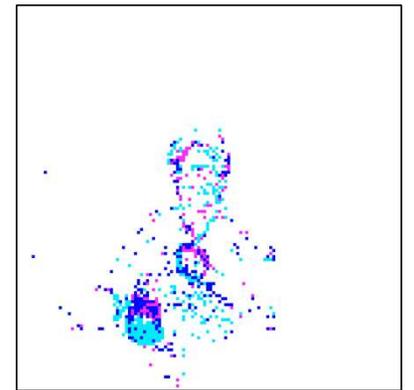# Text-to-Events: Synthetic Event Camera Streams from Conditional Text Input

Joachim Ott, Zuowen Wang & Shih-Chii Liu
Institute of Neuroinformatics,
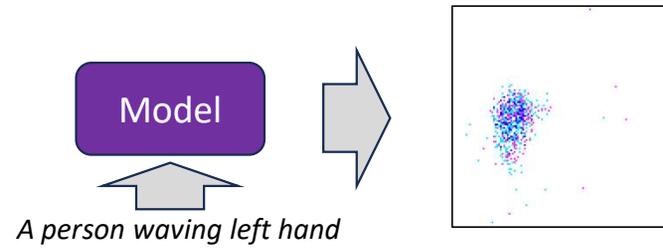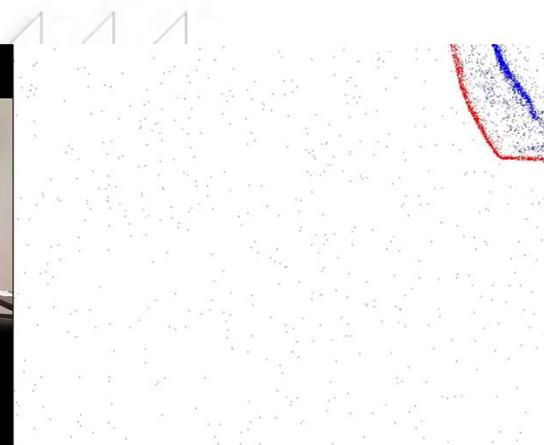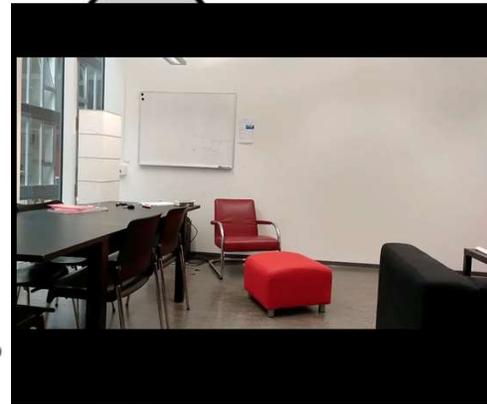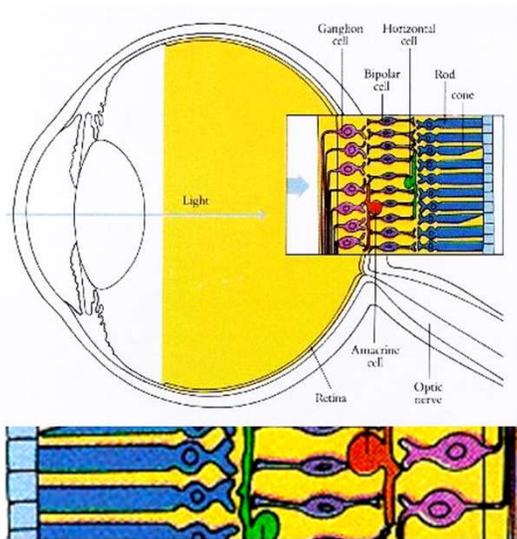University of Zurich & ETH Zurich

# Text2Events



*A person waving left hand*

Ott et al, NICE Conf, SD, 23.04.2024

# Event Camera or Dynamic Vision Sensor



*P. Lichtsteiner et al., JSSC, Feb. 2007*

Ott et al, NICE Conf, SD, 23.04.2024

'A ball bouncing on the grass'
[4]

*'A squirrel eating a burger.'*
[3]

*'A dog driving a car on a suburban street wearing funny sunglasses'*
[2]

*'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.'*
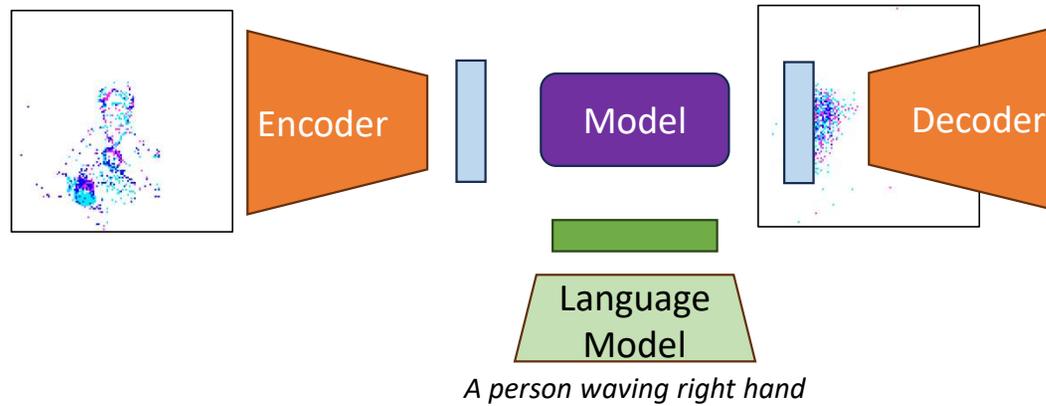
[1]

time

2023

2024

*Prompt texts and videos copied from https://research.nvidia.com/labs/toronto-ai/VideoLDM/, https://lumiere-video.github.io/, and https://openai.com/sora.*
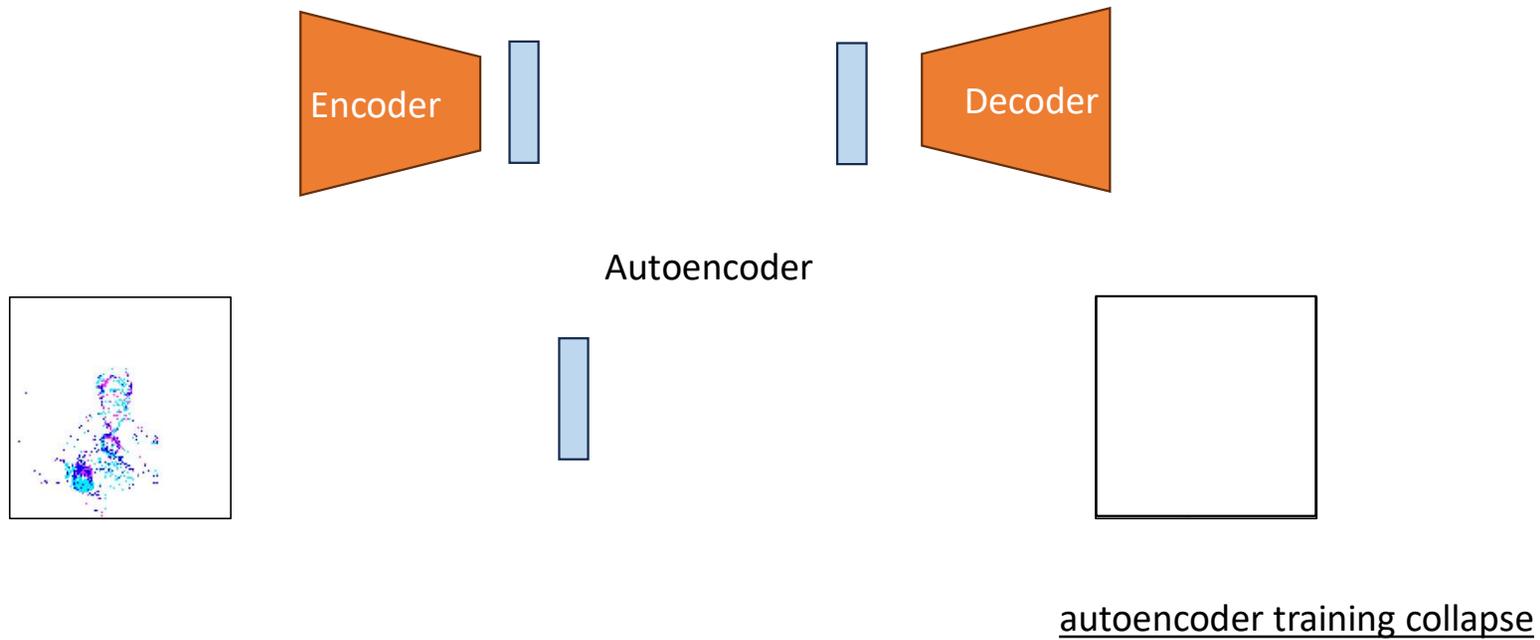*Publications: [1]Brooks, Tim et al. "Video generation models as world simulators", (2024), and [2] Bar-Tal, Omer, et al. "Lumiere: A space-time diffusion model for video generation." (2024), and [3] Blattmann, Andreas, et al. "Align your latents: High-resolution video synthesis with latent diffusion models." 2023, and [4] Luo, Zhengxiong, et al. "VideoFusion: Decomposed Diffusion Models for High-Quality Video Generation." 2023.*
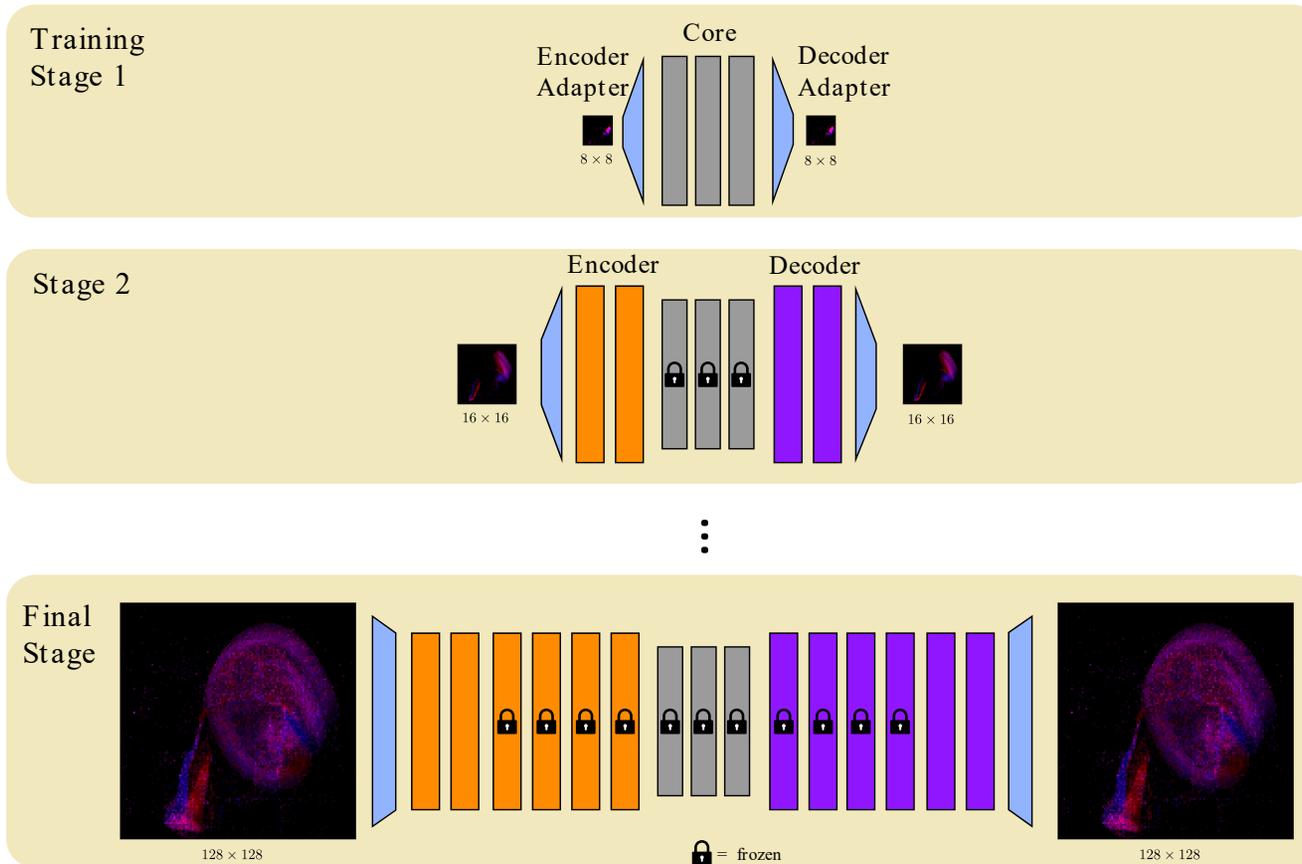
Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Latent Diffusion Model (LDM)



Encoder | Model | Decoder
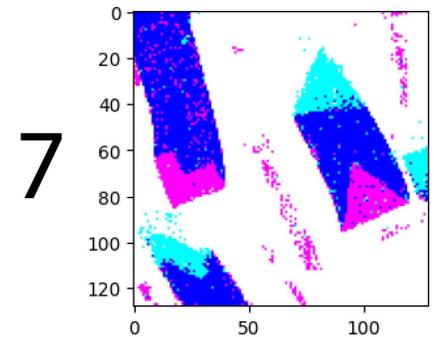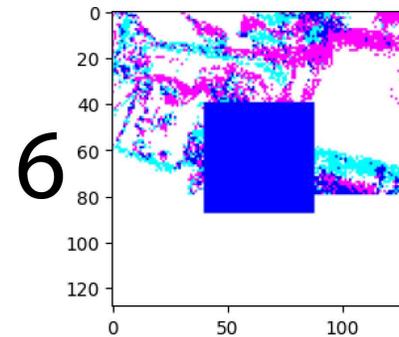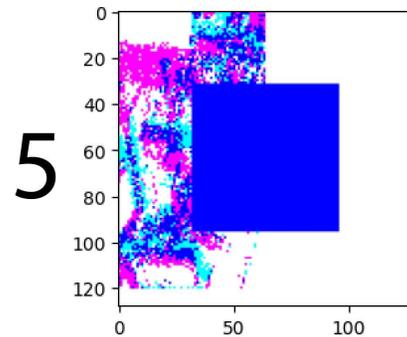
Language Model

*A person waving right hand*

Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Autoencoder
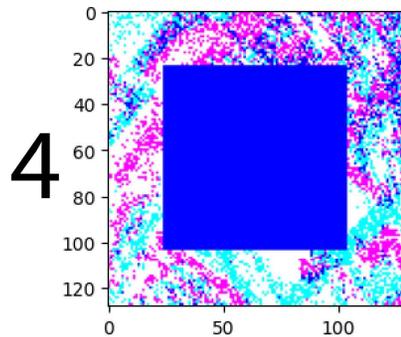


Autoencoder

autoencoder training collapse

Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events - Autoencoder



Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Warm-Up

Epoch



1

2

3

4

5

6

7

Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Latent Diffusion Model



Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Full Model



Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Full Model



**Parameter Count**

Autoencoder:
3.3 million

VLM:
12.3 million

LDM:
120.8 million

Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Experiments

## DVS 128 Dataset



arm roll

left hand clockwise

hand clap

a person is rotating the left arm clockwise

someone is rotating the left arm clockwise

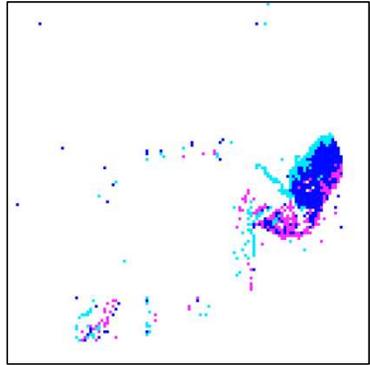this person is rotating their left arm clockwise

→ **36 new labels per class**
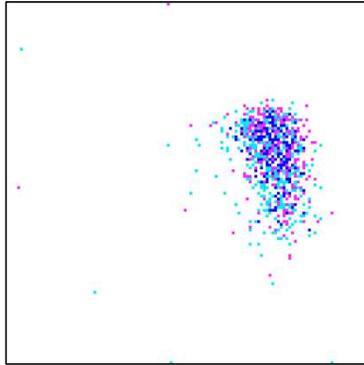
Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events - Results
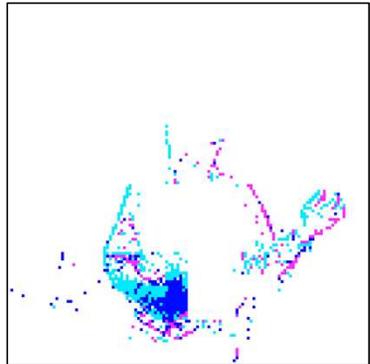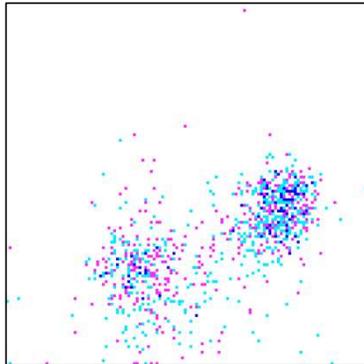
**Real**          **Generated**



left arm clockwise

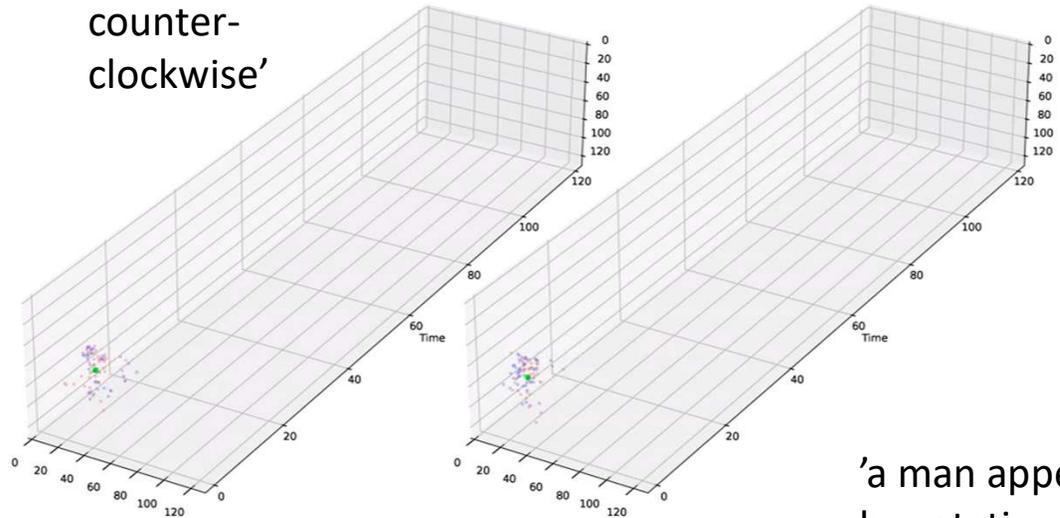'front view of someone rotating their left arm clockwise'



air guitar

'front view of someone playing air guitar'

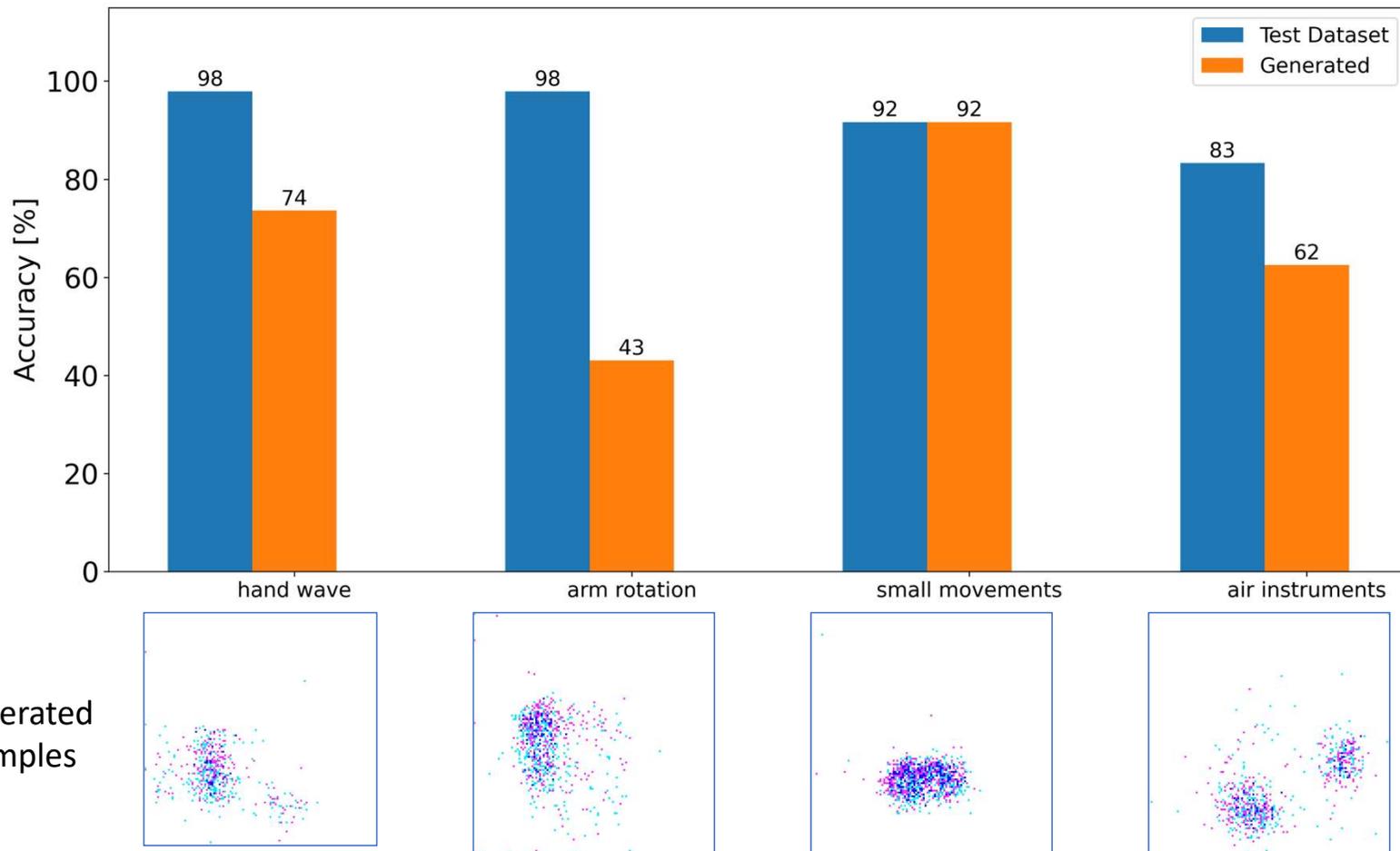'right arm counter-clockwise'



Real          Generated

'a man appears to be rotating his right arm counter-clockwise'
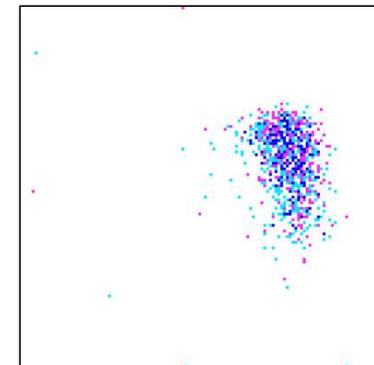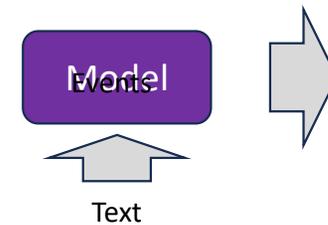
# Text2Events - Results



Test dataset: real event sequences

Generated: Event sequences generated by model

Generated examples

Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events – Summary

- First Text-To-Events model: Method of synthesizing vision event datasets using a latent diffusion model

- Method generates gesture event streams from text prompts instead of using a static text-to-events approach.

- Proposed a progressive training method for the autoencoder to encode from and decode to event frames


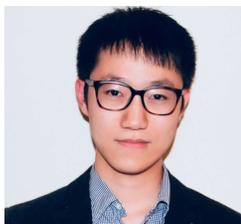
Ott et al, NICE Conf, SD, 23.04.2024

# Text2Events - Outlook

- Further improvements can be done, e.g. by combining intensity frames and events during training such as the frames and DVS event outputs of the DAVIS event camera
- Training the model to capture the behavior of each physical camera
- Universal text-to-events pipeline – can be applied to any sensor that produces events

Ott et al, NICE Conf, SD, 23.04.2024

Joachim Ott

Zuowen Wang

ETH zürich

https://sensors.ini.ch/

sensors
Institute of Neuroinformatics Zurich

University of
Zurich UZH

Ott et al, NICE Conf, SD, 23.04.2024