

# World Models: Background, Applications and Opportunities

Presented by Sakura

October 10, 2025

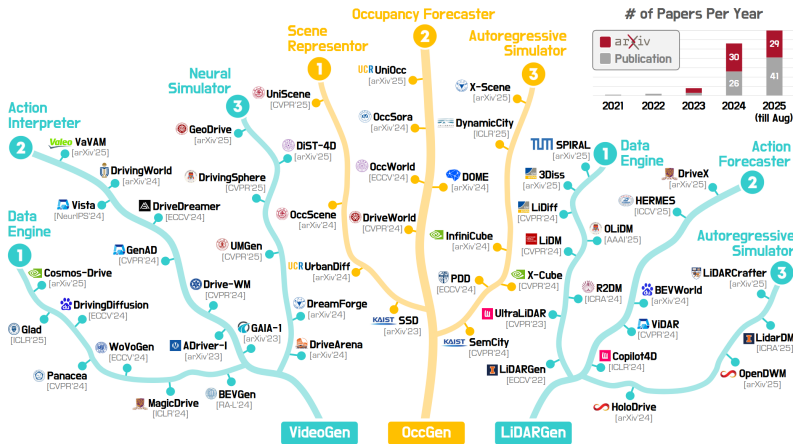
[bili\\_sakura@zju.edu.cn](mailto:bili_sakura@zju.edu.cn)



Teaser: Plato's allegory of the cave. Image Credit: [Wikipedia](#)

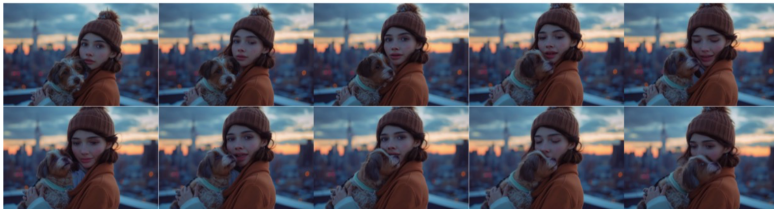
Note: This keynote is still in progress.

# Background and Related Works on World Modeling



**Figure:** Summary of representative video-based generation (VideoGen), occupancy-based generation (OccGen), and LiDARbased generation (LiDARGen) models from existing literature. Image Credit: (Kong et al., 2025)

**Text Prompt:** A girl lowers her head and rubs her face against a puppy, the puppy looks up at the girl



**Text Prompt:** A woman presses a camera shutter, her hair flying

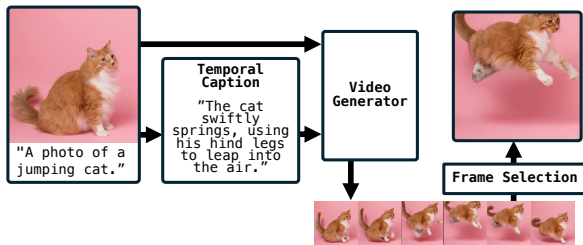


**Figure:** CogVideoX (Yang et al., 2025) for Text-to-Video Generation.

# Using CogVideoX for Image Editing

## Intuition

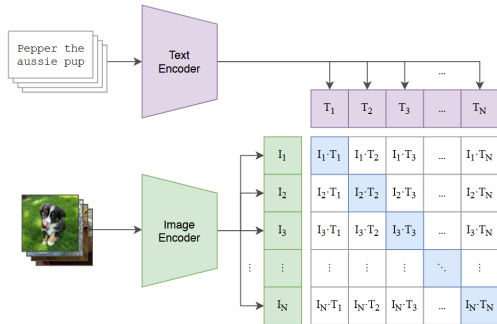
Intuition of using large video generation model is that the sequence of frames in video clip reflect the real-world physics with strong **image consistency**.



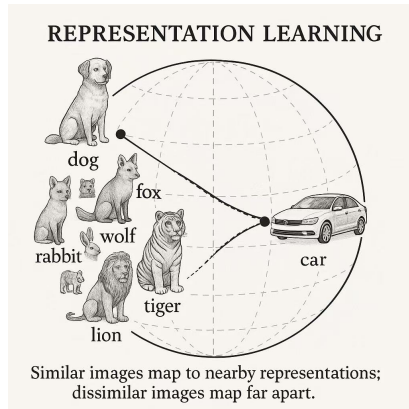
**Figure:** Frame2Frame (Rotstein et al., 2025) is a training-free method that uses a pretrained image-to-video diffusion model to synthesize a sequence of intermediate frames, and then selects the frame that best satisfies the edit.

# Representation Learning

## (1) Contrastive pre-training



CLIP: Visual & Semantic Representation



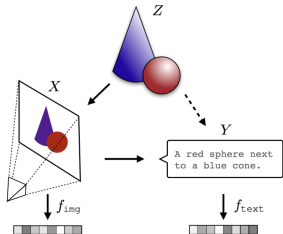
Spherical latent space (Shared embedding)

**Figure:** Overview of CLIP (Radford et al., 2021). A novel example of joint visual and semantic representation learning.

# Platonic Hypothesis

## The Platonic Representation Hypothesis

Neural networks, trained with different objectives on different data and modalities, are converging to a shared statistical model of reality in their representation spaces.



Platonic solids (Abstraction of forms)



Plato's Allegory of the Cave

The Platonic Representation Hypothesis (Huh et al., 2024): Images ( $X$ ) and text ( $Y$ ) are projections of a common underlying reality ( $Z$ ).