

000
001
002
003
004
005
006
007
008054
055
056
057
058
059
060
061
062

Short review of state-of-the-art methods on Video Generation.

009
010
011
012
013
014
015
016
017
018
019
020
021063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Anonymous CVPR submission

Paper ID ****

Abstract

This short review provides a summary of state-of-the-art video generation methods and an algorithm for developing a video generation network, which is presented in the final sections.

1. Literature Review

Conventionally, auto-regressive models are preferred over GANs for generating visual data [19, 20, 26] due to their explicit density modeling and stable training. The outputs of an auto-regressive model are usually generated in a pixel-by-pixel manner. This method requires high computational power. Hence, they are mostly used for low-resolution images. CNN [21], pixel RNN [19], image transformer [23], iGPT [5], and Video Transformer [33] have employed instances of such methods.

Video generation from text differs from video prediction, where only the previous frames are used to generate the forthcoming ones. FitVid [2] from GoogleBrain is an example of a video prediction model. FitVid solved the under/over-fitting problem with the use of data augmentation and achieved state-of-the-art performance in video prediction. It has no complex structure, no attention mechanism, no curriculum learning, and no training scheduling. It comprises two LSTM layers for modeling the dynamics and a few convolutional layers.

[1] uses stochastic methods to model the inherent uncertainty in the task of video prediction by considering motion history. It uses generative models with latent variables for appearance and motion. There are other models e.g. SVG [7] and SRVP [10] that study video prediction. [4] implemented a network for generating videos out of an input image through sampling out of latent variables. This framework is adversarially trained in an unsupervised manner. Fig. 1 shows a network structure that generates videos from an input image. In the following sections, we review some of the state-of-the-art models on text-to-video generation, which stem from the task of video prediction.

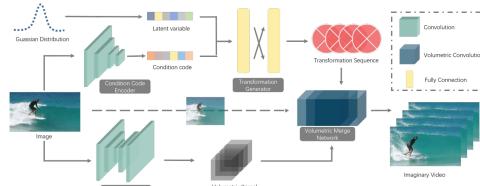


Figure 1. Video Imagination from a Single Image with Transformation Generation [4].

1.1. Models

Much research has recently focused on implementing networks capable of generating videos. [14, 22] combined GANS with 3D convolution layers to generate low-resolution, fixed-length videos. [3] used conditional filters to generate videos of varying lengths. [6] utilized LSTM cells with 2D convolutional networks to retain both temporal coherence and quality of the generated frames. However, the main problem with these works is that they only perform well on small datasets and suffer from low generalization ability. In [12], generating videos from text using transformer-based architectures has been investigated. This architecture receives an input text and tokenizes it. Ultimately, these embedded tokens provided to the transformer-based network create a set of images that form a low-framerate video. The low-quality video gets interpolated recursively by entering into a bidirectional attention network shown in Fig. 2.

1.1.1 Google Imagen [11]

Submitted on 5 October 2020, Imagen is the most recent T2V model. The following procedure describes the input pipeline. First, the input text enters a fully convolutional network, resulting in a low-resolution, low-framerate video. In the subsequent layers, using temporal super-resolution doubles the number of frames resulting in smoother animation. The next stage up-samples the resulting tensor and increases the spatial resolution. Eventually, the final two successive stages double the number of frames and

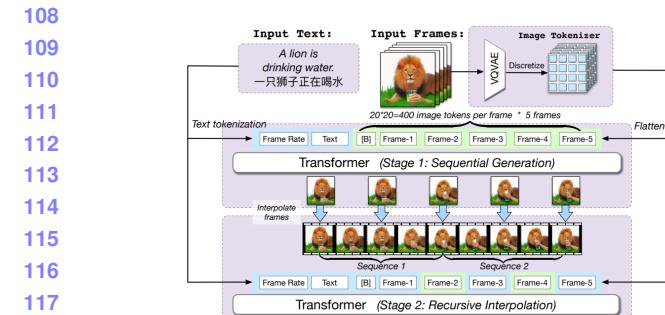


Figure 2. Multi-frame-rate hierarchical generation framework in CogVideo [12].

increase the image's resolution. The structure mentioned above makes use of fully-convolutional temporal and spatial super-resolution and v-parameterization of diffusion models.

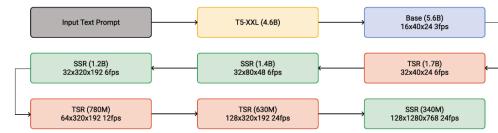


Figure 3. The cascaded sampling pipeline starting from a text prompt input to generating a 5.3-second, $1280 \cdot 768$ video at 24fps. “SSR” and “TSR” denote spatial and temporal super-resolution. The videos are labeled as $\text{frames} \cdot \text{width} \cdot \text{height}$. In practice, the text embedding is injected into all models, not just the base model. [11].

1.1.2 META AI Make-a-Video [28]

This research paper submitted on 29 September 2022 belongs to Meta AI. The input text is first converted into an image. Subsequently, using the created image, a low-framerate movie is generated. At the following stages, there is an improvement in either the framerate by frame interpolation modules or in the quality with the help of spatiotemporal super-resolution convolution and attention layers. Fig. 4 shows the overall structure of the network.

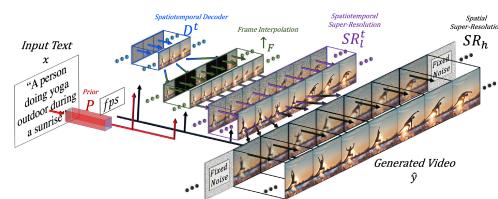


Figure 4. Make-A-Video structure from Meta AI. [28].

1.1.3 Temporal Video Generation [18]

The main component of this architecture is a latent code sequence generator that uses 3D multivariate Gaussian distributions to generate output samples. The generator's structure uses temporal residual connections that extract temporal features both from the past and the future.

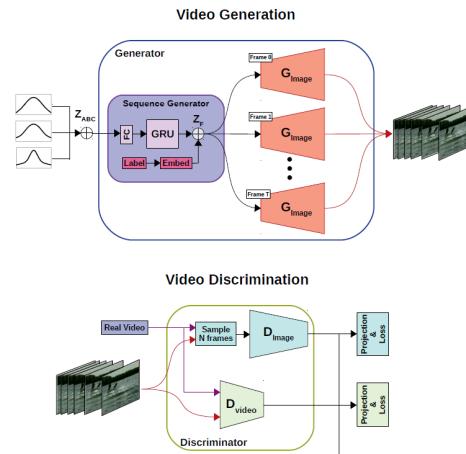


Figure 5. TS-GAN framework.

1.1.4 GODIVA [34]

GODIVA is a pre-trained model for T2V¹ generation, which is a result of a collaboration between Microsoft Research Asia and Duke university. It uses auto-regressive modeling principles and also a 3D sparse attention mechanism. This mechanism decreases the computational complexity. The same collaboration led to the publishing of another model called NUWA, introduced in the next section. GODIVA utilizes VQ-VAEs as the encoder. The idea of employing VQ-VAEs for videos has already been used in the task of video prediction in [24, 27, 31, 32, 37, 38], but GODIVA is the first to use this structure for T2V tasks. According to Fig. 8, which demonstrates the results of training several networks on the Kinetics and Youtube videos datasets, the outputs look blurry in TFGAN, [3] whose network architecture is shown in Fig. 6. The metrics even indicate less promising results in T2V [15]. However, GODIVA's results show higher quality both in temporal and spatial domains. GODIVA is pre-trained on Howto100M, a large-scale text-video dataset with more than 136 million text-video pairs using 64 V100 GPUs. It is also fine-tuned on the MSR-VTT dataset with 8 V100 GPUs comprising 10000 human-annotated text-video pairs.

¹text-to-video

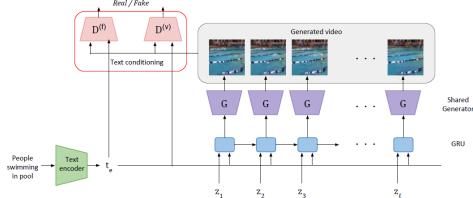


Figure 6. TF-GAN architecture.

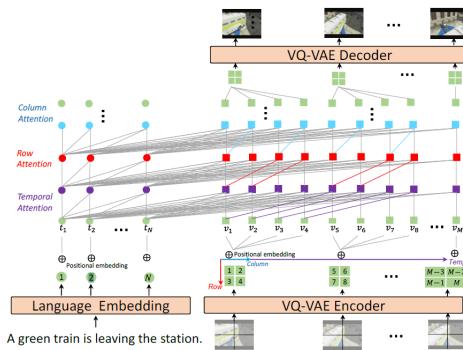


Figure 7. A simple illustration of our GODIVA model with a three-dimensional sparse attention mechanism for text-to-video generation task. The video is auto-regressively predicted with the consideration of four aspects: the input text, the same position of the previously generated frames, the same rows on the same frame, same columns on the same column.



Figure 8. GODIVA visual results

1.1.5 NUWA [35]

NUWA is a unified, multi-modal pre-trained model capable of manipulating existing visual data and generating new ones. NUWA supports texts, images, and videos using an adaptive transformer framework that receives 1D, 2D, or 3D data as the input. The decoder can conduct eight visual synthesis tasks, including T2I, S2I, I2I, TI2I, T2V, S2V, V2V, and TV2V. The transformer uses a 3D Nearby Attention or (3DNA) mechanism capable of considering both

Model	Acc	FID-img	FID-vid	CLIPSIM
T2V(64X64)	42.6	82.13	14.65	0.2853
SC(128X128)	74.7	33.51	7.34	0.2915
TFGAN(128X128)	76.2	31.76	7.19	0.2961
NUWA(128X128)	77.9	28.46	7.05	0.3012
Make-A-Video	-	-	13.17	0.3049
CogVideo (English)	-	-	23.59	0.2631

Table 1. The quantitative results of T2V using four state-of-the-art models trained on the Kinetics dataset. (The last two rows are trained on the MSR-VTT dataset.)

spatial and temporal content information. This mechanism not only reduces the computational complexity but also improves the visual quality of the generated results. NUWA is trained using the Kinetics [13] dataset. The training results are compared to two other models in Tab. 1, verifying that NUWA outperforms other structures. Also, an illustration of the results can be seen in Fig. 9, which indicates that NUWA outperforms other structures in generating unseen videos, such as playing gulf at the swimming pool. This feature is called zero-shot capability.



Figure 9. Results of eight tasks over an unseen text from NUWA, GODIVA, TFGAN, and T2V.

1.1.6 VideoGPT [36]

VideoGPT is a conceptually simple architecture that uses VQ-VAEs to create down-sampled latent representations of video frames with the help of convolutional layers and axial self-attention. It utilizes positional embedding both in the time and spatial domains. The results are comparable to

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
that of GANs and state-of-the-art models. The network has been trained on the BAIR Robot dataset resulting in high-fidelity outputs.

2. Proposed Algorithm

331
332
333
334
335
336
337
338
339
340
The idea behind the proposed algorithm is based on work in [9]. VQ-GAN is a kind of VQ-VAE that exploits a discriminator acting as an extra supervisor; This causes improvements in the quality of the generated data. As represented in Fig. 10, at the first stage of training, images are encoded using the encoder of a VQ-VAE, creating a representation of the frames sampled from the videos of the training set. The input sentence is also embedded into a latent space to assist output generation by defining a loss function measuring the distance between the representations created by the T2I unit.

341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
A codebook is a mapper that quantizes the latent representations created by the encoder. It gets trained during the training phase of the VAE by considering a codebook mapping loss in the total loss term along the GAN loss. We may also consider some extra points in the latent space close to the created representations to retain the ability to create new and unseen examples. These representations can be concatenated to enter the decoder, where a reconstruction of the input is created. The reconstructed images' qualities are assessed by a reconstruction loss which can be a pixel-wise MSE. 3DNA-based attention can be used both for the input text and the input video frames a sequence generator to generate codebook vectors that make sense and go along with the sequences already seen in the codebook. In order to train the structure, first, the GAN and the codebook sections must be trained. The generator gets trained to map the input images into a latent representation with the help of a GAN loss and a reconstruction loss. Then the generator weights are kept constant so that the discriminator can be trained to learn the generator's flaws. This is executed recursively.

361
362
363
364
365
366
367
368
369
370
After freezing weights of the encoder, the decoder, and the codebook, it is time for the transformer to get trained. The proposed transformer uses a three-dimensional sparse attention model to be used both for images and texts, which means that it exploits temporal, column, and row information. 3DNA-based attention enables the network to be used for other types of sequences so that the network can learn to focus on the important regions of the frames. A transformer has an auto-regressive structure and works best with sequences. That is why it is suitable for generating sequences that best match the codebook examples.

371
372
373
374
375
376
377
The transformer in this structure acts as a sequence generator. It gets trained after the VQ-GAN to generate sequences of latent representations comparable to the codebook. Its training phase is as follows: some codebook vectors get randomly modified after the codebook training. These adjusted latent representations enter the transformer

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
399
401
402
403
404
405
399
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
to generate images. The discriminator measures how real the created images are. As a result, the transformer weights get trained to generate sequences that match the ones stored in the codebook.

In other words, the transformer is supposed to correct the randomly adjusted quantized vectors.

A video's frame samples can be generated simultaneously, but if they are created frame-wisely, a temporal coherency loss term must be added to the structure. Briefly, the proposed structure uses two discriminators; one determines how fake the generated video frame is, and the other measures the temporal coherency of the generated frames. Also, a perceptual loss is defined to determine the visual realistic and relative matching of the output according to [34].

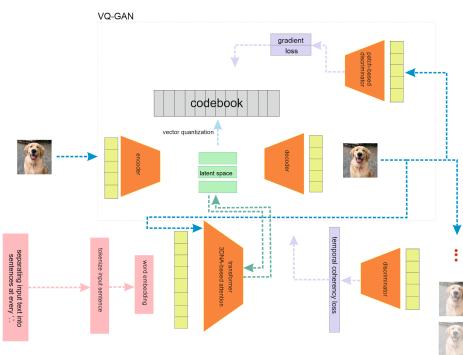


Figure 10. Proposed Structure

3. Advantages of the Proposed Algorithm

410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
Considering previous discussions in the introduction, pixel-by-pixel models [5, 19, 21, 33] have explicit density modeling and stable training, which pushes us to prefer them over GANs [16, 17, 29, 30]. However, pixel-by-pixel data generation is problematic for high-dimensional data synthesis. It is evident that using 3D convolution models to generate videos requires considerable memory. Also, scaling the dimensions up for images of higher resolution is challenging.

422
423
424
425
426
427
428
429
430
431
After introducing VQ-VAEs and auto-regressive models for visual synthesis tasks, the idea of converting images into discrete visual tokens became widespread, as it enabled the researchers to conduct large-scale pre-training e.g. in DALL-E [27] and CogView [8], and for the tasks of video prediction in [25, 36] and video generation in GO-DIVA [34].

422
423
424
425
426
427
428
429
430
431
Also, the study in [35] showed that employing VQ-GANs instead of VQ-VAEs leads to creation of higher-quality videos. According to this paper, the loss function of

432 the encoder in VQ-VAE forces the reconstructed image to
 433 be exactly like the original one, which causes limitations
 434 in the generalization capabilities of VQ-VAEs. Hence by
 435 adding a perceptual loss term to the general loss function
 436 and an extra discriminator to the model, the results
 437 show more satisfactory generalization capabilities and
 438 semantic matching by eliminating exact similarity constraints.
 439

4. Project Schedule

443 The Gantt chart considering the time needed to research
 444 and trying various implementations is shown in Tab. 2.

445 Table 2. Estimation of the RD timeline.

	Oct.	Nov.	Dec.
Research/ Implementing SOTA structures			
Creating a novel structure			

References

- [1] Adil Kaan Akan, Erkut Erdem, Aykut Erdem, and Fatma Güney. Slamp: Stochastic latent appearance and motion prediction, October 2021. 1
- [2] Mohammad Babaeizadeh, Mohammad Saffar, Suraj Nair, Sergey Levine, Chelsea Finn, and Dumitru Erhan. Fitvid: Overfitting in pixel-level video prediction, 06 2021. 1
- [3] Yogesh Balaji, Martin Renqiang Min, Bing Bai, Rama Chellappa, and Hans Peter Graf. Conditional gan with discriminative filter generation for text-to-video synthesis. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 1995–2001. International Joint Conferences on Artificial Intelligence Organization, 7 2019. 1, 2
- [4] Baoyang Chen, Wenmin Wang, and Jinzhuo Wang. Video imagination from a single image with transformation generation, 2017. 1
- [5] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Hee-woo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels, 13–18 Jul 2020. 1, 4
- [6] Kangle Deng, Tianyi Fei, Xin Huang, and Yuxin Peng. Ircgan: Introspective recurrent convolutional gan for text-to-video generation, 7 2019. 1
- [7] Emily Denton and Rob Fergus. Stochastic video generation with a learned prior, 10–15 Jul 2018. 1
- [8] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. Cogview: Mastering text-to-image generation via transformers, 2021. 4
- [9] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis, 2020. 4
- [10] Jean-Yves Franceschi, Edouard Delasalles, Mickael Chen, Sylvain Lamprier, and Patrick Gallinari. Stochastic latent residual video prediction, 13–18 Jul 2020. 1
- [11] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models, 2022. 1, 2
- [12] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022. 1, 2
- [13] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset, 2017. 3
- [14] Yitong Li, Martin Min, Dinghan Shen, David Carlson, and Lawrence Carin. Video generation from text, Apr. 2018. 1
- [15] Yitong Li, Martin Min, Dinghan Shen, David Carlson, and Lawrence Carin. Video generation from text, Apr. 2018. 2
- [16] Gaurav Mittal, Tanya Marwah, and Vineeth N. Balasubramanian. Sync-DRAW. In *Proceedings of the 25th ACM international conference on Multimedia*. ACM, oct 2017.
- [17] Gaurav Mittal, Tanya Marwah, and Vineeth N. Balasubramanian. Sync-draw: Automatic video generation using deep recurrent attentive architectures. In *Proceedings of the 25th ACM International Conference on Multimedia*, MM ’17, page 1096–1104, New York, NY, USA, 2017. Association for Computing Machinery. 4
- [18] Andres Munoz, Mohammadreza Zolfaghari, Max Argus, and Thomas Brox. Temporal shift gan for large scale video generation, 2020. 2
- [19] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks, 2016. 1, 4
- [20] Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu. Conditional image generation with pixelcnn decoders, 2016. 1
- [21] Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu. Conditional image generation with pixelcnn decoders, 2016. 1, 4
- [22] Yingwei Pan, Zhaofan Qiu, Ting Yao, Houqiang Li, and Tao Mei. To create what you tell: Generating videos from captions, 2018. 1
- [23] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Łukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer, 2018. 1
- [24] Ruslan Rakhimov, Denis Volkhonskiy, Alexey Artemov, Dennis Zorin, and Evgeny Burnaev. Latent video transformer, 2020. 2
- [25] Ruslan Rakhimov, Denis Volkhonskiy, Alexey Artemov, Dennis Zorin, and Evgeny Burnaev. Latent video transformer, 2020. 4
- [26] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021. 1

- 540 [27] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray,
541 Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever.
542 Zero-shot text-to-image generation, 18–24 Jul 2021. 2, 4
543 594
544 [28] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An,
545 Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual,
546 Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman.
547 Make-a-video: Text-to-video generation without text-video
548 data, 2022. 2
549 595
550 [29] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan
551 Kautz. Mocogan: Decomposing motion and content for
552 video generation, 2017. 4
553 596
554 [30] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba.
555 Generating videos with scene dynamics, 2016. 4
556 597
557 [31] Jacob Walker, Ali Razavi, and Aäron van den Oord. Predicting
558 video with vqvae, 2021. 2
559 598
560 [32] Jacob Walker, Ali Razavi, and Aäron van den Oord. Predicting
561 video with vqvae, 2021. 2
562 599
563 [33] Dirk Weissenborn, Oscar Täckström, and Jakob Uszkoreit.
564 Scaling autoregressive video models, 2019. 1, 4
565 600
566 [34] Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji,
567 Fan Yang, Guillermo Sapiro, and Nan Duan. GODIVA: generating
568 open-domain videos from natural descriptions, 2021.
569 2, 4
570 601
571 [35] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang,
572 Daxin Jiang, and Nan Duan. NÜwa: Visual synthesis pre-
573 training for neural visual world creation, 2021. 3, 4
574 602
575 [36] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind
576 Srinivas. Videogpt: Video generation using vq-vae and trans-
577 formers, 2021. 3, 4
578 603
579 [37] Yunzhi Zhang, Wilson Yan, Pieter Abbeel, and Aravind
580 Srinivas. Videogen: Generative modeling of videos using
581 {vq}-{vae} and transformers, 2021. 2
582 604
583 [38] Yunzhi Zhang, Wilson Yan, Pieter Abbeel, and Aravind
584 Srinivas. Videogen: Generative modeling of videos using
585 {vq}-{vae} and transformers, 2021. 2
586 605
587 606
588 607
589 608
590 609
591 610
592 611
593 612