

Paper Critique: CogVideoX

Dan Peng
Department of Computer Science
Nov 25, 2024
danpeng@unc.edu

1. Research Problem

1.1. What research problem does the paper address?

This paper addresses the research problem that existing text-to-video generation models fail to produce multi movements with long duration. These models also struggle to generate videos with coherent narratives based on text.

1.2. What is the motivation of the research work?

To solve the challenges above, especially of long-term consistency with coherent dynamic contents, this paper proposes a new large-scale text-to-video generation model based on diffusion transformer which can generate 10s continuous videos with better alignment and fidelity.

2. Technical Novelty

2.1. What are the key technical challenges identified by the authors?

The key technical challenges are:

- How to efficiently consume high-dimension video data?
- How to solve misalignment between videos and texts?
- How to caption video content accurately?
- How to train the model progressively to further enhance the generation performance and stability of CogVideoX?

2.2. How significant is the technical contribution of the paper? If you think that the paper is incremental, please provide references to the most similar work

This paper offers the following contributions:

- Introduced a 3D VAE for video compression, which effectively reduces spatial and temporal dimensions while maintaining high fidelity and temporal consistency

- Developed a video captioning pipeline capable of accurately describing video content
- Proposed an expert Transformer with expert adaptive Layer-Norm to facilitate the fusion between the two modalities
- Adopt and design progressive training techniques, including multi resolution frame pack and resolution progressive training
- Achieved SOTA results across multiple benchmarks for text-to-video generation such as video fidelity, motion dynamics, and text-video alignment

2.3. Identify 1-5 main strengths of the proposed approach.

- The 3D VAE compresses video data across spatial and temporal dimensions, significantly reducing computational costs while preserving video fidelity and minimizing flickering
- The expert transformer design with adaptive LayerNorm enhances the integration of text and video modalities, ensuring improved alignment and coherent narrative generation
- The multi-resolution frame packing technique enables the model to process videos of different lengths and resolutions efficiently
- The developed video captioning and data filtering pipelines generate high-quality text-video pairs, addressing the issue of noisy datasets

2.4. Identify 1-5 main weaknesses of the proposed approach.

- The model's scalability to handle much longer (such as 1 minute) video durations or higher compression ratios is not fully explored, leaving room for improvement in extending its capabilities.

- Qualitative analysis of generated videos, though mentioned, is limited, and additional visual comparisons with other models could be a good supplement to show how well this model performed instead of just showing digits

3. Empirical Results

3.1. Identify 1-5 key experimental results, and explain what they signify.

- CogVideoX-5B achieves the highest scores in 4 out of 7 evaluation metrics, including Human Action, Multiple Objects, Dynamic Quality, and GPT4o-MT score, compared to competing models (Table 3). This result shows that CogVideoX can generate high-quality videos with rich motion dynamics and accurate semantic alignment.
- The proposed 3D VAE shows the lowest flickering score (85.5) and achieves a PSNR of 29.1, outperforming other spatiotemporal compression VAEs like Open-Sora, which scored 92.4 flickering and 28.5 PSNR. This highlights the 3D VAE's ability of reducing visual artifacts and improving video reconstruction quality (Table 2).
- Ablation studies demonstrate that replacing 2D + 1D attention with the proposed 3D full attention significantly improves model stability and reduces training steps (Figure 8b), validating the importance of comprehensive temporal and spatial modeling.

3.2. Are there any weaknesses in the experimental section (i.e., unfair comparisons, missing ablations, etc)?

The comparison with other models relies heavily on quantitative metrics, with limited qualitative visual examples provided to validate the performance arguments. It's unfair to showcase so many video generations by CogVideoX in the end while not producing the videos with same textual prompt produced by other models.

4. Summary

This paper presents a novel text-to-video generation model that addresses challenges like limited motions and text-video alignment in prior methods. By introducing a 3D VAE for efficient video compression, an expert transformer with adaptive LayerNorm for text-video alignment, and progressive training techniques, the model achieves most advancements across benchmarks as SOTAs. It produces coherent, high-quality videos with rich motion and accurate alignment to text prompts. I like CogVideoX because it effectively overcomes limitations in traditional joint-image

training by proposing the mixed-duration training. It allows videos of varying lengths and resolutions to be trained together. By utilizing the Multi-Resolution Frame Pack inspired by Patch'n Pack, the authors ensure consistent batch shapes despite differences in video durations. This innovative solution improves generalization by using diverse datasets.

5. QA Prompt for a Paper Discussion

5.1. Discussion Question

Will this method generate a video longer than 10s if the textual prompt is more than it expected, for example, a chapter of Harry Potter?

5.2. Your Answer

The CogVideoX is unlikely to directly generate a significantly longer video about a full chapter of Harry Potter based on an extended textual prompt. This is because it was trained with limitations designed for shorter video durations while maintaining higher quality and fidelity.