
CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Large-scale pretrained transformers have created milestones in text (GPT-3) and
2 text-to-image (DALL-E and CogView) generation. Its application to video gen-
3 eration is still facing many challenges: The potential huge computation makes
4 it unaffordable for a full training; The scarcity and weak relevance of text-video
5 datasets hinder the model to understand complex movement semantics. In this
6 work, we present 9B-parameter transformer CogVideo, trained by inheriting a
7 pretrained text-to-image model, CogView2. We also propose multi-frame-rate
8 hierarchical training strategy to better align text and video clips. As (probably)
9 the first open-source large-scale pretrained text-to-video model, CogVideo outper-
10 forms all publicly available models at a large margin in both machine and human
11 evaluations.

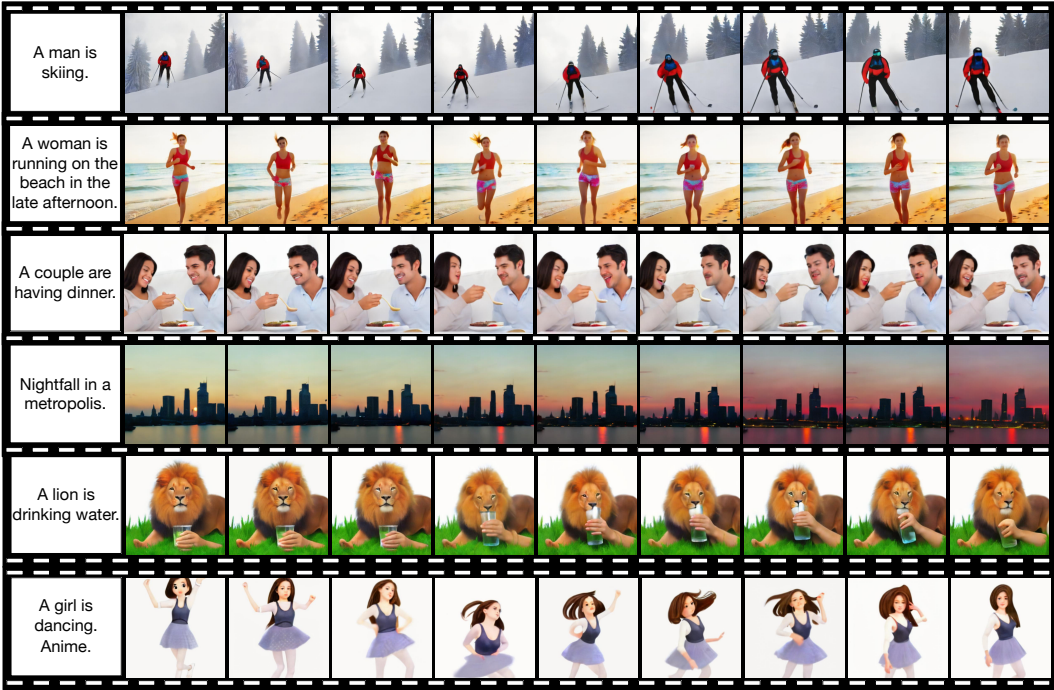


Figure 1: Samples generated by CogVideo. The actual text inputs are in Chinese. Each sample is a 4-second clip of 32 frames, and here we sample 9 frames uniformly for display purpose.

12 1 Introduction

13 Autoregressive transformers, e.g. DALL-E [19] and CogView [5], have revolutionized text-to-image
14 generation recently. It is natural to investigate the potential of autoregressive transformers on text-
15 to-video generation. Previous works followed this basic framework [36, 9], e.g. VideoGPT [37],
16 verifying its superiority over GAN-based methods [4, 27], but are still far from satisfaction.

17 One common challenge is that the generated video frames tend to gradually deviate from the text
18 prompt, making the generated characters hard to perform the desired actions. Vanilla autoregressive
19 models might be good at synthesizing videos with regular (e.g. straightly moving cars) or random
20 patterns (e.g. speaking by randomly moving lips), but fail on text prompt such as “a lion is drinking
21 water”. The main difference between the two cases is that, in the former case the first frame already
22 provides sufficient information for the subsequent changes, while in the latter the model has to
23 precisely understand the action “drink” in order to correctly generate the desired action — the lion
24 lifts the glass to its lip, drinks and then puts down the glass.

25 Why do the autoregressive transformers well understand the text-image relations, but struggle to
26 understand the text-action relations in videos? We hypothesize that the datasets and the way to utilize
27 them are the main reasons.

28 First, it is possible to collect billions of high-quality text-image pairs from Internet [19], but the
29 text-video data are more scarce. The largest annotated text-video dataset, VATEX [32], has only
30 41,250 videos. The retrieval-based text-video pairs, e.g. Howto100M [17], are weakly relevant and
31 most of them only describe the scene without the temporal information.

32 Second, the duration of videos varies a lot. Previous models split the video into many clips with a
33 fixed number of frames for training, which destroys the alignment between the text and its temporal
34 counterparts in the video. If a “drinking” video is split into four individual clips of “holding a glass”,
35 “lifting”, “drinking” and “putting down” with the same text “drinking”, the model will be confused to
36 learn the accurate meaning of drinking.

37 **Present Work.** Here we present a large-scale pretrained text-to-video generative model, CogVideo,
38 which is of 9.4 billion parameters and trained on 5.4 million text-video pairs. We build CogVideo
39 based on a pretrained text-to-image model, CogView2 [6], in order to inherit the knowledge learned
40 from the text-image pretraining. To ensure the alignment between text and its temporal counterparts
41 in the video, we propose the *multi-frame-rate hierarchical training*. The flexibility of the textual
42 condition makes it possible to simply prepend a piece of text describing the frame rate to the original
43 text prompt for modeling different frame rates. To keep the text-video alignment, we choose a proper
44 frame rate description to include the complete action in each training sample. The frame rate token
45 also controls the intensity of the changes throughout continuous frames in generation. Specifically,
46 we train a sequential generation model and a frame interpolation model. The former model generates
47 key frames according to the text, and the latter recursively fill the middle frames by varying the frame
48 rates to make the video coherent. As shown in Figure 1, CogVideo can generate high-resolution
49 (480×480) videos. Human evaluation demonstrates that CogVideo outperforms all publicly available
50 models at a large margin. Our main contributions can be concluded as follows:

- 51 • We present CogVideo, which is the **largest** and **the first open-source** pretrained transformer
52 for text-to-video generation in the general domain.
- 53 • CogVideo elegantly and efficiently finetunes a text-to-video generative model from a pre-
54 trained text-to-image generative model, avoiding the expensive full pretraining from scratch.
- 55 • We propose the multi-frame-rate hierarchical training to better align text-clip pairs, which
56 significantly improves the generation accuracy, in particular for movements of complex
57 semantics. This training strategy endows CogVideo with the capacity of controlling the
58 intensity of changes during the generation.

59 2 Related Work

60 2.1 Video Generation

61 Video generation is a long-standing research topic. Most previous works focus on the next-frame
62 prediction task — forecasting the future frames based on the first video frame. Early works, e.g.

63 CDNA [8] and PredRNN [33], leverage deterministic methods to directly predict the next frame
64 via CNNs or RNNs. However, these deterministic models are unable to capture the stochastic
65 temporal patterns and synthesize coherent complex scenes. Generative models, especially Generative
66 Adversarial Networks [10] (GANs), begin to dominate the area as they can perform unconditional or
67 class-conditional video synthesis without the first frames. VGAN [31] is the first one to use GAN
68 for video generation. It decomposes video to a static background and a moving foreground, and
69 then generates them with 2D and 3D convolutional networks respectively. TGAN[20] proposes
70 to separately generate the temporal latent variables and spatial information, and MoCoGAN [27]
71 similarly decomposes the latent space into context and motion subspaces. DIGAN [38] applies
72 implicit neural representations for video encoding. Recently, text-to-video generation emerges as a
73 promising direction. The framework of VQVAE [29] and autoregressive transformers [30, 1] quickly
74 becomes the mainstream method [35, 36, 9]. Ho et al. [11] proposes video diffusion model along with
75 a gradient method recently for text-to-video generation. The previous methods are basically trained
76 on a specific dataset, e.g. UCF-101 [23], making the trained model domain-specific. Moreover, most
77 of these models are not publicly available.

78 2.2 Autoregressive Transformer

79 Recent years have witnessed the autoregressive transformer emerging as a powerful generative model.
80 The autoregressive models become the most prevalent framework for text generation [24]. With
81 its prominent capacity of fitting, transformer [30] gradually becomes the standard neural structure
82 for text generation. One milestone is GPT-3 [1]. In computer vision, van den Oord et al. [29]
83 first proposes to train a VQVAE to compress the image into a sequence of tokens from a learned
84 dictionary, which can be efficiently handled by autoregressive models. VQ-GAN [7] learns a more
85 semantic-aware dictionary for unconditional image generation. In the text-to-image generation, pre-
86 trained autoregressive transformers such as DALL-E [19] and CogView [5] have shown superiority
87 in open-domain image generation. Besides the pure GPT-style generation, CogView2 [6] proposes a
88 new language model CogLM for infilling in the image generation.

89 Recent autoregressive transformers [18, 37, 35, 36] have also shown their superiority in video
90 generation. Among them, GODIVA [35] and NUWA [36] focus on the open-domain text-to-video
91 generation. However, they simply generate frames or frame blocks one by one in a chronological
92 order, and may suffer from poor text-video alignment (Cf. § 1).

93 3 Method

94 In this section, we first introduce *multi-frame-rate hierarchical training* to better align text and
95 video semantics in § 3.1, and then illustrate an efficient method *dual-channel attention* to inherit
96 the knowledge in pretrained text-image models for video generation in § 3.2. To overcome the
97 large memory and time overhead caused by the large model and long sequence, we refer to Swin
98 Attention [14] and extend it to autoregressive video generation in § 3.3.

99 3.1 Multi-frame-rate Hierarchical Training

100 Here we present the *multi-frame-rate hierarchical training* and generation. We follow the framework
101 of VQVAE [29] and first tokenize each frame into image tokens. Each training sample consists
102 of 5 frame of tokens, but our training method differs in the construction of training sequences and
103 generation process.

104 **Training.** The key design is to add a frame-rate token to the text and sample frames at this frame-rate
105 to compose a fixed-length training sequence. The motivations are two folds:

- 106 (1) Directly separating the long video into clips at a fixed frame-rate often leads to semantic mis-
107 matching. We still use the full text but the truncated clip might only contain incomplete action.
- 108 (2) The adjacent frames are usually very similar. A giant change over the previous frame will
109 probably incur a large loss. This will lead the models less inclined to explore the long-range
110 correlation because to simply copy the previous frame acts like a shortcut.

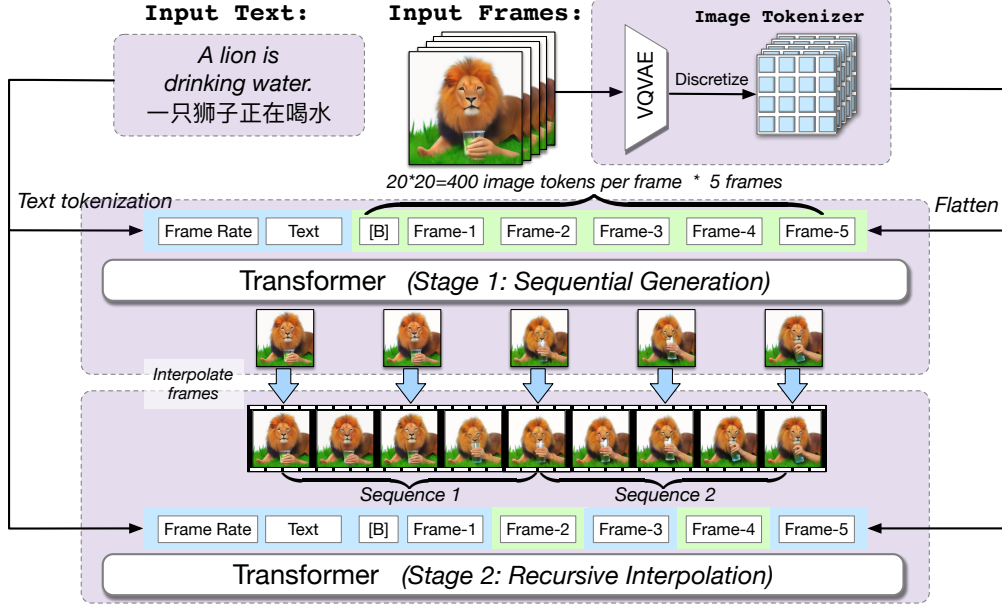


Figure 2: Multi-frame-rate hierarchical generation framework in CogVideo. Input sequence includes frame rate, text, frame tokens. [B] (Begin-of-image) is a separator token, inherited from CogView2. In stage 1, T_s frames are generated sequentially on condition of frame rate and text. Then in stage 2, generated frames are re-input as bidirectional attention regions to recursively interpolate frames. Frame rate can be adjusted during both stages. Bidirectional attention regions are highlighted in blue, and unidirectional regions are highlighted in green.

111 Therefore, in each training sample we want the text and the frames match as possible. We predefined
 112 a series of frame-rates, and select the lowest frame-rate for each text-video pair, as long as we can
 113 sample at least 5 frames at this frame-rate in the video.

114 Although the above method increase the alignment of text and video, the generation at a low frame-
 115 rate could be incoherent. We train another *frame interpolation* model to insert transition frames to the
 116 generated samples of the sequential generation model. Thanks to the generality of CogLM [6], the
 117 two models can share the same structure and training process only with different attention masks.

118 **Generation** The multi-frame-rate hierarchical generation is a recursive process, illustrated in Fig-
 119 ure 2. Specifically, the generation pipeline consists of a sequential generation stage and a recursive
 120 interpolation stage:

- 121 (1) Sequentially generate T_s key frames based on a low frame rate and text. The input sequence
 122 is $[\{\text{Frame Rate}\}\{\text{Text}\} [\text{B}] \{\text{Frame1}\} \dots \{\text{Frame } T_s\}]$. In practice, we always set
 123 $T_s = 5$ and the minimum sampling frame rate to 1 fps.
- 124 (2) Recursively interpolate frames based on the text, frame rate and known frames. In each round
 125 of interpolation, we split generated frames into multiple $\lceil \frac{T_s}{2} \rceil$ -frame blocks overlapping at the
 126 beginning and the end, and interpolate a frame between the successive frames in each block.
 127 The input sequence is $[\{\text{Frame Rate}\}\{\text{Text}\} [\text{B}] \{\text{Frame1}\} \dots \{\text{Frame } T_s\}]$, where
 128 $\text{Frame } 2i (i = 1, 2, \dots, \lfloor \frac{T_s}{2} \rfloor)$ are to be autoregressively generated. By recursively halving $\{\text{Frame}$
 129 $\text{Rate}\}$, we can conduct finer and finer interpolation to generate videos of many frames.

130 **The effect of CogLM.** Tasks such as frame interpolation rely heavily on bidirectional information.
 131 However, most previous works use GPT [35, 37, 36], which is unidirectional. To be aware of the
 132 bidirectional context, we adopt Cross-Modal General Language Model (CogLM) proposed in [6]
 133 which unites bidirectional context-aware mask prediction and autoregressive generation by dividing
 134 tokens into unidirectional and bidirectional attention regions. While bidirectional regions can attend
 135 to all bidirectional regions, unidirectional regions can attend to all bidirectional regions and previous
 136 unidirectional regions. As shown in 2, (1) all frames in stage 1 and the 2nd, 4th frames in stage

137 2 are in the unidirectional region; (2) {Frame Rate}, {Text} and all other frames belong to the
 138 bidirectional region. In this way, bidirectional attention context is fully exploited in text and given
 139 frames without interfering auto-regressive frame prediction.

140 3.2 Dual-channel Attention

141 Large-scale pretraining usually demands a large dataset. For open-
 142 domain text-to-video generation, ideally we need the dataset to
 143 cover sufficient text-video pairs to infer both spatial and
 144 temporal correlation between video and text. However, to collect
 145 high quality text-video pairs is often difficult, expensive and time-
 146 consuming.

147 A natural idea is to make use of the image data to facilitate the
 148 learning of spatial semantics. Video Diffusion Model [11] and
 149 NŪWA [36] try to add text-image pairs into text-video training,
 150 which achieves better results on multiple metrics. However, as
 151 for training a video-only generation model, adding image data
 152 will significantly increase training cost, especially in large-scale
 153 pretraining scenarios.

154 In this paper, we propose to leverage pretrained image generation
 155 models instead of image data. Pretrained text-to-image models,
 156 e.g. CogView2 [6], already have a good command of the text-
 157 image relations. The coverage of the dataset to train these model
 158 is also larger than that of videos.

159 The proposed technique is *dual-channel attention*, where we only
 160 add a new spatial-temporal attention channel to the pretrained CogView2 [6] at each transformer
 161 layer. All the parameters in the CogView2 are frozen in the training, and only the parameters in the
 162 newly added attention layer(See the Attention-plus in Figure 3) are trainable.

163 Here we also emphasize that directly finetuning CogView2 for text-to-video generation cannot well
 164 inherit the knowledge, because the temporal attention follows a different attention pattern and quickly
 165 ruins the pretrained weights during the initial phase of training with large gradients.

166 Specifically, a Transformer layer with dual-channel attention can be computed as

$$\hat{x}_l = \text{LayerNorm}(x_l), \quad (1)$$

$$\tilde{x}_l = \alpha \cdot \text{Attention-base}(\hat{x}_l) + (1 - \alpha) \cdot \text{Attention-plus}(\hat{x}_l), \quad (2)$$

$$x_{l+1} = \text{FFN}(\text{LayerNorm}(x_l + \tilde{x}_l)), \quad (3)$$

167 where x_l denotes input features of layer l ; Attention-base and Attention-plus denote two attention
 168 channels; FFN and LayerNorm represent Feed-Forward Networks and LayerNorm respectively; α
 169 is a vector with length of hidden-size and normalized to (0, 1). The whole structure is the same as
 170 CogView2 when ignoring Attention-plus.

171 Both channels are computed as normal multi-head attention with a certain receptive field formulated
 172 as follows. For token at (t, x, y) in frame block of size (T_s, X, Y) (where (t, x, y) corresponds to
 173 coordination along time, height and width dimension), receptive field RF is a 3D block with extent
 174 $l_t, l_x, l_y \in \mathbb{N}^+$:

$$\text{RF}_{(t,x,y)} = \{(k, i, j) \mid |x - i| < l_x, |y - j| < l_y, |t - k| < l_t, (k, i, j) \notin \text{Mask}_{(t,x,y)}\}, \quad (4)$$

175 where $\text{Mask}_{(t,x,y)}$ represents CogLM attention mask for token (t, x, y) . For Attention-base, we
 176 restrict receptive field to current frame, i.e. $l_x = X, l_y = Y, l_t = 1$, to fully use CogView2’s spatial
 177 modeling ability (therefore referred to as *spatial channel*). For Attention-plus, which is the only
 178 new parameters in CogVideo, we set receptive field to a 3D local block throughout the whole time
 179 dimension, i.e. $l_x = A_x, l_y = A_y, l_t = T_s$ (therefore referred to as *temporal channel*). A_x, A_y
 180 are hyper-parameters satisfying $A_x \leq X, A_y \leq Y$. With A_x and A_y , CogVideo is able to flexibly
 181 trade off between quadratic attention cost and size of receptive field. In practice, we use shifted
 182 window attention [15] as a approximation of 3D block attention and extend it to CogLM scenario, as
 183 illustrated in subsection 3.3.

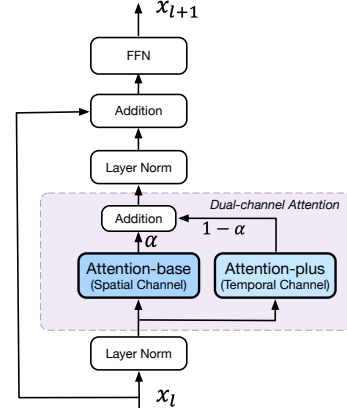


Figure 3: Dual-channel attention. We initialize Attention-plus the same as Attention-base so that the model behaves exactly the same as CogView2 when it is initialized.

184 It is worth noting that two channels are fused and share the same FFN in each layer, because FFN
 185 is a module of heavy parameters containing much vision knowledge. Due to similarity between
 186 images and videos, bringing its knowledge to temporal channel will facilitate video modeling. Finally,
 187 sharing FFN can reduce parameters, thus speed up training and reduce memory overhead.

188 3.3 Shifted Window Attention in Auto-regressive Generation

189 To overcome large time and memory overhead in temporal channel during training and inference,
 190 we refer to Swin Attention proposed in [14] and extend it to auto-regressive scenario by applying
 191 auto-regressive attention mask in shifted windows.

192 Different from non-autoregressive scenario which original Swin
 193 Transformer explores, we propose that Swin Attention can further
 194 accelerate auto-regressive inference because of restricted
 195 receptive field. As shown in Figure 4, receptive field is
 196 restricted by

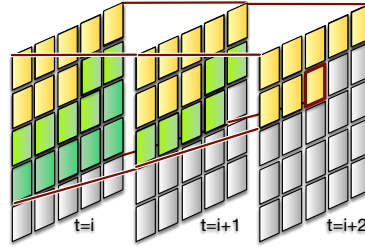


Figure 4: Receptive field (in yellow or green) for the token in red box. Shifted window size is 2×2 in this example.

- 197 • Auto-regressive mask. A token can only attend to pre-
 198 vious frames or tokens before itself in current frame.
- 199 • Shifted window. Only tokens within distance of win-
 200 dow size in both width and height dimension can be
 201 directly attended to.

202 Suppose X, Y is the height and width of each frame, and A_x, A_y
 203 are the height and width of shifted window. For two tokens at (t_1, x_1, y_1) and (t_2, x_2, y_2) , $t_1 < t_2$,
 204 the latter cannot attend to the former either directly or indirectly if

$$(x_1 - x_2)Y + (y_1 - y_2) \geq (t_2 - t_1 + 1)(A_x Y + A_y) \quad (5)$$

205 is satisfied. That is to say, the i -th token in frame t_1 can be generated with the $(i - A_x Y + A_y)$ -th
 206 token in frame $t_1 + 1$ in parallel. In this way, we can generate $\lfloor \frac{XY}{A_x Y + A_y} \rfloor$ tokens in parallel at most,
 207 thus greatly enhance parallelism and accelerate inference compared to auto-regressive with standard
 208 attention which can only generate one token at a time.

209 4 Training

210 Based on methods above, the training details of CogVideo are listed as follows:

211 **Model.** The backbone of CogVideo in both stages is a Transformer with dual-channel attention.
 212 The Transformer has 48 layers, with the hidden size of 3072 in each attention channel, 48 attention
 213 heads and 9.4 billion parameters in total. Among them, 6 billion parameters are fixed to CogView2’s
 214 parameters, which includes Position-wise Feed-Forward Networks (FFN), spatial channel of dual-
 215 channel Attention, first frame’s positional embeddings and all image and text vocabulary embeddings.
 216 The specific implementation of Transformer structure is almost identical to CogView [5] such as
 217 using Sandwich LayerNorm and PB-Relax to stabilize training. Shifted CogLM attention window is
 218 adopted in recursive interpolation model with window size 10×10 .

219 **Dataset.** We pretrain our model on a dataset of 5.4 million captioned videos with a spatial resolution
 220 of 160×160 . For sequential generation model (Stage-1), we adjust frame rate in each sample to
 221 accommodate the whole video, while the minimum frame rate is set to 1 fps. For recursive interpolation
 222 model (Stage-2), we split videos into clips of different length to accommodate prediction on multiple
 223 frame rates including 2, 4, 8 fps.

224 **Pretraining.** The sequence lengths in both stages are 2065, consisting of 64 text tokens, 5 (frames)
 225 \times 400 (per frame) image tokens, and 1 separator token. Both text and images are tokenized with
 226 icetk¹. The parameters are updated by Adam with max learning rate $= 2 \times 10^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.95$,
 227 weight decay $= 1 \times 10^{-2}$. See Appendix for pretraining details.

¹<https://github.com/THUDM/icetk>

Table 1: (Left) Video generation performance on UCF-101. Class labels are used as text inputs. * denotes the model is trained on the training split of UCF-101 only. (Right) Video generation performance on Kinetics-600. Metrics are measured on generated videos of 16 frames priming on first 5 frames, following settings in [18]. ** denotes groundtruth used in FVD testing is blurred with our image tokenizer icetk.

Method	IS (\uparrow)	FVD (\downarrow)	Method	FVD
VideoGPT[37]	24.69	-	Latent Video Transformer[18]	224.73
DVD-GAN[4]	27.38	-	Video Transformer[34]	170
TGANv2[21]*	28.87	1209	DVD-GAN-FP[4]	69.15
MoCoGAN-HD[25]	32.36	838	TriVD-GAN-FP[16]	25.74
DIGAN[38]*	29.71	655	CogVideo (Ours)	109.23
DIGAN[38]	32.70	577	CogVideo (Ours)**	59.55
TATS-base[9]	79.28	332		
CogVideo (Ours)	50.46	626		
CogVideo (Ours)**	-	545		

228 5 Experiments

229 5.1 Machine Evaluation

230 Machine evaluation is conducted on two popular benchmarks for video generation, i.e., UCF101 [23]
 231 and Kinetics-600 [3]. Following Rakhimov et al. [18], Yu et al. [38], we use Fréchet Video Dis-
 232 tance(FVD) [28] and Inception score(IS) [22] as metrics in the evaluation. FVD is calculated based
 233 on I3D model[2] trained on Kinetics-400, and IS is based on C3D model [26] which was first trained
 234 with Sports-1M dataset [12] and then fine-tuned on the UCF101 dataset. Our evaluation code is the
 235 same as the official TGAN-v2 implementation².

236 **UCF-101** is a human action dataset consisted of 13,320 videos annotated with 101 action classes.
 237 Due to the image style and frame rate gap between CogVideo’s training set and UCF-101, we use
 238 class labels as the input text and fine-tune CogVideo on the whole dataset for 10,000 iterations with
 239 batch size = 192. During inference, we sample class labels according to the class distribution. FVD
 240 and IS are evaluated over 2048 and 10,000 samples respectively, following Yu et al. [38]. Results are
 241 shown in Table 1 (Left).

242 **Kinetics-600** dataset contains 600 classes of human action videos, with roughly 350k train and
 243 50k test videos in total. We use the action category as input text, and fine-tune CogVideo on the
 244 training set for 12,000 iterations with batch size of 640. Following the setup of Weissenborn et al.
 245 [34], Rakhimov et al. [18], we center-crop and down-sample each frame to 64x64, and measure with
 246 FVD. Results are shown in Table 1 (Right).

247 5.2 Human Evaluation

248 To further evaluate CogVideo, we invite 90 anonymous evaluators to rate for CogVideo and other open-
 249 source baselines including GAN-based model TGANv2 [21] and GPT-based model VideoGPT [37].
 250 30 classes in UCF101 are randomly picked as text conditions, and several aspects are rated (See
 251 Appendix for details). For VideoGPT, we use the official unconditional pretrained model³ to generate
 252 samples. For TGANv2, we use the official source code to train an unconditional generation model
 253 under the same setting as that in Saito et al. [21]. To assign unconditionally generated samples into
 254 corresponding categories, we choose TSM [13] as the action recognition model and only samples
 255 with confidence >80%. Results in Figure 5 show that CogVideo significantly outperforms baselines
 256 on multiple important aspects including frame texture, motion realism and semantic relevance, and
 257 achieves the top score by overall quality. It can be seen that 49.53% evaluators choose CogVideo as
 258 the best method, and only 15.42% and 5.6% favor VideoGPT and TGANv2, respectively.

²<https://github.com/pfnet-research/tgan2>

³<https://github.com/wilson1yan/VideoGPT>

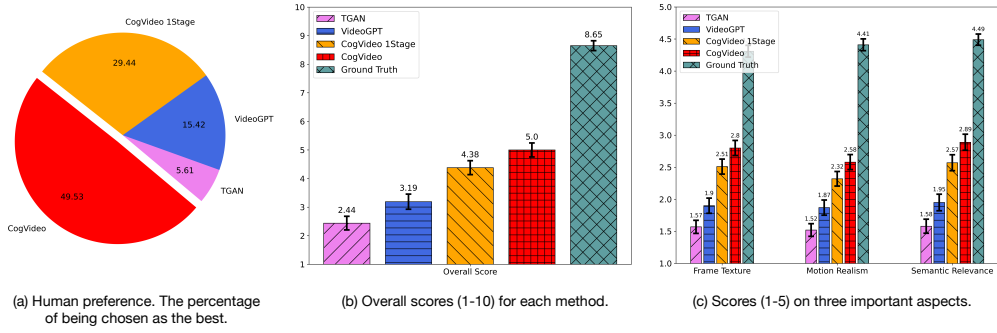


Figure 5: Human evaluation results. "CogVideo 1Stage" refers to the method in ablation study, which generates videos sequentially with CogVideo’s Stage-1 Model only by recursively reinserting last 2 generated frames into input and generate future frames.

Table 2: Ablation study on a 5,000-sample subset of Kinetics-600’s testset. FVD is evaluated on generated 11-frame samples priming on 5 frames and ground-truth blurred by our image tokenizer. The setting column indicates the difference between each method and CogVideo. Models of each setting are trained on Kinetics-600 trainset for 10,000 iterations with batch size of 320.

Method	Setting	FVD (\downarrow)
CogVideo	None	108.27
1-stage Generation($N_{overlap} = 1$)	– hierarchical	137.13
1-stage Generation($N_{overlap} = 2$)	– hierarchical	120.82
Initialized to CogView2	– Pretrain	124.92
Randomly Initialized	– Pretrain – CogView	166.13

259 5.3 Ablation Study

260 To verify the effectiveness of hierarchical multi-frame-rate generation and incorporating CogView2,
 261 we conduct ablation study quantitatively and qualitatively on Kinetics-600 and UCF-101 datasets.

262 **Hierarchical multi-frame-rate generation.** In comparison with CogVideo, we fine-tune a 1-stage
 263 video generation model on Kinetics-600 from the sequential generation model in CogVideo, which
 264 generates long videos by recursively reinserting last $N_{overlap}$ frames into the input to sample next
 265 $N_s - N_{overlap}$ frames. Larger $N_{overlap}$ means more previous frames can be utilized during the
 266 inference, but will increase time overhead.

267 **Dual-channel attention with CogView2’s weights.** We additionally train (1) A randomly initialized
 268 model; (2) A model incorporating CogView2’s weights but leaving temporal channel randomly
 269 initialized and unfixed (equivalent to CogVideo without pretraining on videos) on Kinetics-600.

270 5.3.1 Quantitative Evaluation

271 All aforementioned models have been trained for 11,000
 272 iterations with batch size of 160. Quantitative results are
 273 shown in Table 2. We can see that the hierarchical method
 274 is clearly superior to 1-stage generation with different N_s ,
 275 and model initialized with CogView2’s weights has lower
 276 FVD than randomly initialized one.

277 Figure 6 plots the training loss curve of (1) finetuning
 278 CogVideo; (2) training model from random initialization;
 279 (3) training model initialized to CogView2 and partially
 280 fixed. We can see that CogView2 endows model with a

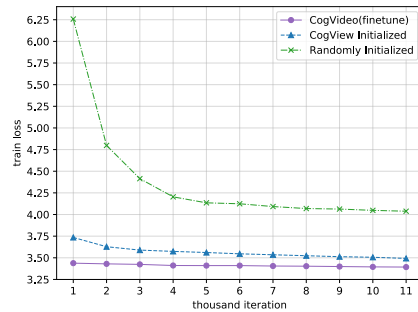


Figure 6: Training loss in ablation study.

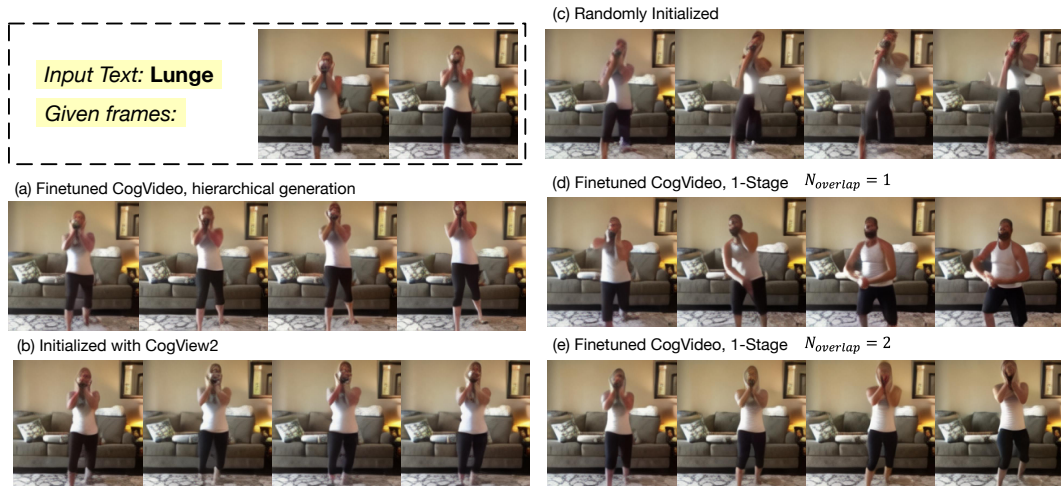


Figure 7: Video samples in ablation study, which are generated priming on class label and first 5 frames in Kinetics-600. All samples are down sampled by extracting one in every three frames for display purpose. (a) Use fine-tuned CogVideo to hierarchically generate samples. (b) Train a model on Kinetics-600 which is initialized as and partially fixed to CogView2, and hierarchically generate samples. (c) Train a model on Kinetics-600 which is randomly initialized, and hierarchically generate samples. (d)(e) Use fine-tuned CogVideo to generate frames in 1 stage with different $N_{overlap}$.

281 good initialization point from which the loss function can converge faster to a lower value. Also,
 282 fixing part of the parameters to CogView2 reduce optimization cost, which gains more than 2x
 283 acceleration when using optimization CPU-offload mode in deepspeed.

284 5.3.2 Qualitative Evaluation

285 Qualitative comparison is shown in Figure 7. While model trained from random initialization tends to
 286 produce irrational deformation, model incorporating CogView2 is able to model objects better. And
 287 samples generated hierarchically performs better on content consistency and motion rationalization.

288 We also conduct human evaluation between 1-stage and hierarchical video generation model under
 289 the same setting as 5.2. As shown in 5, hierarchical model, i.e. CogVideo, outperforms 1-stage model
 290 on semantic relevance, motion realism as well as texture quality. This is probably because 1-stage
 291 model tends to constantly generate small movements which make the whole video unrealistic, and if
 292 one generated frame collapses, the subsequent frames often suffer from severe degradation.

293 6 Conclusion

294 We present CogVideo, to the best of our knowledge, the largest and the first open-source pretrained
 295 transformer for text-to-video generation for the general domain. CogVideo is also the first attempt
 296 to efficiently leverage pretrained text-to-image generative model to text-to-video generation model
 297 without hurting its image generation capacity. With the proposed multi-frame-rate hierarchical
 298 training framework, CogVideo is endowed with better understanding of text-video relation and ability
 299 to control the intensity of changes during generation. We extend swin attention to CogLM, which
 300 achieves acceleration in both training and inference. There are still some limitations in CogVideo, e.g.
 301 restriction on length of the input sequence still exists due to the large scale of model and limitation of
 302 GPU memory, and we leave them for future work.

303 **Broader Impact.** This paper aims to advance the open-domain text-to-video generation, which
 304 will ease the effort of short video and digital art creation. The efficient training method transfers
 305 knowledge from text-to-image models to text-to-video models, which helps avoid training from
 306 scratch, and thus reduce the energy consumption and carbon emission. A negative impact is the risk
 307 of misinformation. To alleviate it, we can train an additional classifier to discriminate the fakes. We
 308 believe the benefits outweigh the downsides.

References

- 309
- 310 [1] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan,
311 P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *arXiv preprint*
312 *arXiv:2005.14165*, 2020.
- 313 [2] J. Carreira and A. Zisserman. Quo vadis, action recognition? a new model and the kinetics
314 dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
315 pages 6299–6308, 2017.
- 316 [3] J. Carreira, E. Noland, A. Banki-Horvath, C. Hillier, and A. Zisserman. A short note about
317 kinetics-600. *arXiv preprint arXiv:1808.01340*, 2018.
- 318 [4] A. Clark, J. Donahue, and K. Simonyan. Adversarial video generation on complex datasets.
319 *arXiv preprint arXiv:1907.06571*, 2019.
- 320 [5] M. Ding, Z. Yang, W. Hong, W. Zheng, C. Zhou, D. Yin, J. Lin, X. Zou, Z. Shao, H. Yang,
321 et al. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural*
322 *Information Processing Systems*, 34, 2021.
- 323 [6] M. Ding, W. Zheng, W. Hong, and J. Tang. Cogview2: Faster and better text-to-image generation
324 via hierarchical transformers. *arXiv preprint arXiv:2204.14217*, 2022.
- 325 [7] P. Esser, R. Rombach, and B. Ommer. Taming transformers for high-resolution image synthesis.
326 *arXiv preprint arXiv:2012.09841*, 2020.
- 327 [8] C. Finn, I. Goodfellow, and S. Levine. Unsupervised learning for physical interaction through
328 video prediction. *Advances in neural information processing systems*, 29, 2016.
- 329 [9] S. Ge, T. Hayes, H. Yang, X. Yin, G. Pang, D. Jacobs, J.-B. Huang, and D. Parikh. Long
330 video generation with time-agnostic vqgan and time-sensitive transformer. *arXiv preprint*
331 *arXiv:2204.03638*, 2022.
- 332 [10] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville,
333 and Y. Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014.
- 334 [11] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet. Video diffusion models.
335 *arXiv preprint arXiv:2204.03458*, 2022.
- 336 [12] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video
337 classification with convolutional neural networks. In *Proceedings of the IEEE conference on*
338 *Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- 339 [13] J. Lin, C. Gan, and S. Han. Tsm: Temporal shift module for efficient video understanding. In
340 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7083–7093,
341 2019.
- 342 [14] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer:
343 Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF*
344 *International Conference on Computer Vision*, pages 10012–10022, 2021.
- 345 [15] Z. Liu, J. Ning, Y. Cao, Y. Wei, Z. Zhang, S. Lin, and H. Hu. Video swin transformer. *arXiv*
346 *preprint arXiv:2106.13230*, 2021.
- 347 [16] P. Luc, A. Clark, S. Dieleman, D. d. L. Casas, Y. Doron, A. Cassirer, and K. Simonyan.
348 Transformation-based adversarial video prediction on large-scale data. *arXiv preprint*
349 *arXiv:2003.04035*, 2020.
- 350 [17] A. Miech, D. Zhukov, J.-B. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic. Howto100m: Learning
351 a text-video embedding by watching hundred million narrated video clips. In *Proceedings of*
352 *the IEEE/CVF International Conference on Computer Vision*, pages 2630–2640, 2019.
- 353 [18] R. Rakhimov, D. Volkhonskiy, A. Artemov, D. Zorin, and E. Burnaev. Latent video transformer.
354 *arXiv preprint arXiv:2006.10704*, 2020.

- 355 [19] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever.
356 Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*, 2021.
- 357 [20] M. Saito, E. Matsumoto, and S. Saito. Temporal generative adversarial nets with singular
358 value clipping. In *Proceedings of the IEEE international conference on computer vision*, pages
359 2830–2839, 2017.
- 360 [21] M. Saito, S. Saito, M. Koyama, and S. Kobayashi. Train sparsely, generate densely: Memory-
361 efficient unsupervised training of high-resolution temporal gan. *International Journal of*
362 *Computer Vision*, 128(10):2586–2606, 2020.
- 363 [22] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved
364 techniques for training gans. In *Proceedings of the 30th International Conference on Neural*
365 *Information Processing Systems*, pages 2234–2242, 2016.
- 366 [23] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from
367 videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- 368 [24] I. Sutskever, J. Martens, and G. Hinton. Generating text with recurrent neural networks. In
369 *ICML'11*, page 1017–1024, 2011.
- 370 [25] Y. Tian, J. Ren, M. Chai, K. Olszewski, X. Peng, D. N. Metaxas, and S. Tulyakov. A good image
371 generator is what you need for high-resolution video synthesis. *arXiv preprint arXiv:2104.15069*,
372 2021.
- 373 [26] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features
374 with 3d convolutional networks. In *Proceedings of the IEEE international conference on*
375 *computer vision*, pages 4489–4497, 2015.
- 376 [27] S. Tulyakov, M.-Y. Liu, X. Yang, and J. Kautz. Mocogan: Decomposing motion and content
377 for video generation. In *Proceedings of the IEEE conference on computer vision and pattern*
378 *recognition*, pages 1526–1535, 2018.
- 379 [28] T. Unterthiner, S. van Steenkiste, K. Kurach, R. Marinier, M. Michalski, and S. Gelly. To-
380 wards accurate generative models of video: A new metric & challenges. *arXiv preprint*
381 *arXiv:1812.01717*, 2018.
- 382 [29] A. van den Oord, O. Vinyals, and K. Kavukcuoglu. Neural discrete representation learning. In
383 *Proceedings of the 31st International Conference on Neural Information Processing Systems*,
384 pages 6309–6318, 2017.
- 385 [30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and
386 I. Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- 387 [31] C. Vondrick, H. Pirsivash, and A. Torralba. Generating videos with scene dynamics. *Advances*
388 *in neural information processing systems*, 29, 2016.
- 389 [32] X. Wang, J. Wu, J. Chen, L. Li, Y.-F. Wang, and W. Y. Wang. Vatex: A large-scale, high-
390 quality multilingual dataset for video-and-language research. In *Proceedings of the IEEE/CVF*
391 *International Conference on Computer Vision*, pages 4581–4591, 2019.
- 392 [33] Y. Wang, M. Long, J. Wang, Z. Gao, and P. S. Yu. Predrnn: Recurrent neural networks for
393 predictive learning using spatiotemporal lstms. *Advances in neural information processing*
394 *systems*, 30, 2017.
- 395 [34] D. Weissenborn, O. Täckström, and J. Uszkoreit. Scaling autoregressive video models. *arXiv*
396 *preprint arXiv:1906.02634*, 2019.
- 397 [35] C. Wu, L. Huang, Q. Zhang, B. Li, L. Ji, F. Yang, G. Sapiro, and N. Duan. Godiva: Generating
398 open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021.
- 399 [36] C. Wu, J. Liang, L. Ji, F. Yang, Y. Fang, D. Jiang, and N. Duan. N\ " uwa: Visual synthesis
400 pre-training for neural visual world creation. *arXiv preprint arXiv:2111.12417*, 2021.

- 401 [37] W. Yan, Y. Zhang, P. Abbeel, and A. Srinivas. Videogpt: Video generation using vq-vae and
402 transformers. *arXiv preprint arXiv:2104.10157*, 2021.
- 403 [38] S. Yu, J. Tack, S. Mo, H. Kim, J. Kim, J.-W. Ha, and J. Shin. Generating videos with dynamics-
404 aware implicit generative adversarial networks. *arXiv preprint arXiv:2202.10571*, 2022.

405 Checklist

- 406 1. For all authors...
- 407 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
408 contributions and scope? [Yes]
- 409 (b) Did you describe the limitations of your work? [Yes]
- 410 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 411 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
412 them? [Yes]
- 413 2. If you are including theoretical results...
- 414 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 415 (b) Did you include complete proofs of all theoretical results? [N/A]
- 416 3. If you ran experiments...
- 417 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
418 mental results (either in the supplemental material or as a URL)? [No] We will release
419 code later.
- 420 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
421 were chosen)? [Yes]
- 422 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
423 ments multiple times)? [No]
- 424 (d) Did you include the total amount of compute and the type of resources used (e.g., type
425 of GPUs, internal cluster, or cloud provider)? [No]
- 426 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 427 (a) If your work uses existing assets, did you cite the creators? [Yes] See footnotes.
- 428 (b) Did you mention the license of the assets? [No]
- 429 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 430 (d) Did you discuss whether and how consent was obtained from people whose data you’re
431 using/curating? [No]
- 432 (e) Did you discuss whether the data you are using/curating contains personally identifiable
433 information or offensive content? [No]
- 434 5. If you used crowdsourcing or conducted research with human subjects...
- 435 (a) Did you include the full text of instructions given to participants and screenshots, if
436 applicable? [Yes] See supplemental material.
- 437 (b) Did you describe any potential participant risks, with links to Institutional Review
438 Board (IRB) approvals, if applicable? [N/A]
- 439 (c) Did you include the estimated hourly wage paid to participants and the total amount
440 spent on participant compensation? [Yes]