

Diffusion (and more) text-to-video generation methods

**Denis Dimitrov** 

Managing director in data science, Sber Al Scientific advisor, AIRI



# Agenda



«Time-lapse of a flower blooming: growth, beauty, and the passage of time» by Kandinsky Video 1.1

#### **01** Task history

- Diffusion generation of images from text, Kandinsky 3
- **03** Diffusion generation of videos from text
- **04** Models Application & Conclusions



# Task history

## **Brief history of visual generative models**

- 2014: the emergence of GAN was a breakthrough in the development of image generation models ٠
- 2016: the first mention of image generation from textual descriptions was made at S. Reed et al. (ICML 2016 and NeurIPS 2016)
- Since 2017, no ML/DL/CV conference A/A\* passes without publications on the topic of **image** synthesis from textual descriptions
- 2018: one of the first works on **video** generation from text descriptions, Y. Li et al.
- Since 2020 the most active research and development of **text-to-image** architectures
- Since 2023 the most active research and development of **text-to-video** architectures





Learning What and Where to Draw

Video Generation From Text									
<b>Yitong Li</b> Duke University Durham, NC 27708	Martin Rengiang Min NEC Labs America Princeton, NJ 08540	Dinghan Shen Duke University Durham, NC 27708	David Carlson Duke University Durham, NC 27708	Lawrence Carin Duke University Durham, NC 27708					
Generating videos fi lenge for existing g by training a conditi and dynamic inform brid framework, em and a Generative A tures, called "gist", ground color and o are considered by tr To obtain a lange an model, we develop a internal newlo show	Abstract rom text has proven to be a significat enerative model. We tackle this pr onal generative model to extract both aion from text. This is manifested ipolying a Variational Autoencoder diversarial Network (GAN). The sta- are used to sketch text-conditioned bject layout structure. Dynamic for ansforming input text into an image out of data for training the deeple methodise to automative and the state of the methodise to automative and the state of the methodise of the state of the state of the state methodise of the state of th	Tulyakov tive Adv oblem static diversity, static diversity, is tatic diversity, is faito a hy- (VAE) generatic back- generatic atures a good v aming Eliman M ucheon simply r entres to provi	et al. 2017). Both of the ci rsrarial Network (GAN) (C s shown encouraging result er, in contrast with these p n, here we conditionally text in requires both a good c ideo generator. There are or text-to-image generat ideo generator. There are or text-to-image generat sainov and Stakhutdin eplacing the image genera	ted works use a Genera- ioodfellow et al. 2014), s on sample fidelity and revious works on video generate videos baede captions. Text-to-video onditional scheme and a number of existing on (Reed et al. 2016; v 2016); unfortunately, tor by a video genera- severe mode collapse).					

imental results show that the proposed framework generates



#### **Generative architectures**





#### Advantages

- Show better performance
- No adversarial training
- No mode collapse

#### Disadvantages

Large inference time



### Text-to-visual models of Sber & AIRI





Demo ruDALLE A demo text-to-image model, 1.3B



text

Kandinsky 2.0 Multilanguage diffusion model to generate images from



Kandinsky 2.1 A new diffusion model to generate images from text



Kandinsky 2.2 Multilanguage diffusion model to generate photorealistic images from text













02

Diffusion generation of images from text, Kandinsky - 3



## **Denoising Diffusion Probabilistic Models (DDPM)**

The diffusion process consists of two parts:

- Forward diffusion process iterative noise addition
- Reverse diffusion process iterative noise removal

Forward diffusion process (fixed)





Noise

Reverse denoising process (generative)



#### **Forward diffusion process**



## **Training of a reverse diffusion process**



$$p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$$
$$p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

$$L_{\text{simple}}(\theta) \coloneqq \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}} \Big[ \big\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \big\|^2 \Big]$$

# Base architectures – U-Net, DiT



#### Unconditional U-Net architecture



Conditional U-Net architecture



The Diffusion Transformer (DiT) architecture



## General inference scheme of diffusion text-to-image models



# Why DiT?

#### Architectures

U-Net	Mostly applied to image generation tasks
Transformer	<ul> <li>Compatible with LLM distributed learning frameworks</li> <li>Clear scaling ability</li> <li>More memory efficient than U-Net</li> <li>Sora-like models use this approach</li> </ul>

#### **Training type**

Diffusion	<ul> <li>High quality of output generations</li> <li>Sora-like models use exactly this approach</li> </ul>
Autoregression	LLM architectures can be used



# Kandinsky 3.0 / 3.1

#### Inference pipeline





#### Proprietary U-Net architecture, learning with cosine scheduling



# Kandinsky 3.0 training

# Amount of data for training (including filters):

256x256	1.1B
384x384	780M
512x512	450M
768x758	224M
Mixed	280M

#### **Datasets:**

#### Not so good:

- LAION 800M
- COYO 700M
- TIGER 80M

#### Good:

- Depositphotos 100M
- RussianThemed 236K



15

«Oriental tea party, baklava, tea, cups, honey, hospitality, harmonious serving, oriental sweets, anime style»



«Image in the golden circle of the Kul Sharif Mosque, Kazan, photorealism, maximum detail»

«Image of a black headdress with oriental gold ornament, gold thread embroidery, detailing»

### **Kandinsky Flash**



Kandinsky

Flash











A Pikachu with an angry expression and red eyes, with lightning around it, hyper realistic style



A panda is playing a guitar



A Pomeranian is sitting on the Kings throne wearing a crown. Two tiger soldiers are standing next to the throne

#### Speedup by more than x10 0.4s for generation

of 1 image



# Kandinsky 3.1 = Kandinsky 3.0 + Kandinsky Flash

### Kandinsky 3.1 vs Kandinsky 3.0



30 annotators23 812 for text24 741 for visual



## Kandinsky 3.1 vs DALL-E 3

### Kandinsky 3.1 vs DALL-E v3



Text



Visual



# Kandinsky 3.1 vs Midjourney v5.2

## SBS, Kandinsky 3.1 vs Midjourney v5.2





# 03

# Diffusion generation of videos from text



# The main challenges in the field of video generation, which are tackled by the world's most advanced AI Labs

To develop a high-quality Sora-level long video generation model: 1) Up to 1 minute duration; 2) **1920x1080** resolution, **24 fps** (Full HD 1080p or even better format)

At the same time, it is necessary to train intermediate versions of the model with a lower spatial resolution of 1280x720 (HD or 720p) and a shorter duration of 10 seconds per video generation. This will demonstrate the ability to scale the developed architecture.

The final model should support the ability to generate video by text, by the initial frame, by the initial video fragment, as well as the ability to supplement an arbitrary "incomplete" video fragment (including solving the problem of merging two videos). In addition, this model must adapt to the tasks of instructional image and video editing, video-inpainting, as well as to any other reasonable additional conditions (ControlNet, IP Adapter).

#### **R** runway

Prompt: A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about.

@ProperPrompter





#### **VGM development**

#### In the field of video generation (VGM), there is competition similar so that of LLM

On February 15, 2024, Open AI announced its video generation model called Sora, which significantly surpasses all previous models. Since then, more than 5 companies have presented their video generation models of comparable quality

Company	Model	Country	Status	Date	Computing resources	Data (presumably)	Architecture	Output
Open Al	Sora	USA	Preview	15.02.2024			DiT	1080p video up to 60 seconds
ShengShu	Vidu	China	Preview	27.04.2024			U-ViT	1080p video up to 16 seconds
DeepMind, Google	Veo	USA	Preview	14.05.2024	Cluster	Hundreds of millions of pairs of «text-video»,	DiT	1080p video up to 60 seconds
Luma Labs	Dream Machine	USA	Open demo	12.06.2024	~2 000—10 000 H100	~2 000—10 000 H100 «text-image», as well as just images and videos	DiT (?)	1080p video up to 120 seconds
Runway Al	Gen-3	USA	Open demo	17.06.2024			DIT (?)	1080p video up to 10 seconds
Kuaishou Technology	Kling	China	Open demo	24.07.2024			DiT (?)	1080p video up to 120 seconds

MovieGen	Meta*	USA	Preview	04.10.2024	6144 H100	Around 400 millons		DIT	768p up to 16 seconds (16 fps)
					'	'	24	KAND	INSKY

. . .

# **Difficulties of the video generation task**

- New field without established architectures
- No high-quality Open Source solutions
- The computational complexity is higher than in the LLM field

	Large Language Models (LLM)	Video Generation Models (VGM, VPT, GWM)
Start of intensive model development	2017	2023
Established architectures that are explicitly described (or coded) and have been verified to scale	GPT, BERT, T5 etc	— (most likely it would be a special version of DiT)
Open Source	LLaMA (Meta), Mistral (Mistral AI), Gemma (Google, DeepMind), Grok (X.ai) etc	<b>Open-Sora</b> , <b>Open-Sora-Plan</b> (community enthusiasts) <b>CogVideoX</b> (Zhipu AI)
Model context length	128k (GigaChat)	864k (for 10 sec video with 720p resolution)



# VGM development



Промпт: Reflections in the window of a train traveling through the Tokyo suburbs. KLING



SORA





Prompt: Drone view of waves crashing against the rugged cliffs along Big Sur's Garay Point beach. The crashing blue waters create white-tipped waves, while the golden light of the setting sun illuminates the rocky shore. A small island with a lighthouse sits in the distance, and green shrubbery covers the cliff's edge. The steep drop from the road down to the beach is a dramatic feat, with the cliff's edges jutting out over the sea. This is a view that captures the raw beauty of the coast and the rugged landscape of the Pacific Coast Highway.

Source: OpenAl, Sora

## VGM development



#### Pika 1.5 (Pikaffects)



## But it's not perfect yet....



We're sharing our research progress early to start working with and getting feedback from people outside of OpenAl and to give the public a sense of what Al capabilities are on the horizon.





**Промпт:** Five gray wolf pups frolicking and chasing each other around a remote gravel road, surrounded by grass. The pups run and leap, chasing each other, and nipping at each other, playing. Weakness: Animals or people can spontaneously appear, especially in scenes containing many entities.

https://openai.com/index/sora/

https://openai.com/index/video-generation-models-as-world-simulators/



# General inference scheme of diffusion text-to-video models

#### **Pipeline inference** Historical footage of California Historical footage during the Gold Rush of California during the Gold Rush Video-VAE Text Encoder **Historical footage** of California during the Gold Rush Diffusion Video-VAE Transformer (Encoder) (Latent) Latent Video-VAE (Decoder) Video-VAE (Decoder)





#### **Training approaches**

#### Current algorithm: Noise



Time

#### Autoregressive algorithm:

Noise



**Different degrees of noise level** 



Don't make some patches noisy



**Completely noising some patches** 



#### **Loss functions**

#### Problem

Now diffusion is taught with MSE Loss, which penalizes the image pixel by pixel and does not focus on important details.

#### Diffusion Model with **MSE Loss** (No CFG)



Diffusion Model with Self-Perceptual Loss (No CFG)

The solution is to use more advanced loss functions:

- → Perceptual loss (there are studies on images)
- Adversarial loss (there is an insight from Kandinsky 3.0 Flash – the discriminator allows to increase the detail)
- → Physics inspired loss (still at the hypothesis level)

Kandinsky 3.0

Kandinsky 3.1 Flash











## **Compute (Sora)**



Base compute

4x compute

32x compute

Same seed(!) Same data for training (!!)



Kandinsky 4.0 XL



If we don't have enough resources to train one big model, we can complicate the pipeline a little bit, but simplify each part of this pipeline

## Kandinsky Video 1.0/1.1 pipeline

#### Base Frames generation model + Frames interpolation model



# Kandinsky Video 1.0

**Base Frames:** 

Shape: (bs, t, c, w, h)

## Interpolation:

(base frames are concatenated with noise)

- Shape: (bs, t+2, c, w, h)
- Shape: (bs, t, c\*3, w, h)







# Kandinsky Video 1.1

17.5B parameters



Prompt: Large white octopus sitting on top of a building

Prompt: A ladybug flies away from a mushroom



#### https://evalcrafter.github.io



Prompt: Red Car Drift



Prompt: The Russian flag is waving in the wind



Generation from image. Prompt: Prompta: Cartoon But if you have the enough GPU (computing resources), a single model shows the best quality!

#### Main directions of development



#### Data



#### Data



#### «Raw»:

- ~400,000 hours
- ~150 million «text-video» pairs

#### After processing:

~170,000 hours ~60 million «text-video» pairs

Approximately 40% of the data is included in the training



#### Pipeline

#### «Raw» videos table

dataset_name	obs_endpoint	obs_path	status	process_start	width	height	duration	fps	id	process_node		109
youtube	obs.ru-moscow- 1.hc.sbercloud.r u	gigaeye- data/youtube/2 024-02- 1/1.mp4	process	2024-07-12 11:15:32.847311 +03	1920	1080	00:01:11	24	1	job-3cc09c65- a72d-5b6e05155	$\left  \right\rangle$	, <b>LU</b> <sup>s</sup> rows

#### «Processed» videos table

	id	raw_id	obs_endpoint	text prompt	width	height	duration	dynam icity	aesthetic_score		
-	10	1	gigaeye-kandinsky- spark/db_pipeline/scenes/1/0:01:16.m p4	The scene gives off a casual and relaxed atmosphere, with the woman seemingly enjoying	1920	1080	12	0.976	0.812	1 rc	208 2005
	11	1	gigaeye- <del>kandinsky-</del> spark/db_pipeline/scenes/1/0:01:42.m p4	The scene is in black and white, giving it a classic feet	1920	1080	13	<del>0.312</del>	0.115		



**Optical flow** 



# Video captioning

#### CogVLM2-Video



#### Есть проблемы

20%	0,5%	Main
descriptions	descriptions	bottleneck
are inaccurate	are incorrect	(2.2 svsc)



A group of individuals engaging in various outdoor activities on a paved pathway surrounded by lush greenery. The scenes depict people walking, **riding bicycles**, and **skateboarding**,

with some individuals in red and white attire. The pathway is well-maintained, with long shadows cast by the trees, suggesting it's either early morning or late afternoon.



# **Comparison with captions from LAION**



Revolights-eclipse

(600, 329)



Grosjean's car was ripped onto two pieces after careering through the guardrails

(768, 512) (768, 512)



Shooting stars elec ireland (640, 360) (640, 320)



2013 DACIA DUSTER 1.5DCI AMBIANCE 4X2 \*\*\*\*OWN THIS CAR FROM £33 PER WEEK

(682, 512) (640, 512)



#### Main part: Diffusion Transformer



## Working with a diffusion transformer



# **DiT: Baseline**





## **MMDiT, Sparse Attention**



#### Sparse Attention for resolutions > 256:



- Learning speedup **x4 times**
- Inference speedup **x6 times**
- Memory consumption decrease x1.3 times
- x1.5 more parameters

# **Training Pipeline: Baseline**

The main factors of memory consumption: model weights, optimizer state, activations, attention



# **Training Pipeline: Engineering optimizations (FSDP)**



# **Training Pipeline: Engineering optimizations (Flash Attention)**



#### Memory consumption decrease more than 15 times:

Strategy 96 GPU, 3B DiT	DDP	DDP + Flash Attn	FSDP+ Flash Attn	FSDP + Checkpointing + Flash Attn 2
Memory 256x256, 2.5s, 24fps	OOM	76 Gb	20 Gb	5 Gb
Step Speed 256x256, 2.5s, 24fps	inf	0.6 s	0.4 s	0.5 s
I				

# **Training Pipeline: Engineering optimizations** (Ring Flash Attention)



# **MMDiT: training on videos**

Number of parameters	Resolution	fps	Duration	Video Dataset	GPU
2.7B	256x256	24	<10 sec	31M 120B Tokens	64



SORA: base compute



A panda, dressed in a small, red jacket and a tiny hat, sits on a wooden stool in a serene bamboo forest...



A knight riding on a horse through the countryside



Confident teddy bear surfer rides the wave in the tropics



«Mandala, with elements of oriental carpet, kaleidoscope style movement, clear lines, calibrated graphics»

0

I

P

«Beautiful Tatar girl on the background of the Kazan Kremlin, elements of Tatar culture, Karina in art style, maximum camera zoom, facial details»

1.00

0001

«Bright postcard with the image of the mosque of Kazan, in the foreground detailed frame of blue flowers with golden patterns on a green background, with elements of Tatar floral patterns, dynamics of movement, revitalization of the postcard»

## Main directions of development



# Video-VAE

#### High computational complexity



#### Open Source solutions do not have a high enough compression ratio

VAE model	Compression (TxHxW)	Number of tokens
—	Without compression	55 296k
Image VAE	1x8x8	864k
Open-Sora v1.2 Open-Sora Plan v1.2 CV-VAE CogVideoX	4x8x8	216k
Cascade Video-VAE (Ours)	8x8x8	108k
Goal	8x16x16	27k

## Video-VAE



- Small context 17 frames
- Low compression ratio 4
- Requires lots of resources 32 GPU



- Large context 128 frames
- Higher compression ratio 8
- Requires less resources 32 GPU

H

W





## Video-VAE



Original video



Reconstruction - x512 times smaller dimension!



04

# Models Application & Conclusions



## **General World Models**

#### **R** runway

A world model is an Al system that builds an internal representation of an environment, and uses it to simulate future events within that environment.

We believe the next major advancement in Al will come from systems that understand the visual world and its dynamics, which is why we're starting a new long-term research effort around what we call general world models.

#### 🜀 Sora

Sora serves as a foundation for models that can understand and simulate the real world, a capability we believe will be an important milestone for achieving AGI. We believe the capabilities Sora has today demonstrate that continued scaling of video models is a promising path towards the development of capable simulators of the physical and digital world, and the objects, animals and people that live within them.



World Models – unlimited source of data for training Foundation Models

# General World Models: симуляторы







Using video generation models to simulate our reality (or any other virtual one) and to adjust the robot's actions will be significant improvement in this area

**Cons** - long and computationally expensive



## **General World Models: downstream-tasks**

A fundamental model of computer vision that can enhance other computer vision tasks

- Synthetic data generation
- Pretrain for perception tasks: depth map, detection, segmentation (example Marigold)
- Pretrain for generative tasks: 3D object/scene generation





## **General World Models: visual content**



- Entertainment, advertising, cinema, art
- Video generation by text query
- Image animation



# CONTACTS



#### **Denis Dimitrov**

Scientific advisor, AIRI Managing director in data science, Sber AI



@DENDI\_MATH\_AI



@DENDIMITROV