a_{<t}) \neq p(z_{t,i} | z_{t,<i}, c_{<t}, a_{<t}), \quad (3)$$

where $\hat{z}_{t,<i}$ contains tokens that may differ from the ground truth $z_{t,<i}$. Consequently, small errors propagate through the sequence, amplifying discrepancies in later tokens and potentially degrading entire frames.

During inference, the model generates a video sequence by sequentially sampling tokens from these categorical distributions. The choice of sampling strategy is therefore critical: greedy decoding

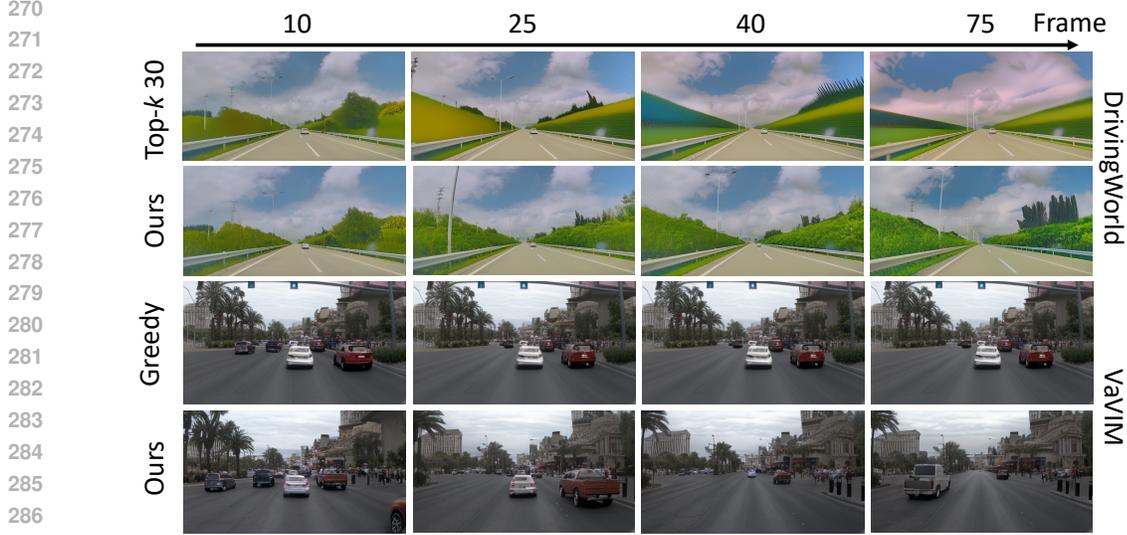


Figure 4: Visual results of DrivingWorld and VaVIM models with our strategy..

Algorithm 1 ENkG Sampling**Require:** probability distribution $\mathbf{p} \in \mathbb{R}^V$; hyperparameters $(\alpha, \beta, p_{\text{low}}, p_{\text{high}}, k_g)$ **Ensure:** Sampled token indices $\mathbf{y} \in \mathbb{Z}$ 1: Compute normalized entropy: $\mathcal{H} \leftarrow -\frac{1}{\log V} \sum_i p_i \log p_i$

2: Map entropy to nucleus probability via affine and clip:

$$p_{t,i} \leftarrow \text{clip}(\alpha \widehat{\mathcal{H}}_{t,i} + \beta, p_{\text{low}}, p_{\text{high}})$$

3: Let $\{q_{(i)}\}_{i=1}^V$ be probabilities sorted in descending order4: Find cutoff $c = \min\{j : \sum_{i=1}^j q_{(i)} \geq p\}$ 5: Set $c \leftarrow \max(c, k_g)$ 6: Define truncated distribution $\tilde{q}_i = \frac{q_{(i)}}{\sum_{j=1}^c q_{(j)}}$ for $i \leq c$, 0 otherwise7: Sample token $y \sim \tilde{q}$ 8: **return** \mathbf{y}

305 often produces blurry frames or repetitive collapse, while excessively random sampling amplifies
306 noise and disrupts temporal coherence. To address this limitation, we introduce an **Entropy-Guided**
307 **k-Guard (ENkG) sampling strategy** that dynamically adjusts sampling diversity according to the
308 model’s predictive confidence, balancing structural fidelity and the richness of textures.

4.3 ENTROPY-GUIDED K-GUARD SAMPLING

312 To quantify token-level uncertainty, we consider the predicted categorical distribution $q_{t,i}$ for each
313 token $z_{t,i}$ at image token site (t, i) :

$$314 \quad q_{t,i}(v) := p(z_{t,i} = v \mid z_{<t}, z_{t,<i}, c_{\leq t}, a_{<t}), \quad v \in \mathcal{V}, \quad (4)$$

316 where \mathcal{V} denotes the discrete codebook. The uncertainty associated with this prediction is measured
317 by its Shannon entropy:

$$318 \quad \mathcal{H}_{t,i} = - \sum_{v \in \mathcal{V}} q_{t,i}(v) \log q_{t,i}(v). \quad (5)$$

321 To obtain a standardized measure on the unit interval, we normalize the entropy by the maximum
322 possible value $\log |\mathcal{V}|$:

$$323 \quad \widehat{\mathcal{H}}_{t,i} = \frac{\mathcal{H}_{t,i}}{\log |\mathcal{V}|} \in [0, 1]. \quad (6)$$

Table 1: Quantitative results on Saturn and Nuplan. * *Cosmos uses a fixed 33-frame generation window; hence its metrics are computed on the first 33 frames (vs. 75 for others).*

Model	DiverseDrive					Nuplan				
	FVD ₇₅ ↓	FID ₇₅ ↓	LPIPS↓	PSNR↑	SSIM↑	FVD ₇₅ ↓	FID ₇₅ ↓	LPIPS↓	PSNR↑	SSIM↑
DrivingWorld(top- <i>k</i> 30)	696	61.78	0.401	14.03	0.43	583	37.80	0.380	14.22	0.39
DrivingWorld(+Ours)	489	26.61	0.350	15.87	0.45	565	31.34	0.360	14.96	0.40
VaVIM(greedy)	1473	91.75	0.396	16.46	0.50	927	65.26	0.315	14.82	0.44
VaVIM(+Ours)	1055	46.76	0.426	14.76	0.46	1031	41.60	0.327	14.43	0.42
Cosmos(top- <i>p</i> 0.8)*	1260	87.82	0.48	16.56	0.54	814	80.45	0.29	17.52	0.54
Cosmos(+Ours)*	1132	84.67	0.47	16.61	0.53	801	75.01	0.29	17.81	0.55

The normalized entropy $\widehat{\mathcal{H}}_{t,i}$ serves as a direct indicator of the model’s confidence in predicting token $z_{t,i}$. Low values of $\widehat{\mathcal{H}}_{t,i}$ correspond to sharply peaked distributions, indicating high confidence in a dominant token, whereas high values indicate flatter distributions and thus greater uncertainty. Specifically, $\widehat{\mathcal{H}}_{t,i} \approx 0$ corresponds to a highly confident, nearly deterministic prediction, while $\widehat{\mathcal{H}}_{t,i} \approx 1$ corresponds to a nearly uniform, highly uncertain prediction.

Entropy-guided adaptive nucleus. To adaptively control sampling diversity, the normalized predictive entropy $\widehat{\mathcal{H}}_{t,i}$ is mapped to a target cumulative probability $p_{t,i} \in [p_{\text{low}}, p_{\text{high}}]$ via an affine transformation with clipping:

$$p_{t,i} = \text{clip}\left(\alpha \widehat{\mathcal{H}}_{t,i} + \beta, p_{\text{low}}, p_{\text{high}}\right). \quad (7)$$

where $\alpha = \frac{p_{\text{high}} - p_{\text{low}}}{\widehat{\mathcal{H}}_{\text{high}} - \widehat{\mathcal{H}}_{\text{low}}}$, $\beta = p_{\text{low}} - \alpha \widehat{\mathcal{H}}_{\text{low}}$.

where $\text{clip}(x, a, b) = \min(\max(x, a), b)$. In the experiments, $p_{\text{low}} = 0.65$, $p_{\text{high}} = 0.9$, and $\widehat{\mathcal{H}}_{\text{low}} = 0.25$, $\widehat{\mathcal{H}}_{\text{high}} = 0.6$.

Based on $p_{t,i}$, the adaptive nucleus set $\mathcal{S}_p(\mathbf{q}_{t,i})$ is defined as the minimal subset of tokens whose cumulative probability meets or exceeds $p_{t,i}$:

$$\mathcal{S}_p(\mathbf{q}_{t,i}) = \arg \min_{\mathcal{S} \subseteq \mathcal{V}} \left\{ |\mathcal{S}| \mid \sum_{v \in \mathcal{S}} q_{t,i}(v) \geq p_{t,i} \right\}. \quad (8)$$

Tokens with low entropy correspond to high-confidence predictions, allowing a small nucleus and near-greedy sampling that preserves fine structures such as edges and boundaries. High-entropy tokens indicate greater uncertainty; a larger nucleus in these cases encourages exploration, enhances diversity, and mitigates the compounding of errors in sequential autoregressive generation. This entropy-guided adaptation provides a principled mechanism to balance fidelity and diversity in token-level sampling.

***k*-Guard for robust exploration.** Direct entropy-guided adaptive sampling can become overly greedy in high-confidence (low-entropy) regions, where the nucleus set may contain only a few tokens. To preserve minimal exploration without compromising stability, the nucleus is augmented with the k_g most probable tokens, forming a *k*-guard:

$$\mathcal{S}_{t,i} = \mathcal{S}_p(\mathbf{q}_{t,i}) \cup \mathcal{G}_k(\mathbf{q}_{t,i}, k_g), \quad (9)$$

where $\mathcal{S}_p(\mathbf{q}_{t,i})$ denotes the nucleus (top-*p*) candidate set, and $\mathcal{G}_k(\mathbf{q}_{t,i}, k_g)$ returns the indices of the k_g tokens with the highest probabilities under $\mathbf{q}_{t,i}$. Typical choices for k_g are small integers such as 3, 5, or 10. If the nucleus already contains these tokens, the union leaves the set unchanged. To limit computational cost, a maximum size n_{max} can be enforced by retaining only the top n_{max} tokens sorted by probability.

Token selection is then performed by sampling from the renormalized distribution over $\mathcal{S}_{t,i}$. This combined framework leverages token-wise uncertainty to adaptively regulate sampling diversity, while the *k*-guard ensures minimal exploration in highly confident regions, enhancing the robustness and stability of sequential autoregressive generation.



Figure 5: Visual results of Cosmos model with our strategy.

Table 2: Ablation study with DrivingWorld model on self-collected dataset.

Method	FVD↓	FID↓	LPIPS↓	PSNR↑
Full Strategy	489	26.61	0.350	15.87
w/o Entropy	532	41.43	0.591	13.96
w/o k-Guard	552	39.76	0.421	15.18

5 EXPERIMENTS

5.1 EXPERIMENTAL SETTING

Models. Since ENkG is a plug-and-play solution, we integrate it with existing AR-based video world model for experiments, including DrivingWorld (Hu et al., 2024), VaVIM (Bartoccioni et al., 2025) and Cosmos (NVIDIA, 2025). We keep generation parameters consistent with the original model for fair comparison. Specifically, DrivingWorld adopts top- k ($k = 30$), VaVIM adopts greedy sampling, while Cosmos uses a top- p ($p = 0.8$).

Evaluation Dataset. We conduct evaluations on two datasets: DiverseDrive and nuPlan. DiverseDrive is a self-collected high-quality driving dataset, which consists of 50 video clips. Compared to nuPlan datasets, DiverseDrive contains more scenarios and a richer variety of plants. These characteristics promote stronger generalization, making DiverseDrive a closer match to the requirements of world-model evaluation.

Metrics. To assess the quality of generated videos, we report the Fréchet Video Distance (FVD) as a measure of video-level realism, and the Fréchet Inception Distance (FID) to evaluate per-frame image fidelity. In addition, we include low-level metrics such as LPIPS, PSNR, and SIIM as supplementary evaluations, though these metrics are not well-suited for the video generation task.

5.2 MAIN RESULTS

Quantitative Comparison. As shown in Table 1, integrating ENkG consistently yields substantial gains across different architectures. On DiverseDrive, our method reduces FVD and FID by an average of 22.8% and 36.5% respectively, while also lowering LPIPS and improving PSNR/SSIM, indicating both perceptual and structural benefits. DrivingWorld also benefits from ENkG on NuPlan, despite being sufficiently trained on this dataset. Even Cosmos, which has a relatively weak AR backbone, achieves modest improvements on DiverseDrive. Notably, VaVIM tends to generate repeated frames, which artificially yields relatively lower FVD values; our strategy, by contrast, effectively alleviates this frame-freezing issue.

Qualitative Comparison. As shown in Fig. 4 and Fig. 5, the existing sampling techniques frequently result in textural degradation, where crucial details in road markings, such as crosswalks, and surrounding vegetation become indistinct and blurry. Furthermore, these approaches are susceptible to color distortion, leading to unnatural and washed-out hues that compromise the scene’s realism. The generated sequence of VaVIM collapses into a static or near-static frame, failing to capture the inherent dynamics of the driving environment and rendering vehicles motionless. In contrast, our entropy-guided approach, which dynamically adjusts the size of candidates to prevent overconfidence, demonstrates substantial improvements in perceptual quality. Our strategy effectively mitigates the aforementioned issues, producing videos with sharp, well-defined textures and accurate color fidelity.



442 Figure 6: Entropy-adaptive guidance prevents collapse in the video model.



454 Figure 7: The k -guard design prevents frame-freezing in the video model.

455
456 5.3 ABLATION STUDY

457
458 **Effect of entropy-adaptive guidance.** The core contribution of the entropy-adaptive guidance is
459 the dynamic adjustment of the sampling nucleus based on the model’s predictive uncertainty. As
460 shown in Figure 1 and Figure 6, it effectively mitigates issues of textural decay and color shifting
461 commonly seen in baseline methods. By allowing a wider range of tokens when uncertainty is high
462 and narrowing it when the model is confident, our method preserves high-frequency details, resulting
463 in significantly sharper textures on surfaces like road markings and vegetation. This enhancement
464 in per-frame visual fidelity translates to lower FID and LPIPS scores, indicative of more realistic
465 and perceptually similar generated frames. Consequently, entropy-adaptive guidance significantly
466 improves the visual quality and realism of the generated videos.

467
468 **Effect of k -guard.** The k -guard mechanism ensures a minimum level of diversity in candidate to-
469 kens. Without the k -guard, the model can, even with entropy-adaptive guidance, become overly
470 confident in certain contexts. In Figure 7, this leads to significant temporal artifacts, such as vehi-
471 cles that remain nearly stationary when they should be in motion, a physically implausible scenario.
472 The introduction of k -guard directly addresses this failure mode, leading to more fluid and realistic
473 motion dynamics, as quantitatively reflected in the substantial reduction of the FVD score, which
474 measures temporal consistency. Therefore, the k -guard is crucial for maintaining temporal coher-
475 ence and preventing the generation of degenerate, static sequences.

476 6 CONCLUSION

477
478 In this work, we investigated the challenge of error accumulation in autoregressive video genera-
479 tion and highlighted the overlooked role of the sampling process. We proposed *Uncertainty-aware*
480 *Adaptive Sampling*, a simple yet effective strategy that dynamically modulates token diversity based
481 on predictive entropy with a minimal k -guard. Unlike prior approaches that require architectural
482 changes or retraining, our method operates purely at inference time, making it broadly applicable
483 to existing large-scale models. Extensive experiments demonstrate that our approach significantly
484 improves temporal coherence, preserves fine-grained details, and extends the effective generation
485 horizon. These results suggest that inference-time uncertainty-aware strategies provide a practical
and generalizable path toward more robust and high-fidelity video world models.

REFERENCES

- 486
487
488 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,
489 Shijie Wang, Jun Tang, Humen Zhong, Yanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,
490 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,
491 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
492
- 493 Florent Bartoccioni, Elias Ramzi, Victor Besnier, Shashanka Venkataramanan, Tuan-Hung Vu,
494 Yihong Xu, Loick Chambon, Spyros Gidaris, Serkan Odabas, David Hurych, Renaud Marlet,
495 Alexandre Boulch, Mickael Chen, Eloi Zablocki, Andrei Bursuc, Eduardo Valle, and Matthieu
496 Cord. Vavim and vavam: Autonomous driving through video generative modeling. *arXiv preprint*
497 *arXiv:2502.15672*, 2025.
- 498 Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled Sampling for Sequence
499 Prediction with Recurrent Neural Networks. *arXiv preprint arXiv:1506.03099*, 2015.
500
- 501 Paul Benjamin and Anna Smith. Measuring compounding errors in sequential prediction models.
502 *NeurIPS*, 2018.
- 503 Alex Clark and Sanja Fidler. Adversarial video generation on large datasets. In *ICCV*, 2019.
504
- 505 Jingtao Ding, Yunke Zhang, Yu Shang, Yuheng Zhang, Zefang Zong, Jie Feng, Yuan Yuan,
506 Hongyuan Su, Nian Li, Nicholas Sukiennik, et al. Understanding world or predicting future?
507 a comprehensive survey of world models. *ACM Computing Surveys*, 2024.
- 508 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
509 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
510 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
511 *arXiv:2010.11929*, 2020.
512
- 513 Hao Feng and Jie Li. Sequence-level error accumulation in autoregressive video models. *TPAMI*,
514 2021.
- 515 Raghav Goyal and Sung Lee. Non-markovian effects and memory bottlenecks in ar video models.
516 *ICML*, 2022.
517
- 518 Yuchao Gu, Weijia Mao, and Mike Zheng Shou. Long-context autoregressive video modeling with
519 next-frame prediction. *arXiv preprint arXiv:2503.19325*, 2025.
- 520 David Ha and Jürgen Schmidhuber. World models. *CoRR*, 2018.
521
- 522 Haoran He, Yang Zhang, Liang Lin, Zhongwen Xu, and Ling Pan. Pre-trained video generative
523 models as world simulators. *arXiv preprint arXiv:2502.07825*, 2025.
- 524 Jonathan Ho and Ajay Jain. Video diffusion models. In *NeurIPS*, 2022.
525
- 526 Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shot-
527 ton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving. *arXiv*
528 *preprint arXiv:2309.17080*, 2023.
- 529 Xiaotao Hu, Wei Yin, Mingkai Jia, Junyuan Deng, Xiaoyang Guo, Qian Zhang, Xiaoxiao Long, and
530 Ping Tan. Drivingworld: Constructingworld model for autonomous driving via video gpt. *arXiv*
531 *preprint arXiv:2412.19505*, 2024.
532
- 533 X. Huang, Z. Li, G. He, M. Zhou, and E. Shechtman. Self forcing: Bridging the train-test gap in
534 autoregressive video diffusion. 2025.
- 535 Zhewei Kang, Xuandong Zhao, and Dawn Song. Scalable best-of-n selection for large language
536 models via self-certainty. *arXiv preprint arXiv:2502.18581*, 2025.
537
- 538 Lingdong Kong, Wesley Yang, Jianbiao Mei, Youquan Liu, Ao Liang, Dekai Zhu, Dongyue Lu,
539 Wei Yin, Xiaotao Hu, Mingkai Jia, et al. 3d and 4d world modeling: A survey. *arXiv preprint*
arXiv:2509.07996, 2025.

- 540 Chenhao Li, Andreas Krause, and Marco Hutter. Robotic world model: A neural network simulator
541 for robust policy optimization in robotics, 2025. URL [https://arxiv.org/abs/2501.](https://arxiv.org/abs/2501.10100)
542 10100.
- 543 Xiaoxiao Long, Qingrui Zhao, Kaiwen Zhang, Zihao Zhang, Dingrui Wang, Yumeng Liu, Zhengjie
544 Shu, Yi Lu, Shouzheng Wang, Xinzhe Wei, Wei Li, Wei Yin, Yao Yao, Jia Pan, Qiu Shen, Ruigang
545 Yang, Xun Cao, and Qionghai Dai. A survey: Learning embodied intelligence from physical
546 simulators and world models, 2025. URL <https://arxiv.org/abs/2507.00917>.
- 547 Xiaoxiao Ma, Feng Zhao, Pengyang Ling, Haibo Qiu, Zhixiang Wei, Hu Yu, Jie Huang, Zhixiong
548 Zeng, and Lin Ma. Towards better faster autoregressive image generation: From the perspective
549 of entropy, 2025. URL <https://arxiv.org/abs/2510.09012>.
- 550 Sicheng Mo, Ziyang Leng, Leon Liu, Weizhen Wang, Honglin He, and Bolei Zhou. Dream-
551 land: Controllable world creation with simulator and generative models. *arXiv preprint*
552 *arXiv:2506.08006*, 2025.
- 553 Georgy Noarov, Soham Mallick, Tao Wang, Sunay Joshi, Yan Sun, Yangxinyu Xie, Mengxin Yu,
554 and Edgar Dobriban. Foundations of top- k decoding for language models, 2025. URL <https://arxiv.org/abs/2505.19371>.
- 555 NVIDIA. Cosmos world foundation model platform for physical ai. *arXiv preprint*
556 *arXiv:2501.03575*, 2025. URL <https://arxiv.org/abs/2501.03575>.
- 557 OpenAI. Video generation models as world simulators. [https://openai.com/research/
558 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators), 2024.
- 559 Raghu Parthipan, Mohit Anand, Hannah M Christensen, and J Scott Hosking. Defining error accu-
560 mulation in ML atmospheric simulators. *arXiv preprint arXiv:2405.14714*, 2024.
- 561 Zexuan Qiu, Zijing Ou, Bin Wu, Jingjing Li, Aiwei Liu, and Irwin King. Entropy-based decoding for
562 retrieval-augmented large language models, 2025. URL [https://arxiv.org/abs/2406.](https://arxiv.org/abs/2406.17519)
563 17519.
- 564 Shauli Ravfogel, Yoav Goldberg, and Jacob Goldberger. Conformal nucleus sampling. In *Findings*
565 *of the Association for Computational Linguistics: ACL 2023*, pp. 27–34, 2023.
- 566 S. Ren, Q. Yu, J. He, X. Shen, A. Yuille, and L.-C. Chen. Beyond next-token: Next-x prediction for
567 autoregressive visual generation. 2025a.
- 568 Shuhuai Ren, Shuming Ma, Xu Sun, and Furu Wei. Next block prediction: Video generation via
569 semi-autoregressive modeling, 2025b. URL <https://arxiv.org/abs/2502.07737>.
- 570 Ankit Saxena and Rohan Kumar. Error propagation and memory bottlenecks in ar video generation.
571 *ICLR*, 2024.
- 572 Florian Schmidt. Generalization in generation: A closer look at exposure bias. *EMNLP-IJCNLP*
573 *2019*, pp. 157, 2019.
- 574 Toby Simonds. Entropy adaptive decoding: Dynamic model switching for efficient inference. *arXiv*
575 *preprint arXiv:2502.06833*, 2025.
- 576 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
577 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- 578 Alessandro Stolfo, Ben Wu, Wes Gurnee, Yonatan Belinkov, Xingyi Song, Mrinmaya Sachan, and
579 Neel Nanda. Confidence regulation neurons in language models. *Advances in Neural Information*
580 *Processing Systems*, 37:125019–125049, 2024.
- 581 Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang,
582 Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable
583 length video generation from open domain textual description. *arXiv preprint arXiv:2210.02399*,
584 2022.

- 594 Brian Walker and Abhinav Gupta. Uncertain predictions in sequential generative models. *NeurIPS*,
595 2016.
- 596
- 597 Xiaofeng Wang, Zheng Zhu, Guan Huang, Boyuan Wang, Xinze Chen, and Jiwen Lu. World-
598 dreamer: Towards general world models for video generation via predicting masked tokens. *arXiv*
599 *preprint arXiv:2401.09985*, 2024a.
- 600 Xiaofeng Wang, Zheng Zhu, Guan Huang, Boyuan Wang, Xinze Chen, and Jiwen Lu. World-
601 dreamer: Towards general world models for video generation via predicting masked tokens.
602 *CoRR*, 2024b.
- 603
- 604 Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit. Scaling autoregressive video models. In
605 *International Conference on Learning Representations*.
- 606 Jialong Wu, Shaofeng Yin, Ningya Feng, Xu He, Dong Li, Jianye Hao, and Mingsheng Long.
607 ivideogpt: Interactive videogpts are scalable world models. *Advances in Neural Information*
608 *Processing Systems*, 37:68082–68119, 2024.
- 609
- 610 Kaixuan Yu and Lei Chen. Magi: Mitigating accumulated generative inference error in ar video
611 models. *ICML*, 2024.
- 612 Hangjie Yuan, Weihua Chen, Jun Cen, Hu Yu, Jingyun Liang, Shuning Chang, Zhihui Lin, Tao
613 Feng, Pengwei Liu, Jiazheng Xing, Hao Luo, Jiasheng Tang, Fan Wang, and Yi Yang. Lumos-
614 1: On autoregressive video generation from a unified model perspective, 2025. URL <https://arxiv.org/abs/2507.08801>.
- 615
- 616 Kaiwen Zhang, Zhenyu Tang, Xiaotao Hu, Xingang Pan, Xiaoyang Guo, Yuan Liu, Jingwei Huang,
617 Li Yuan, Qian Zhang, Xiao-Xiao Long, et al. Epona: Autoregressive diffusion world model for
618 autonomous driving. *arXiv preprint arXiv:2506.24113*, 2025.
- 619
- 620 Shimao Zhang, Yu Bao, and Shujian Huang. Edt: Improving large language models’ generation
621 by entropy-based dynamic temperature sampling, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2403.14541)
622 [2403.14541](https://arxiv.org/abs/2403.14541).
- 623
- 624 D. Zhou, Q. Sun, Y. Peng, K. Yan, R. Dong, D. Wang, N. Duan, X. Zhang, H.-Y. Ni, and H.-Y.
625 Shum. Taming teacher forcing for masked autoregressive video generation. 2025.
- 626
- 627
- 628
- 629
- 630
- 631
- 632
- 633
- 634
- 635
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647

A APPENDIX

A.1 EXPERIMENTAL PROTOCOLS (CONCISE)

We report FVD, FID, LPIPS, PSNR, and SSIM under identical preprocessing and frame sampling across methods unless noted. **Cosmos note:** Cosmos uses a fixed 33-frame generation window; thus all Cosmos metrics are computed on the first 33 frames of each sequence, while others use 75 (FVD₇₅, FID₇₅).

A.2 USE OF LLMs (DISCLOSURE)

LLMs were used for language polishing and minor code refactoring suggestions only; authors verified all outputs. No dataset labeling or claim-critical content was delegated without human validation.

A.3 ETHICS (BRIEF)

We comply with dataset licenses and anonymization policies; no personally identifiable information was used. Potential risks (e.g., misuse in safety-critical scenarios) are discussed and bounded by research-only usage.

B HYPERPARAMETER SENSITIVITY AND BASELINE TUNING

This section provides additional analyses on the sensitivity of ENkG to its hyperparameters and on the tuning of static sampling baselines. Unless otherwise noted, all experiments are conducted on DrivingWorld.

B.1 SENSITIVITY OF ENkG HYPERPARAMETERS

Robustness to entropy and probability thresholds. We first study the robustness of ENkG with respect to the entropy thresholds ($H_{\text{low}}, H_{\text{high}}$) and the corresponding probability range ($p_{\text{low}}, p_{\text{high}}$) used for entropy-guided truncation. Starting from the default configuration (*Mid*), we construct two extreme variants: a conservative configuration (*Left*), which prefers smaller p and narrower entropy range, and an aggressive configuration (*Right*), which allows larger p and a wider entropy range.

Table 3 shows that all metrics remain stable across these settings. Although the extreme variants cause moderate degradation in FVD/FID, the overall differences are small, indicating a broad performance plateau. This suggests that ENkG does not require delicate hand-tuning of entropy or probability thresholds to remain effective.

Table 3: Sensitivity to entropy thresholds ($H_{\text{low}}, H_{\text{high}}$) and probability band ($p_{\text{low}}, p_{\text{high}}$). The default configuration (*Mid*) achieves the best trade-off, while both conservative (*Left*) and aggressive (*Right*) shifts only moderately affect the performance.

Setting	$H_{\text{low/high}}$	$p_{\text{low/high}}$	FVD ↓	FID ↓	LPIPS ↓	SSIM ↑	PSNR ↑
Left	0.0 / 0.5	0.60 / 0.90	522.05	30.93	0.35	0.49	16.26
Mid (Default)	0.25 / 0.60	0.65 / 0.90	489.00	26.61	0.35	0.45	15.87
Right	0.40 / 0.90	0.80 / 0.95	497.52	29.91	0.35	0.48	15.98

Insensitivity to guard size k_g . Next, we fix all other ENkG hyperparameters and vary the guard size k_g , which controls how many top candidates are preserved by the k -guard at each step. As shown in Table 4, the performance is highly stable for any $k_g \in [2, 15]$, and only the degenerate case $k_g = 1$ (which effectively disables the guard) leads to a clear degradation. This indicates that while the *presence* of the k -guard is crucial to prevent collapse, its exact value is not sensitive in a wide range.

Table 4: Sensitivity to guard size k_g (other hyperparameters fixed). ENkG remains stable for a broad range of k_g , and only the degenerate case $k_g = 1$ (no guard) leads to noticeable degradation.

k_g	1	2	3 (Default)	7	15
FVD ↓	552.00	503.98	489.00	510.96	510.64
FID ↓	39.76	29.62	26.61	27.55	29.67

Table 5: Static top- p baselines on DrivingWorld. ENkG (default configuration from Table 3) is shown for reference. Even the best static top- p setup remains substantially worse than ENkG in FVD/FID.

Metric	$p = 0.5$	$p = 0.7$	$p = 0.8$	$p = 0.9$	$p = 1.0$
FVD ↓	836.49	680.62	642.97	625.44	530.22
FID ↓	52.64	46.86	40.03	47.26	43.73
LPIPS ↓	0.40	0.36	0.37	0.38	0.38
SSIM ↑	0.46	0.50	0.46	0.43	0.34
PSNR ↑	13.50	15.76	14.75	16.46	14.66

B.2 TUNING OF STATIC SAMPLING BASELINES

To further exclude the possibility that our gains come from under-tuned baselines, we perform a systematic grid search over *static* top- p , *static* top- k , and combined pk sampling on DrivingWorld, while keeping all other settings (including temperature) fixed. For each baseline family, we report the *best* configuration found in the grid and compare it against ENkG.

Static top- p baselines. We sweep over $p \in \{0.5, 0.7, 0.8, 0.9, 1.0\}$. The quantitative results are summarized in Table 5. Even under the best static configuration, the FVD remains above 530 and FID around 40, which are clearly worse than ENkG (FVD = 489.00, FID = 26.61).

Static top- k baselines. We further sweep over $k \in \{30, 60, 90, 120, 150, 500\}$. Table 6 reports the results. Although FID is minimized around $k = 90$ and FVD around $k = 150$, even these best-performing static top- k configurations still lag behind ENkG in both FVD and FID.

Combined pk baselines. Finally, we explore combined pk sampling (first top- k , then top- p within the truncated set). Among the tested configurations, we find the best performance around ($p = 0.8, k = 1000$):

$$\text{FVD} = 595.93, \quad \text{FID} = 43.76, \quad \text{LPIPS} = 0.357, \quad \text{SSIM} = 0.50, \quad \text{PSNR} = 15.93.$$

While the perceptual metrics (LPIPS/SSIM/PSNR) are comparable to ENkG, the FVD and FID are still notably worse.

Discussion. Across all these sweeps, we always compare ENkG against the *best* static configuration of each baseline family (top- p , top- k , and pk). ENkG consistently achieves substantially better FVD and FID, indicating that its advantage does not come from weak or under-tuned baselines, but from the proposed *entropy-guided dynamic truncation with k -guard*, which provides a strictly stronger sampling strategy than any single static choice of (p) or (k).

Interestingly, we also observe that for very large candidate sets (e.g., $p = 1.0$ in top- p or $k > 100$ in top- k), the generated videos can exhibit visibly fragmented or “shattered” structures, yet FVD does not necessarily increase and can even improve. This behavior is consistent with known limitations of FVD in penalizing spatial incoherence, and explains why DrivingWorld’s original implementation adopts relatively defaults (e.g., $k = 30$) to preserve vehicle structural consistency rather than aggressively minimizing FVD under heavily fragmented scenes.

Table 6: Static top- k baselines on DrivingWorld. The best static top- k settings (e.g., $k = 90$ for FID, $k = 150$ for FVD) remain inferior to ENkG.

Metric	$k = 30$	$k = 60$	$k = 90$	$k = 120$	$k = 150$	$k = 500$
FVD ↓	696.14	661.52	615.37	569.10	554.74	564.10
FID ↓	61.78	41.54	34.50	37.86	39.02	39.10
LPIPS ↓	0.40	0.39	0.39	0.37	0.38	0.38
SSIM ↑	0.44	0.45	0.44	0.48	0.45	0.45
PSNR ↑	14.04	14.32	14.05	15.73	15.34	15.08

C QUALITATIVE RESULTS

C.1 GENERAL-DOMAIN MODELS

To validate the effectiveness of ENkG in general domains, we evaluate it on Lumos-1 (Yuan et al., 2025) and NBP (Ren et al., 2025b). Our method mitigates error accumulation during long-horizon video generation and produces more temporally consistent results with reduced color drift and collapse artifacts. Representative qualitative results are shown in Figure 8 and Figure 9.

C.2 ADDITIONAL COMPARISONS

We further provide more qualitative results on DrivingWorld in Figure 11 and VaVim in Figure 12.

C.3 LONG-HORIZON GENERATION ON DRIVINGWORLD

To further evaluate ENkG under long-range autoregressive rollout, we conduct a 200-frame generation experiment on DrivingWorld. This setting is particularly prone to error accumulation and low-entropy collapse. As shown in Figure 10, ENkG substantially mitigates visual drift and maintains scene stability over very long horizons, whereas the baseline model exhibits background smearing and global color shift.

D ANALYSIS OF ENTROPY.

D.1 EMPIRICAL EXAMINATION OF ENTROPY IN AR VIDEO MODELS

To better understand the entropy dynamics underlying low-entropy collapse, we compare (a) the probability distributions of the top-20 tokens between a large language model (Qwen2.5 Bai et al. (2025)) and an autoregressive video model (DrivingWorld Hu et al. (2024)), and (b) the average token entropy across generation timesteps on the NuPlan dataset. As illustrated in Figure 13a, AR video models exhibit significantly sharper and more rapidly collapsing distributions, making them especially vulnerable to trajectory locking and deterministic failure modes.

D.2 CONSEQUENCES OF THE LOW-ENTROPY TRAP

In autoregressive (AR) video models, the predictive distribution at step t can be written as $p_\theta(x_t | x_{<t})$ with entropy $H_t = -\sum_x p_\theta(x | x_{<t}) \log p_\theta(x | x_{<t})$. We refer to a *low-entropy trap* as the regime where H_t collapses prematurely and p_θ becomes pathologically overconfident around a small set of locally consistent tokens, even though the corresponding trajectory is globally suboptimal.

Local overconfidence and trajectory locking. Once the model enters such a regime, the effective candidate set under common sampling schemes (top- p , top- k , or p^k) often shrinks to one or two tokens with probability mass close to 1. From that point on, the model repeatedly feeds its own highly deterministic predictions back into the context, reinforcing the same local pattern step after step. This behavior is analogous to language models falling into repetitive loops (e.g., “the the the”), persisting in a wrong chain-of-thought branch, or doubling down on an early but incorrect inference.

810 **Visual manifestations in video space.** In video generation, the low-entropy trap does not merely
811 lead to trivially frozen frames. More subtly, it manifests as: (i) **background smearing**, where large
812 regions of the background collapse into blurry blobs and lose meaningful structural detail; (ii) **global**
813 **color shift**, where the entire frame drifts toward an unnatural color cast, making consecutive frames
814 look consistently tinted or over-/under-exposed; and (iii) **texture freezing**, where fine-grained ap-
815 pearance patterns (e.g., grass, water, sky) become unnaturally static and appear glued to the camera
816 or object rather than evolving with the motion. These effects are illustrated in Figure 11 and 12,
817 where the model settles into visually coherent yet clearly undesirable trajectories.

818 **Distinction from high-entropy noise.** It is important to distinguish the low-entropy trap from
819 the more familiar high-entropy failure mode. High entropy typically corresponds to excessive ran-
820 domness, producing noisy or chaotic frames that visibly violate short-term consistency. In contrast,
821 low-entropy collapse yields *overly deterministic* behavior: short-term consistency can even look
822 improved, while global realism, long-horizon dynamics, and scene plausibility deteriorate.

823 **Effect of k -guard.** Our ENkG sampler directly targets this overconfidence. When entropy falls
824 below the lower threshold H_{low} , conventional truncation schemes would typically select only the
825 single most probable token. In contrast, ENkG enforces a minimum guard size k_g : even in very
826 low-entropy regimes, at least k_g candidates are preserved, preventing the effective distribution from
827 collapsing to a point mass. This mechanism maintains a controlled level of uncertainty and allows
828 the model to explore alternative continuations that can recover from early mistakes. Empirically,
829 we observe that enabling k -guard reduces the frequency and severity of texture locking, background
830 freezing, and unnatural motion, leading to more coherent long-horizon trajectories and higher overall
831 generative quality.
832

833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

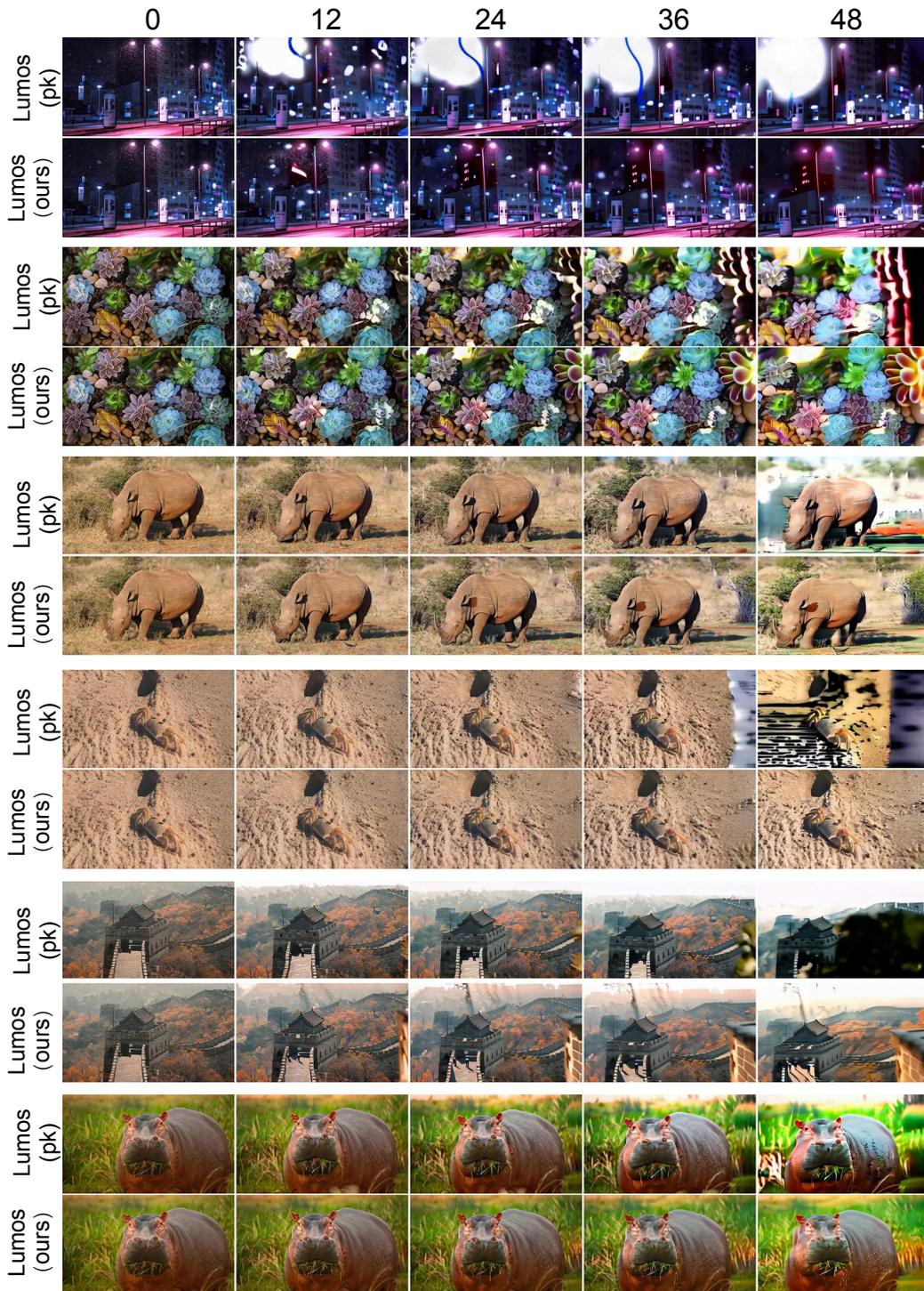


Figure 8: Comparative experiments on the *Lumos-1* model.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

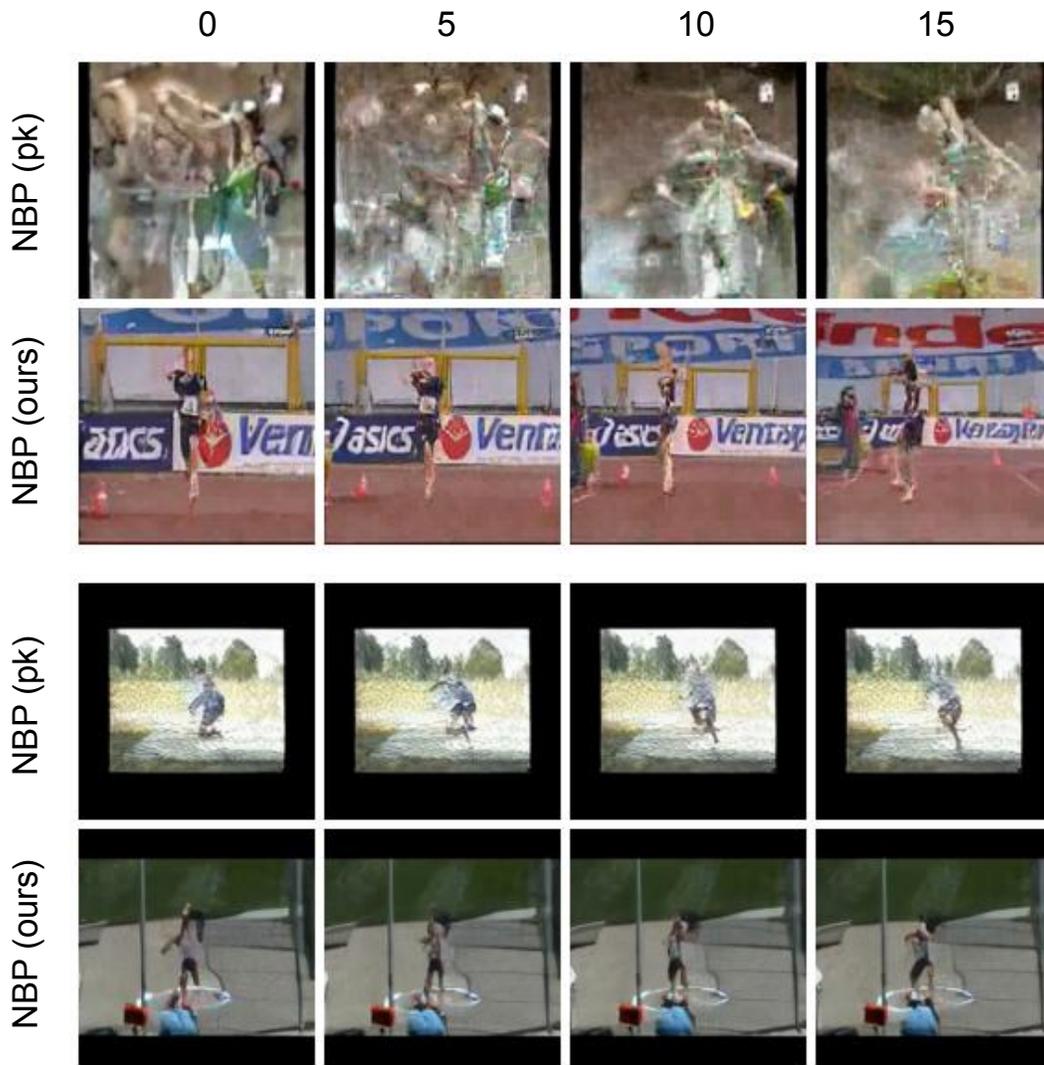


Figure 9: A Qualitative Demonstration of NBP on the UCF-101 Dataset.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

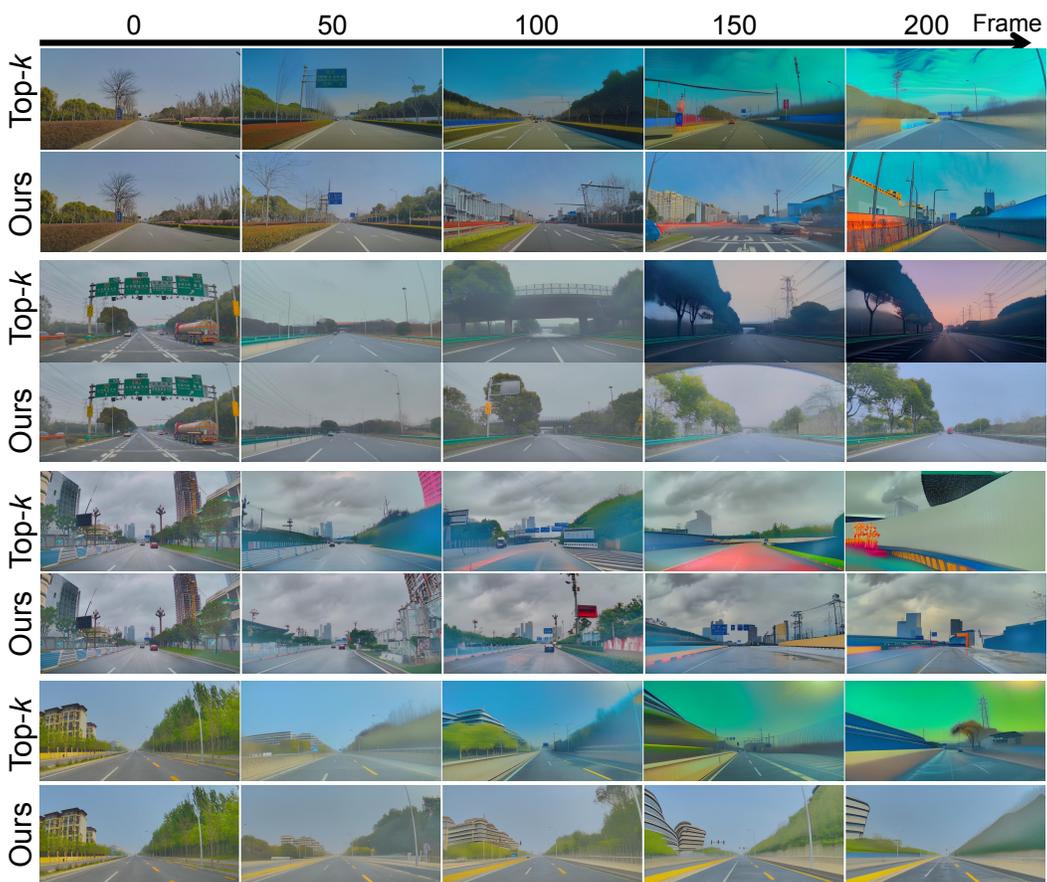


Figure 10: Long-term generation results of Drivingworld.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

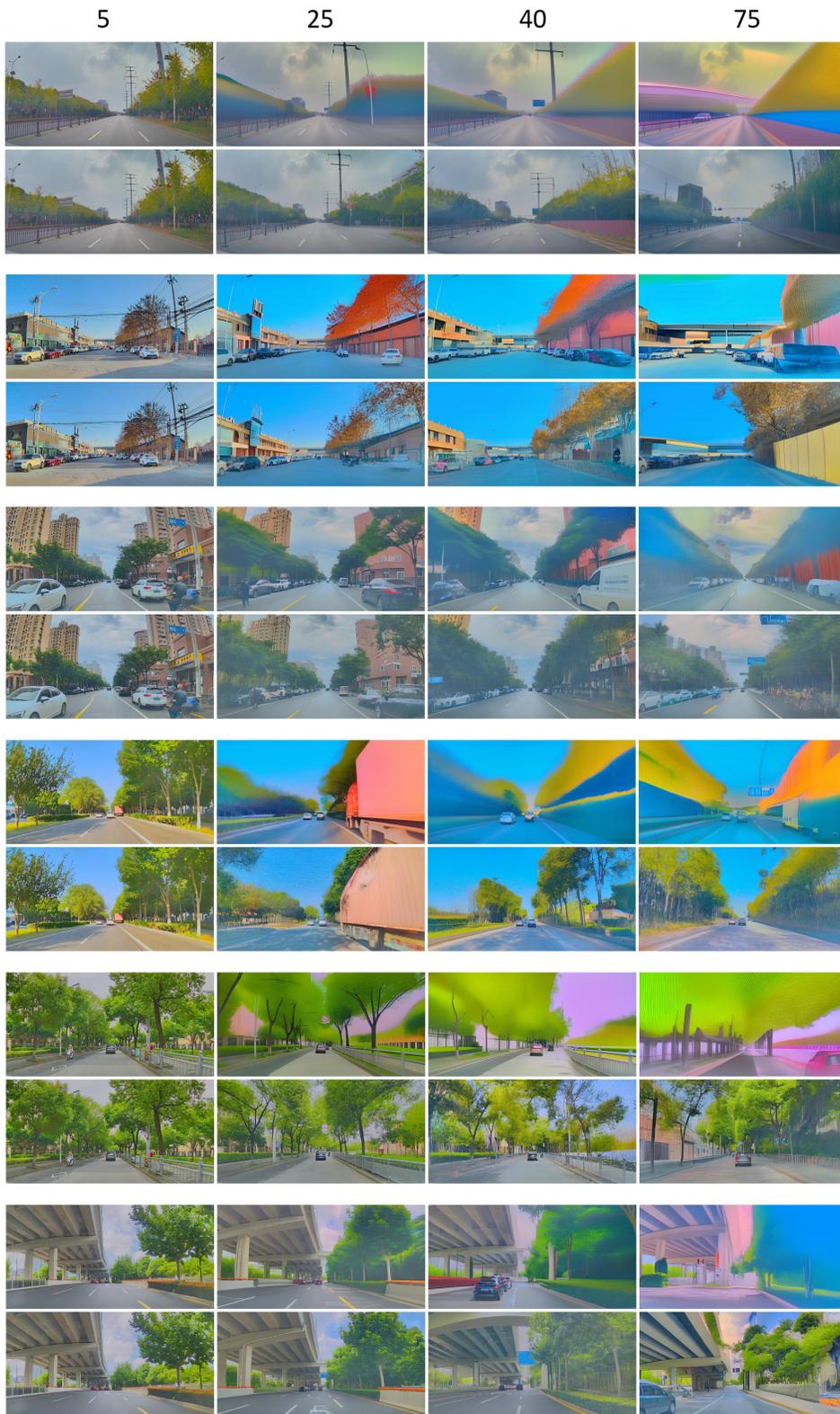


Figure 11: Additional DrivingWorld comparisons between the baseline and ENkG. (Part1)

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

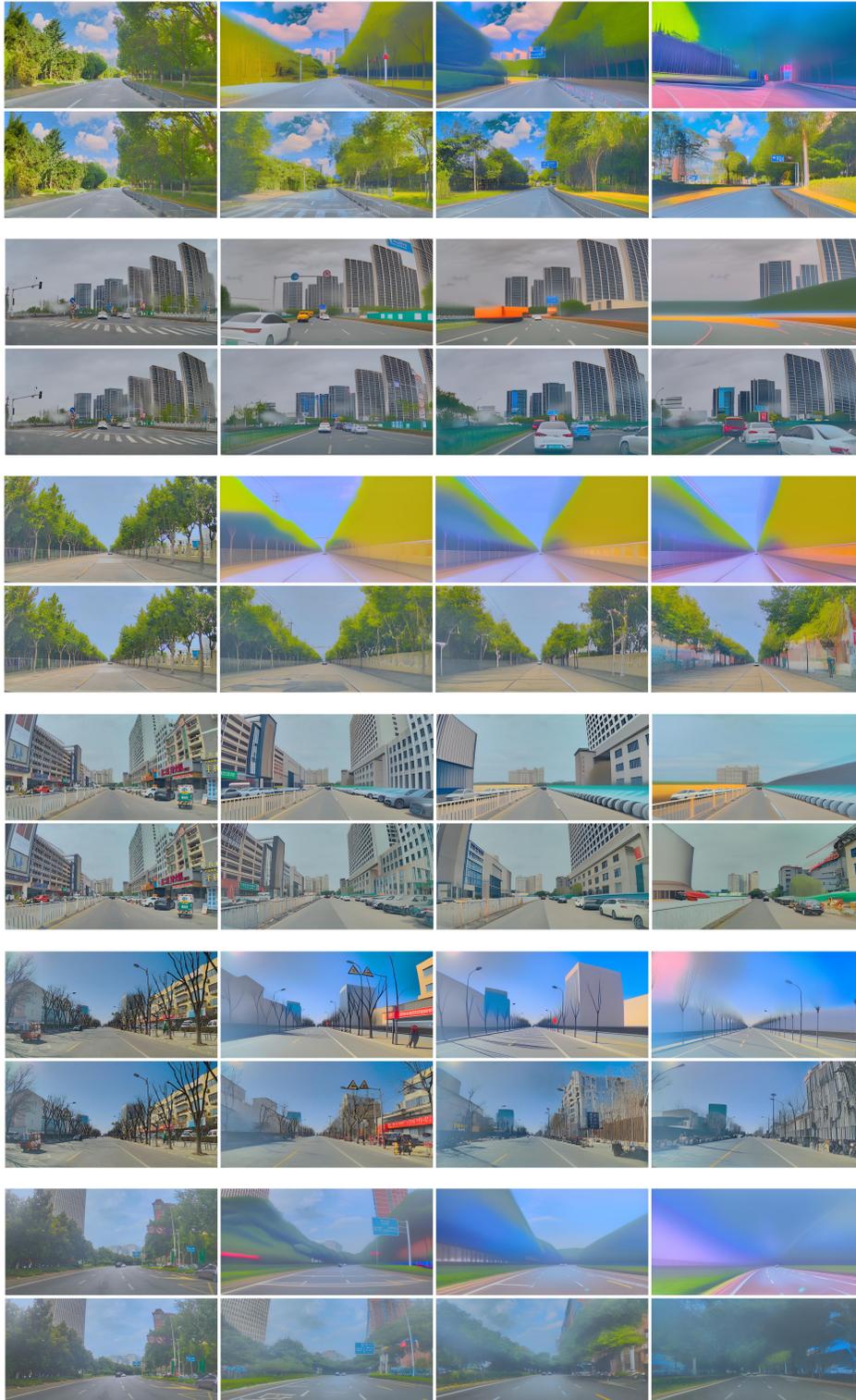


Figure 11: Part2

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

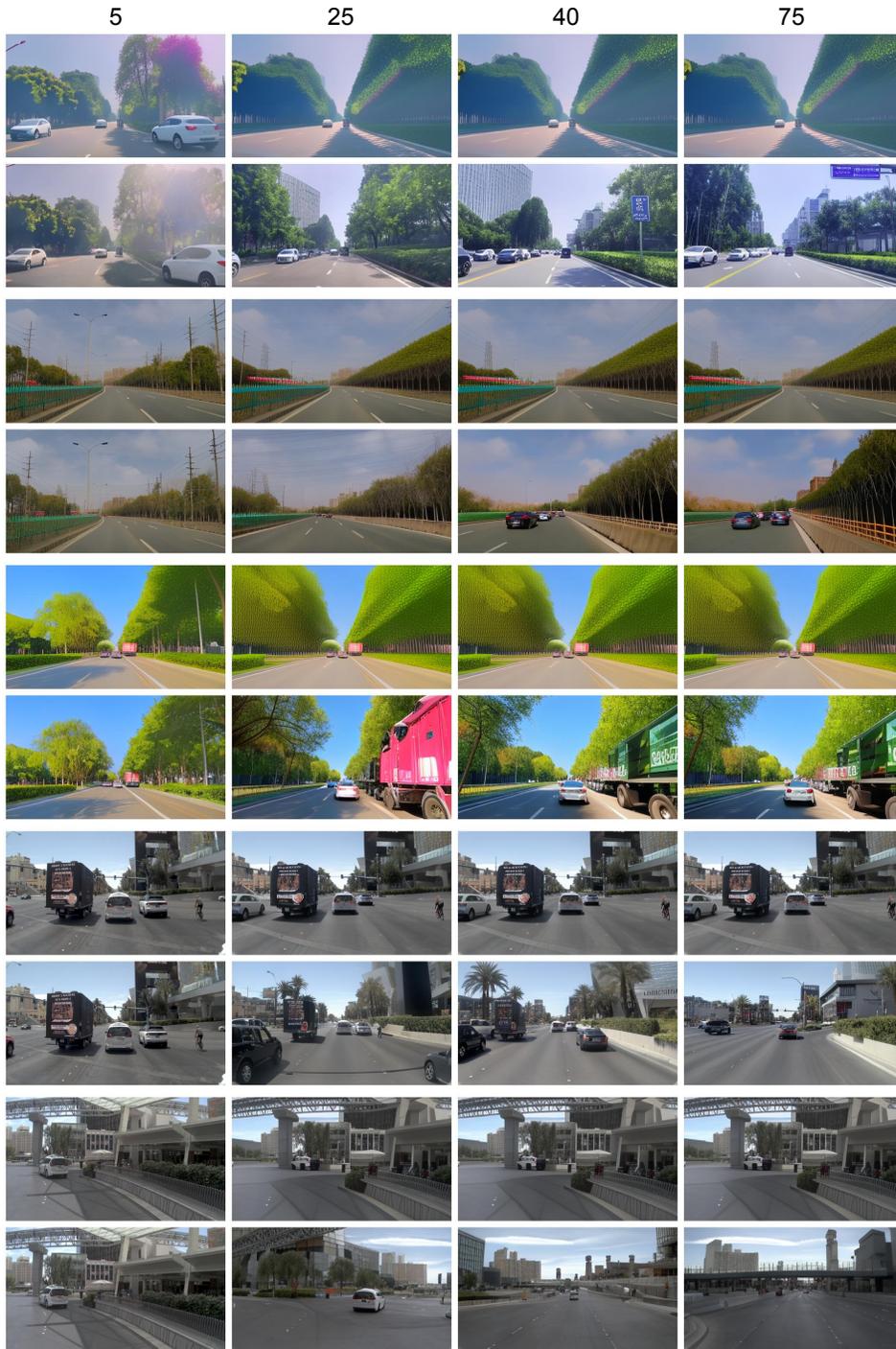


Figure 12: Additional Vavim comparisons between the baseline and ENkG.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

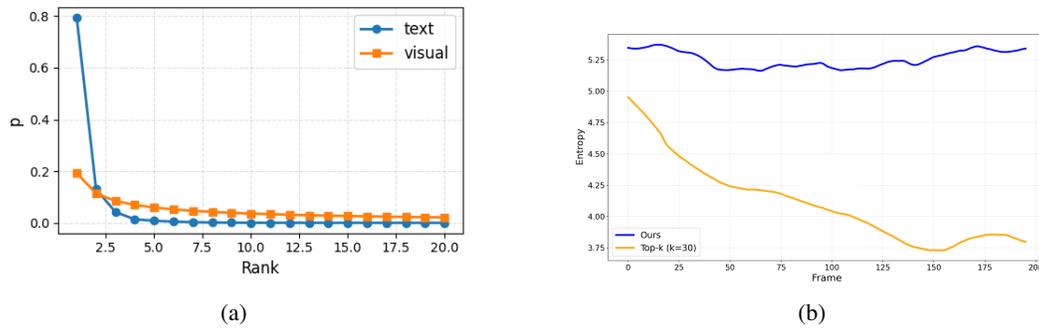


Figure 13: (a).Comparison of the probabilities of the top-20 tokens between LLMs (Qwen2.5 Bai et al. (2025)) and video AR model (DrivingWorld Hu et al. (2024)). (b).Average token entropy of DrivingWorld at each frame as a function of generation timestep on the NuPlan dataset.