

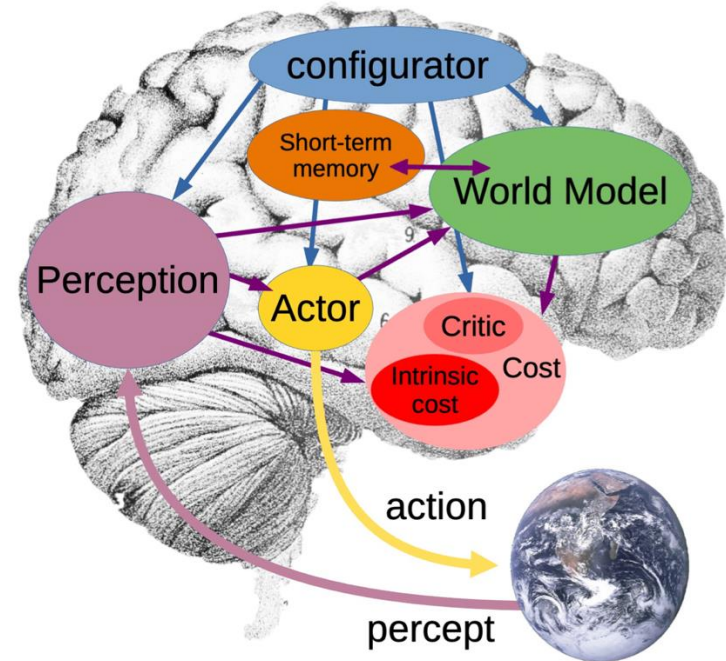
# Building Generalizable World Models for Autonomous Driving

Shenyuan Gao

# LeCun's AGI System



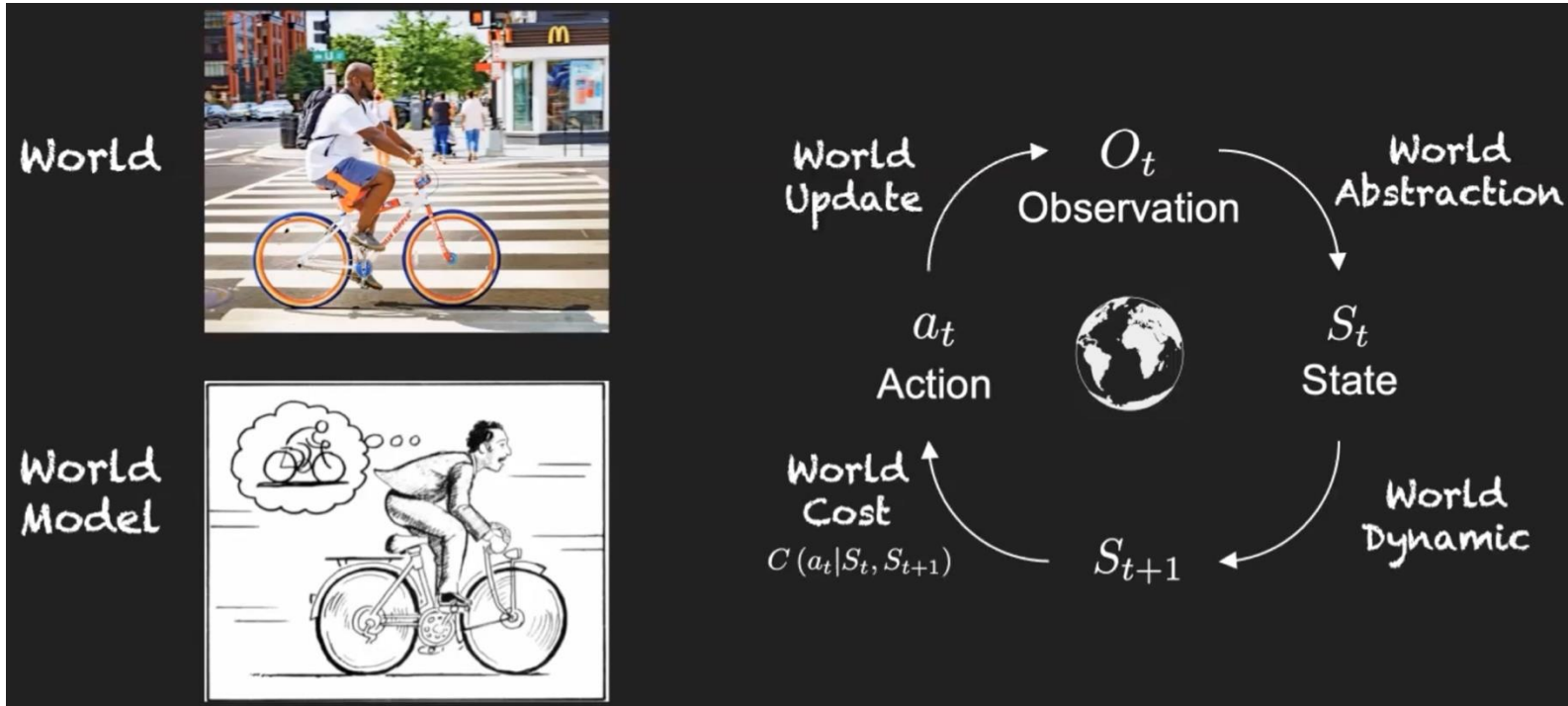
- Perception
  - Observe and estimate the world state
- Actor
  - Propose actions
- World model (core)
  - Predict plausible future state
- Cost
  - Evaluate the state



# What is “World Model”

- A substitute that simulates the real world
- Given previous states and actions, predict the future

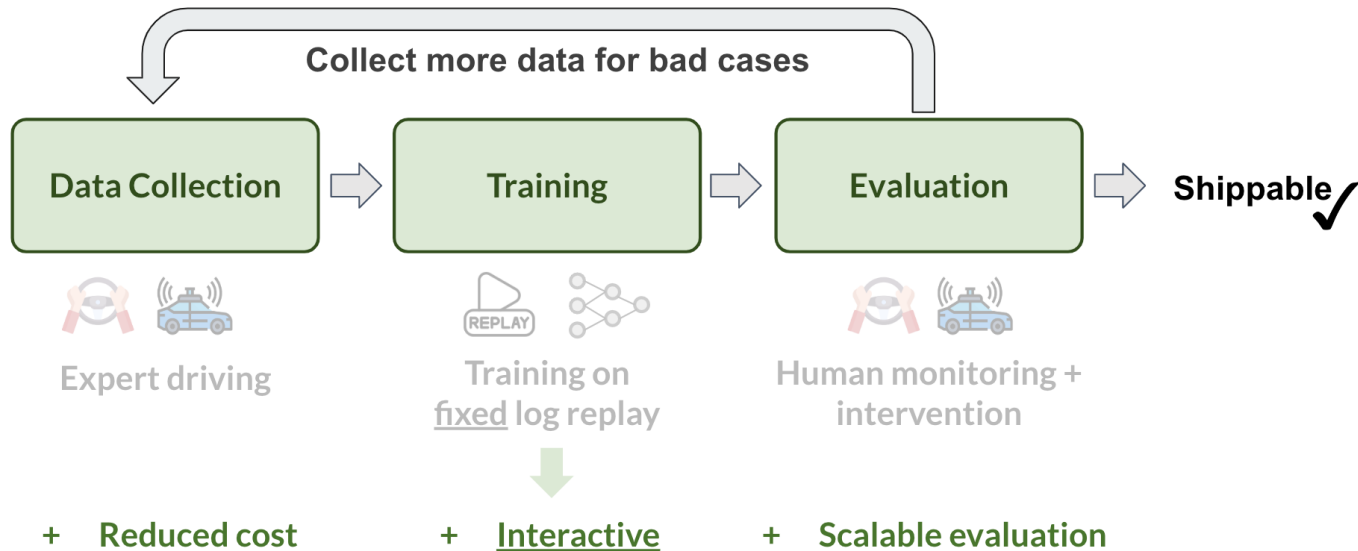
$$p(s_{t+1} | s_t, a_t)$$



# Introducing World Models to Autonomous Driving



- **Expectation**



# Limitations of Existing Driving World Models

## Generalization

- Limited data scale and geographic coverage



5 hours only Singapore & Boston



## Representation Capacity

- Inferior fidelity
- Low resolution and frame rate



## Action Control

- Single action modality
- No zero-shot action controllability



## Application

- Underexplored



- Not applicable for diverse real-world scenarios

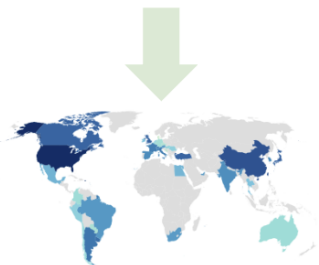
# Our Goal

## Generalization

- Limited data scale and geographic coverage



5 hours only Singapore & Boston



Worldwide  
Generalization

## Representation Capacity

- Inferior fidelity
- Low resolution and frame rate



High Fidelity  
High Resolution

## Action Control

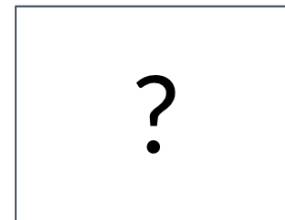
- Single action modality
- No zero-shot action controllability



Versatile Action  
Controllability

## Application

- Underexplored



Generalizable  
Reward Function

- How to build such a driving world model?







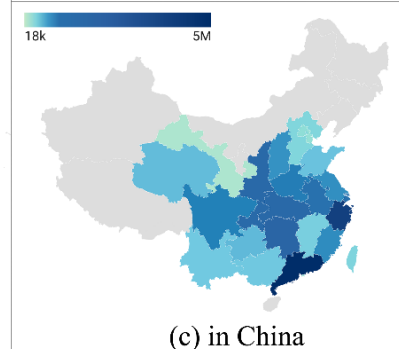
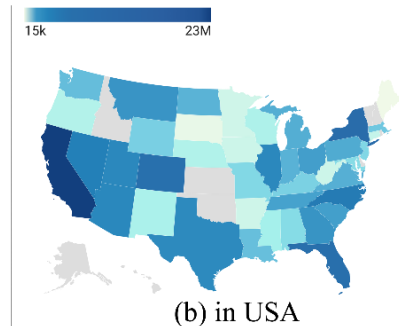
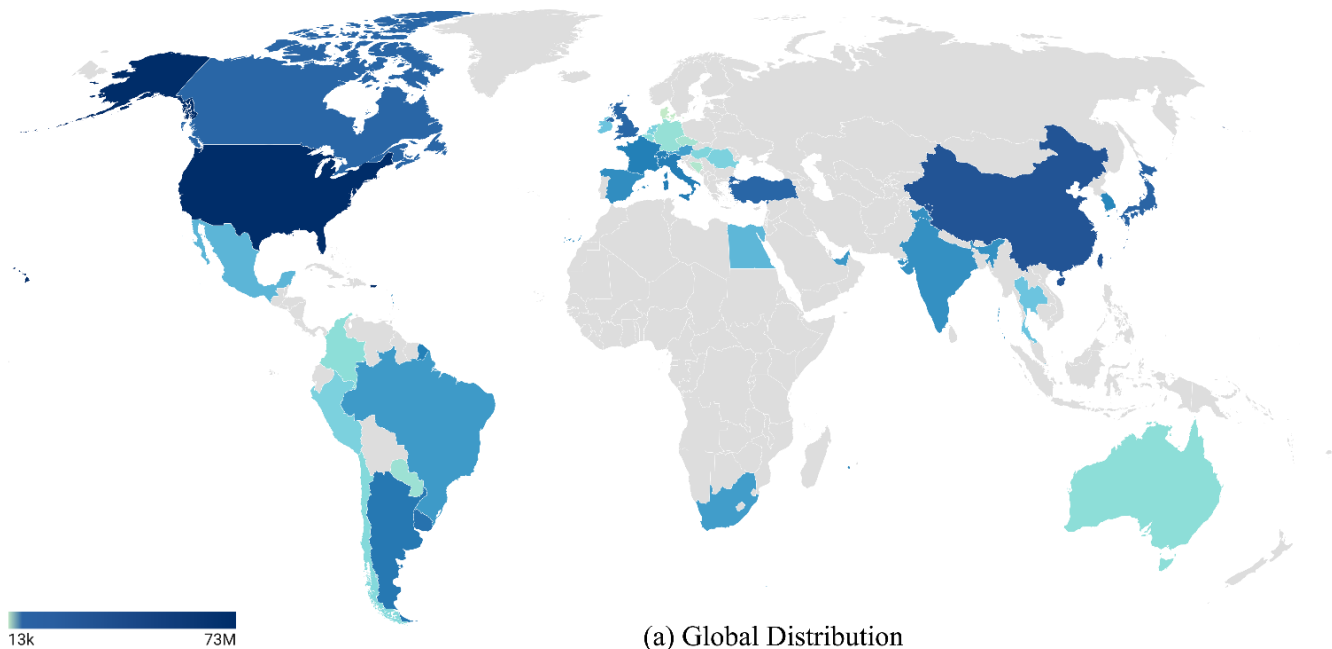
# OpenDV Dataset

Wind Walk Travel Videos

# Largest Public Driving Video Dataset

- Sourced from YouTube 
- 370x larger than nuScenes 
- 2059 hours, 40+ nations, 244+ cities, 65M+ frames
- Hours-long video samples

<https://github.com/OpenDriveLab/DriveAGI>





# From Generation Model to Prediction Model

- Base model: SVD (image2video generation, non-predictive)



- Our objective

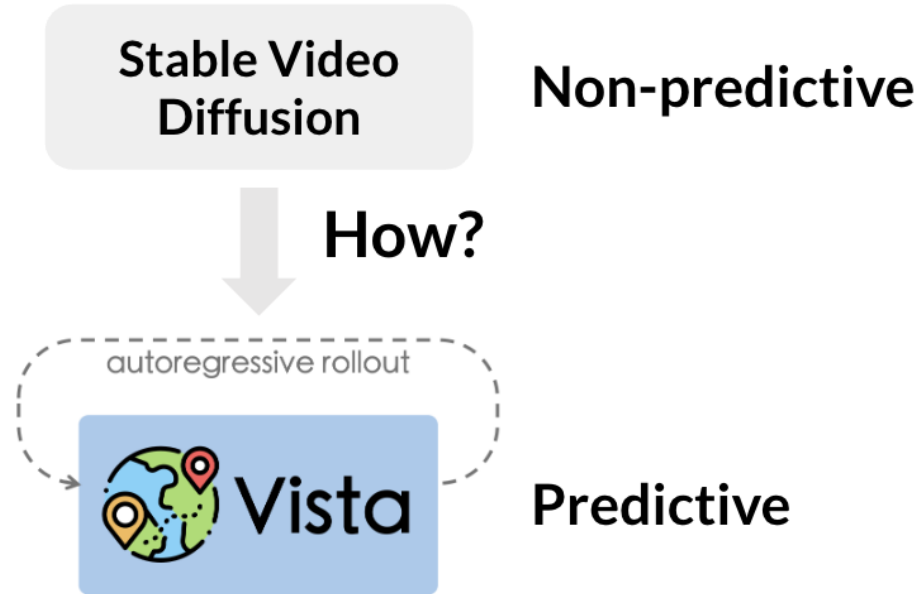


Coherent Long-Horizon Rollout

Blattmann, Andreas, et al. "Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets." arXiv preprint arXiv:2311.15127 (2023).

Gao, Shenyuan, et al. "Vista: A Generalizable Driving World Model with High Fidelity and Versatile Controllability." Advances in Neural Information Processing Systems 38 (2024).

# From Generation Model to Prediction Model

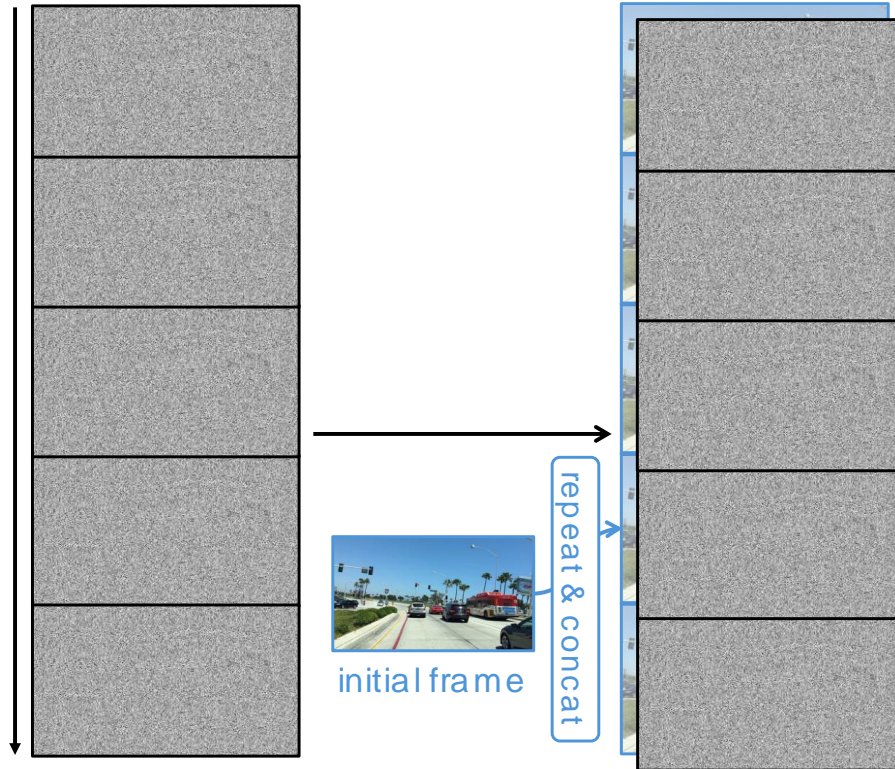


Blattmann, Andreas, et al. "Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets." arXiv preprint arXiv:2311.15127 (2023).

Gao, Shenyuan, et al. "Vista: A Generalizable Driving World Model with High Fidelity and Versatile Controllability." Advances in Neural Information Processing Systems 38 (2024).

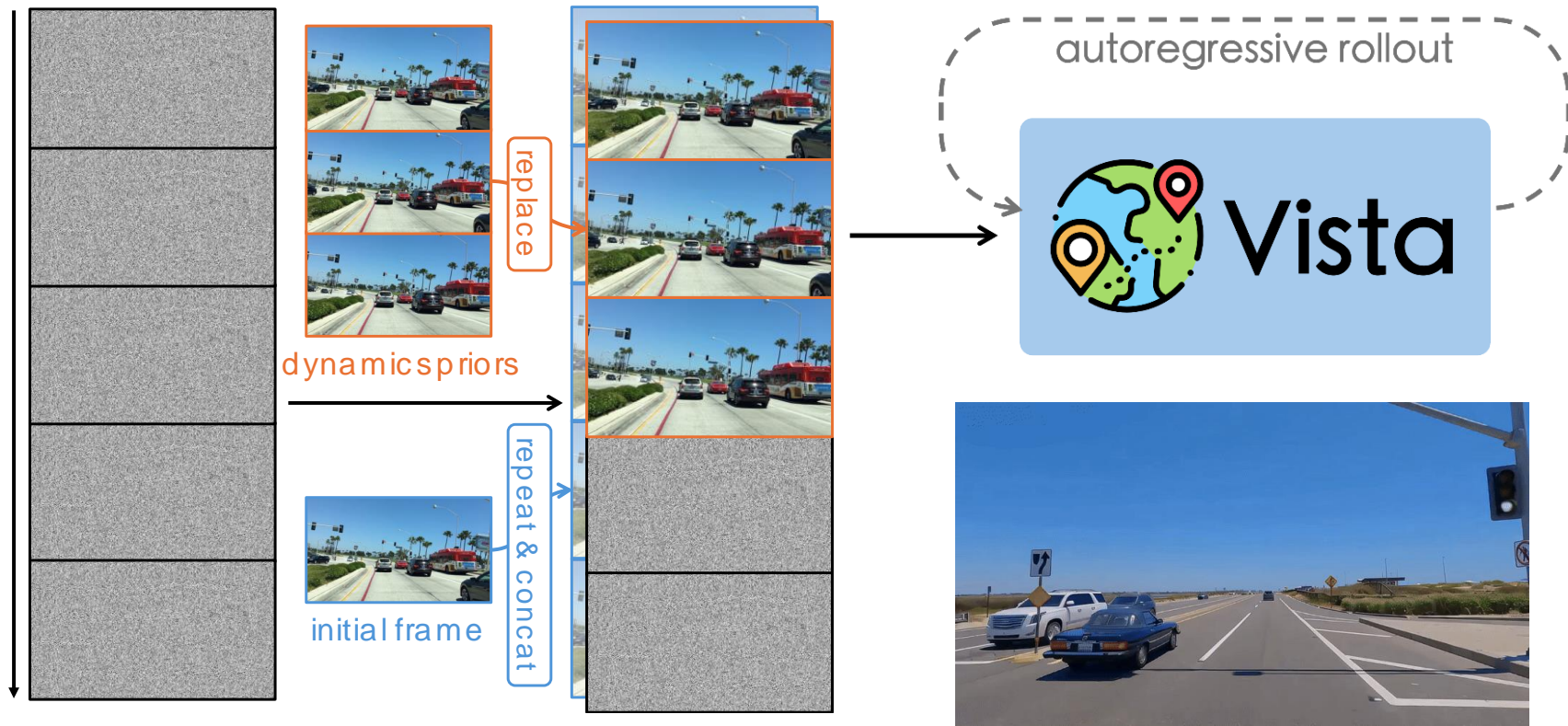
# From Generation Model to Prediction Model

- Impose the first frame as the condition image



# From Generation Model to Prediction Model

- Inject previous frames as dynamic priors to enable coherent rollout

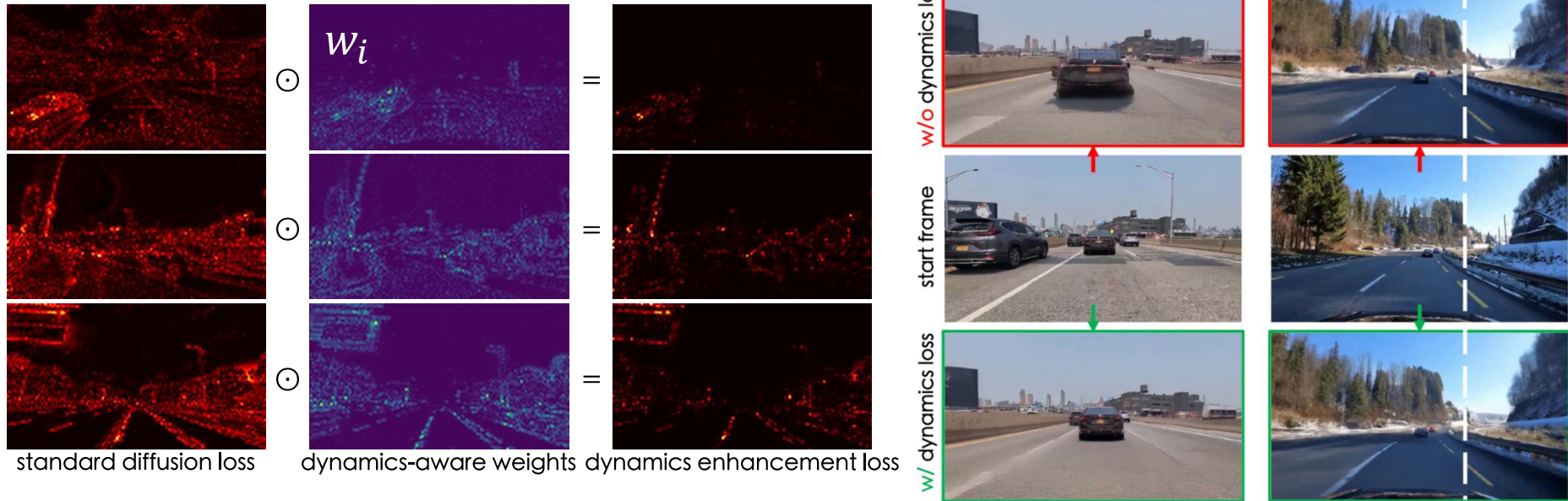




# Improving Fidelity (for Driving Scenarios)

- Dynamics-critical regions only occupy a rather small area
- **How to emphasize?**
  - Highlight motion disparities between prediction and ground truth

$$w_i = \|(D_\theta(\hat{n}_i; \sigma) - D_\theta(\hat{n}_{i-1}; \sigma)) - (z_i - z_{i-1})\|^2$$

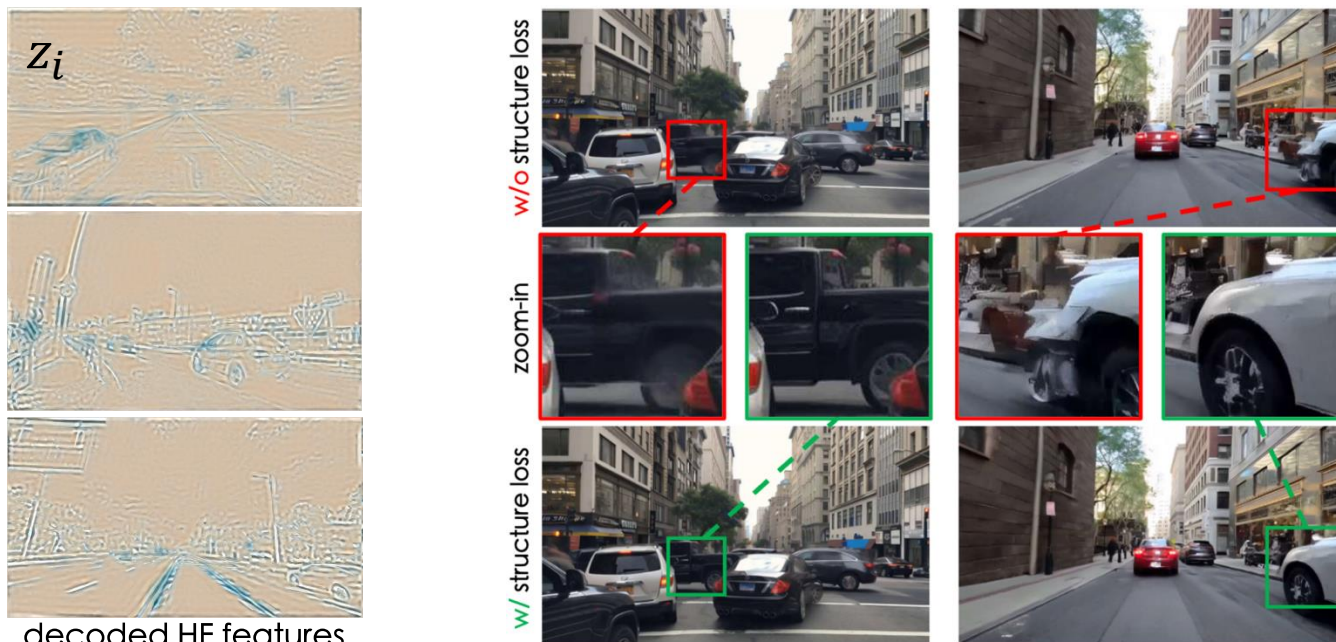




# Improving Fidelity (for Driving Scenarios)

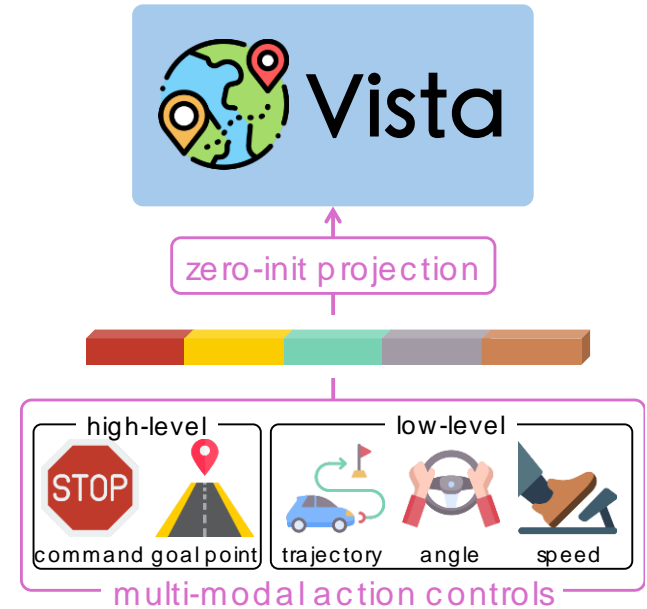
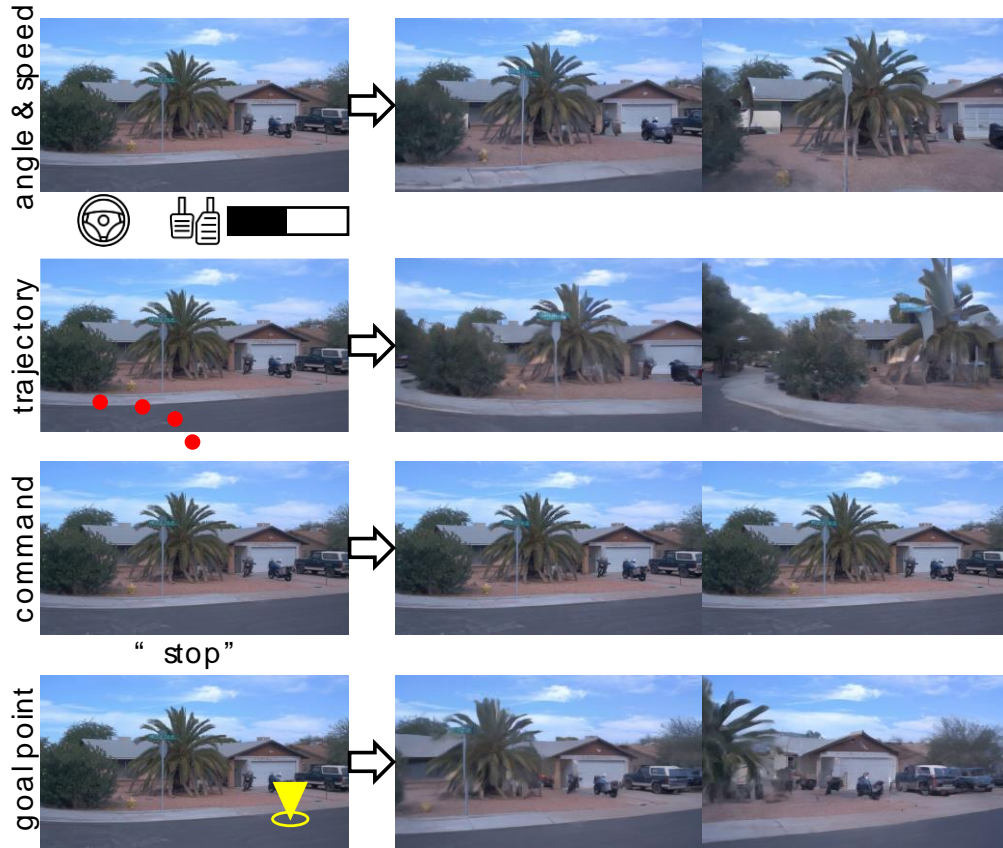
- Generating fast-moving objects tend to result in corruptions
- **How to alleviate?**
  - Enhance high-frequency structural information (e.g., edges and lanes)

$$z'_i = \mathcal{F}(z_i) = \text{IFFT}(\mathcal{H} \odot \text{FFT}(z_i))$$



# Multi-Modal Action Controllability

- Unified control interface for a versatile action suit



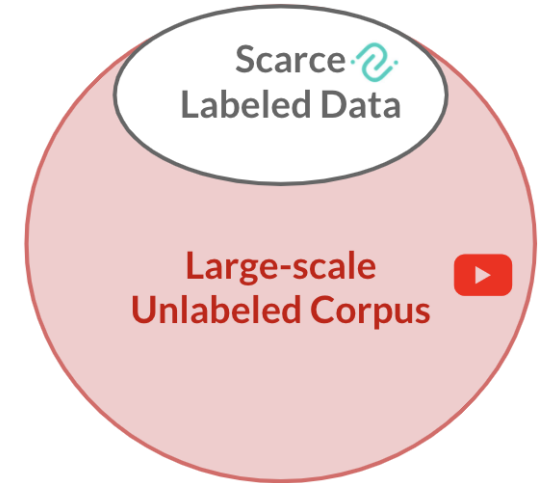
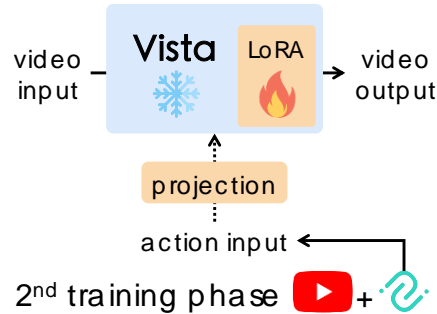
# Multi-Modal Action Controllability

- Two training phases

phase1: + generalization



phase2: + controllability



- Effective learning strategies
  - Multi-stage (low-res → high-res)
  - Parameter-efficient LoRA
  - Action independence constraint

# World Model as A Reward Function

- Our insights
  - OOD action condition will increase prediction uncertainty
- Generalizable reward estimation
  - **Conditional variance** as a source of reward
    - Without referring to ground truth trajectories
    - Inherit the strong generalization ability of Vista
    - Could be used to evaluate actions in the wild



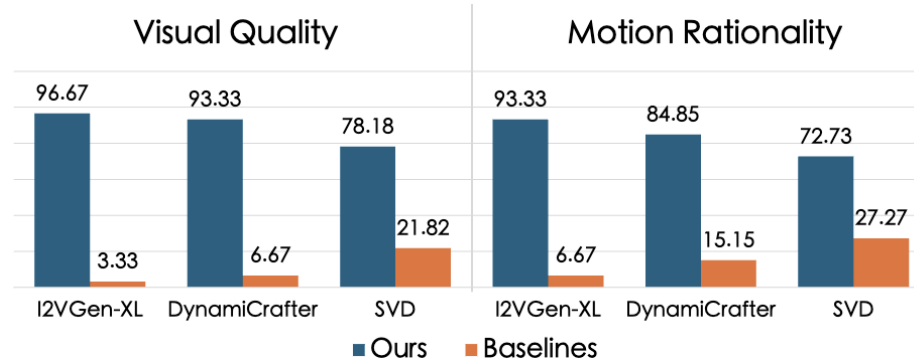
# Quantitative Results

ours

- Compared to existing driving world models

Metric	DriveGAN [100]	DriveDreamer [123]	WoVoGen [88]	Drive-WM [125]	GenAD [134]	Vista (Ours)
FID ↓	73.4	52.6	27.6	15.8	15.4	<b>6.9</b>
FVD ↓	502.3	452.0	417.7	122.7	184.0	<b>89.4</b>

- Compared to prominent video generators



Yang, Jiazhi, et al. "Generalized Predictive Model for Autonomous Driving." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

Gao, Shenyuan, et al. "Vista: A Generalizable Driving World Model with High Fidelity and Versatile Controllability." Advances in Neural Information Processing Systems 38 (2024).



# Visualization: Open-World Generalization

- Generalization to novel scenarios



- Coherent long-horizon rollout (up to 15s)



# Visualization: Action Controllability

turn left



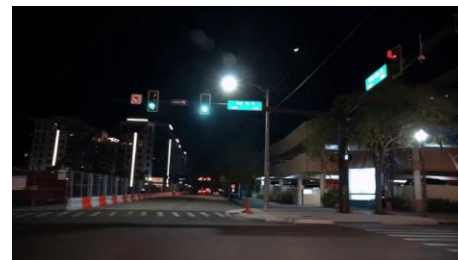
go straight



turn right



stop





# Visualization: Counterfactual Reasoning Ability



- Could be used for close-loop simulation

# Our Code and Model are Open-Sourced!

- OpenDV Dataset

<https://github.com/OpenDriveLab/DriveAGI>

- Vista Code & Model

<https://github.com/OpenDriveLab/Vista>

- Video Demo

<https://opendriveab.com/Vista>

Yang, Jiazhi, et al. "Generalized Predictive Model for Autonomous Driving." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2024.

Gao, Shenyuan, et al. "Vista: A Generalizable Driving World Model with High Fidelity and Versatile Controllability." Advances in Neural Information Processing Systems 38 (2024).