# Our Story

# Three Missions

Conduct world-class breakthrough research in AI, putting Vietnam on the AI world map.

Train the top future AI talents for Vietnam.
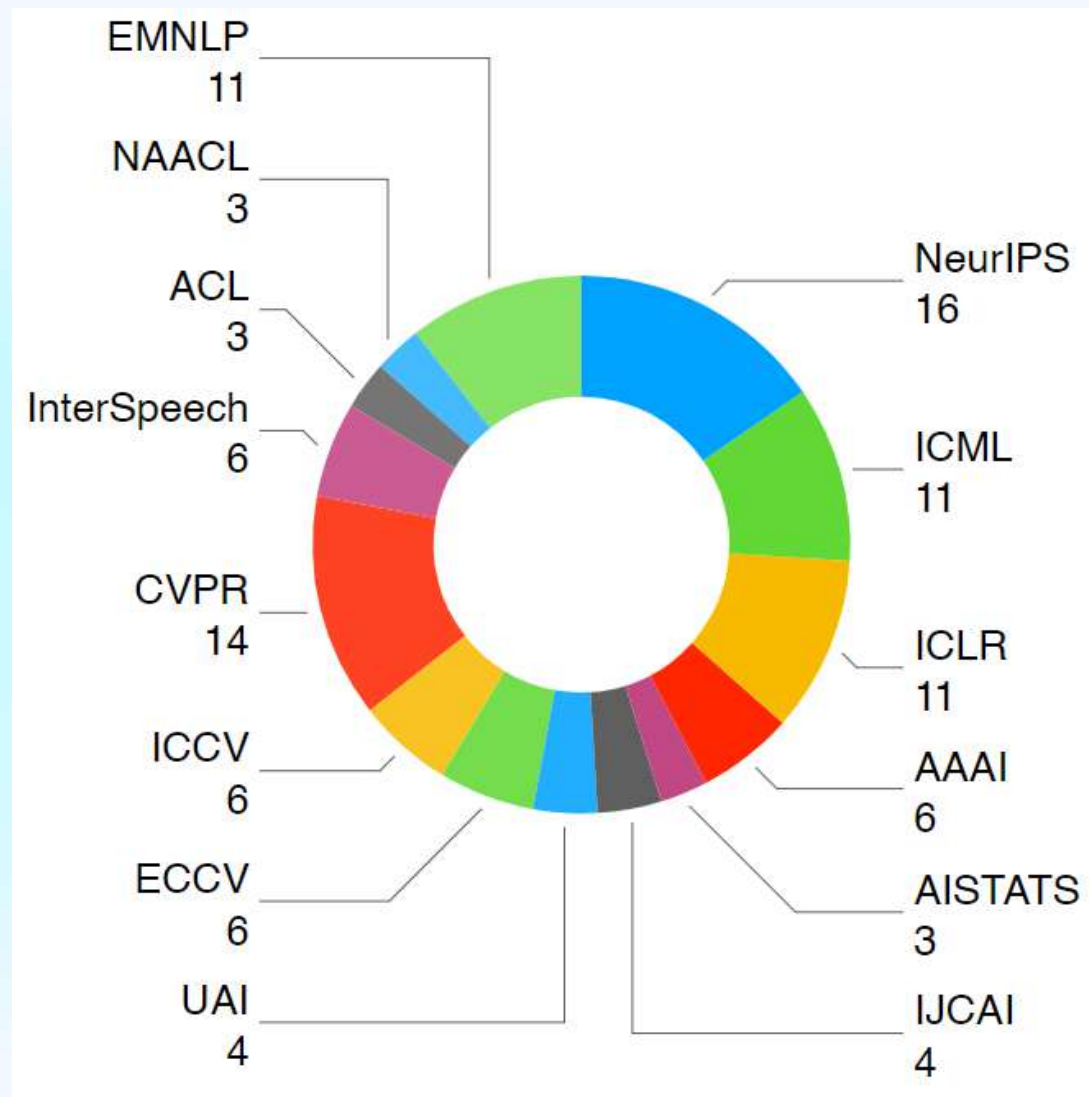
Build AI-powered products that offer the best customer's value.

- Founded **2019**, headquartered in Hanoi, Vietnam & USA, Australia.

- **150+ employees total**

- Key investor is **Vingroup**

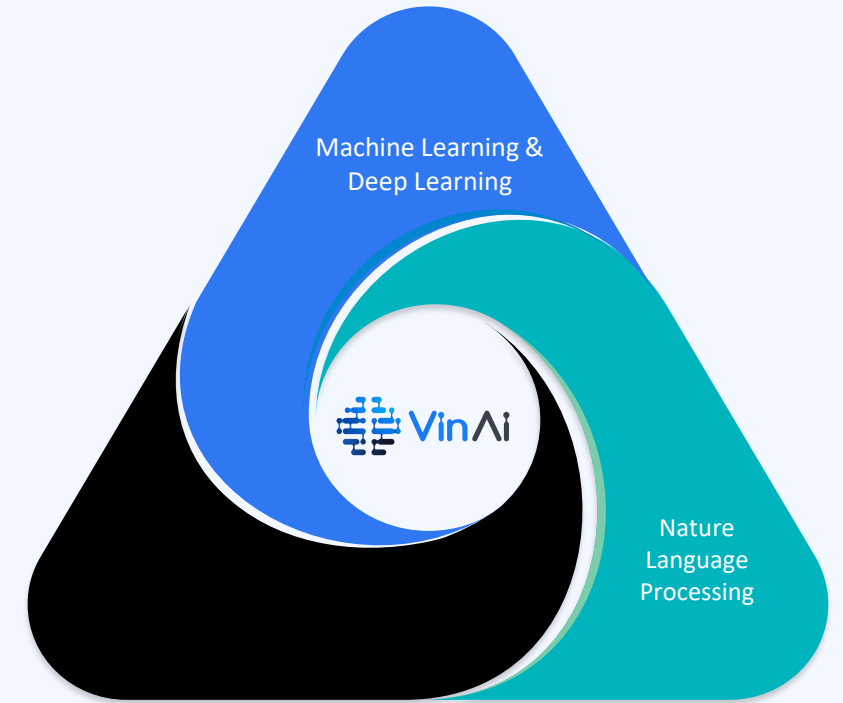# Putting Vietnam on the Global AI Map

**2022**

**104** top tier papers

**2021**

**68** top tier papers

**2020**

**23** top tier papers

**2019**

**6** top tier papers



EMNLP 11
NAACL 3
ACL 3
InterSpeech 6
CVPR 14
ICCV 6
ECCV 6
UAI 4
NeurIPS 16
ICML 11
ICLR 11
AAAI 6
AISTATS 3
IJCAI 4

# Our World Class AI Research

1. Google (USA) - 200.2
2. Microsoft (USA) - 79.3
3. Facebook (USA) - 54.9
4. Amazon (USA) - 26.5
5. IBM (USA) - 26.3
6. Huawei (China) - 21.8
7. Alibaba (China) - 13.1
8. NVIDIA (USA) - 12.5
9. Tencent (China) - 10.2
10. Samsung (South Korea) - 10.0
11. Baidu (China) - 9.7
12. NTT (Japan) - 7.5
13. Apple (USA) - 7.0
14. OpenAI (USA) - 6.7
15. Intel (USA) - 6.7
16. Adobe (USA) - 6.2
17. Salesforce (USA) - 6.0
18. Yendex (Russia) - 6.0
19. NEC (Japan) - 5.0
20. VinAI (Vietnam) - 4.5
21. Bosch (Germany) - 4.2
22. Criteo (France) - 3.6
23. Bytedance (China) - 3.5
24. JD (China) - 3.5
25. Kuaishou Technology (China) - 3.2
26. Megvii (China) - 3.0
27. SenseTime (China) - 2.9
28. Naver (South Korea) - 2.8
29. AITRICS (Korea) - 2.7
30. Ant (China) - 2.5

VinAI ranked 20th in the **top Global AI Research Companies** in 2022

Source: Thundermark Capital https://thundermark.medium.com/ai-research-rankings-2022-sputnik-moment-for-china-64b693386a4

## 3 Research Pillars

Machine Learning & Deep Learning

Nature Language Processing

# Home grown talents in VN

# Talent - AI Residency Program

## The First and Most Prestigious AI Residency Program in Vietnam

80 Residents

World-class research training & Mentorship

43 accepted Top-tiers papers

66 Ph.D. Scholarships (from the world's TOP 20 CS universities)

45 Filed Patents

Our Current Global Alumni Network

# Speaker



**Anh Tran**

- Research Scientist at VinAI Research
- Research field: Computer Vision
- Former Amazon
- PhD degree from University of Southern California (USA)
- VEF 2012
- 1st prize at Vietnam Talent 2010

Website : https://sites.google.com/site/anhttranusc/
Google Scholar: https://scholar.google.com/citations?user=FYZ5ODQAAAAJ

# Outline

1. Diffusion Models

2. Text-to-image models

3. Other text-to-X models

# Image Generation



Map an input $z \sim N(0, \mathbb{I})$ to an output image x in the desired domain

➤ Unconditional generation: $\quad\quad\quad\quad\quad\quad\quad\quad x = G(z)$

➤ Conditional generation: additional target attributes **c**: $\quad x = G(z, c)$

# Typical Approaches



Generative Adversarial Networks

2016
2017
2018

High Quality Samples

Fast Sampling

Mode Coverage / Diversity

Variational Autoencoders,

Denoising Diffusion Models

(b) A cute corgi lives in a house made out of sushi.

(c) A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.

Encoder $q_\phi(z|x)$

Decoder $p_\theta(x|z)$

Others:
- Autoregressive models (ARM)
- Energy-based models (EBM)
- Normalizing flow models

# 1. Diffusion Models

# High-level Ideas

➤ Noise adding process (forwards diffusion)



➤ Denoising process (reverse diffusion)
- Can generate image from random noise!!!



J. Ho, A. Jain, and P. Abbeel. "Denoising diffusion probabilistic models". *In NeurIPS 2020*, *33*, 6840-6851.

# Details

- Noise adding steps for training: $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1-\bar{\alpha}_t}\epsilon$
  - In training, we want to predict the noise $\epsilon \sim N(0, I)$ given $x_t$ using a UNet:



- Denoising strategy (sampling): $x_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t, t)\right) + \sigma_t\epsilon_t$

K. Kreis, R. Gao, and A. Vahdat. Denoising diffusion-based generative modeling: foundations and applications. CVPR Tutorial. 2022.

# Image Quality: Diffusion models beat GANs



**BigGAN**

**ADM (diff. model)**

P. Dhariwal and A. Nichol. Diffusion models beat GANs on image synthesis. In *NeurIPS* 2021.

# Latent Difussion Models (LDMs)

R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. "High-resolution image synthesis with latent diffusion models". In *CVPR* 2022 (pp. 10684-10695).

# Diversity: Much Better Than GANs

- Text: most common

- Sketches, bounding-box layouts, scene graph, semantic segmentation map, depth map, style image, human poses
  ➔ More control over the image generation!

- Representation from other modalities: audio, fMRI signal…

# Text Conditioning: Text-to-image



R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. "High-resolution image synthesis with latent diffusion models". In *CVPR* 2022 (pp. 10684-10695).

# Image Conditioning

**Image/Semantic maps** conditional:

Image translation, Inpainting, Super-resolution/Restoration...



R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. "High-resolution image synthesis with latent diffusion models". In *CVPR* 2022 (pp. 10684-10695).

# Various Types of Conditioning



"A doll in the shape of letter 'A'"

"A car with flying wings"

"Two girls"

"A cool man in the room"

"Two fluffy rabbit ears"

"A magic world, bright stars in sky"

"A skier, high quality"

Caption: "A woman sitting in a restaurant with a pizza in front of her"
Grounded text: table, pizza, person, wall, car, paper, chair, window, bottle, cup

Caption: "Elon Musk and Emma Watson on a movie poster"
Grounded text: Elon Musk, Emma Watson; Grounded style image: blue inset

Mou, Chong, et al. "T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models." *arXiv* 2023.

# Conditioning from Other Modalities



(a) Single Sound

(b) Mixed Sound

(c) Sound-Text Mix

"MinD-Vis"

Qin, Can, et al. "GlueGen: Plug and Play Multi-modal Encoders for X-to-image Generation." *arXiv* 2023.

Chen, Zijiao, et al. "Seeing beyond the brain: Conditional diffusion model with sparse masked modeling for vision decoding." *arXiv 2022*

# 2. State-of-the-art Text-to-image Models

# Preliminary: CLIP

## Contrastive Language-Image Pretraining (CLIP)

➢ Pretrained **aligned** encoders using contrastive loss

- ✓ Image
- ✓ Text

➢ Only employ the text encoder in text-to-image models

A. Radford, J.W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, and G. Krueger. "Learning transferable visual models from natural language supervision". In *ICML* 2021.

# Text-to-image Models



| Stability AI | OpenAI | Adobe | Midjourney Lab | Google/NVIDIA |
|:---:|:---:|:---:|:---:|:---:|
| Open-source | Query via website | Query via website | Query via Discord | Closed to public |

# Text-to-image Models



MIDJOURNEY        DALL-E 2        STABLEDIFFUSION        FIREFLY

film still, portrait of an old man, wrinkles, dignified look, grey silver hair, peculiar nose, wise, eternal wisdom and beauty, incredible lighting and camera work, depth of field, bokeh, screenshot from a hollywood movie

# MidJourney versions



| V1 | V2 | V3 | V4 | V5 | V5.1 |
|----|----|----|----|----|------|
| Released February 2022 | Released April 12, 2022 | Released July 25, 2022 | Released November 5, 2022 | Released March 15, 2022 | Released May 3, 2022 |

# **Stable Diffusion**

➢ LDM + CLIP

➢ Open-source. Many tools

➢ Various applications



Text-guided img2img

# Text-to-image Applications

# AI Art

## Art Made by AI Wins Fine Arts Competition

AI-generated artwork won a recent art competition in the US, sparking controversy and fury among artists

by **Belinda Teoh** — September 13, 2022 in **Art**, **Culture**, **Society**, **Tech**



Share on Facebook    Share on Twitter

An artwork made by Artificial Intelligence (AI) won first place at the Colorado State Fair's fine arts competition last week, sparking controversy about whether AI-generated art can be used to compete in competitions.

# Path for future VFX

# Path for future VFX

https://wonderdynamics.com/

# Story Synthesis

# Story Synthesis



https://onceuponabot.com/story

# Personalization



Input images — in the Acropolis — swimming — sleeping — in a doghouse — in a bucket — getting a haircut

N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein and K. Aberman. "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation". *arXiv preprint arXiv:2208.12242*.

# Personalization

Input images

Vincent Van Gogh  Michelangelo  Rembrandt

Johannes Vermeer  Pierre-Auguste Renoir  Leonardo da Vinci

N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein and K. Aberman. "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation". *arXiv preprint arXiv:2208.12242*.

# Personalization

N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein and K. Aberman. "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation". *arXiv preprint arXiv:2208.12242*.

# Text-guided Image Editing



Wang, Qian, et al. "MDP: A Generalized Framework for Text-Guided Image Editing by Manipulating the Diffusion Path." *arXiv* 2023

# Text-guided Image Editing



Wang, Qian, et al. "MDP: A Generalized Framework for Text-Guided Image Editing by Manipulating the Diffusion Path." *arXiv* 2023
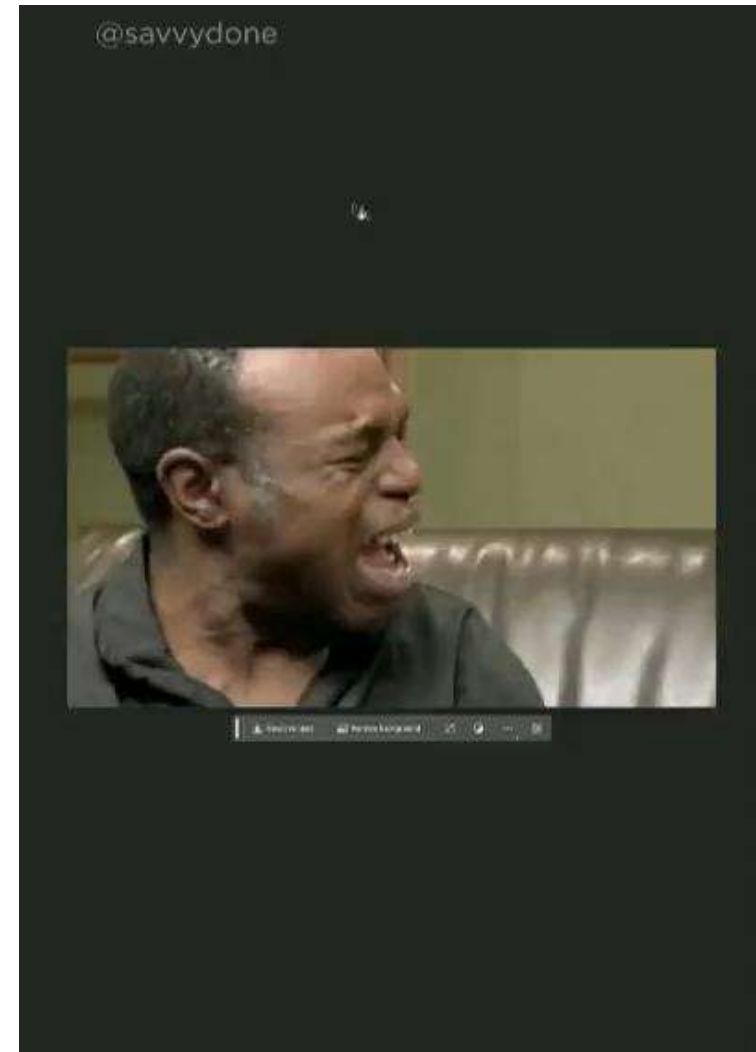
# Outpainting

Generative Fill (Photoshop Beta)

# Outpainting

**Generative Fill (Photoshop Beta)**

# Outpainting

**Zoom-out (Midjourney 5.2)**

# 3. State-of-the-art Text-to-X Models

# Text-to-Video



Melting ice cream dripping down the cone.



Campfire at night in a snowy forest with starry sky in the background.



Wooden figurine surfing on a surfboard in space



A happy elephant wearing a birthday hat walking under the sea

https://makeavideo.studio          https://imagen.research.google/video/                    https://phenaki.video/

# Text-guided Video Editing

**+ Van Gogh Style Painting**

**+ Watercolor Painting**

**Bear ➡ A Red Tiger**
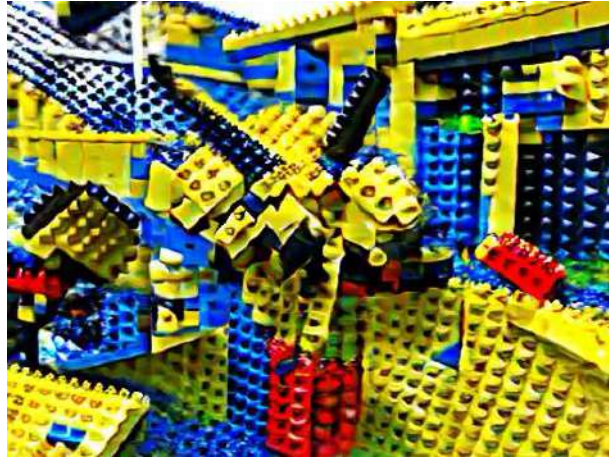
**Swan ➡ White Duck***

# Text-to-3D

A car made out of sushi.

A peacock on a surfboard.

https://dreamfusion3d.github.io/gallery.html

Magic3D: High-Resolution

# Text-guided 3D Editing

make it Lego blocks



convert it to a mechanical flower made of silver metal



make it chocolate



What if it was made of diamonds?

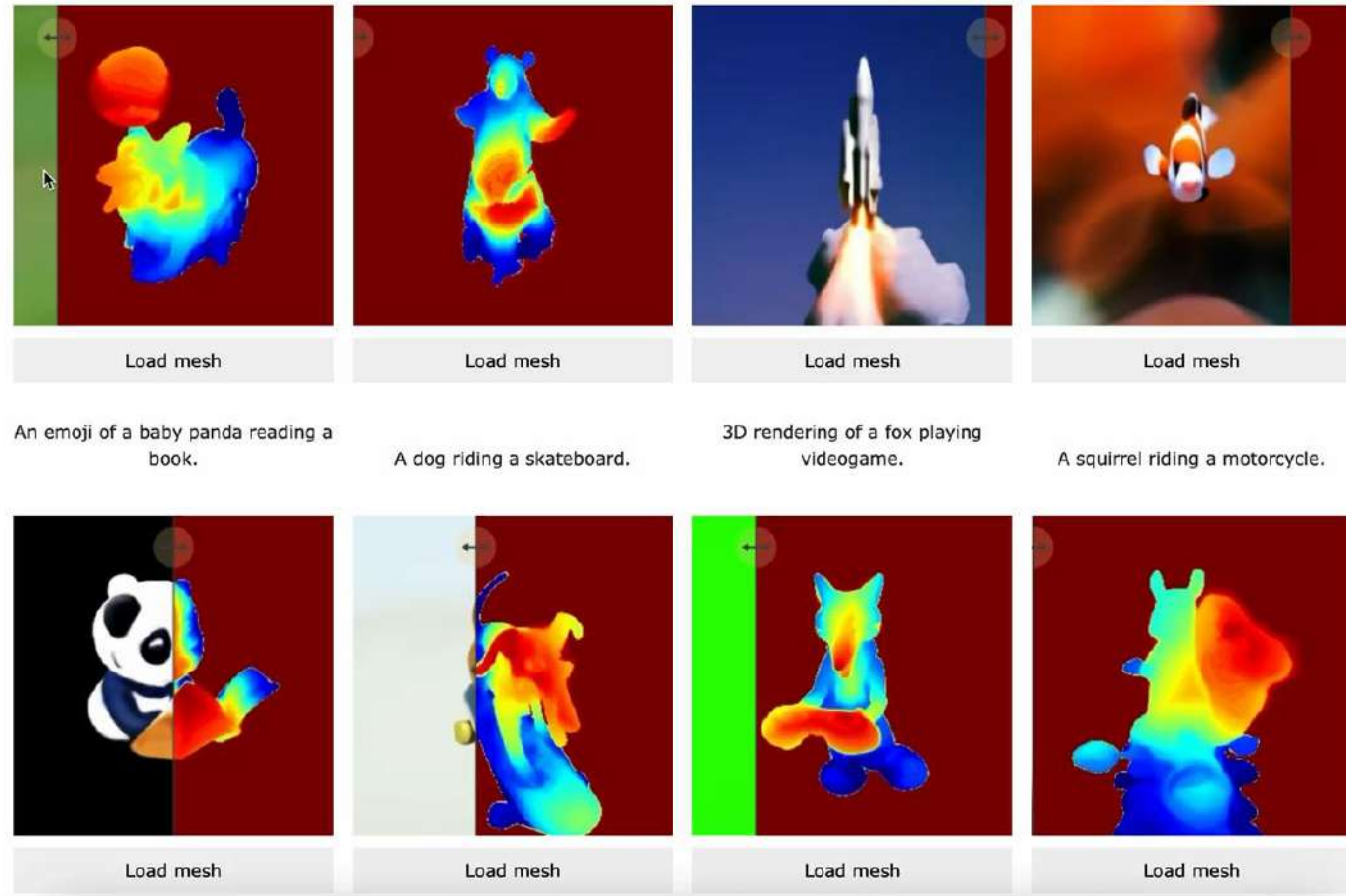# 3D Models from a Single 2D Image



GeNVS (nvlabs.github.io)

# Text-to-4D



Text-To-4D Dynamic Scene Generation (make-a-video3d.github.io)

# Computer Vision is picking up steam!!!

Thanks for listening