

# Supplementary material: OSEG: Improving Diffusion sampling through Orthogonal Smoothed Energy Guidance

Masud An Nur Islam Fahim  
University of Vaasa,  
Finland.  
masud.fahim@uwasa.fi

Nazmus Saqib  
Jeju National University,  
Korea.  
nsaqib1995@gmail.com

Joon-Min Gil  
Jeju National University,  
Korea.  
jmgil@jejunu.ac.kr

## 1. Proof

### Orthogonal Query Leads to Uniform Attention

Here, we present complete proof for the uniform attention when the query appears to be orthogonal to the key vectors. If the query vector  $\mathbf{Q} \in \mathbb{R}^d$  is orthogonal to all key vectors  $\mathbf{K}_i \in \mathbb{R}^d$  (i.e.,  $\mathbf{Q}\mathbf{K}_i^\top = 0; \forall i$ ), then the attention weights  $a_i$  are uniform:

$$A_i = \frac{1}{n}; \forall i$$

*Proof.* Given  $\mathbf{A}$  is the attention weights,  $\mathbf{Q}$  is the query vector,  $\mathbf{K}$  is the key vector, and  $d$  is the dimensionality of the vectors. For a single instance  $A_i$  in the attention weights given by:

$$A_i = \frac{\exp\left(\frac{\mathbf{Q}\mathbf{K}_i^\top}{\sqrt{d}}\right)}{\sum_{j=1}^n \exp\left(\frac{\mathbf{Q}\mathbf{K}_j^\top}{\sqrt{d}}\right)} \quad (1)$$

Given the orthogonality condition  $\mathbf{Q}\mathbf{K}_i^\top = 0$  for all  $i$ , substitute into the exponent:

$$\frac{\mathbf{Q}\mathbf{K}_i^\top}{\sqrt{d}} = \frac{0}{\sqrt{d}} = 0.$$

Thus, the upper part of the 1 simplifies to:

$$\exp\left(\frac{\mathbf{Q}\mathbf{K}_i^\top}{\sqrt{d}}\right) = \exp(0) = 1.$$

Substitute this into 1 and we get:

$$A_i = \frac{1}{\sum_{j=1}^n \exp\left(\frac{\mathbf{Q}\mathbf{K}_j^\top}{\sqrt{d}}\right)}.$$

For the denominator, the summation becomes:

$$\sum_{j=1}^n \exp\left(\frac{\mathbf{Q}\mathbf{K}_j^\top}{\sqrt{d}}\right) = \sum_{j=1}^n 1 = n.$$

Therefore, the attention weight for each  $i$  is:

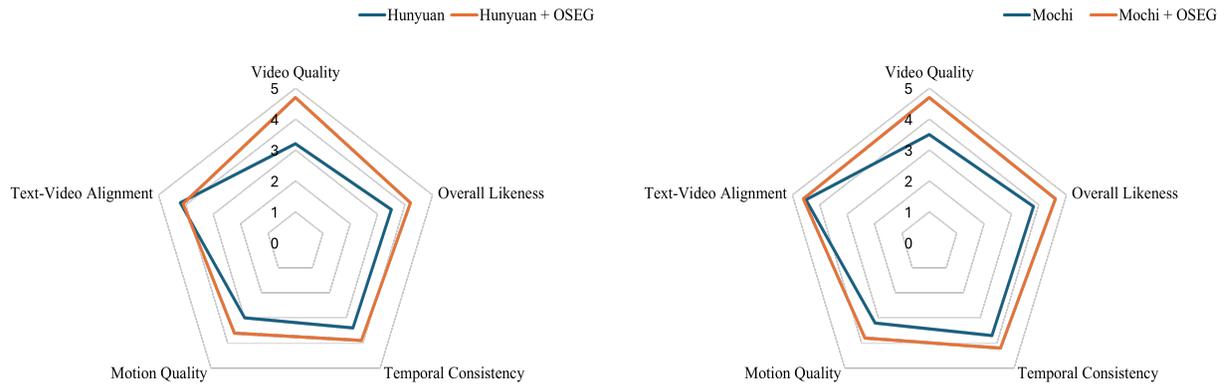
$$A_i = \frac{1}{n}.$$

Thus, we conclude the proof.

## 2. Human Evaluation

We follow the same procedure of [6] to evaluate subjective opinions from human interpretation which are designed across five key aspects: a) Video quality, where the higher scores indicates that the generated video consumes less blurriness, noise or unwanted artifacts; b) Text-video alignment, which highlights the alignment between the input text prompt and the generated video; c) Motion quality, which evaluates the correctness and realism of the motions depicted in the video; d) Temporal Consistency, measuring coherence between the frame-to-frame of the generated video to assess the smoothness of the movement, and e) Subjective Likeness which is an aesthetic score for human preferences.

For each metric of this experiment, we have collected feedback from five users who are privileged to rate the generalized video between scale 1 to 5, where higher scores represent better alignment. We use 100 prompts from Evalcrafter for T2V generation with Mochi [7] and Hunyuan [5].



(a) Human evaluation using Hunyuan [5]

(b) Human evaluation using Mochi [7]

Figure 1. User study results for OSEG on Mochi Hunyuan [5], and [7] using 100 prompts from EvalCrafter [6]. The results demonstrate that incorporating OSEG leads to improved video quality across all evaluated aspects.

### 3. Additional qualitative results

Figure 2 provides more visual comparison results of T2I generation.

#### 3.1. Conditional Image generation

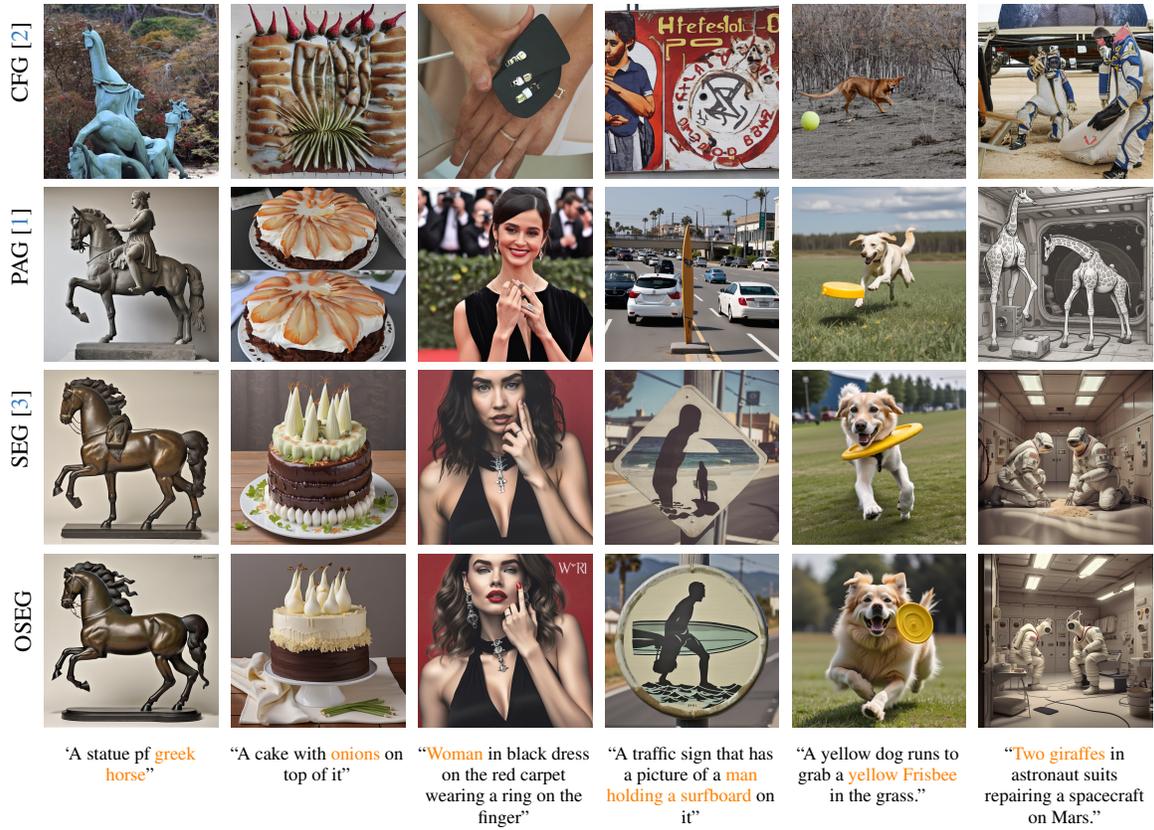


Figure 2. Qualitative comparison under conditional generation between CFG, PAG, SEG and OSEG.

## 4. Conditional Video generation

Figure 3 and Figure 4 provide additive T2V results generated by Mochi and Hunyuan, respectively, with the guidance of OSEG. Each figure provides a single frame of the videos, clicking on which will direct you to the original video link.

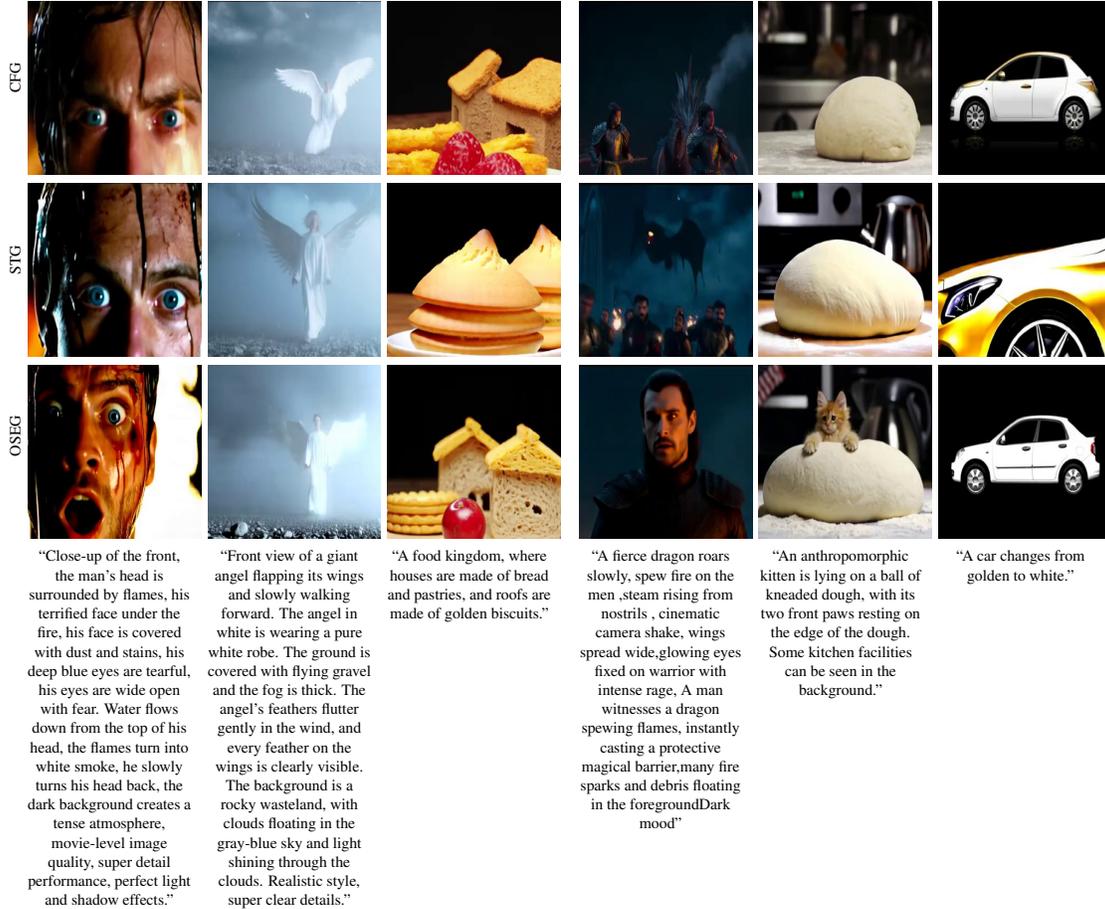


Figure 3. Conditional video generation comparison between CFG[2], STG[4], and OSEG by Hunyuan[5]. The first video produces a better scale of the human face, while the second video provides a clearer representation of the giant angel in thick fog. The third video demonstrates improved content preservation of golden biscuits. Similarly, OSEG generates the cat on the dough with better prompt alignment. Finally, both OSEG and CFG show similar appearance in the car video.



Figure 4. Conditional video generation comparison between CFG[2], STG[4], and OSEG by Mochi[7]. Compared to the existing SOTA guidance methods, incorporating OSEG alongside CFG noticeably improves spatial resolution with natural semantic coherence of the structures within the samples. This combination effectively strengthens fine-grained details and overall harmony, leading to high-resolution generations such as the modern microscopic view (first video) and better exposure (second video). In the third video, OSEG creates the structure of **ZE**, which addresses better prompt alignment. The remaining videos also provide better background details with diversity. For better visualization, please visit the website: <https://oseg-guidance.github.io/>

## References

- [1] Donghoon Ahn, Hyoungwon Cho, Jaewon Min, Wooseok Jang, Jungwoo Kim, SeonHwa Kim, Hyun Hee Park, Kyong Hwan Jin, and Seungryong Kim. Self-rectifying diffusion sampling with perturbed-attention guidance. In *European Conference on Computer Vision*, pages 1–17. Springer, 2024. 3
- [2] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3, 4, 5
- [3] Susung Hong. Smoothed energy guidance: Guiding diffusion models with reduced energy curvature of attention. *Advances in Neural Information Processing Systems*, 37:66743–66772, 2024. 3
- [4] Junha Hyung, Kinam Kim, Susung Hong, Min-Jung Kim, and Jaegul Choo. Spatiotemporal skip guidance for enhanced video diffusion sampling. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 11006–11015, 2025. 4, 5
- [5] 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. 2, 4
- [6] Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tiejong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large video generation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22139–22149, 2024. 2
- [7] Genmo Team. Mochi 1. <https://github.com/genmoai/models>, 2024. 2, 5