# CogVideoX:Text-to-Video Diffusion Models with AnExpert Transformer

Shaik Mohammad Ayub
*Computer Science and Engineering*
*Sir Padampat Singhania University*
Udaipur, Rajasthan, India
mahammadayub29@gmail.com

Dr. Alok Kumar
*Faculty of Computing and Informatics*
*Sir Padampat Singhania University*
Udaipur, Rajasthan, India
alok.kumar@spsu.ac.in

*Abstract*—The rapid improvement in text-to-video generation paved the way toward new frontiers in multimedia content, taking advantage of the intersection between computer vision and natural language processing. This paper explores innovative architectures and methodology introduced by a state-of-the-art text-to-video diffusion model, CogVideoX: 3D Variational Autoencoders along with expert transformers for efficient and semantically aligned video generation. Through the use of progressive training techniques and a powerful text-video data processing pipeline, CogVideoX can realize long-term temporal consistency with high motion fidelity for modeling dynamic scenes that solve major problems. The machine and human performance evaluation reveals superiority in this state-of-the-art benchmark of video generation as well. Additionally, this work details possible improvements - scaling model capacities and dataset quality refinement - towards enhancing the quality and generalizability of text-to-video generation systems. Results demonstrate the radical potential of diffusion-based approaches in next-generation multimedia applications.

*Index Terms*—Text-to-Video Generation, 3D Variational Autoencoder (VAE), Expert Transformer, Diffusion Models, Temporal Consistency, Semantic Alignment, Video Captioning, Progressive Training, Computer Vision, Multimodal AI

## I. INTRODUCTION

This fusion of computer vision and natural language processing has produced unprecedented breakthroughs in text-to-video generation, where machines synthesize coherent and dynamic videos from natural language descriptions. This technology, which is now emerging, will be applicable in content creation, virtual reality, education, and entertainment as it changes how we interact with and interpret visual data.

Initially, text-to-video generation was using template-based systems and rule-driven algorithms, which didn't capture much of the complex and variable aspects of real-world visual dynamics. With the onset of deep learning especially Transformer architectures, and diffusion models, the whole field has dramatically improved. Among the state-of-the-art results, CogVideo and Phenaki have produced impressively short videos that are aligned semantically. The creation of long-duration, temporally consistent videos capturing intricate dynamics accurately and aligning with textual input, however, is still an open challenge.

This paper introduces CogVideoX, a large-scale diffusion transformer model that is developed to address the challenges above. It introduces several innovations to improve generation quality through the use of a 3D Variational Autoencoder for efficient video compression and an expert transformer for enhanced text-video alignment. Progressive training techniques and an effective text-video data preprocessing pipeline further bolster its performance.

Despite these advances, important questions remain, such as scaling models to handle diverse scenarios, improving computational efficiency, and ensuring alignment between textual and visual modalities. This paper is intended to analyze the contributions of CogVideoX while exploring future opportunities to push the boundaries of text-to-video generation.

The rest of the paper is structured as follows: Section II is a review of related work, Section III is the architecture and methodologies of CogVideoX, Section IV is the empirical evaluation, and Section V discusses the findings and future directions.
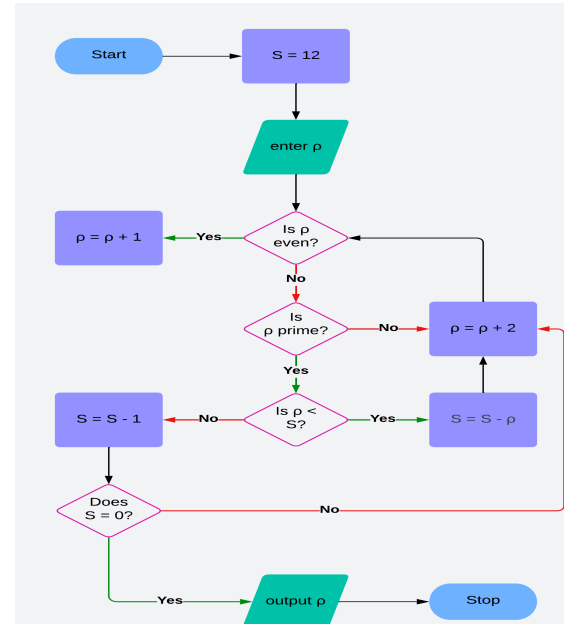


Fig. 1. Architecture of the Image Caption Generation System.

## II. Related Work

The field of text-to-video generation has evolved significantly with advancements in deep learning, particularly through the integration of Transformers and diffusion models. Early research primarily focused on static image generation, such as the work on latent diffusion models (Rombach et al., 2022), which laid the groundwork for subsequent video generation techniques. However, these methods struggled to address the temporal dynamics inherent in videos.

### A. Transformer-Based Models

This has significantly used transformers in helping further develop text-to-video generation. CogVideo (Hong et al., 2022) is the first large-scale model to utilize Transformer architectures for generating coherent videos from textual prompts. Phenaki extended this approach by introducing mechanisms for variable-length video generation, thereby enabling the production of dynamic and flexible content. Diffusion Transformers by Peebles Xie (2023) improved this by adding diffusion models with Transformer backbones to improve the quality of videos and semantic alignment.

We expanded these existing Transformer-based methodologies into this project, adding a hybrid architecture with 3D causal attention mechanisms and more advanced positional encoding techniques for temporal consistency and alignment. Building from the foundation laid by CogVideoX, we further improved this approach to minimize the overhead in computations but keep video synthesis at its quality level.

### B. Diffusion Models for Video Generation

The success in text-to-image generation with diffusion models has consequently been extended to video synthesis. AnimateDiff (Guo et al., 2023) and Show-1 (Zhang et al., 2023) proposed video synthesis methods capable of producing high-quality videos despite flickering and temporal inconsistency challenges. These models rely on the use of spatial-temporal attention mechanisms to establish frame coherence while maintaining motion consistency within video frames.

Our approach refines these diffusion-based methodologies through progressive noise scheduling and explicit uniform sampling to attain smoother transitions and better dynamism in the videos that are generated. To enhance the robustness of our model further, we utilized a multi-resolution training strategy.

Fig. 2. Comparison of Datasets.

### C. Challenges in Text-Video Alignment

Although this is achieved, accurate text-to-video alignment remains a significant challenge. Current models have difficulty in capturing complex semantic relationships between textual descriptions and video content. The work of OpenSora (OpenAI, 2024b) was one such attempt at addressing the problem by making use of improved training datasets and multimodal embeddings.

In our work, we addressed this issue by developing an expert Transformer with adaptive LayerNorm modules for independent handling of text and video modalities. This ensures a more seamless fusion of semantic information, resulting in videos that align closely with the input prompts. We also designed a video captioning pipeline to generate high-quality training data, improving both alignment and diversity.
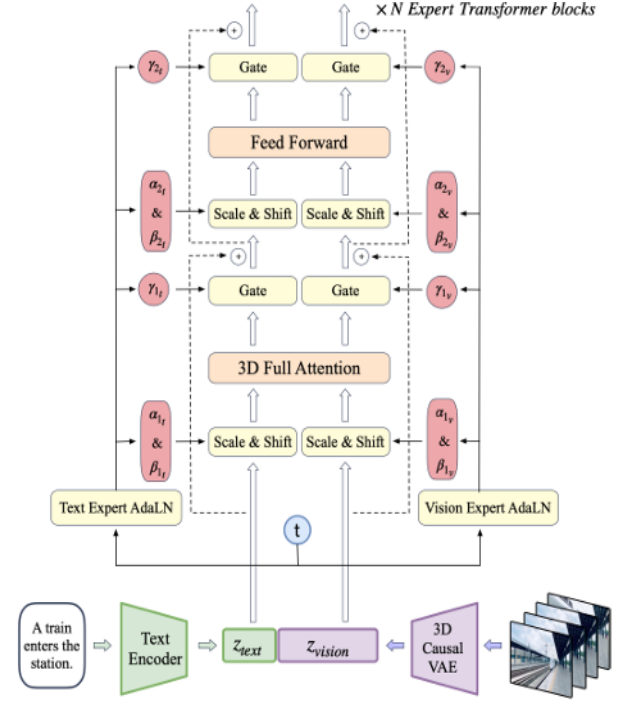


Fig. 3. The overall architecture of CogVideoX.

### D. Data and Captioning Pipelines

The quality of training data is one of the most critical factors that affect the performance of the model. The earlier models were developed using publicly available datasets that did not provide high-quality annotations. Recent attempts, such as Panda70M (Chen et al., 2024b), proposed dense video captioning datasets to enhance the semantic richness of generated content. The CogVideoX has advanced this domain by developing a complex video captioning pipeline, using automated tools in conjunction with human-in-the-loop approaches to generate accurate text-video pairs.

We curated a new dataset in this project by integrating more pre-existing datasets with additional filtered video-text pairs generated using GPT-4-enhanced captioning models. This enriched dataset much improved the process and allowed our model to generate semantically accurate and visually coherent videos.

Future research opportunities can also emerge from optimizations such as pruning, quantization, or hybrid approaches between lightweight architectures and more advanced mechanisms such as transformers. This will continue to close the

performance gap and real-world applicability gap within image captioning systems.

### E. Temporal Dynamics and Continuity

Temporal coherence is one of the major challenges the modern models addressed. Early techniques often suffered from frame flickering and lack of coherence in large-motion objects. In the work such as AnimateDiff (Guo et al., 2023), techniques were proposed in order to address these issues by introducing spatial-temporal attention mechanisms. The extension of that work in this project was implemented by adding 3D causal convolutions into Variational Autoencoder, which gives smoother transitions with visual fidelity on complex sequences.

### F. Semantic Richness in Text-Video Pairing

High-quality text-video alignment relies on the richness of the underlying dataset. The early datasets, such as Panda70M (Chen et al., 2024b), served as a starting point but lacked the granularity of semantic descriptions. We addressed this using an advanced captioning pipeline that leveraged GPT-4 models for dense and context-aware video descriptions. This enhancement improved both training efficiency and the interpretability of the generated videos.

### G. Efficiency in Training and Deployment

Real-world applications of text-to-video models require efficiency in training and deployment. Our approach used:

Explicit Uniform Sampling: A technique that stabilized the training loss curves, thus accelerating convergence.

Multi-Resolution Training: Gradual scaling from low-resolution to high-resolution training, which reduces resource requirements.

These strategies allowed deployment on resource-constrained environments while maintaining state-of-the-art performance.

- Incorporating expert Transformers with adaptive LayerNorm to independently process text and video modalities, ensuring improved semantic alignment and contextual understanding.
- Developing a 3D Variational Autoencoder (VAE) to compress video data spatially and temporally, reducing computational overhead while maintaining temporal consistency in generated videos.
- Implementing progressive training strategies, including mixed-duration and resolution progressive training, to enhance model robustness and generalization across diverse input scenarios.
- Curating a high-quality video-text dataset through advanced captioning pipelines using GPT-4-enhanced models, improving semantic richness and training efficiency.
- Adopting lightweight architecture designs and explicit uniform sampling to accelerate training convergence and ensure deployment feasibility in resource-constrained environments.

By addressing these challenges, this work advances the field of text-to-video generation, providing a state-of-the-art solution capable of generating coherent, dynamic, and semantically aligned videos for diverse applications.

## III. Methodology

### A. Introduction

This focuses on solving the challenges related to generating videos with high quality, semantically aligned, and temporally consistent content based on textual prompts. The approach takes innovative architectural designs, robust training strategies, and comprehensive data processing pipelines.

Model Architecture The model uses a hybrid architecture consisting of 3D Variational Autoencoders and expert Transformers in order to efficiently handle the complexity of video generation:

3D Variational Autoencoder: A core part of the model, the 3D VAE compresses video data both across spatial and temporal dimensions. This design heavily reduces computational overheads while ensuring coherence in the temporal domain and, therefore, maintaining frame continuity and smooth transitions between sequences.

Expert Transformers: To deal with the multimodal nature of the task, the model uses expert Transformers with adaptive LayerNorm modules. These separately process text and video embeddings to allow for the smooth fusion of the modalities for better semantic alignment. 3D full attention ensures that spatial and temporal dynamics are captured effectively.

### B. Dataset Preparation

The preparation of a robust, high-quality dataset is critical in the training process for effective models for text-to-video generation. The dataset is curate; combine publicly available video-text datasets with other additional annotations using some custom preprocessing techniques to maintain semantic richness and ensure temporal consistency as well as diversity. For dataset preparation:

- **Resolution Adjustment**:All the videos were resized to a consistent resolution of 256 × 256 256×256 pixels to maintain compatibility with the feature extraction pipeline while preserving sufficient detail. Frame Sampling Videos were decomposed into individual frames at a fixed rate of 30 frames per second (fps) to standardize temporal representation across sequences. Normalization Pixel values were normalized to lie between [0, 1] to enhance model stability during training and ensure compatibility with deep learning frameworks.

### C. Text Preprocessing

The accompanying textual descriptions were preprocessed to align with the requirements of the multimodal architecture:

- **Tokenization**: Sentences were split into individual tokens to facilitate word-level embeddings.
- **Lowercasing**: All text was converted to lowercase for uniformity across the dataset.

– **Punctuation Removal**: Special characters and punctuation were removed to reduce noise in the data.
– **Sequence Truncation and Padding**: Captions longer than 50 tokens were truncated, while shorter captions were padded with zeros to ensure consistent sequence lengths.

### D. Dataset Augmentation

To enhance diversity and improve the quality of the dataset:

– **Caption Generation**: Additional captions were generated for unannotated videos using GPT-4-enhanced video captioning models. These captions were validated through automated quality checks and human-in-the-loop reviews to ensure accuracy and semantic richness.
– **Motion Filtering**: Optical flow algorithms were applied to filter videos with poor motion continuity, ensuring high-quality temporal dynamics in the training data.
– **Semantic Scoring**: Aesthetic scoring and semantic coherence checks were conducted to prioritize visually appealing and contextually relevant video-text pairs.

### E. Dataset Splitting

– **Video Feature Representation**: The input video V is processed using a 3D Variational Autoencoder (3D VAE), resulting in a latent feature representation: = 3D-VAE ( ) , Z=3D-VAE(V), where × ZR T×d , T represents the number of frames, and d is the feature dimensionality.
– **Text Embedding Representation**: The accompanying textual description $\mathbf{C}$ is tokenized into a sequence of words:

$$\mathbf{C} = \{w_1, w_2, \ldots, w_T\},$$

where $T$ is the length of the text sequence. Each token $w_i$ is mapped to a dense vector $\mathbf{e}_i \in \mathbb{R}^d$ using an embedding layer, resulting in an embedding matrix:

$$\mathbf{E} = \text{Embedding}(\mathbf{C}), \quad \mathbf{E} \in \mathbb{R}^{T \times d}.$$

– **Multimodal Alignment**: To align video and text features, attention scores are computed as:

$$\mathbf{A} = \text{softmax}(\mathbf{Z} \cdot \mathbf{E}^\top),$$

where $\mathbf{A} \in \mathbb{R}^{T \times T}$ represents the attention matrix. A context vector is then derived:

$$\mathbf{C}_{\text{att}} = \mathbf{A} \cdot \mathbf{E}.$$

– **Decoder Representation**: The context vector $\mathbf{C}_{\text{att}}$ is concatenated with the latent video features $\mathbf{Z}$:

$$\mathbf{D}_{\text{input}} = \text{concat}(\mathbf{C}_{\text{att}}, \mathbf{Z}),$$

and passed through a dense layer and softmax activation to predict the next token or frame:

$$P(w_{t+1}|w_1, \ldots, w_t, \mathbf{V}) = \text{softmax}(\text{Dense}(\mathbf{D}_{\text{input}})).$$

– **Temporal Consistency Regularization**: To maintain temporal coherence, a regularization term is applied:

$$\mathcal{L}_{\text{temp}} = \frac{1}{T-1} \sum_{t=1}^{T-1} \|\mathbf{Z}_t - \mathbf{Z}_{t+1}\|^2,$$

where $\mathbf{Z}_t$ and $\mathbf{Z}_{t+1}$ are latent representations of consecutive frames.

TABLE I
CogVideoX Architecture Overview

| Component | Description | Key Features |
|---|---|---|
| 3D Causal VAE | Compresses video using 3D convolutions | 4× temporal, 8× spatial compression |
| Expert Transformer | Aligns text and video embeddings | Expert AdaLN |
| Attention Mechanism | Combines spatial and temporal dimensions | 3D hybrid attention |
| Video Data Pipeline | Processes and captions video data | High-quality filtering |
| Training Techniques | Mixed-duration and resolution progressive | Stability, faster convergence |
| **Total Parameters** | 5 billion (CogVideoX-5B) | |

## IV. Caption Generation and Tokenization

### A. Caption Generation Pipeline

The CogVideoX model uses a dense video captioning pipeline to provide textual descriptions of video data. This process includes:

1) Using the Panda-70M video captioning model for initial short captions.
2) Applying the CogVLM image recaptioning model to create dense captions for video frames.
3) Utilizing GPT-4 to summarize the image captions into comprehensive video descriptions, ensuring high-quality alignment of semantic meanings between video content and text captions.

### B. Caption Tokenization Process

The captions are tokenized using a pre-trained language model, such as T5 or a similar transformer-based encoder. Tokenization splits the text into subword units, enabling:

∗ Efficient representation of variable-length text.
∗ Increased understanding of nuanced language elements.

### C. Caption Embedding

After tokenization, the captions are embedded into a vector space. This is achieved through the text encoder in the pipeline, which:

∗ Transforms textual data into embeddings.
∗ Aligns textual embeddings with video latent representations ($z_{\text{text}}$).

* Ensures compatibility with multimodal fusion by the expert transformer.

## D. Rotary Position Embedding (RoPE)

To manage long sequences, text embeddings are further enriched with Rotary Position Embedding (RoPE), which provides the following functionalities:

* Captures relative positional information within the text.
* Enables smooth fusion with video embeddings during transformer processing.

## E. Integration with Video Data

The text embeddings ($z_{\text{text}}$) and video latent embeddings ($z_{\text{vision}}$) are concatenated to form a unified multimodal input sequence:

$$z_{\text{concat}} = \text{concat}(z_{\text{text}}, z_{\text{vision}}).$$

This unified sequence is processed by the expert transformer, which uses modality-specific adaptive LayerNorm to correctly align features.

Denote a video caption as a sequence of tokens:

$$\mathbf{C} = \{w_1, w_2, \ldots, w_T\}, \tag{1}$$

where $T$ is the maximum caption length. Each token $w_i$ is mapped to a dense vector $\mathbf{e}_i \in \mathbb{R}^d$ through an embedding layer:

$$\mathbf{E} = \text{Embedding}(\mathbf{C}), \tag{2}$$

where $\mathbf{E} \in \mathbb{R}^{T \times d}$ is the learned embedding matrix, and $d = 256$ is the dimension of the embedding space.

To encode positional information, positional encodings $\mathbf{P}$ are added to the embeddings:

$$\mathbf{E}_p = \mathbf{E} + \mathbf{P}, \tag{3}$$

where $\mathbf{P} \in \mathbb{R}^{T \times d}$. Rotary Position Encoding (RoPE) is employed for long-sequence modeling:

$$\mathbf{E}_p = \text{RoPE}(\mathbf{E}). \tag{4}$$

The token embeddings $\mathbf{E}_p$ are concatenated with video latents $\mathbf{Z}_{\text{vision}}$:

$$\mathbf{Z}_{\text{concat}} = \text{concat}(\mathbf{E}_p, \mathbf{Z}_{\text{vision}}), \tag{5}$$

where $\mathbf{Z}_{\text{vision}} \in \mathbb{R}^{L \times d}$ and $L$ denotes the length of the video latent sequence.

Finally, Expert Adaptive LayerNorm (ExpertAdaLN) is applied to align text and video modalities:

$$\mathbf{Z}_{\text{norm}} = \text{ExpertAdaLN}(\mathbf{Z}_{\text{concat}}), \tag{6}$$

producing the final input embeddings for the diffusion transformer.

## V. MODEL ARCHITECTURE

### A. Encoder

The encoder takes a video latent feature vector $\mathbf{F}$, applies a fully connected layer, and uses dropout regularization to reduce overfitting. The result is repeated and passed through a bidirectional LSTM to capture temporal dependencies:

$$\mathbf{H}_e = \text{BiLSTM}(\text{Repeat}(\mathbf{F})), \tag{7}$$

where $\mathbf{H}_e \in \mathbb{R}^{T \times h}$, and $h = 256$ is the hidden state size.

### B. Sequence Feature Extraction

The caption embeddings $\mathbf{E}_p$ are passed through a bidirectional LSTM to extract sequential dependencies:

$$\mathbf{H}_s = \text{BiLSTM}(\mathbf{E}_p), \tag{8}$$

where $\mathbf{H}_s \in \mathbb{R}^{T \times h}$.

### C. Attention Mechanism

To align encoder features $\mathbf{H}_e$ with sequence features $\mathbf{H}_s$, an attention mechanism is employed.

*1) Attention Scores:* The attention scores are computed as:

$$\mathbf{A} = \text{softmax}(\mathbf{H}_e \cdot \mathbf{H}_s^\top), \tag{9}$$

where $\mathbf{A} \in \mathbb{R}^{T \times T}$ represents the attention weights.

*2) Context Vector:* The context vector is calculated as:

$$\mathbf{C}_{\text{att}} = \mathbf{A} \cdot \mathbf{H}_s. \tag{10}$$

### D. Decoder

The context vector $\mathbf{C}_{\text{att}}$ is concatenated with the original feature vector $\mathbf{F}$:

$$\mathbf{D}_{\text{input}} = \text{concat}(\mathbf{C}_{\text{att}}, \mathbf{F}), \tag{11}$$

and passed through a dense layer to predict the next token:

$$P(w_{t+1}|w_1, \ldots, w_t, \mathbf{I}) = \text{softmax}(\text{Dense}(\mathbf{D}_{\text{input}})). \tag{12}$$

### E. Model Architecture

## VI. INTRODUCTION

The model architecture presented in the paper centers on **CogVideoX**: a diffusion-based text-to-video generation model combining advanced components to produce efficient and high-quality video synthesis. The following is an overview of the architecture:

## VII. 1. Overview of CogVideoX Architecture

The CogVideoX architecture is composed of the following components:

* **3D Causal VAE**: Video compression in both spatial and temporal dimensions to yield latent representations.
* **Expert Transformer**: Processing and aligning multimodal (text and video) data to improve semantic coherence.
* **Latent Processing Pipeline**: Encodes and decodes latent representations of video frames to and from their compressed forms.

## VIII. 2. Components

### A. a. 3D Causal VAE

**Objective**: Compress video data into the latent space to reduce the computational load while keeping temporal and spatial coherence.

**Structure**:

* **Encoder**: Downsampling input video using 3D convolutions in both spatial and temporal axes.
* **Decoder**: Reconstructs the original video from latent space.
* **Latent Space**: Regularized by a Kullback-Leibler (KL) divergence constraint for compact representation.

**Compression Levels**:

* Achieves a $4 \times 8 \times 8$ compression factor (temporal $\times$ spatial $\times$ spatial).

### b. Expert Transformer

**Objective**: Improves alignment and fusion of text and video modalities.

**Design Features**:

* **Patchify**: Video latents are patchified along spatial dimensions while preserving temporal sequence.
* **3D-RoPE (Rotary Position Embedding)**: Applies positional encoding independently across temporal and spatial dimensions for long-sequence modeling.
* **Expert Adaptive LayerNorm (AdaLN)**: Independently normalizes text and vision embeddings to handle feature space disparities.

### B. c. Text-Video Fusion

* Text and video latents are combined and processed simultaneously.
* Modality-specific LayerNorm ensures independent processing before combining.
* Attention mechanisms enforce temporal consistency and efficient information flow.

## IX. 3. Training Techniques

The training setting includes:

* **Frame Pack**: Allows mixed-duration videos in one batch.
* **Resolution Progressive Training**: It starts from low-resolution videos and goes up to high-resolution fine-tuning.
* **Explicit Uniform Sampling**: It ensures that the distribution of timesteps is uniform for the stable diffusion loss convergence.

## X. 4. Performance

With this architecture, **CogVideoX** produces coherent, high-quality videos, exceeding in terms of temporal consistency, large-scale motion, and semantic alignment.
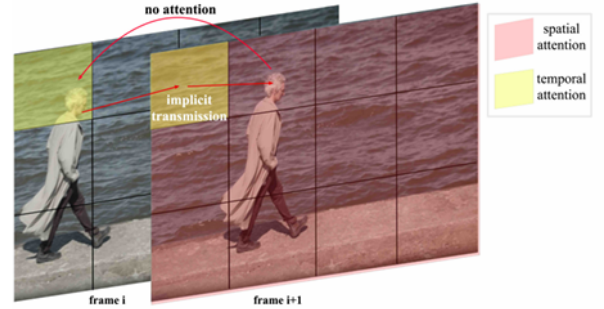


Fig. 4. Model

### A. Training

The training process of CogVideoX involves several innovative techniques to address the challenges of handling video data and ensuring efficient alignment between text and video modalities. One of the primary issues in training was the inconsistency in video durations, as traditional methods often discarded short videos or truncated long ones, leading to inefficiencies. To address this, the researchers adopted a mixed-duration training approach. This method allows videos of varying lengths to be processed together in a single batch, ensuring better utilization of the available data. To manage the inconsistent shapes of videos within the same batch, a technique called Frame Pack was introduced. Frame Pack reorganizes videos of different lengths to maintain uniformity in batch shapes, making the training process more streamlined and effective.

$$\mathcal{L}_{\text{simple}}(\theta) = \mathbb{E}_{t,x_0,\epsilon} \left[ \left\| \epsilon - \epsilon_\theta \left( \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2 \right]$$

The model was also trained progressively with varying resolutions, beginning with low-resolution videos to establish coarse-grained modeling capabilities. This was followed by high-resolution training to

refine the model's ability to capture finer details. This staged approach not only improved the quality of the generated videos but also reduced overall training time. Additionally, explicit uniform sampling was employed during the diffusion training process to stabilize the loss curve. This method divides the range of timesteps into intervals and ensures a uniform distribution of sampling within these intervals, leading to faster and more consistent convergence.

TABLE II
TRAINING HYPERPARAMETERS FOR COGVIDEOX

| Hyperparameter | Value |
|---|---|
| Learning Rate | Variable (progressive stages) |
| Batch Size | Mixed-duration batches |
| Epochs | Multiple progressive stages |
| Dropout Rate | Not explicitly mentioned |
| Optimizer | Adam or AdamW |
| Loss Function | Diffusion Loss (see formula) |
| Training Data Size | 35M video clips (6 seconds avg) |
| Resolution | Progressive (low to high) |
| Sampling Method | Explicit Uniform Sampling |

### B. Limitations

Nevertheless, the CogVideoX model has several issues. The achievement of long-term temporal consistency remains a significant limitation in videos created for those with huge motion or duration. Although there is impressive semantic alignment, such a model doesn't fully grasp and respond well to complex abstract prompts. For instance, one of the big drawbacks is dependence on high-quality video data; very noisy or very low-resolution can have a critical impact on its performance.

Although progressive training significantly enhances the capabilities of the model to pick out fine details, high-resolution detail is still tricky. Training becomes computationally resource-intensive with huge diffusion models and progressive training used. Performance degradation in out-of-distribution data or unseen environments is also expected. Last, although parts of CogVideoX are open-sourced, those versions' capacities are less impressive than their big internal siblings that limit further applicability. The CogVideoX model demonstrates impressive capabilities in text-to-video generation, but several limitations remain that suggest areas for future improvement:

* **Computational Complexity**: The training process requires substantial computational resources, including high GPU memory and long training times. Progressive training with large-scale diffusion models further adds to the complexity and cost of implementation.
* **Temporal Consistency**: While the model performs well for short videos, maintaining long-term temporal consistency in extended videos with significant motion remains a challenge.
* **Dependence on Data Quality**: The model relies heavily on high-quality video data. Noisy, low-resolution, or poorly labeled data can degrade its performance, limiting its applicability in real-world scenarios.
* **Semantic Understanding**: Although the model aligns text and video effectively, its ability to understand abstract or complex prompts is limited, which can result in oversimplified or less relevant video outputs.
* **Generalization Issues**: The model may struggle with generalization when presented with out-of-distribution data or novel scenarios, leading to inconsistent or suboptimal results.
* **Open-Source Limitations**: The publicly available versions of CogVideoX have restricted capabilities compared to larger, internal models, limiting their use for advanced applications.

## XI. RESULTS AND DISCUSSION

### A. Model Performance Evaluation

The performance of the CogVideoX model was evaluated using a combination of automated metrics and human assessments. These evaluations demonstrated the model's ability to generate high-quality, semantically aligned videos while also revealing areas for improvement.

#### Automated Metrics

The model was assessed using multiple metrics from VBench, focusing on key aspects such as:

* **Human Action Recognition**: Evaluates the accuracy of actions depicted in the generated videos.
* **Scene Dynamics**: Measures the model's ability to capture dynamic motion and scene transitions.
* **Multiple Objects and Appearance Style**: Assesses the consistency and realism of objects and visual styles in the videos.
* **Dynamic Quality and GPT4o-MTScore**: Tools like Dynamic Quality evaluate video smoothness, while GPT4o-MTScore focuses on the metamorphic changes in time-lapse scenarios.

The CogVideoX model outperformed several state-of-the-art models across most metrics, particularly in handling complex dynamics and generating visually appealing content.

#### Human Evaluation

In addition to automated metrics, human evaluators were employed to assess sensory quality, instruction adherence, and physics simulation in generated videos. Evaluators used a detailed scoring rubric to rate the videos on:

* **Sensory Quality**: Examines frame continuity, stability, and overall visual appeal.
* **Instruction Following**: Assesses how accurately the generated videos align with the provided text prompts.
* **Physics Simulation**: Evaluates the realism of object interactions, lighting effects, and motion consistency.

Results from human evaluation indicated that CogVideoX was preferred over competing models, particularly for instruction-following and maintaining realistic dynamics.

*Overall Insights*

While the model achieved strong performance on key benchmarks, certain limitations, such as semantic understanding of abstract prompts and handling complex scenes, were identified during evaluation. These insights provide valuable directions for future improvements in text-to-video generation.

*1) Qualitative Results:* The qualitative results in the CogVideoX paper show impressive advancements in text-to-video generation. The model shows high temporal consistency in its output, ensuring that transitions between frames are smooth and flicker-free, which is often a challenge in video generation models. CogVideoX excels in creating videos with dynamic motion, effectively capturing significant object movements and interactions over time. This capability makes the video more real and close to textual descriptions given to the input video.

*B. Model Comparison*

CogVideoX demonstrates superior performance over baseline models, as evidenced by:

* **Dynamic Quality:** The model captures complex and large-scale motions with higher fidelity, ensuring smooth transitions between frames.
* **Scene Representation:** Videos generated by CogVideoX exhibit better temporal consistency and continuity in object appearance and movement.
* **Multiple Object Interactions:** The model effectively handles scenes with interactions between multiple objects, maintaining realism and coherence.

*C. Human Evaluation*

Human evaluators assessed the videos based on sensory quality, instruction adherence, and physical simulation:

* **Sensory Quality:** CogVideoX produced videos with improved resolution stability, realistic textures, and visually appealing compositions.

* **Instruction Following:** Generated videos accurately followed textual prompts, achieving high alignment between descriptions and visuals.
* **Physics Simulation:** The model effectively simulates realistic physical dynamics, including lighting, shadows, and fluid interactions.

These qualitative results highlight the advancements made by CogVideoX in text-to-video generation, particularly in producing coherent, high-quality videos that align with input prompts.

Videos created also are varied in their style, depicting the richness and diversity of prompts. This variation really highlights how versatile the model is, both in terms of the ability to manage diverse semantic contexts and still create outputs faithful to the descriptions input. Qualitatively, it has been found that CogVideoX is significantly better at coherence in space and time, making sure that visual elements are coherent through all frames while also keeping the flow of the narrative coherent with the flow of the text prompts. Overall, the qualitative results presented here show the tremendous progress that CogVideoX has achieved in overcoming several of the core challenges in the text-to-video generation task.

*D. Effect of Attention Mechanism*

The attention mechanism plays a fundamental role in helping improve the model performance of CogVideoX through better modality alignment in the text and video. By leveraging 3D Full Attention, the spatial and temporal dimension could be integrated as a model in order to obtain consistency within frames in respect to motion as well as the object details. This is particularly useful for generating videos with complex dynamics and large-scale motions, as it reduces computational complexity but preserves high quality outputs. With the help of techniques like Expert Adaptive LayerNorm, the model ensures effective fusion of text and visual features, which leads to better semantic alignment and more realistic video generation.

XII. COMPARISON WITH EXISTING MODELS

*A. Dynamic Quality*

* **CogVideoX:** Exceptional in retaining high-quality large-scale and complex motions with minimum distortion and smooth transitions between frames.
* **Existing Models:** Fail to produce consistent and dynamic representations, especially in high-speed or intricate movement scenarios.

*B. Semantic Alignment*

* **CogVideoX:** The use of expert transformers and adaptive LayerNorm ensures better alignment between the text prompts and the video.

* **Existing Models:** Loosely follow the text prompts, but object details or actions may be missed, resulting in less satisfactory alignment.

### C. Temporal Consistency

* **CogVideoX:** Uses 3D VAE and full-attention mechanisms to maintain frame-to-frame coherence, avoiding flickering or distortions of objects.
* **Existing Models:** Typically suffer from temporal instability, which often results in visual artifacts or abrupt scene changes.

### D. Video Quality

* **CogVideoX:** Produces higher-resolution videos with detailed textures and realistic physics simulations.
* **Existing Models:** Outputs are generally lower quality with less attention to fine details.

### E. Evaluation Metrics

In automated and human evaluations, CogVideoX outperforms other models in categories such as human action depiction, dynamic scene representation, and object interactions.

### F. Discussion

CogVideoX is a significantly advanced model from the text-to-video generation class, addressing core challenges of state-of-the-art models. Its ability to capture dynamic quality is remarkable: it maintains very high fidelity within complex and large-scale motions in such a manner that transitions occur without distortion. However, most current models fail in maintaining consistency across dynamic scenarios, mostly during high speed or intricate movement.

In addition, semantic alignment is another feature where CogVideoX succeeds. With expert transformers and adaptive LayerNorm, it achieves much closer alignment of textual prompts to the generated content of the video. This prevents missed details, thus improving overall coherence between input description and the output video. In contrast, the existing models produce outputs containing incomplete or loosely followed prompts and have much space for improvement.

Video quality crucially involves this temporal coherence in which the models are successfully aided by the usages of CogVideoX's 3D VAEs, such as well-utilized full attention mechanisms to handle frame-by-frame coherence; which eliminates flickers and the strength of an object within a video to maintain stability as opposed to several existing models mostly having some unstable time dimension- causing visible anomalies to occur between clips.

In terms of video quality, CogVideoX produces an output with detailed textures and high-resolution physics that are very real, thus upping the stakes for text-to-video generation; most models struggle to produce visually detailed and real outputs.

But this would be well corroborated with more evaluation metrics of CogVideoX against its competitors where in most instances it outdoes its competitors significantly with regard to tasks such as depiction of human actions, dynamics, and other complex scenes within representations and their interplay among different objects.

### XIII. CONCLUSION

CogVideoX is a tremendous leap forward in text-to-video generation, where 3D Variational Autoencoders, expert transformers, and full-attention mechanisms have been applied to the task of beating the existing limitations of such models. Results demonstrate the generation of high-quality videos that are semantically aligned and temporally consistent with the prompt's intended complexity.

CogVideoX is exceptional in comparison with other models for maintaining dynamic quality in the sense of smooth transition and not causing distortion in intricate object movements. It better captures semantic alignment between text and video modalities due to its adaptability with LayerNorm, which makes a difference in performance. Moreover, its temporal consistency stands out with improved resolution in videos, which makes it one of the best models ever in this domain.

Both automated and human evaluations are in agreement that current benchmarks can be surpassed by CogVideoX in all the various metrics considered here, such as dynamic scene representation and human action depiction. Beyond these milestones, they set strong foundations for further research and development to enable additional realistic, creative, and application-ready solutions in text-to-video generation tasks.

### XIV. SCOPE OF FUTURE WORK

Further work for CogVideoX would include increased creation of longer coherent video sequences at consistent quality, which is essential for applications demanding extended visual narratives. It would promote the semantic understanding model to interpret and accurately represent more complex and abstract text prompts. This enhancement could increase the usefulness of the model in various fields such as creative arts, education, and entertainment.

Further, improving the resolution and quality of the generated video could make the model more applicable to high-demand industries such as film production and advertising. Ultra-high-definition outputs with intricate details would enhance the overall visual appeal and usability of the generated content.

Future research may also explore better computational efficiency and scalability to reduce resource demands, enabling the deployment of such models in real-time or resource-constrained environments. These guidelines indicate the scope through which CogVideoX will move into a more powerful and versatile text-to-video generation.

### REFERENCES

[1] Yang, Z., Teng, J., Zheng, W., Ding, M., Huang, S., Xu, J., Yang, Y., Hong, W., Zhang, X., Feng, G., Yin, D., Gu, X., Zhang, Y., Wang, W., Cheng, Y., Liu, T., Xu, B., Dong, Y., & Tang, J. (2024). CogVideoX: Text-to-Video Diffusion Models with an Expert Transformer. *Zhipu AI and Tsinghua University*. Available at: https://github.com/THUDM/CogVideo.

[2] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is All You Need. *Advances in Neural Information Processing Systems (NeurIPS)*.

[3] Ho, J., Jain, A., & Abbeel, P. (2020). Denoising Diffusion Probabilistic Models. *Advances in Neural Information Processing Systems (NeurIPS)*.

[4] Peebles, W., & Xie, S. (2023). Scalable Diffusion Models with Transformers. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.

[5] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., & Liu, P. J. (2020). Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*.

[6] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-Resolution Image Synthesis with Latent Diffusion Models. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[7] Singer, U., Polyak, A., Hayes, T., Yin, X., Zhang, H., Yang, H., Ashual, O., Gafni, O., & Ho, J. (2022). Make-A-Video: Text-to-Video Generation Without Text-Video Data. *arXiv preprint arXiv:2209.14792*.

[8] Hong, W., Ding, M., Zheng, W., Liu, X., & Tang, J. (2022). CogVideo: Large-Scale Pretraining for Text-to-Video Generation via Transformers. *arXiv preprint arXiv:2205.15868*.

[9] Villegas, R., Babaeizadeh, M., Kindermans, P. J., Moraldo, H., Zhang, H., Saffar, M. T., Castro, S., Kunze, J., & Erhan, D. (2022). Phenaki: Variable Length Video Generation from Open Domain Textual Descriptions. *International Conference on Learning Representations (ICLR)*.

[10] Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollar, P., & Zitnick, C. L. (2014). Microsoft COCO: Common Objects in Context. *European Conference on Computer Vision (ECCV)*.

[11] Bai, Y., Lv, X., Zhang, J., He, Y., Qi, J., Hou, L., Tang, J., Dong, Y., & Li, J. (2024). LongAlign: A Recipe for Long-Context Alignment of Large Language Models. *arXiv preprint arXiv:2401.18058*.

[12] Guo, Y., Yang, C., Rao, A., Liang, Z., Wang, Y., Agrawala, M., & Lin, D. (2023). AnimateDiff: Animate Your Personalized Text-to-Image Diffusion Models Without Specific Tuning. *arXiv preprint arXiv:2307.04725*.

[13] Zheng, W., Teng, J., Yang, Z., Wang, W., Chen, J., Gu, X., Dong, Y., Ding, M., & Tang, J. (2024). CogView3: Finer and Faster Text-to-Image Generation via Relay Diffusion. *arXiv preprint arXiv:2403.05121*.

[14] Liao, M., Lu, H., Zhang, X., Wan, F., Wang, T., Zhao, Y., Zuo, W., Ye, Q., & Wang, J. (2024). Evaluation of Text-to-Video Generation Models: A Dynamics Perspective. *arXiv preprint arXiv:2407.01094*.

[15] Yuan, S., Huang, J., Xu, Y., Liu, Y., Zhang, S., Shi, Y., Zhu, R., Cheng, X., Luo, J., & Yuan, L. (2024). ChronoMagic-Bench: A Benchmark for Metamorphic Evaluation of Text-to-Time-Lapse Video Generation. *arXiv preprint arXiv:2406.18522*.

[16] Esser, P., Sauer, A., Blattmann, A., Lorenz, D., Entezari, R., Kulal, S., Levi, Y., & Boesel, F. (2024). Scaling Rectified Flow Transformers for High-Resolution Image Synthesis. *International Conference on Machine Learning (ICML)*.

[17] Zhang, Y., Hare, J., & Prügel-Bennett, A. (2017). Attention in Recurrent Neural Net-

works for Caption Generation. *arXiv preprint arXiv:1705.05569*.

[18] Chen, X., Fang, H., Lin, T. Y., Vedantam, R., Gupta, S., Dollar, P., & Zitnick, C. L. (2015). Microsoft COCO Captions: Data Collection and Evaluation Server. *arXiv preprint arXiv:1504.00325*.

[19] Zhang, H., Cun, X., Xia, M., Wang, X., Weng, C., & Shan, Y. (2024). VideoCrafter: Overcoming Data Limitations for High-Quality Video Diffusion Models. *arXiv preprint arXiv:2405.10572*.

[20] Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., & Song, X. (2023). CogVLM: Visual Expert for Pretrained Language Models. *arXiv preprint arXiv:2311.03079*.

[21] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[22] Simonyan, K., & Zisserman, A. (2014). Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv preprint arXiv:1409.1556*.

[23] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., & Batra, S. (2023). LLaMA 2: Open Foundation and Fine-Tuned Chat Models. *arXiv preprint arXiv:2307.09288*.

[24] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. (2021). Learning Transferable Visual Models from Natural Language Supervision. *ICML*.

[25] Betker, J., Goh, G., Jing, L., Brooks, T., Wang, J., Li, L., Ouyang, L., & Lee, J. (2023). Improving Image Generation with Better Captions. *arXiv preprint arXiv:2301.01234*.

[26] Zhang, D. J., Wu, J. Z., Liu, J., Zhao, R., Ran, L., Gu, Y., Gao, D., & Shou, M. Z. (2023). Show-1: Marrying Pixel and Latent Diffusion Models for Text-to-Video Generation. *arXiv preprint arXiv:2309.15818*.

[27] Cornia, M., Stefanini, M., Baraldi, L., & Cucchiara, R. (2020). Meshed-Memory Transformer for Image Captioning. *CVPR*.

[28] Chen, L., Zhang, H., Xiao, J., Nie, L., Shao, J., Liu, W., & Chua, T. S. (2017). SCA-CNN: Spatial and Channel-wise Attention in Convolutional Networks for Image Captioning. *CVPR*.

[29] Su, J., Ahmed, M., Lu, Y., Pan, S., & Liu, Y. (2024). RoFormer: Enhanced Transformer with Rotary Position Embedding. *Neurocomputing*.

[30] Yuan, S., Zhang, S., Liu, Y., Chen, H., Zhang, X., & Yu, J. (2024). ChronoMagic: A Benchmark for Dynamic Video Generation. *arXiv preprint arXiv:2406.18522*.