

DS-GA.3001 Embodied Learning and Vision

Mengye Ren

NYU

Spring 2026

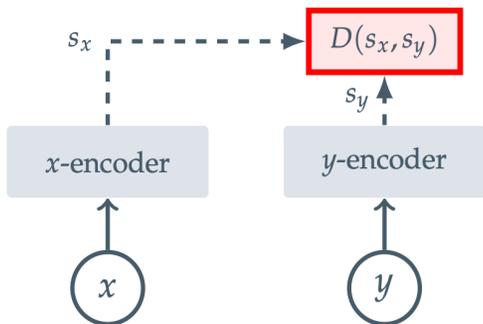
elvcourse.org





Joint Embedding Predictive Architecture (JEPA)

- Predict what changes from one view to another.

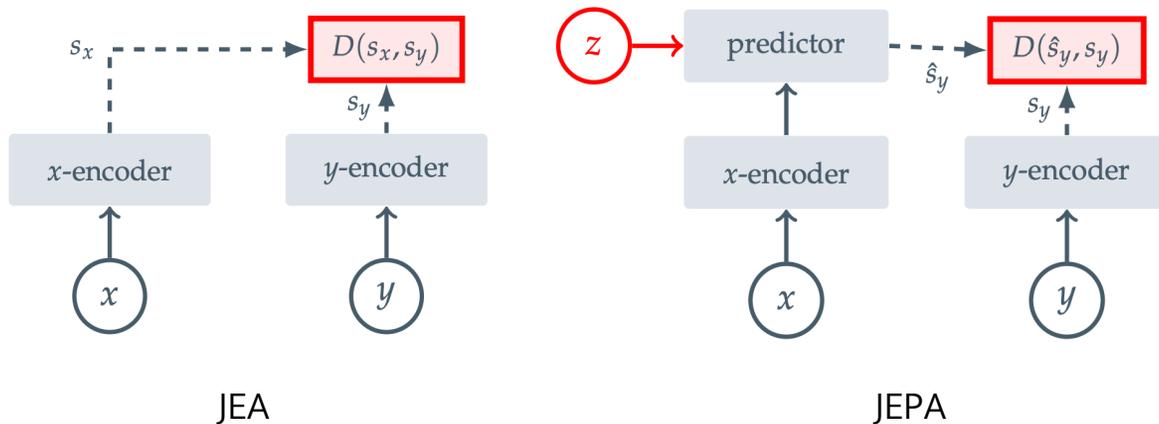


JEPA



Joint Embedding Predictive Architecture (JEPA)

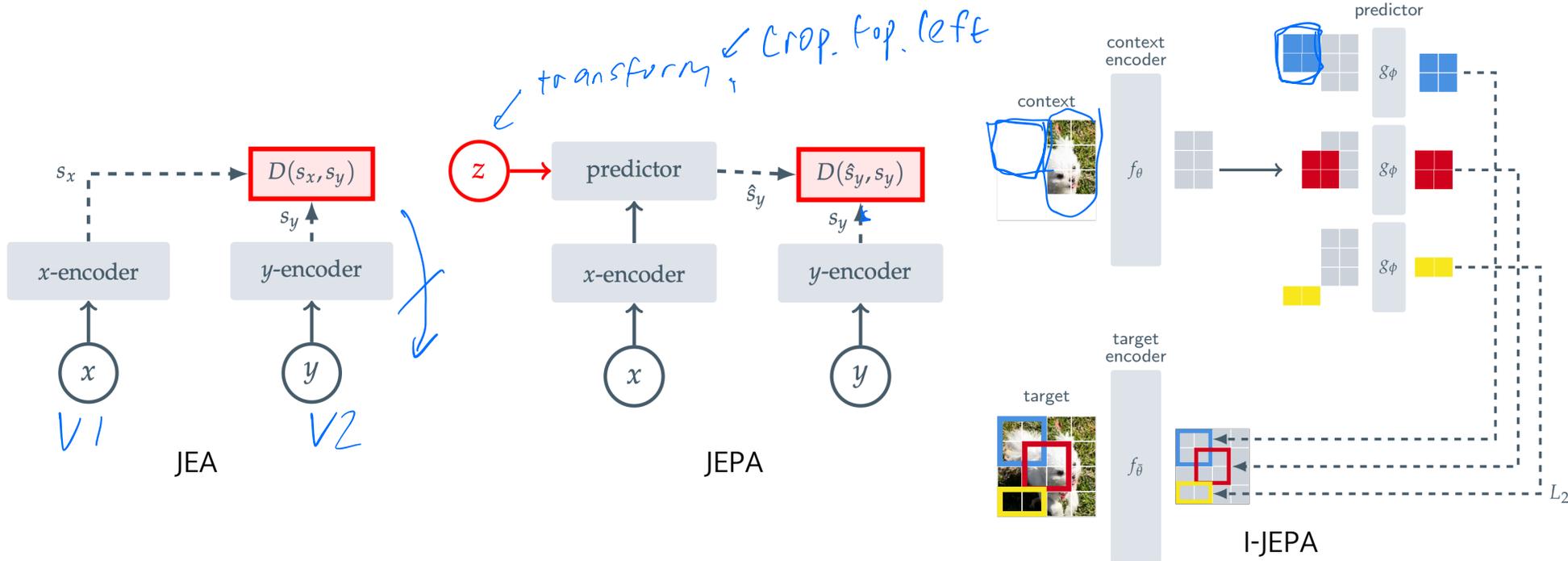
- Predict what changes from one view to another.





Joint Embedding Predictive Architecture (JEPA)

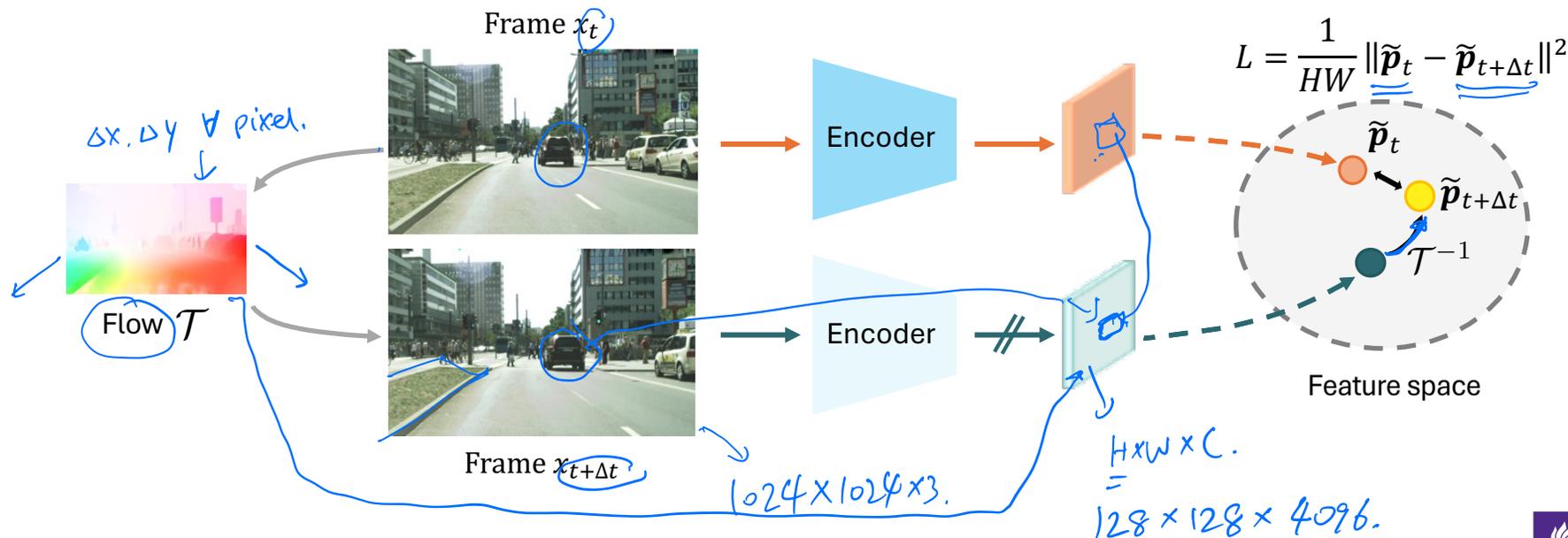
- Predict what changes from one view to another.



Joint Embedding Prediction with Dense Motion

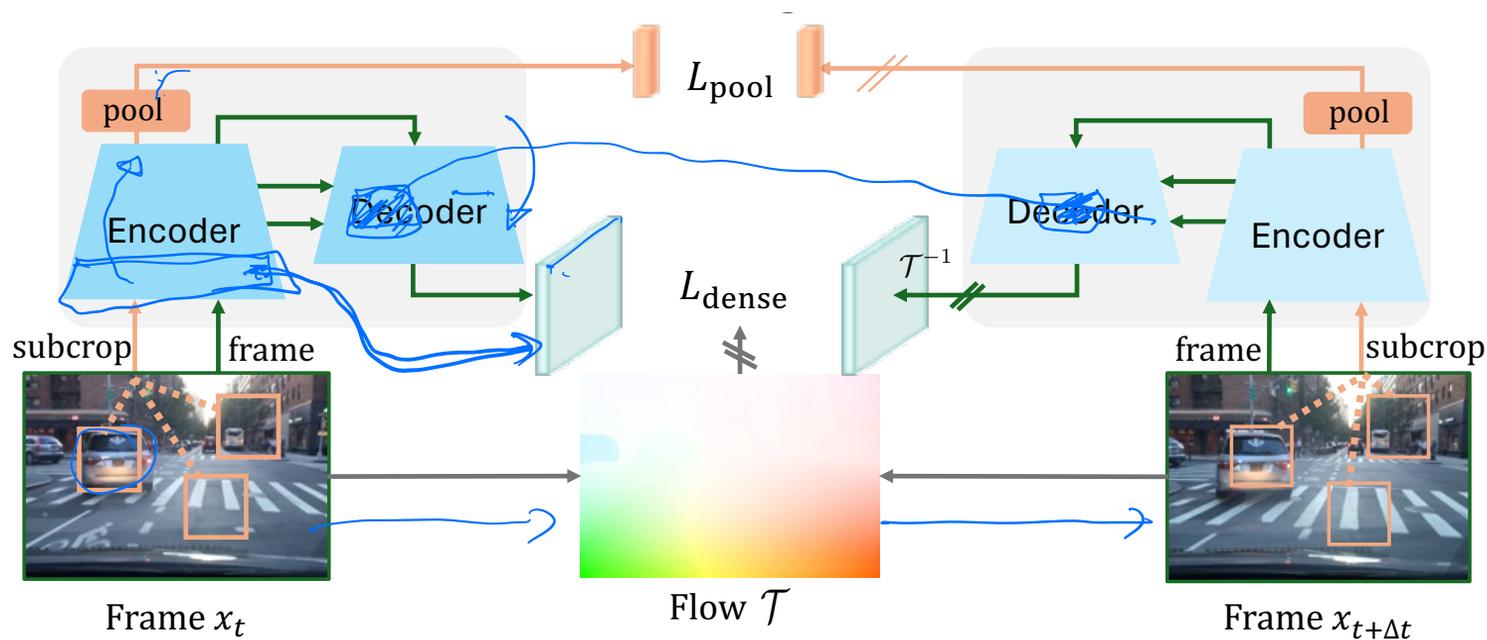
- Can we use adjacent video frames as self-supervision?
- Objects move densely throughout the image.

motion equivariance.



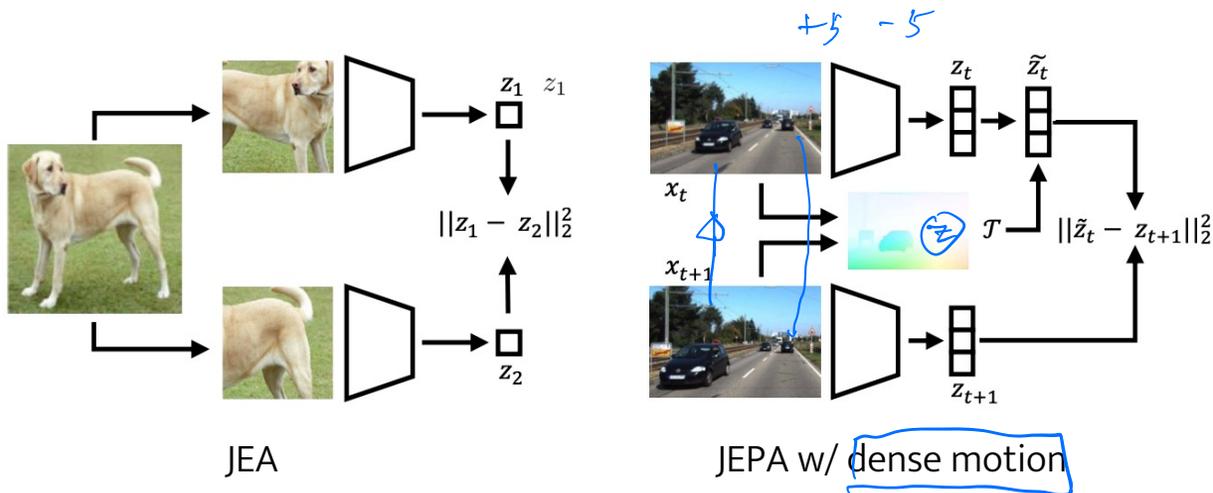
Joint Embedding Prediction with Dense Motion

- Perform SSL in multiple scales (small objects vs. big regions).



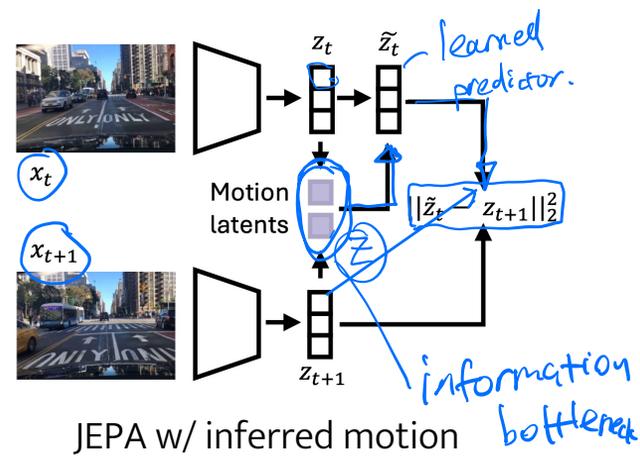
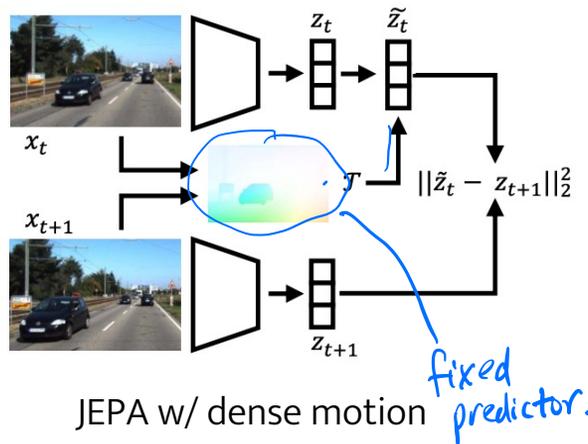
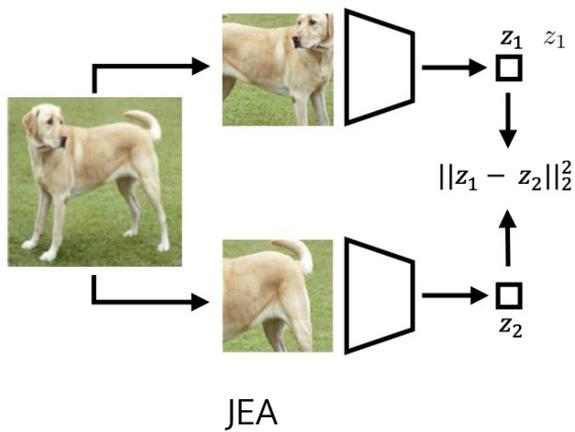
Joint Motion and Semantic Learning

- We either create or assume some external motion / action information to send to predictors.
- Additional motion needs to be inferred.



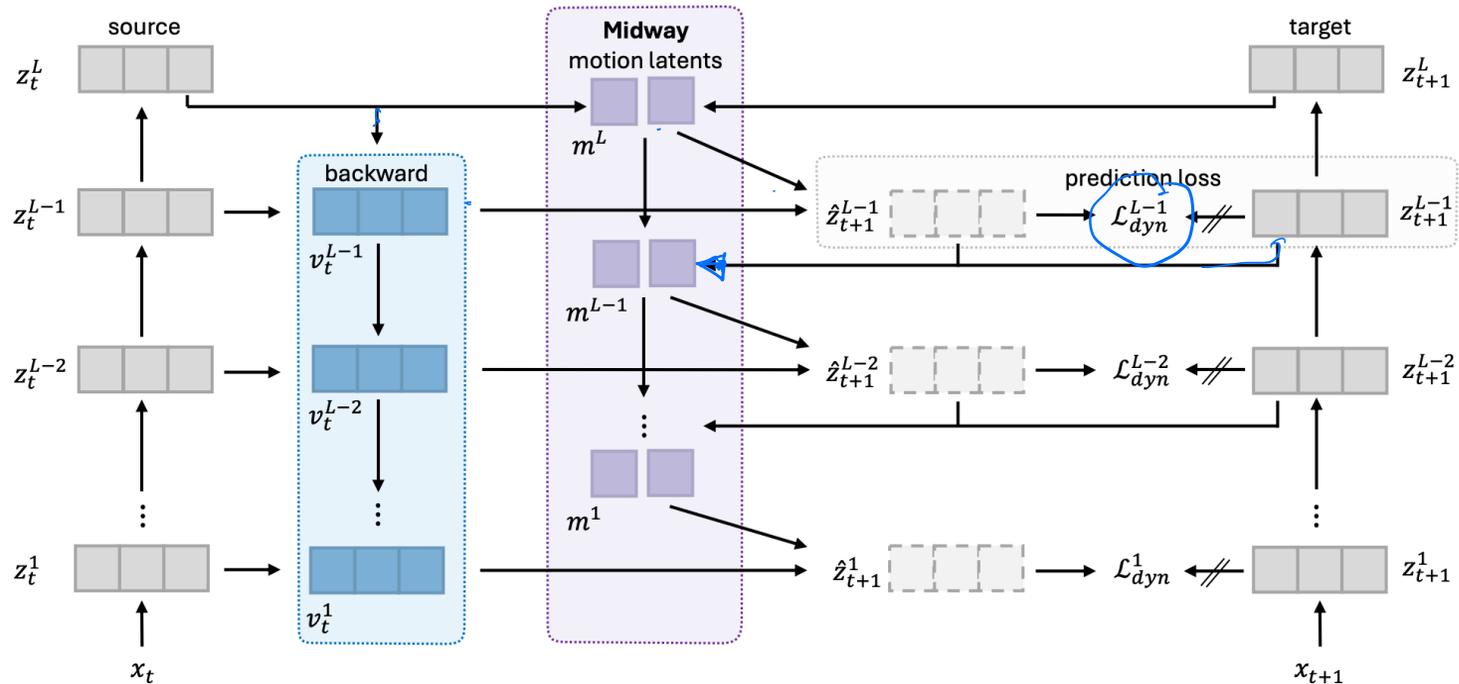
Joint Motion and Semantic Learning

- We either create or assume some external motion / action information to send to predictors.
- Additional motion needs to be inferred.



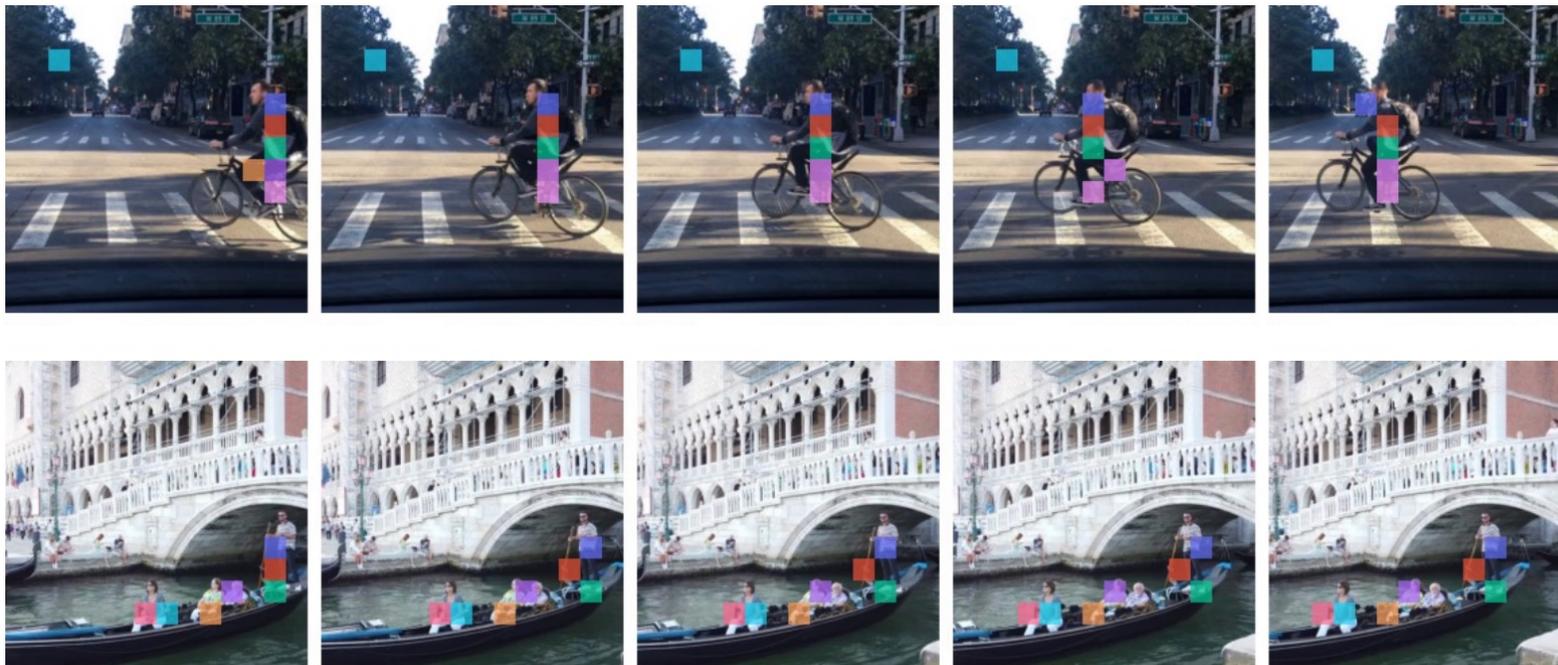
Midway Networks

- Using backward refinement from Optical Flow networks.



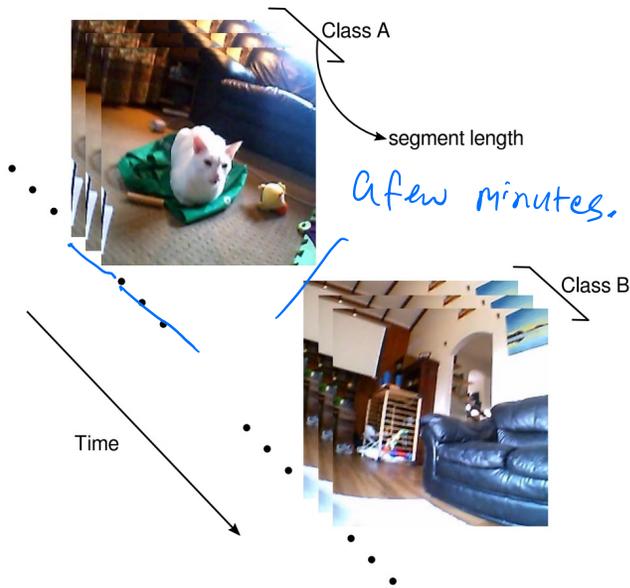
Midway Networks

- Use perturbation to visualize the motion learned from motion tokens.

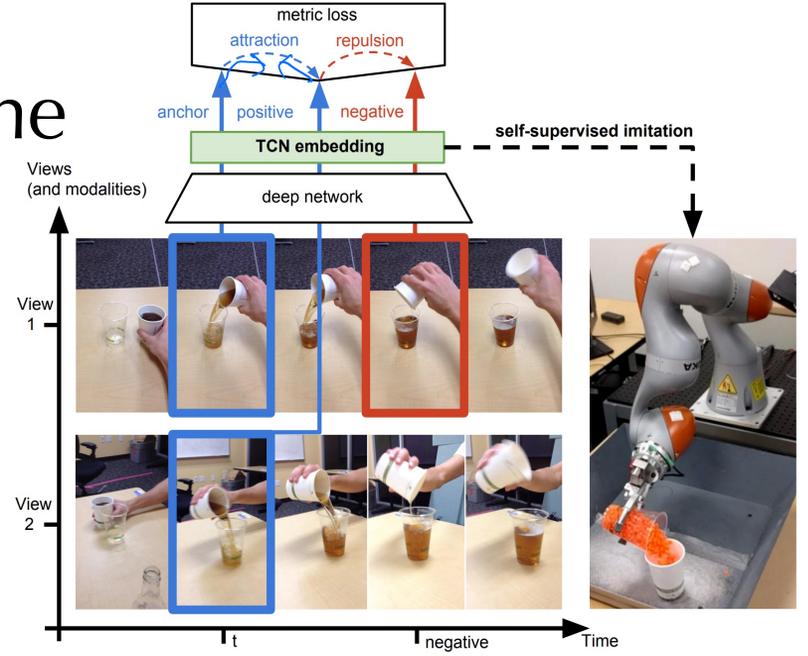
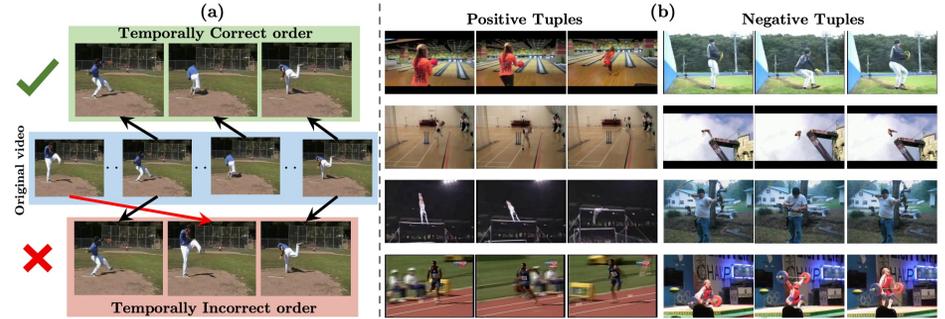


SSL with Time

- Use time as an additional source of supervision.

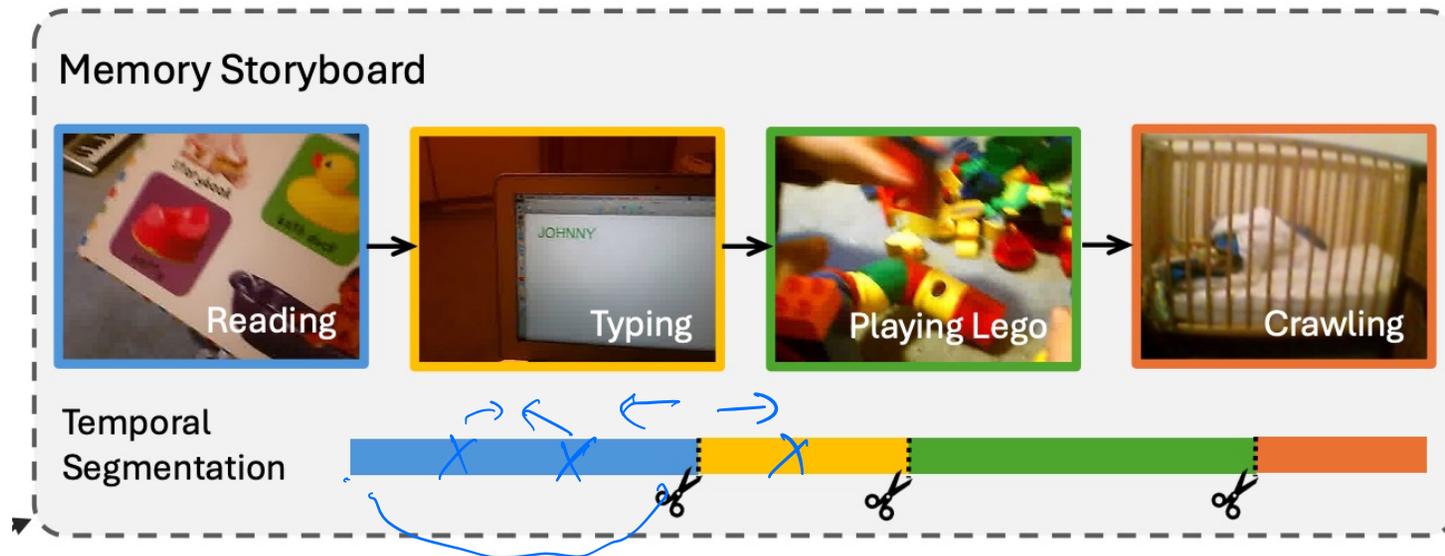


Similar adjacent frames.

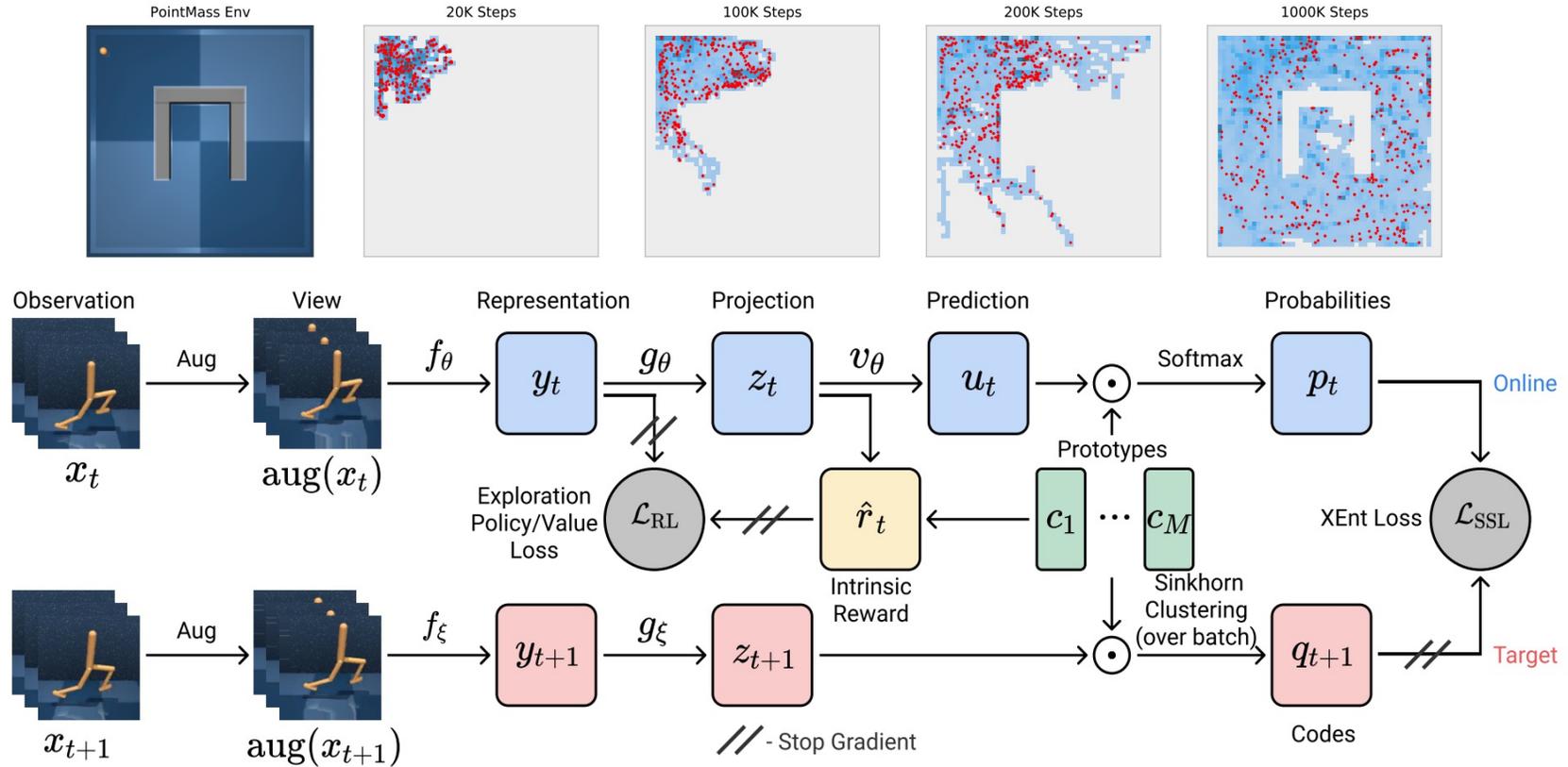


SSL with Time

- We can segment videos into meaningful events.
- Leverage the spatiotemporal continuity structure.



SSL for Visual Control

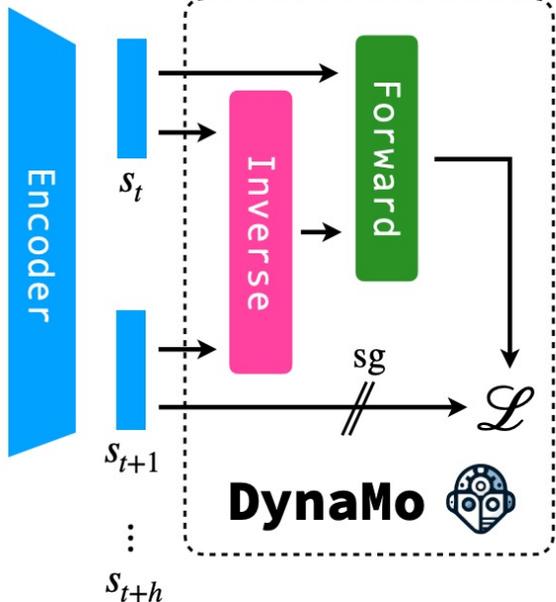


SSL for Visual Control

Observations

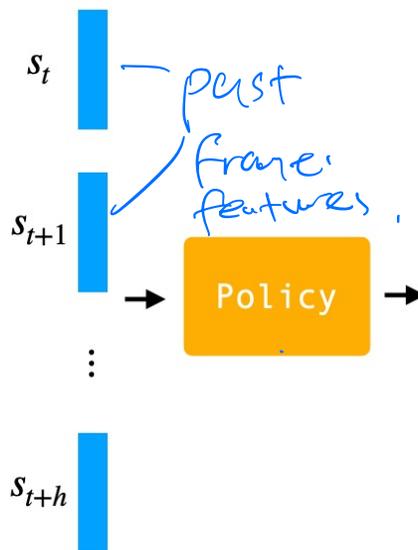


\vdots
 O_{t+h}

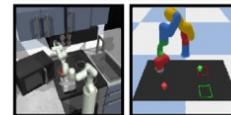


(a) Representation learning

Embeddings



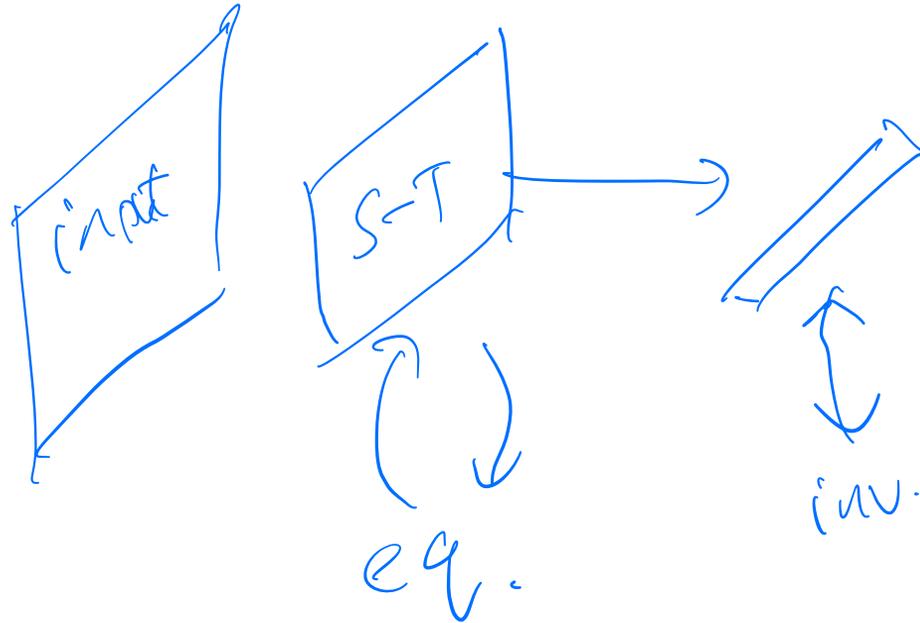
Environments



(b) Policy on pretrained representations

Summary

- Representation learning leverage the information in unlabeled data.
- Inductive biases matter
 - spatial-time invariance
 - motion equivariance



Summary

- Representation learning leverage the information in unlabeled data.
- Inductive biases matter
 - spatial-time invariance
 - motion equivariance
- A foundation for sensorimotor learning – pretrained representations as a starter for policy learning.

Emergent Attention, Object Discovery

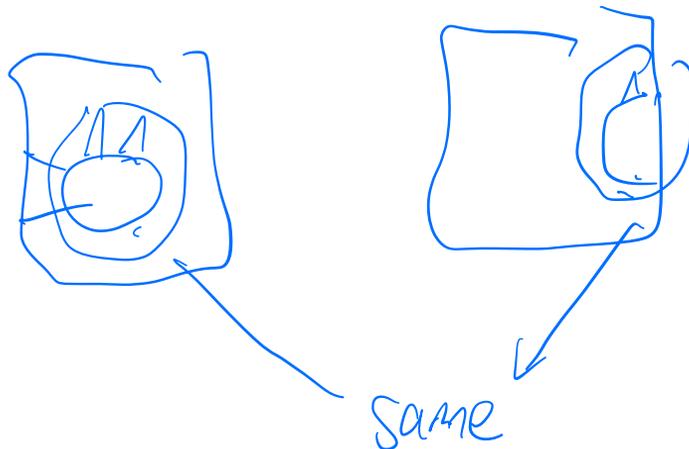
- SSL representations show awareness of object classes and instance identities.

Emergent Attention, Object Discovery

- SSL representations show awareness of object classes and instance identities.
- Why does attention show awareness of objects?

Emergent Attention, Object Discovery

- SSL representations show awareness of object classes and instance identities.
- Why does attention show awareness of objects?
- The network is encouraged to associate different parts of the objects together in order to identify whether two inputs belong to the same image or not.



Emergent Attention, Object Discovery

- SSL representations show awareness of object classes and instance identities.
- Why does attention show awareness of objects?
- The network is encouraged to associate different parts of the objects together in order to identify whether two inputs belong to the same image or not.
- Attending to semantically similar parts facilitates the process.

Emergent Attention, Object Discovery

- SSL representations show awareness of object classes and instance identities.
- Why does attention show awareness of objects?
- The network is encouraged to associate different parts of the objects together in order to identify whether two inputs belong to the same image or not.
- Attending to semantically similar parts facilitates the process.
- The network is a hierarchical information processing pipeline – Lower layers integrate more granular and smaller neighborhood.

Weak-to-Strong Supervision

- General idea: Use self-supervised learning to learn good features, which allow us to generate low-quality masks.

Weak-to-Strong Supervision

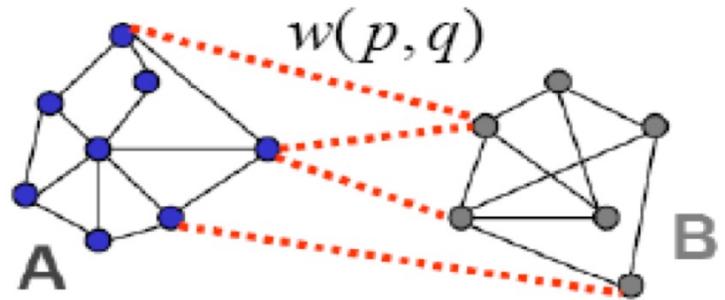
- General idea: Use self-supervised learning to learn good features, which allow us to generate low-quality masks.
- Then use these masks as pseudo labels and supervise the network to predict these low-quality masks.

Weak-to-Strong Supervision

- General idea: Use self-supervised learning to learn good features, which allow us to generate low-quality masks.
- Then use these masks as pseudo labels and supervise the network to predict these low-quality masks.
- Question: how do we come up with masks? What loss is used to supervise the network?

Graph Cut

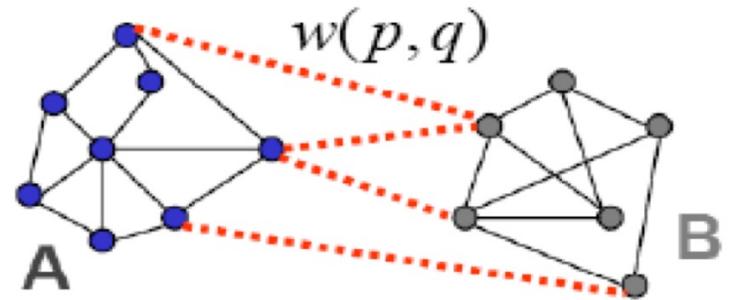
- Segmentation is essentially a clustering problem.



$$\text{cut}(A, B) = \sum_{p \in A, q \in B} w(p, q)$$

Graph Cut

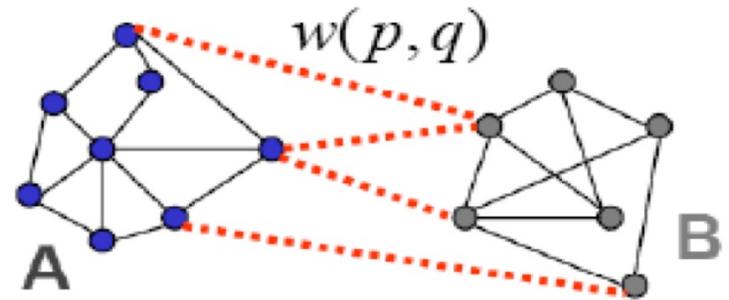
- Segmentation is essentially a clustering problem.
- We can transform the clustering problem with the graph cut problem.



$$\text{cut}(A, B) = \sum_{p \in A, q \in B} w(p, q)$$

Graph Cut

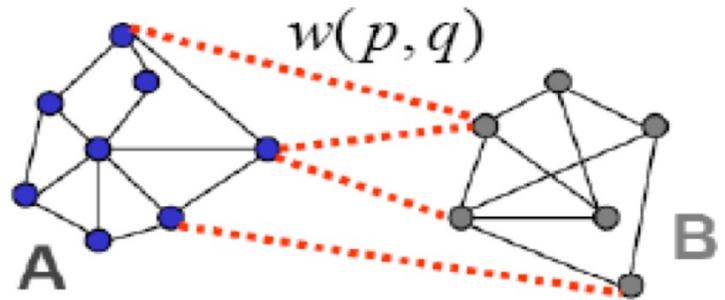
- Segmentation is essentially a clustering problem.
- We can transform the clustering problem with the graph cut problem.
- Pixel = node.



$$\text{cut}(A, B) = \sum_{p \in A, q \in B} w(p, q)$$

Graph Cut

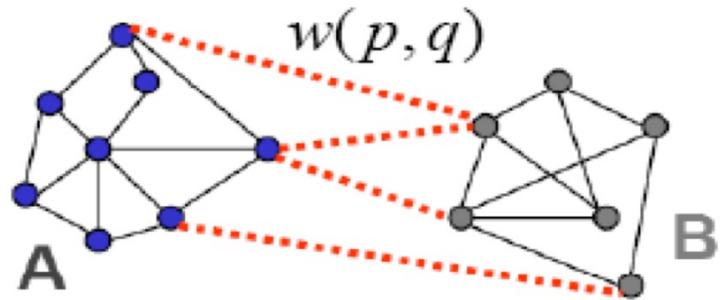
- Segmentation is essentially a clustering problem.
- We can transform the clustering problem with the graph cut problem.
- Pixel = node.
- Affinity between the two pixels = edge value (flow).



$$cut(A, B) = \sum_{p \in A, q \in B} w(p, q)$$

Graph Cut

- Segmentation is essentially a clustering problem.
- We can transform the clustering problem with the graph cut problem.
- Pixel = node.
- Affinity between the two pixels = edge value (flow).
- Objective: Cut the graph into disconnected components with a minimum sum of edge values.



$$\text{cut}(A, B) = \sum_{p \in A, q \in B} w(p, q)$$

NCut

- Sort the eigenvectors from the smallest to the largest.



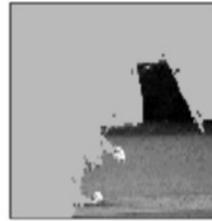
(a)



(b)



(c)



(d)



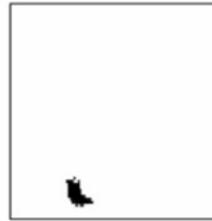
(e)



(f)



(g)



(h)

NCut

- Sort the eigenvectors from the smallest to the largest.
- This was a classic image segmentation technique operating directly on image intensity.



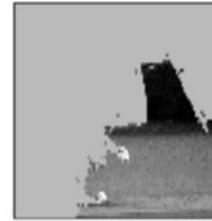
(a)



(b)



(c)



(d)



(e)



(f)



(g)



(h)

NCut

- Sort the eigenvectors from the smallest to the largest.
- This was a classic image segmentation technique operating directly on image intensity.
- Now, instead of segmenting pixels, we can directly segment semantically meaningful representations from self-supervision.



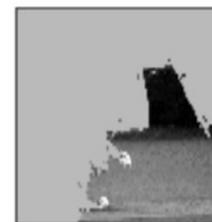
(a)



(b)



(c)



(d)



(e)



(f)



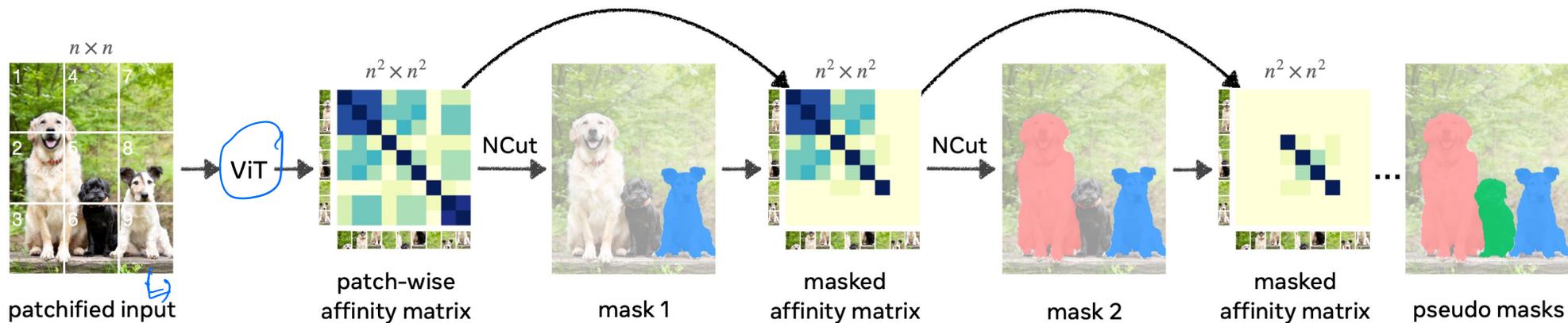
(g)



(h)

MaskCut

- Use a pretrained DINO ViT network.



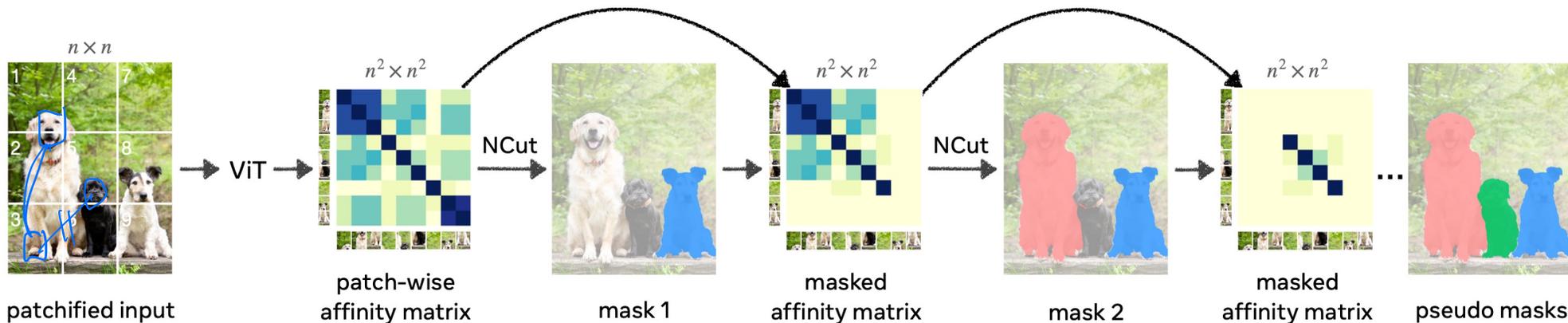
MaskCut

v_i token  f_{oken} 

- Use a pretrained DINO ViT network.

- Use the “key” features from the last attention layer:

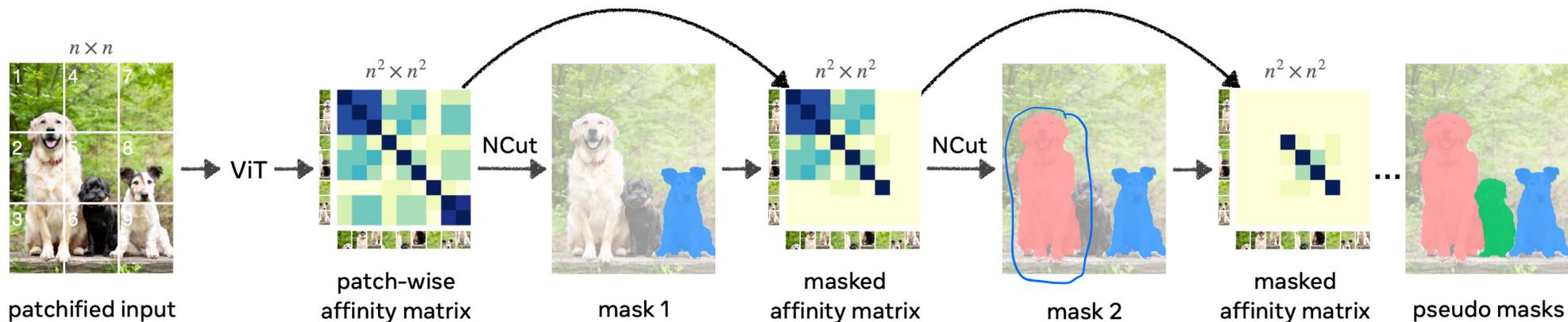
$$W_{ij} = \frac{K_i K_j}{\|K_i\|_2 \|K_j\|_2}$$



MaskCut

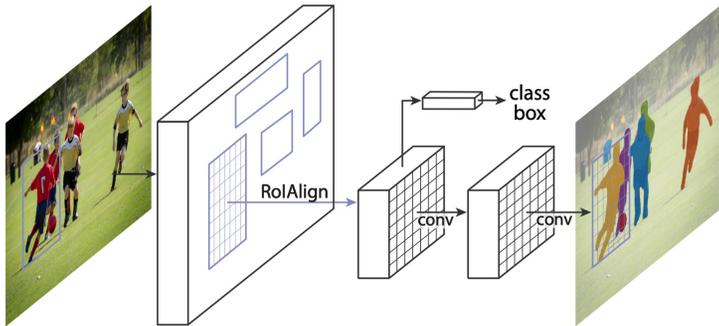
graph cut.

- Use a pretrained DINO ViT network.
- Use the “key” features from the last attention layer: $W_{ij} = \frac{K_i K_j}{\|K_i\|_2 \|K_j\|_2}$
- Iterative NCut on the pairwise matrix by masking out the regions from previous stages.



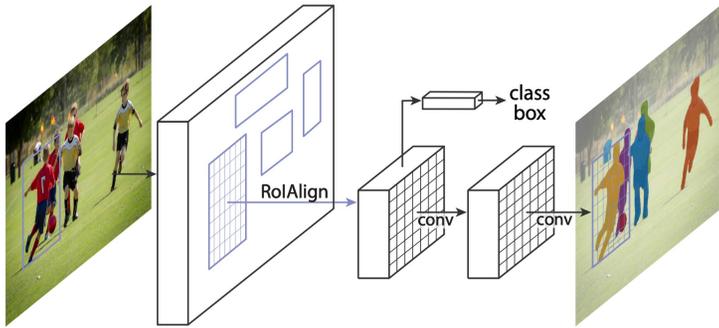
Iterative Self-Training

- Now add a MaskRCNN structure on top of the pretrained network.



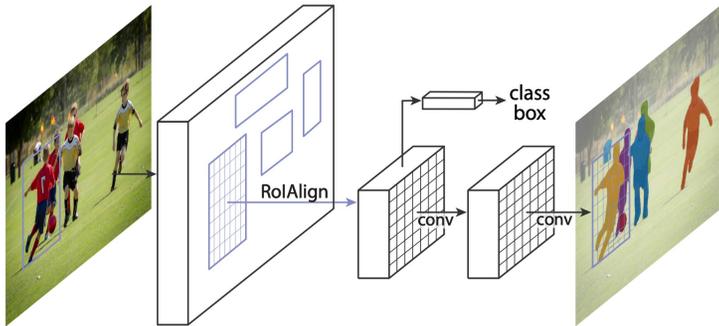
Iterative Self-Training

- Now add a MaskRCNN structure on top of the pretrained network.
- Select the predictions with the highest confidence score and use them as labels.

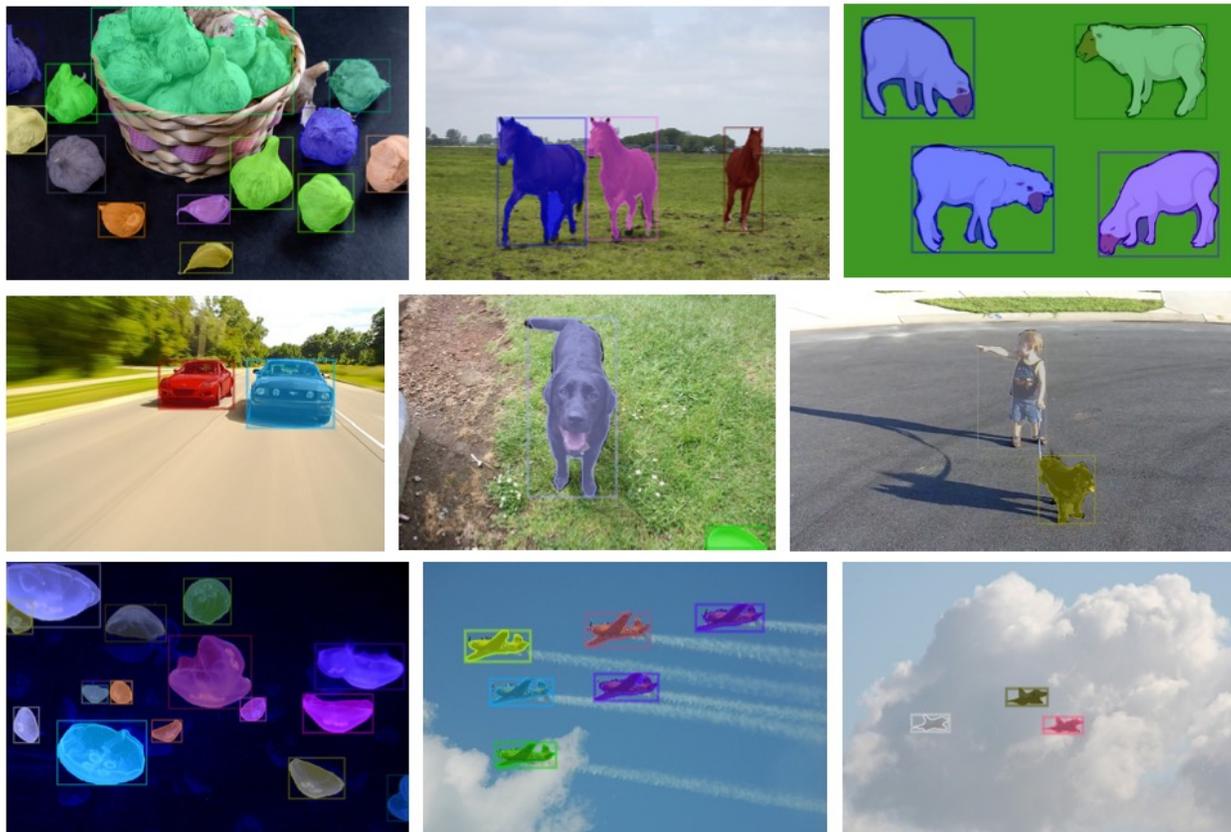


Iterative Self-Training

- Now add a MaskRCNN structure on top of the pretrained network.
- Select the predictions with the highest confidence score and use them as labels.
- Neural networks can learn from the noisy labels and output smoother predictions.

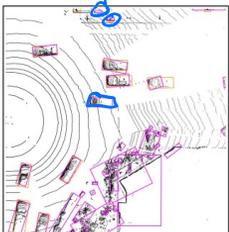


More Visualization

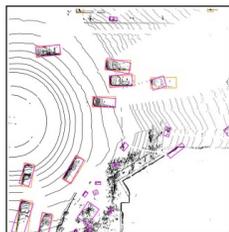


Pseudo Labels in 3D

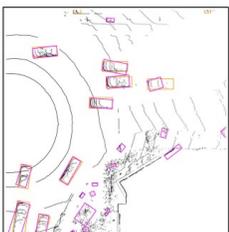
1 Point clustering pseudo-labels



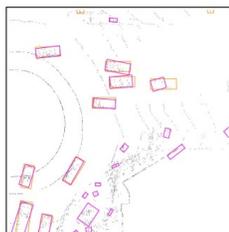
2 Filter out temporally inconsistent tracks



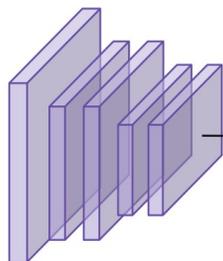
3 Randomly drop lidar beams



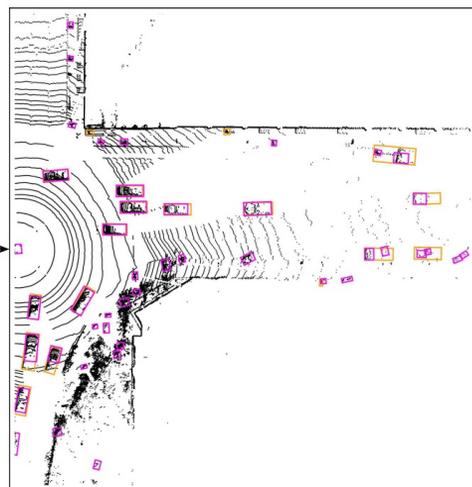
4 Randomly drop spherical rows/cols



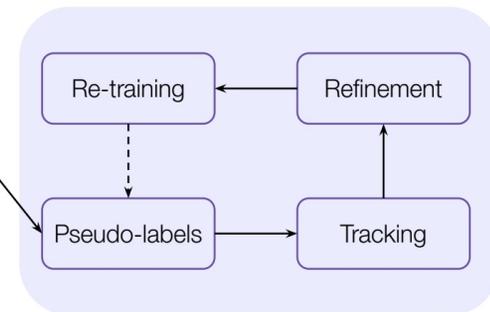
5 Train CNN in short range



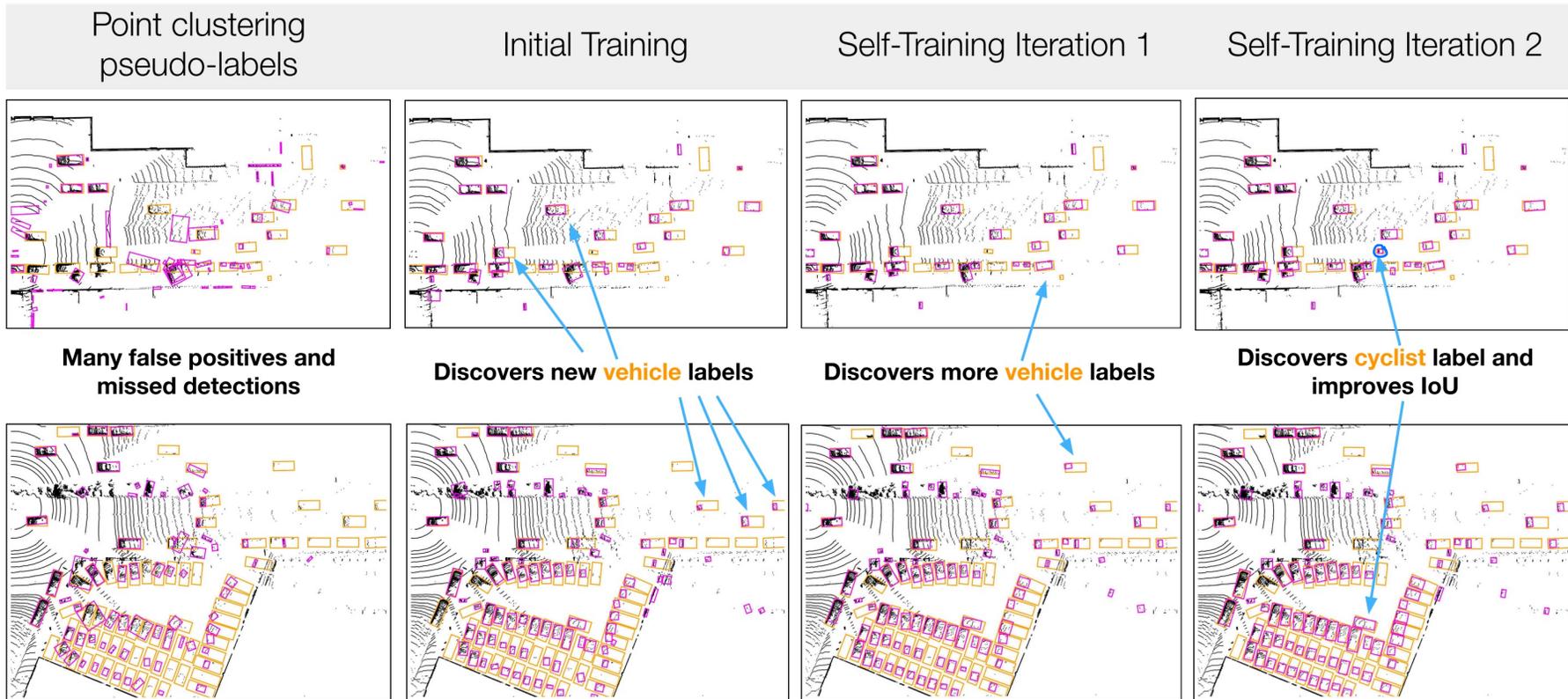
6 Zero-shot generalization to long-range



7 Self-training in long-range

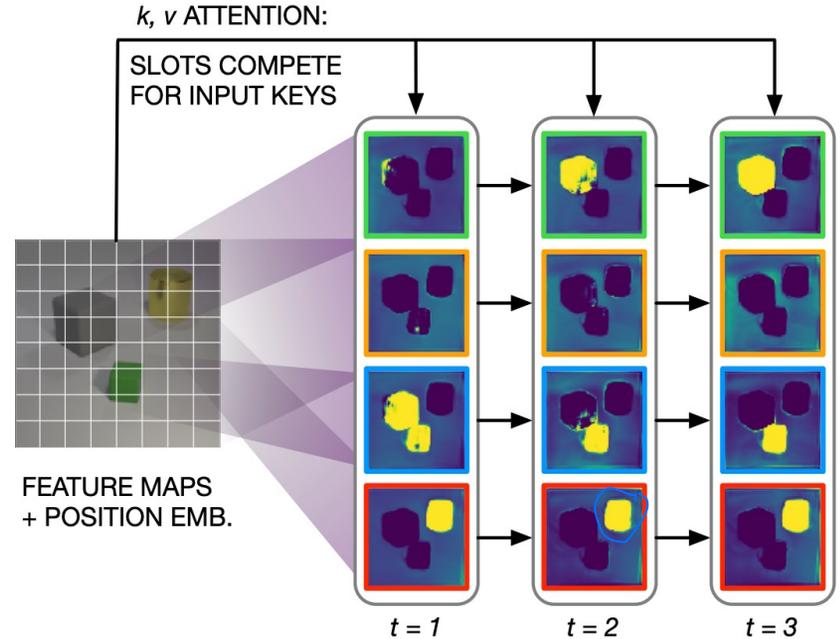


Iterative Refinement of Pseudo Labels



Slot Attention Networks

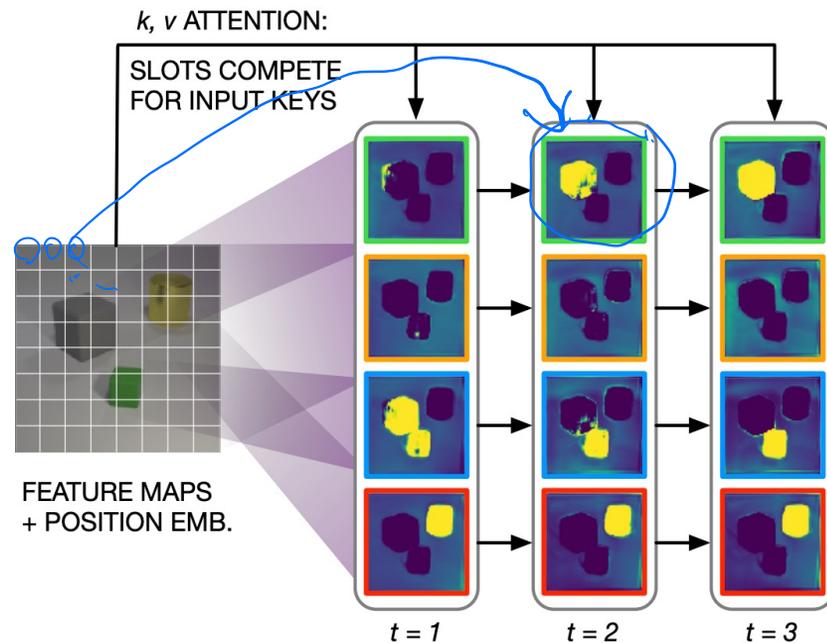
- Can we learn clustering as an end-to-end operation?



(a) Slot Attention module.

Slot Attention Networks

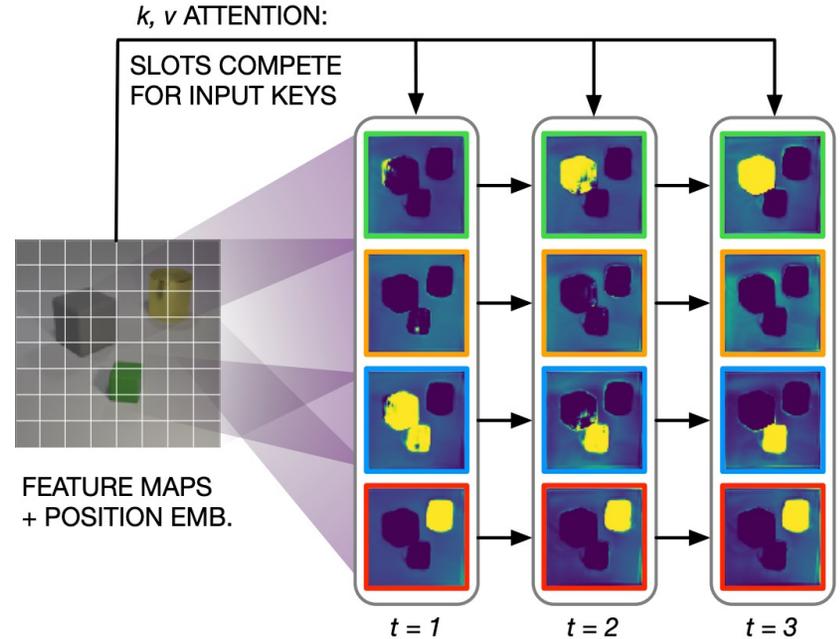
- Can we learn clustering as an end-to-end operation?
- Slot attention is inspired by the success of the attention mechanism.



(a) Slot Attention module.

Slot Attention Networks

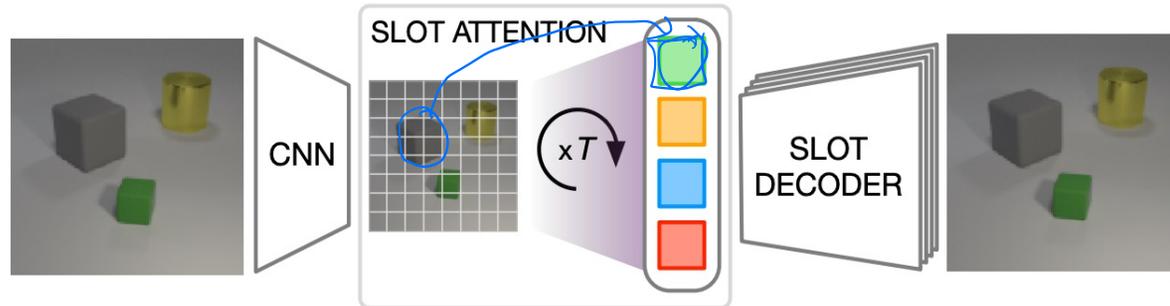
- Can we learn clustering as an end-to-end operation?
- Slot attention is inspired by the success of the attention mechanism.
- Each “slot” attends to a region of the image and stores an object centric representation.



(a) Slot Attention module.

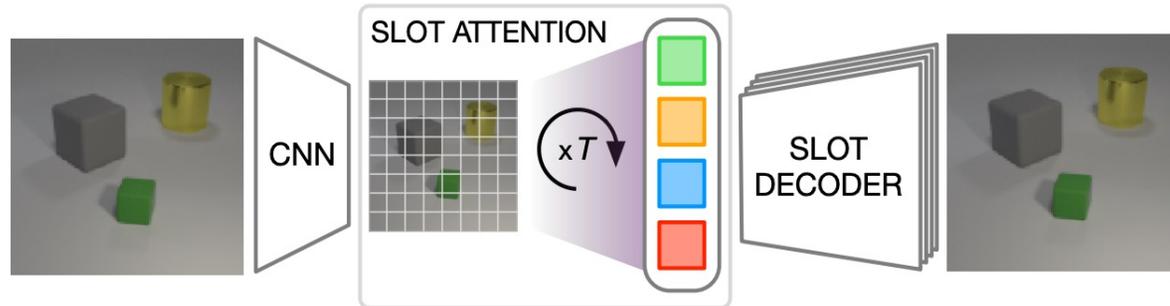
Slot Attention Networks

- Goal: Reconstruct the image with a concise slot-based representation.



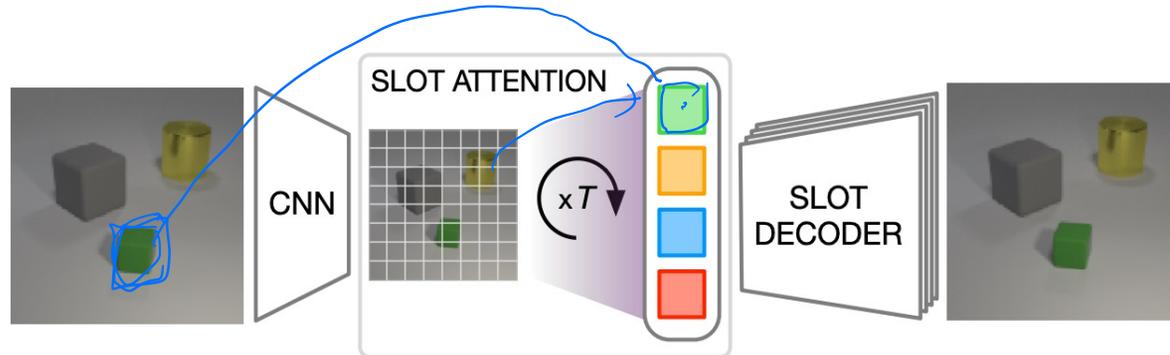
Slot Attention Networks

- Goal: Reconstruct the image with a concise slot-based representation.
- Input: $x \in \mathbb{R}^{N \times D}$ (after encoder), Slots: $m \in \mathbb{R}^{M \times D}$. Normalize:
 $\tilde{m}_{t-1} = LN(m_{t-1})$.
total objects.



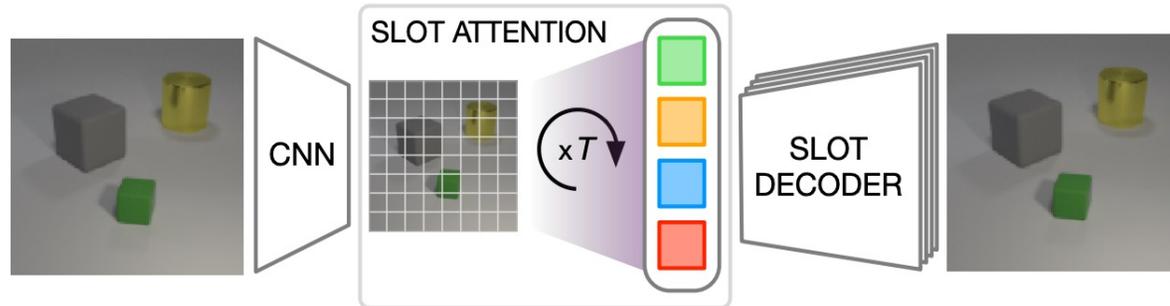
Slot Attention Networks

- Goal: Reconstruct the image with a concise slot-based representation.
- Input: $x \in \mathbb{R}^{N \times D}$ (after encoder), Slots: $m \in \mathbb{R}^{M \times D}$. Normalize:
 $\tilde{m}_{t-1} = LN(m_{t-1})$.
- Attention over slots: $a_{t,i,j} = \frac{\frac{1}{\sqrt{D}} k(x_i) \cdot q(\tilde{m}_j)^\top}{\sum_j \frac{1}{\sqrt{D}} k(x_i) \cdot q(\tilde{m}_j)^\top}$.



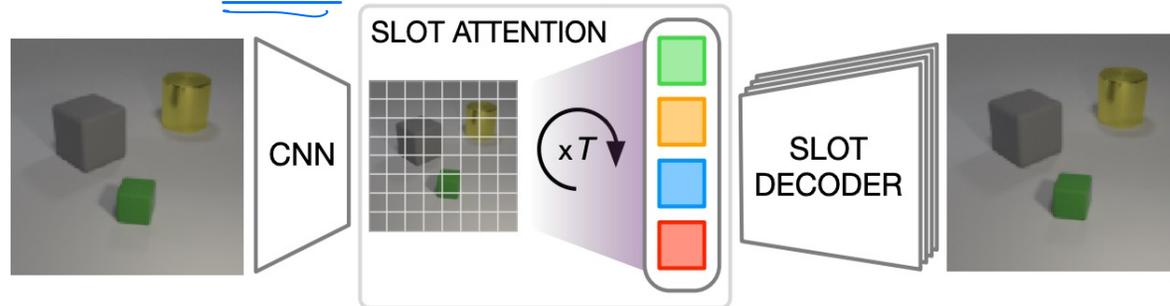
Slot Attention Networks

- Goal: Reconstruct the image with a concise slot-based representation.
- Input: $x \in \mathbb{R}^{N \times D}$ (after encoder), Slots: $m \in \mathbb{R}^{M \times D}$. Normalize: $\tilde{m}_{t-1} = LN(m_{t-1})$.
- Attention over slots: $a_{t,i,j} = \frac{\frac{1}{\sqrt{D}} k(x_i) \cdot q(\tilde{m}_j)^\top}{\sum_j \frac{1}{\sqrt{D}} k(x_i) \cdot q(\tilde{m}_j)^\top}$.
- Updates: $u_{tj} = \sum_i a_{tij} v(x_i)$.

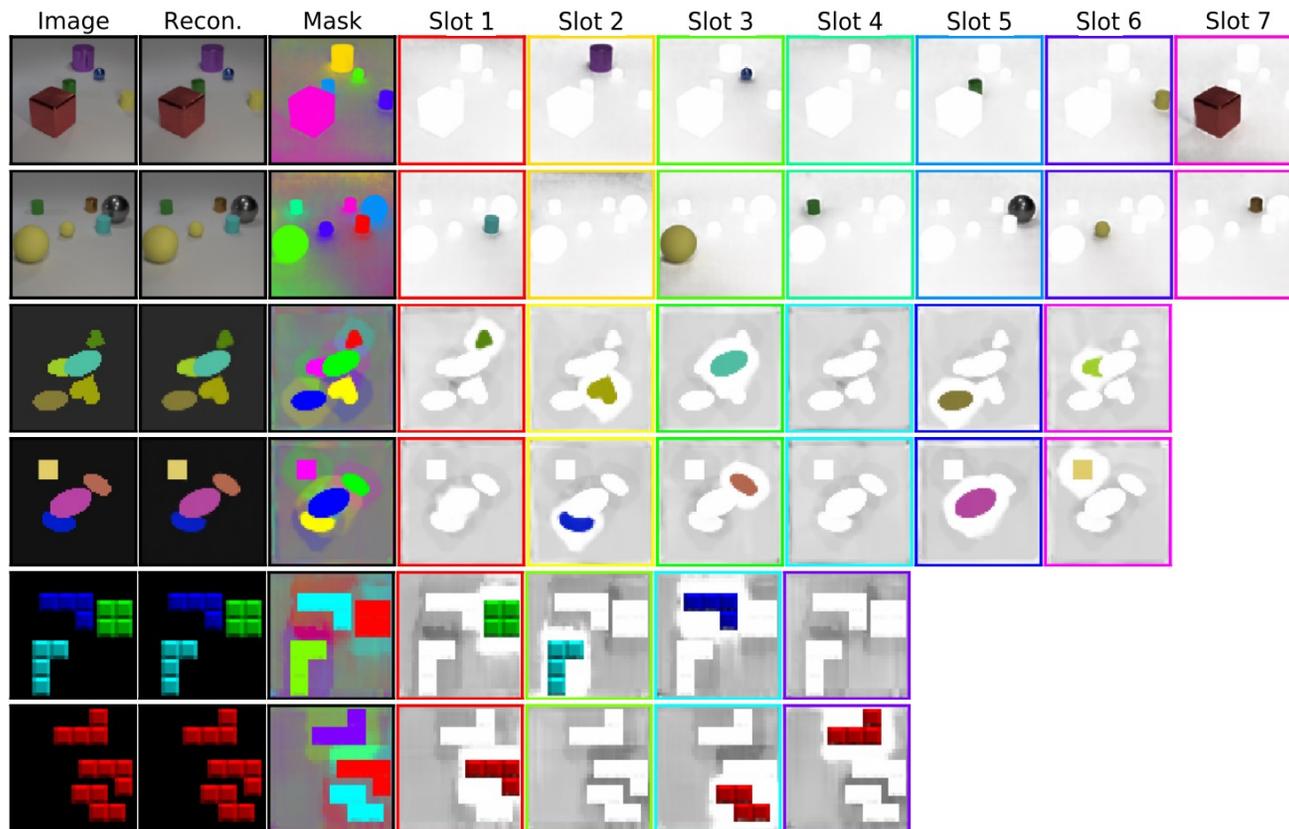


Slot Attention Networks

- Goal: Reconstruct the image with a concise slot-based representation.
- Input: $x \in \mathbb{R}^{N \times D}$ (after encoder), Slots: $m \in \mathbb{R}^{M \times D}$. Normalize: $\tilde{m}_{t-1} = LN(m_{t-1})$.
- Attention over slots: $a_{t,i,j} = \frac{\frac{1}{\sqrt{D}} k(x_i) \cdot q(\tilde{m}_j)^\top}{\sum_j \frac{1}{\sqrt{D}} k(x_i) \cdot q(\tilde{m}_j)^\top}$.
- Updates: $u_{tj} = \sum_i a_{tij} v(x_i)$.
- Write into slots: $m_t = \underline{GRU}(m_{t-1}, u_t) + MLP(\tilde{m}_{t-1})$.

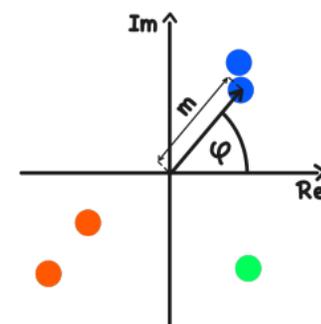
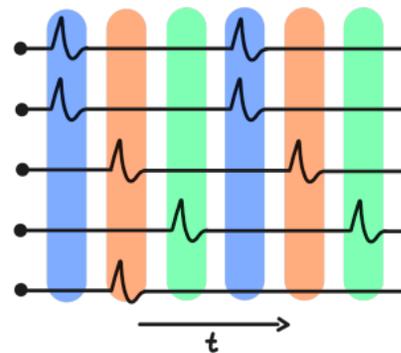
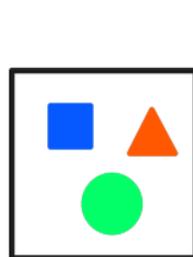
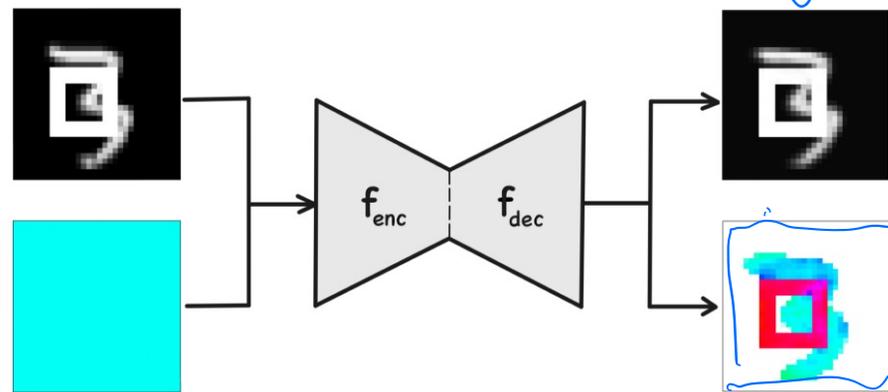


Slot Attention Networks



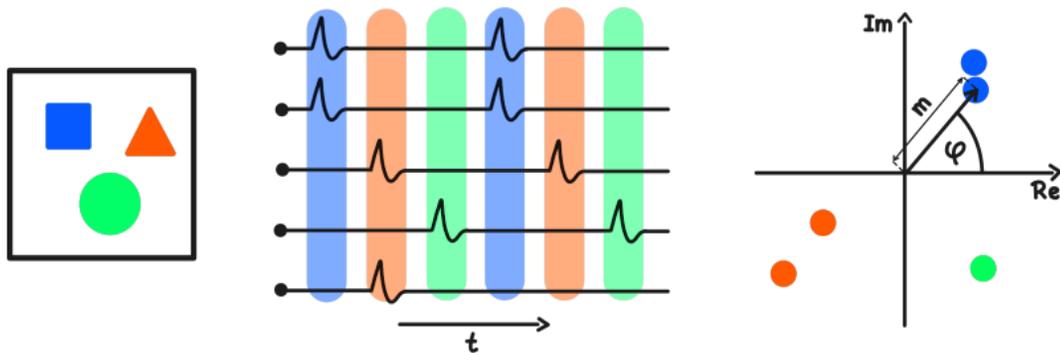
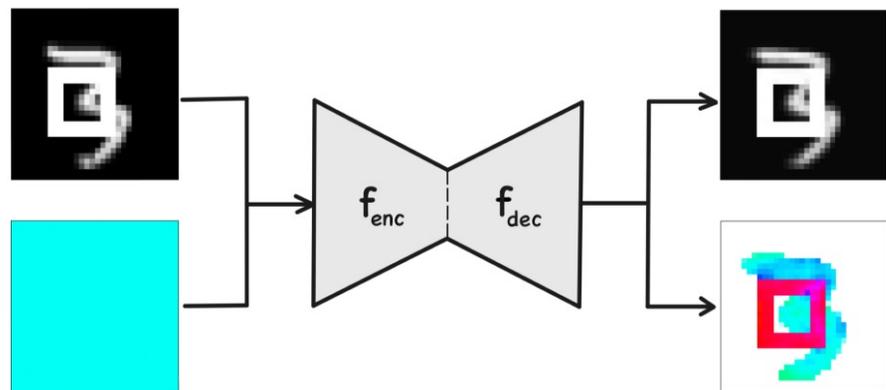
Complex-Valued Autoencoders (CAEs)

- The complex number can represent magnitude and phase: $z = m \cdot e^{i\varphi} \in \mathbb{C}$.



Complex-Valued Autoencoders (CAEs)

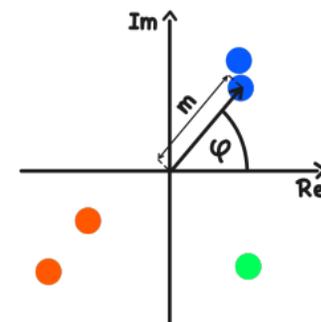
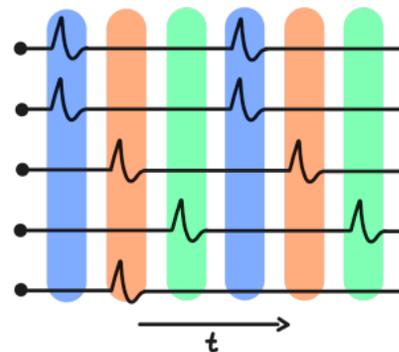
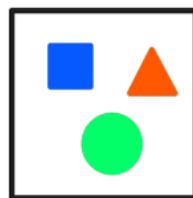
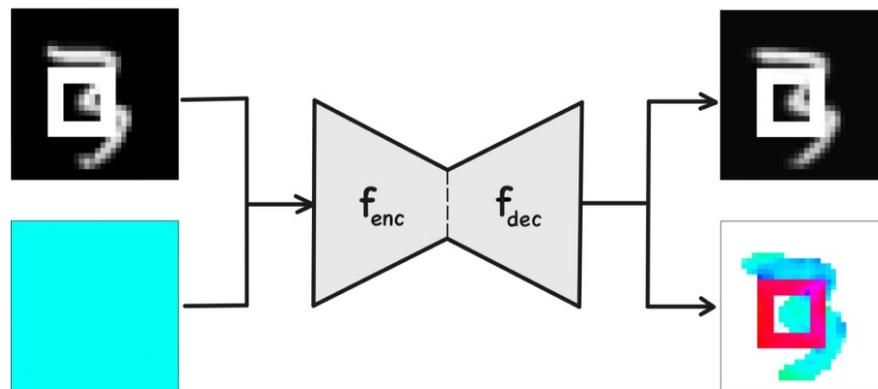
- The complex number can represent magnitude and phase: $z = m \cdot e^{i\varphi} \in \mathbb{C}$.
- Each pixel starts with an initial phase 0.



Complex-Valued Autoencoders (CAEs)

- The complex number can represent magnitude and phase: $z = m \cdot e^{i\varphi} \in \mathbb{C}$.
- Each pixel starts with an initial phase 0.

$$\hat{\mathbf{z}} = \underline{\underline{f_{dec}(f_{enc}(\mathbf{x}))}} \in \mathbb{C}^{h \times w}.$$

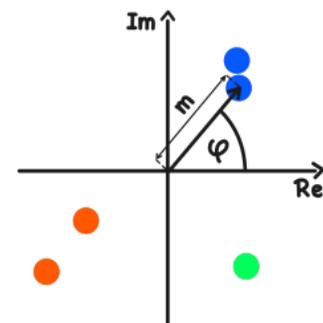
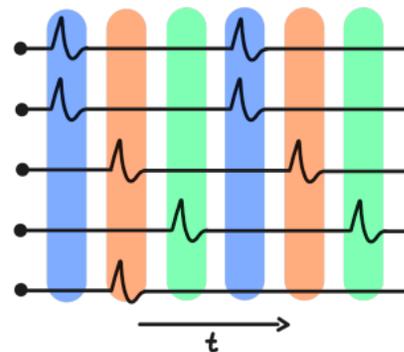
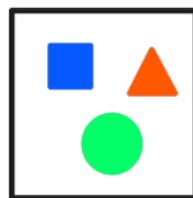
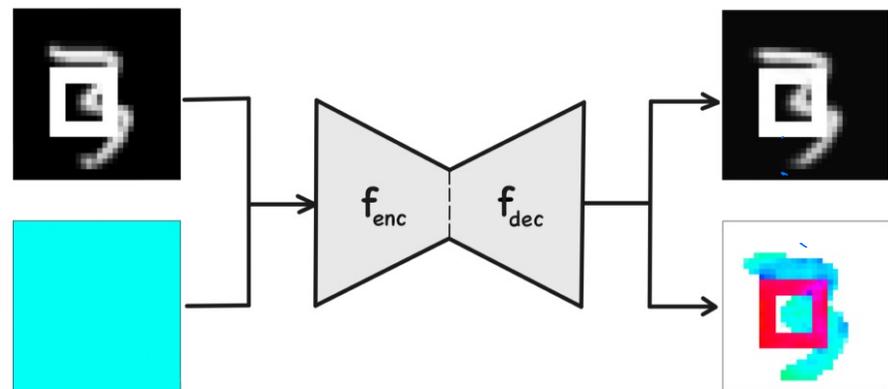


Complex-Valued Autoencoders (CAEs)

- The complex number can represent magnitude and phase: $z = m \cdot e^{i\varphi} \in \mathbb{C}$.
- Each pixel starts with an initial phase 0.

$$\hat{\mathbf{z}} = f_{\text{dec}}(f_{\text{enc}}(\mathbf{x})) \in \mathbb{C}^{h \times w}.$$

$$\hat{\mathbf{x}} = f_{\text{out}}(\hat{\mathbf{z}}) \in \mathbb{R}^{h \times w}.$$



CAE: More Details

- Apply weights separately to real and imaginary:

$$\psi = f_{\mathbf{w}}(\mathbf{z}) = \underbrace{f_{\mathbf{w}}(\operatorname{Re}(\mathbf{z}))}_{\text{shared network}} + \underbrace{f_{\mathbf{w}}(\operatorname{Im}(\mathbf{z}))}_{\text{shared network}} \cdot i \in \mathbb{C}_{d_{\text{out}}}$$

shared network.

CAE: More Details



- Apply weights separately to real and imaginary:

$$\psi = f_{\mathbf{w}}(\mathbf{z}) = f_{\mathbf{w}}(\text{Re}(\mathbf{z})) + f_{\mathbf{w}}(\text{Im}(\mathbf{z})) \cdot i \in \mathbb{C}_{d_{\text{out}}}$$

- Bias on magnitude and phase:

$$m_{\psi} = |\psi| + \mathbf{b}_m \in \mathbb{R}^{d_{\text{out}}} \quad \varphi_{\psi} = \arg(\psi) + \mathbf{b}_{\varphi} \in \mathbb{R}^{d_{\text{out}}}$$

CAE: More Details

- Apply weights separately to real and imaginary:

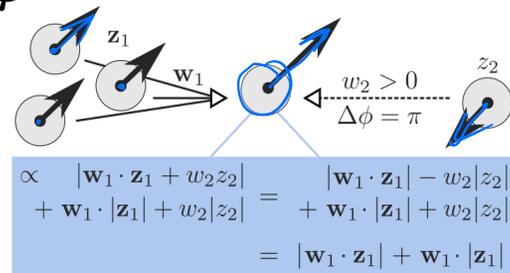
$$\psi = f_{\mathbf{w}}(\mathbf{z}) = f_{\mathbf{w}}(\text{Re}(\mathbf{z})) + f_{\mathbf{w}}(\text{Im}(\mathbf{z})) \cdot i \in \mathbb{C}_{d_{\text{out}}}$$

- Bias on magnitude and phase:

$$\mathbf{m}_{\psi} = |\psi| + \mathbf{b}_m \in \mathbb{R}^{d_{\text{out}}} \quad \varphi_{\psi} = \arg(\psi) + \mathbf{b}_{\varphi} \in \mathbb{R}^{d_{\text{out}}}$$

- Gating: $\chi = f_{\mathbf{w}}(|\mathbf{z}|) + \mathbf{b}_m \in \mathbb{R}^{d_{\text{out}}}$

$$\mathbf{m}_{\mathbf{z}} = 0.5 \cdot \mathbf{m}_{\psi} + 0.5 \cdot \chi \in \mathbb{R}^{d_{\text{out}}}$$



CAE: More Details

- Apply weights separately to real and imaginary:

$$\psi = f_{\mathbf{w}}(\mathbf{z}) = f_{\mathbf{w}}(\text{Re}(\mathbf{z})) + f_{\mathbf{w}}(\text{Im}(\mathbf{z})) \cdot i \in \mathbb{C}^{d_{\text{out}}}$$

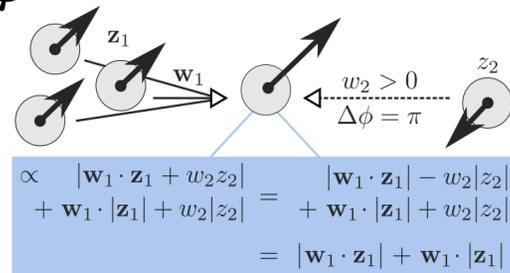
- Bias on magnitude and phase:

$$\mathbf{m}_{\psi} = |\psi| + \mathbf{b}_m \in \mathbb{R}^{d_{\text{out}}} \quad \varphi_{\psi} = \arg(\psi) + \mathbf{b}_{\varphi} \in \mathbb{R}^{d_{\text{out}}}$$

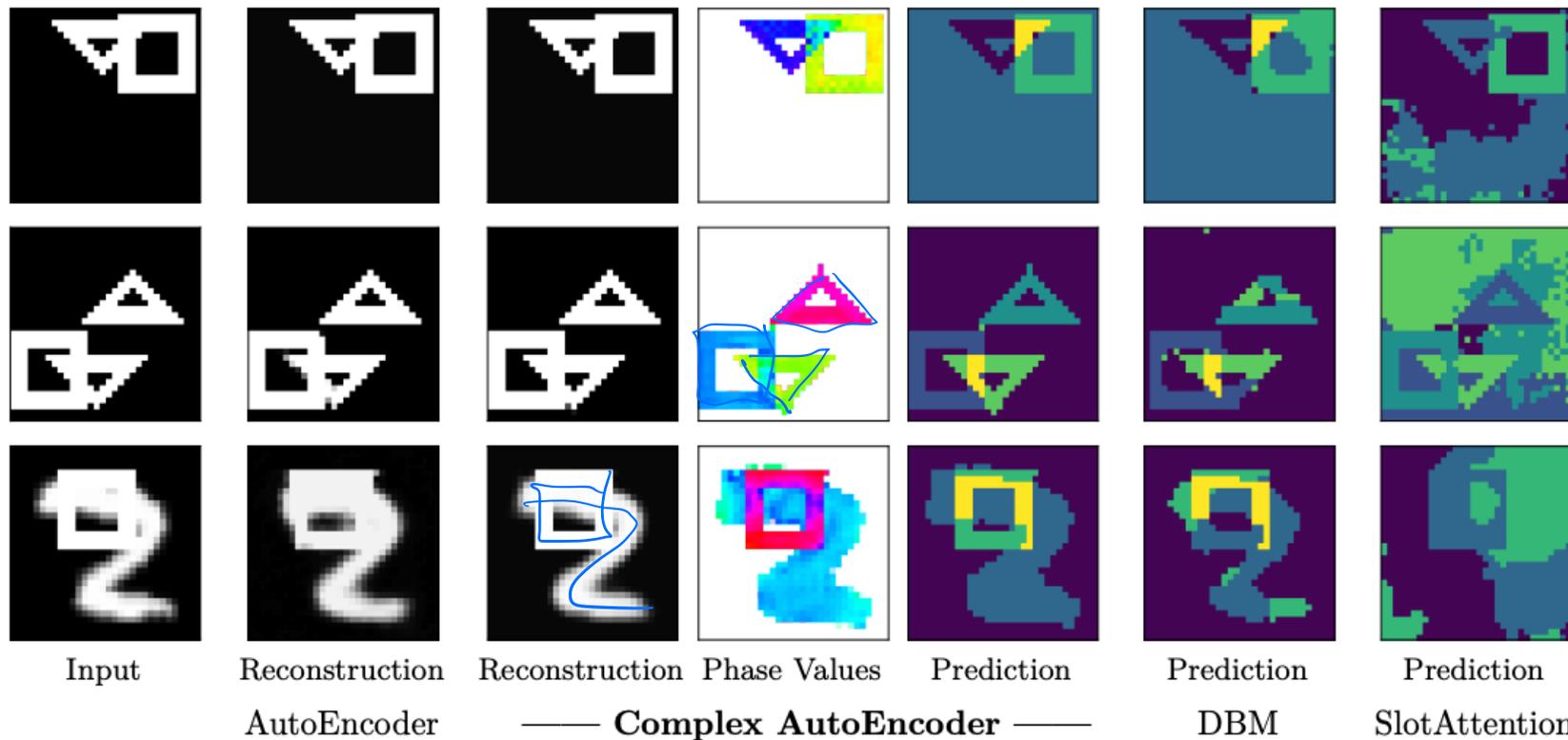
- Gating: $\chi = f_{\mathbf{w}}(|\mathbf{z}|) + \mathbf{b}_m \in \mathbb{R}^{d_{\text{out}}}$

$$\mathbf{m}_{\mathbf{z}} = 0.5 \cdot \mathbf{m}_{\psi} + 0.5 \cdot \chi \in \mathbb{R}^{d_{\text{out}}}$$

- Activation: $\mathbf{z}' = \text{ReLU}(\text{BatchNorm}(\mathbf{m}_{\mathbf{z}})) \cdot e^{i\varphi_{\psi}} \in \mathbb{C}^{d_{\text{out}}}$



Complex-Valued Autoencoders



Summary: Object Discovery

- Combine deep features with clustering algorithms.

Summary: Object Discovery

- Combine deep features with clustering algorithms.
- Pseudo-labels to train detector networks.

Summary: Object Discovery

- Combine deep features with clustering algorithms.
- Pseudo-labels to train detector networks.
- Creative end-to-end learning-based solutions exist, but there are still plenty room for improvement.
 - Possible to train from scratch!

Summary: Object Discovery

- Combine deep features with clustering algorithms.
- Pseudo-labels to train detector networks.
- Creative end-to-end learning-based solutions exist, but there are still plenty room for improvement.
 - Possible to train from scratch!
- What do we make use of the discovered objects? Is it better to keep the awareness in the latent space?

Module 4: World Models and End-to-End Planning

World Models: Predicting Future

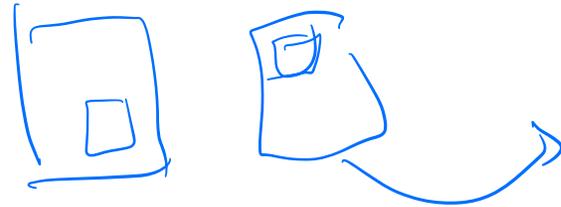
3D, latent, \rightarrow predict future \leftarrow planning.

World Models: Predicting Future

- World models: Models that predict the future states based on the past sequences of states and actions.

World Models: Predicting Future

- World models: Models that predict the future states based on the past sequences of states and actions.
- Predicting future is a natural self-supervised learning objective.



World Models: Predicting Future

- World models: Models that predict the future states based on the past sequences of states and actions.
- Predicting future is a natural self-supervised learning objective.
- Why?

World Models: Predicting Future

- World models: Models that predict the future states based on the past sequences of states and actions.
- Predicting future is a natural self-supervised learning objective.
- Why?
 - Reward is sparse, rollout trajectories are abundant.

World Models: Predicting Future

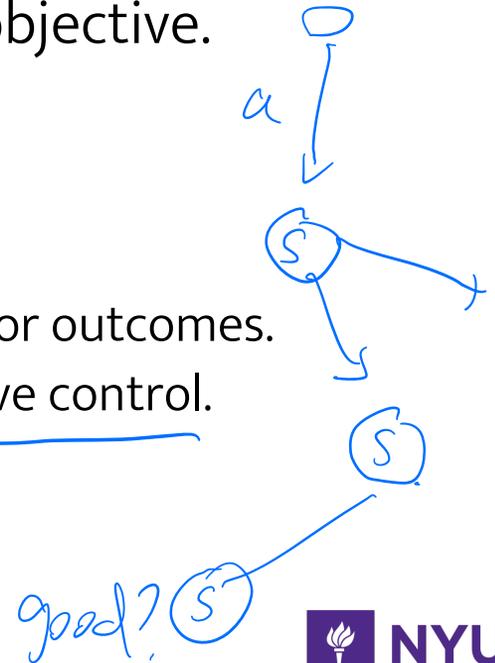
- World models: Models that predict the future states based on the past sequences of states and actions.
- Predicting future is a natural self-supervised learning objective.
- Why?
 - Reward is sparse, rollout trajectories are abundant.
 - Learning task /planning aware representations.

World Models: Predicting Future

- World models: Models that predict the future states based on the past sequences of states and actions.
- Predicting future is a natural self-supervised learning objective.
- Why?
 - Reward is sparse, rollout trajectories are abundant.
 - Learning task /planning aware representations.
 - Predicting future can “simulate” rollouts without querying for outcomes.

World Models: Predicting Future

- World models: Models that predict the future states based on the past sequences of states and actions.
- Predicting future is a natural self-supervised learning objective.
- Why?
 - Reward is sparse, rollout trajectories are abundant.
 - Learning task /planning aware representations.
 - Predicting future can “simulate” rollouts without querying for outcomes.
 - Can help planning by predicting at test time. Model predictive control.



Model-Predictive Control (MPC)

- A classic example of world models.
- Analytical forms, complete knowledge of the dynamical system.

Model-Predictive Control (MPC)

- A classic example of world models.
- Analytical forms, complete knowledge of the dynamical system.

General Form:

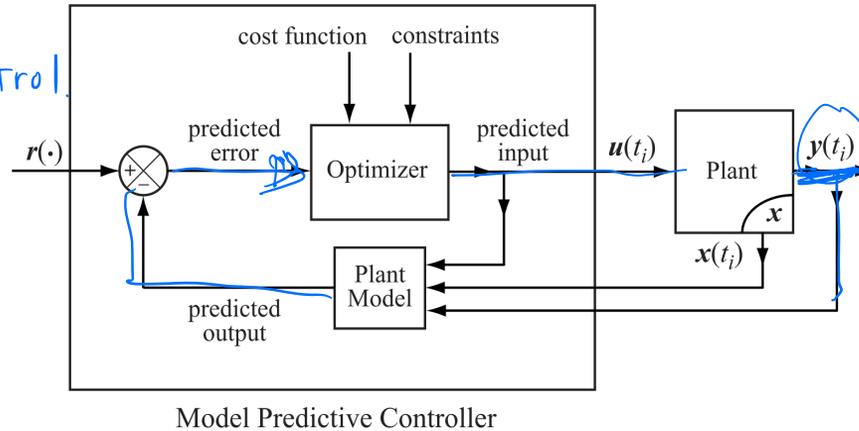
dynamical system. $\frac{dx}{dt} = f(\mathbf{x}, \mathbf{u}, t)$

state, control.

$y = h(\mathbf{x}, \mathbf{u}, t)$

$\mathbf{x}(t_0) = \mathbf{x}_0$

Observed \rightarrow *reference traj.*
 $r(x_T)$



Model-Predictive Control (MPC)

- A classic example of world models.
- Analytical forms, complete knowledge of the dynamical system.

General Form:

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \mathbf{u}, t)$$

$$\mathbf{y} = h(\mathbf{x}, \mathbf{u}, t)$$

$$\mathbf{x}(t_0) = \mathbf{x}_0$$

Linear Form:

$$\frac{dx}{dt} = A\mathbf{x} + B\mathbf{u}$$

$$\mathbf{y} = C\mathbf{x} + D\mathbf{u}$$

$$\mathbf{x}(t_0) = \mathbf{x}_0$$

Model-Predictive Control (MPC)

- A classic example of world models.
- Analytical forms, complete knowledge of the dynamical system.

General Form:

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \mathbf{u}, t)$$

$$\mathbf{y} = h(\mathbf{x}, \mathbf{u}, t)$$

$$\mathbf{x}(t_0) = \mathbf{x}_0$$

Linear Form:

$$\frac{dx}{dt} = A\mathbf{x} + B\mathbf{u}$$

$$\mathbf{y} = C\mathbf{x} + D\mathbf{u}$$

$$\mathbf{x}(t_0) = \mathbf{x}_0$$

Optimize Value/Cost Function:

$$\min_{\mathbf{u}} J(\mathbf{x}(0), \mathbf{u})$$

(Note: Blue arrows in the original image point to the initial state $\mathbf{x}(0)$ and the control input \mathbf{u} in the cost function.)

Model-Predictive Control (MPC)

- A classic example of world models.
- Analytical forms, complete knowledge of the dynamical system.

General Form:

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \mathbf{u}, t)$$

$$\mathbf{y} = h(\mathbf{x}, \mathbf{u}, t)$$

$$\mathbf{x}(t_0) = \mathbf{x}_0$$

Linear Form:

$$\frac{d\mathbf{x}}{dt} = A\mathbf{x} + B\mathbf{u}$$

$$\mathbf{y} = C\mathbf{x} + D\mathbf{u}$$

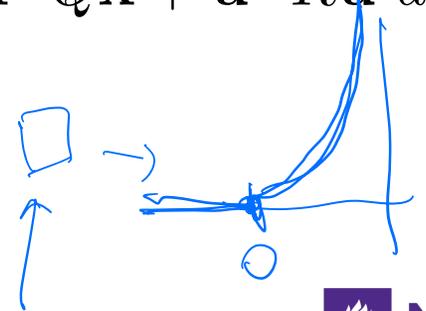
$$\mathbf{x}(t_0) = \mathbf{x}_0$$

Optimize Value/Cost Function:

$$\min_{\mathbf{u}} J(\mathbf{x}(0), \mathbf{u})$$

Quadratic Form (LQR):

$$J = \int_0^{\infty} \mathbf{x}^T Q \mathbf{x} + \mathbf{u}^T R \mathbf{u} dt.$$



Model-Predictive Control (MPC)

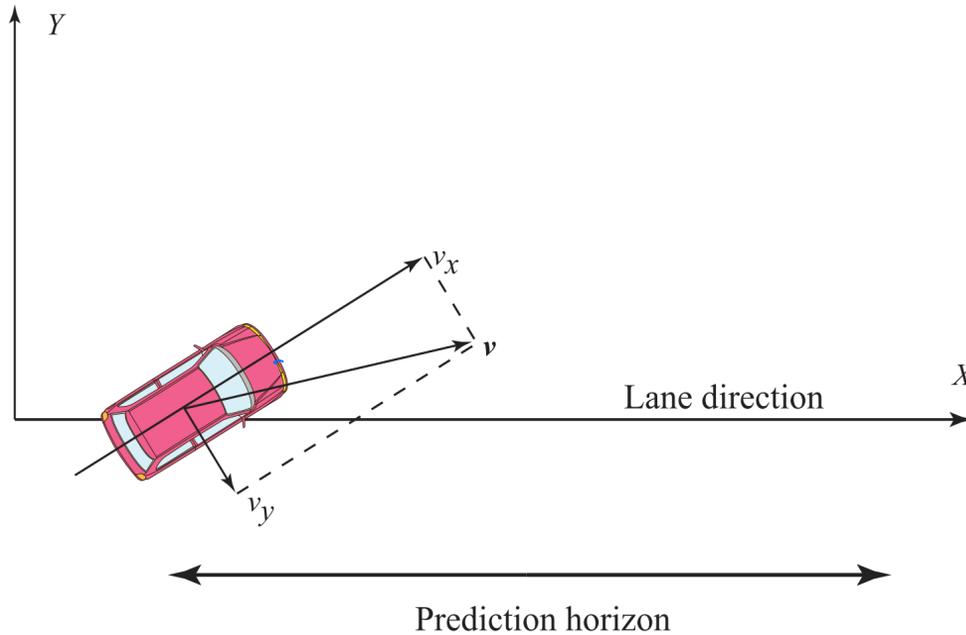
- A classic example
- Analytical for

General Form:

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}, \mathbf{u}, t)$$

$$\mathbf{y} = h(\mathbf{x}, t)$$

$$\mathbf{x}(t_0) = \mathbf{x}_0$$



al system.

value/Cost Function:

$$J(\mathbf{x}(0), \mathbf{u})$$

Form (LQR):

$$\int_0^{\infty} \mathbf{x}^T Q \mathbf{x} + \mathbf{u}^T R \mathbf{u} dt.$$

- Example: Cars: \mathbf{x} : position velocity angle angular velocity; \mathbf{u} : jerk and angular accel.

Several Types of World Models

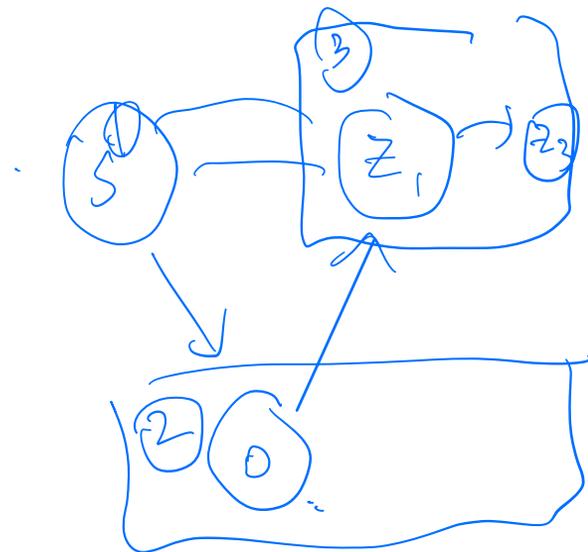
- State
 - Explicitly defined state space
 - Needs a separately trained perception network

Several Types of World Models

- State
 - Explicitly defined state space
 - Needs a separately trained perception network
- Raw Input
 - Generate future camera video frames / point clouds

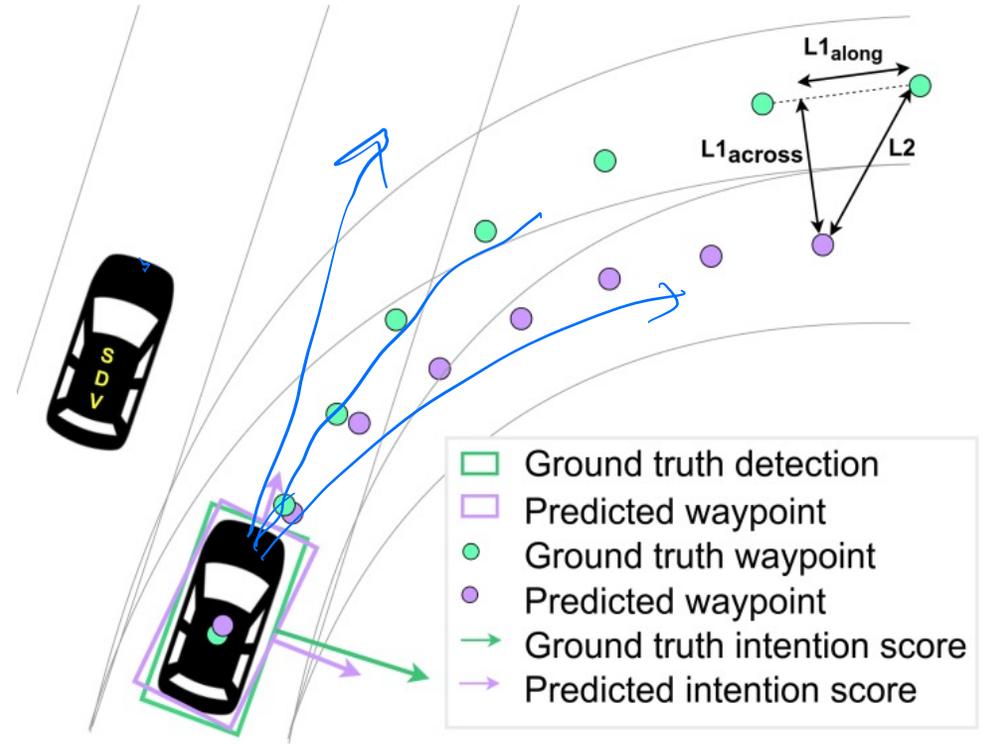
Several Types of World Models

- State
 - Explicitly defined state space
 - Needs a separately trained perception network
- Raw Input
 - Generate future camera video frames / point clouds
- Latent
 - Generate future latent states
 - May go outside of the latent manifold



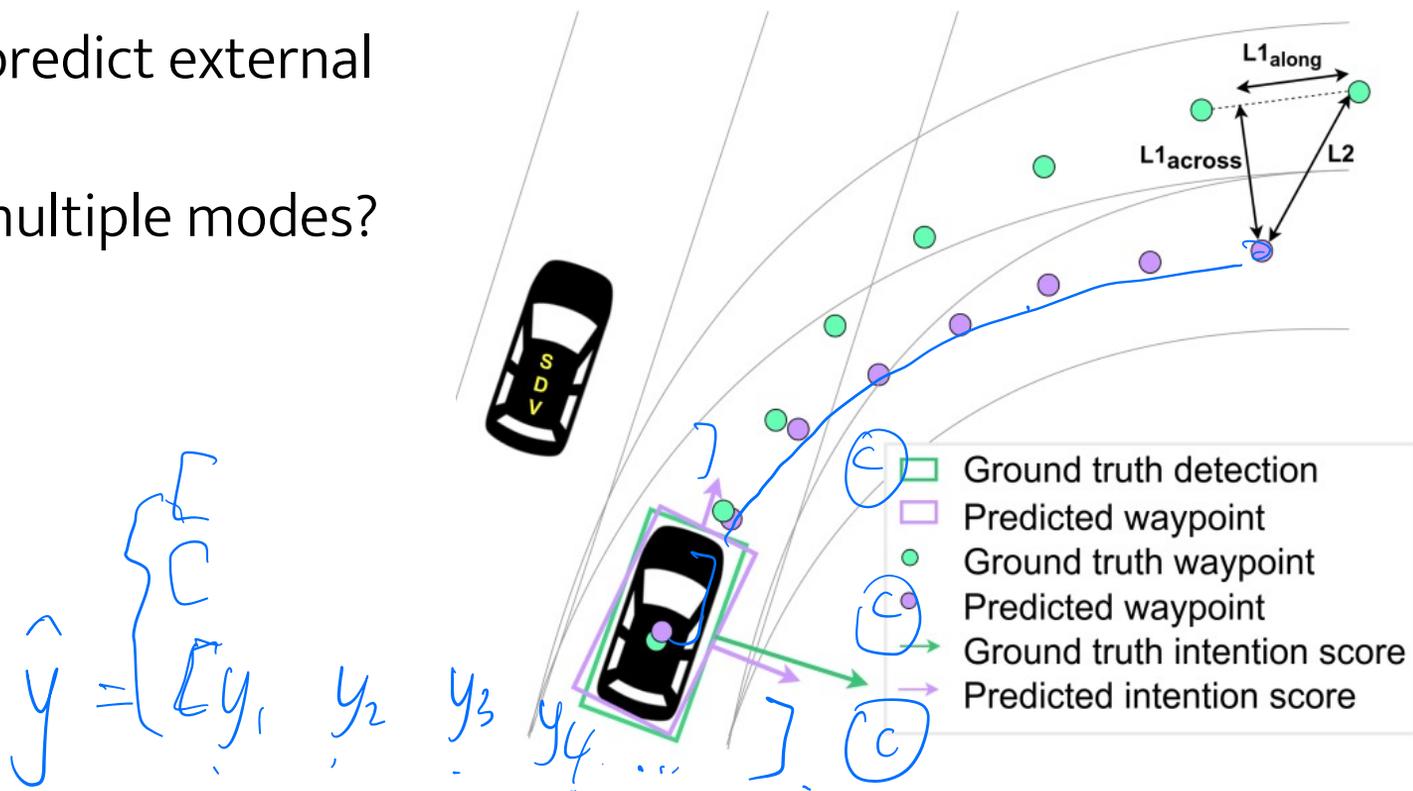
Trajectory Prediction as Object Detection

- We also need to predict external dynamic objects.



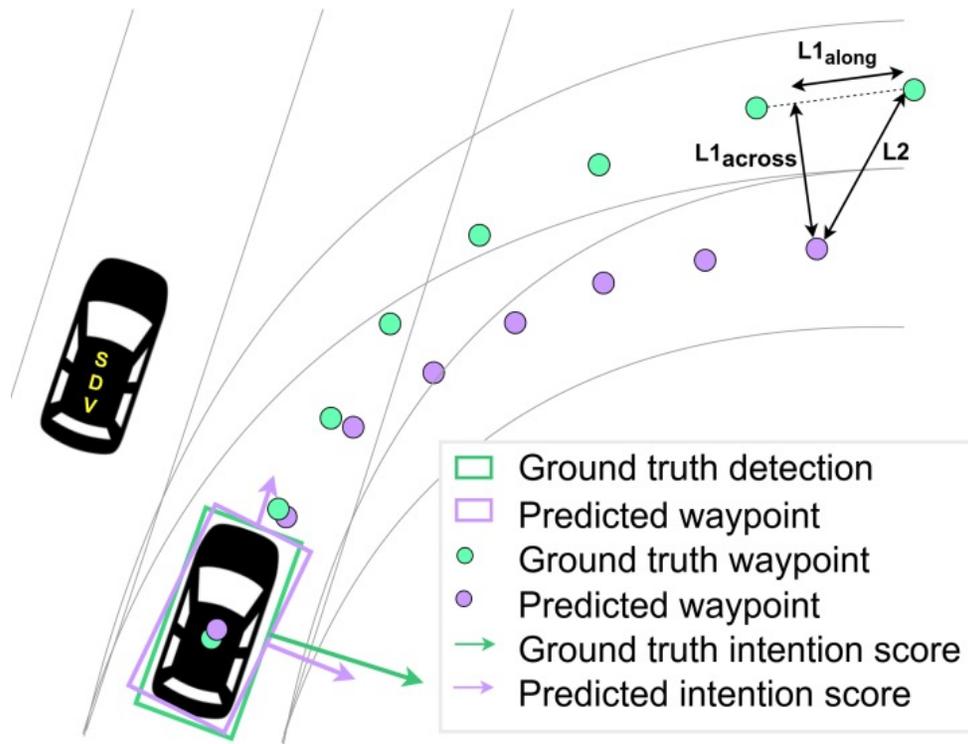
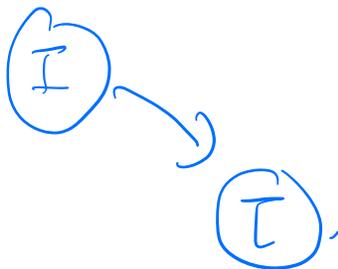
Trajectory Prediction as Object Detection

- We also need to predict external dynamic objects.
- How to capture multiple modes?



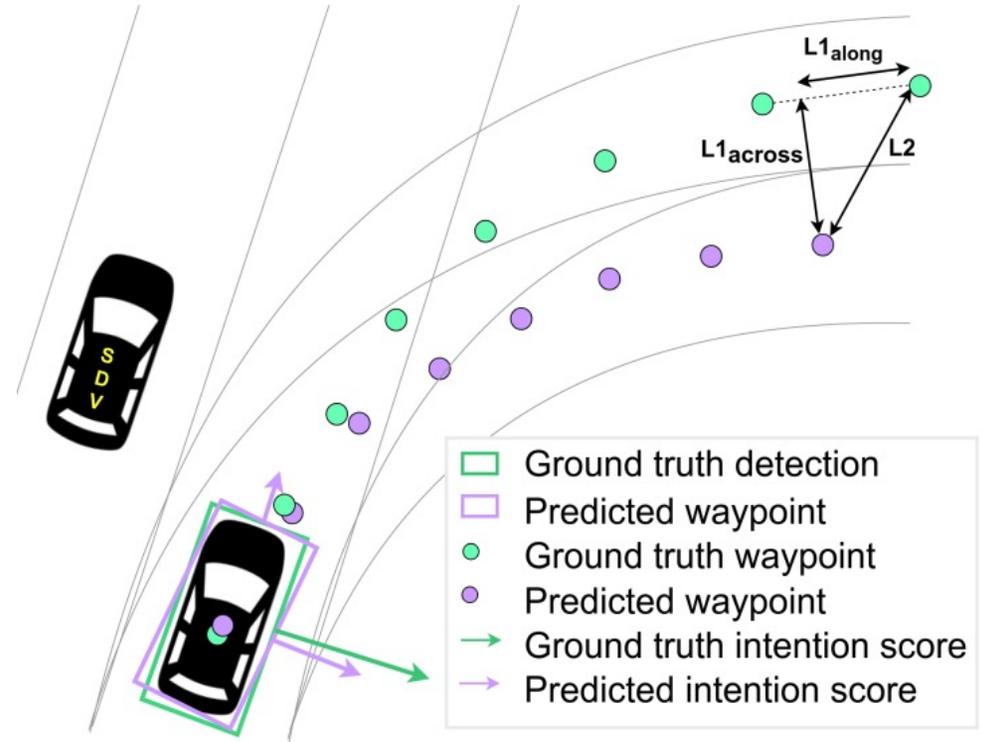
Trajectory Prediction as Object Detection

- We also need to predict external dynamic objects.
- How to capture multiple modes?
- Discrete intention prediction problem:
 - keep lane, turn left/right, left/right change lane, stopping, parked, etc.



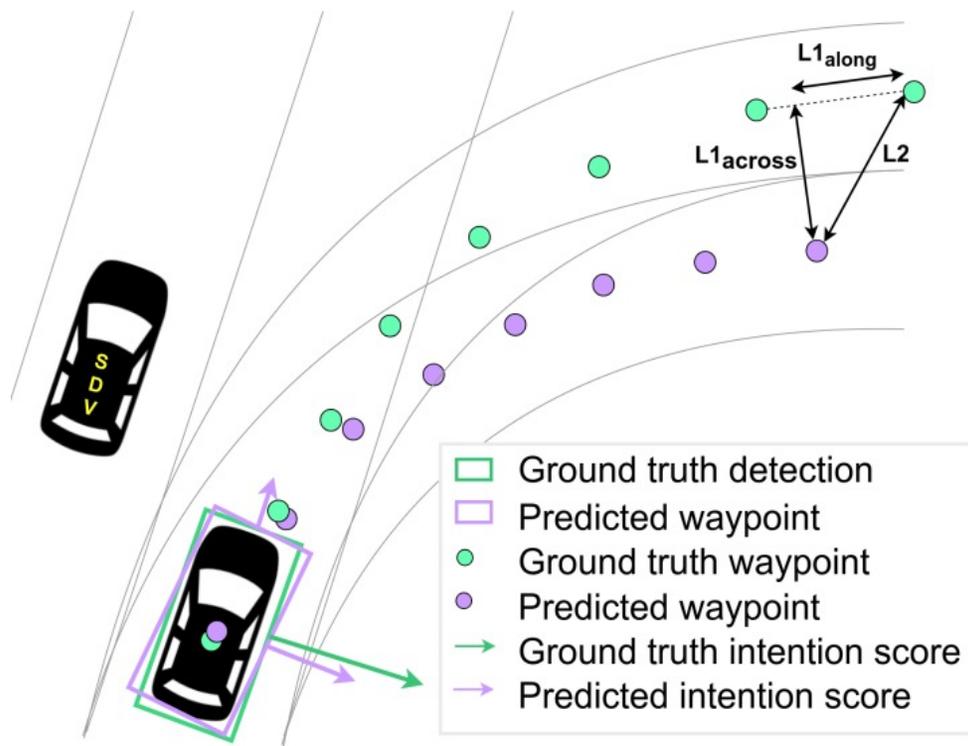
Trajectory Prediction as Object Detection

- We also need to predict external dynamic objects.
- How to capture multiple modes?
- Discrete intention prediction problem:
 - keep lane, turn left/right, left/right change lane, stopping, parked, etc.
- Output multiple trajectories



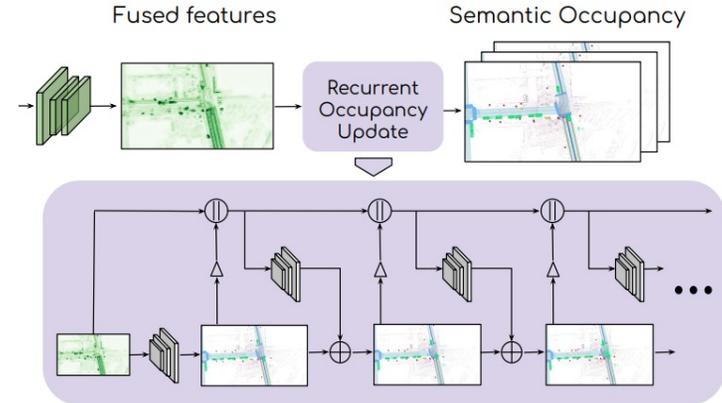
Trajectory Prediction as Object Detection

- We also need to predict external dynamic objects.
- How to capture multiple modes?
- Discrete intention prediction problem:
 - keep lane, turn left/right, left/right change lane, stopping, parked, etc.
- Output multiple trajectories
- Requires high-level action labels.



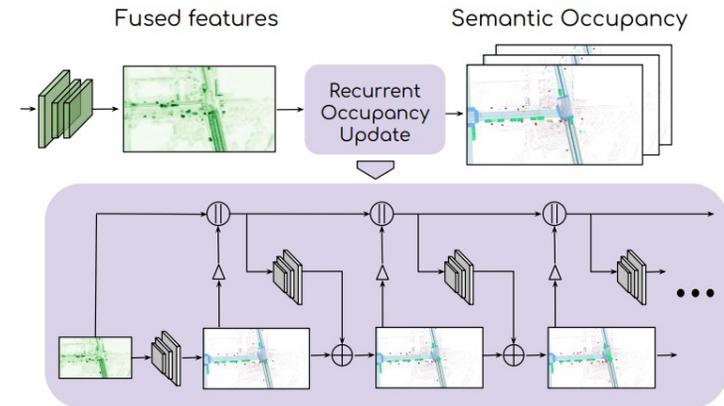
Semantic Occupancy Volume Prediction

- 4-D volume: H, W, T, C (class) $o_{i,j}^{t,c}$
num. time steps.



Semantic Occupancy Volume Prediction

- 4-D volume: H, W, T, C (class) $o_{i,j}^{t,c}$
- Doesn't grow with the increasing number of actors.

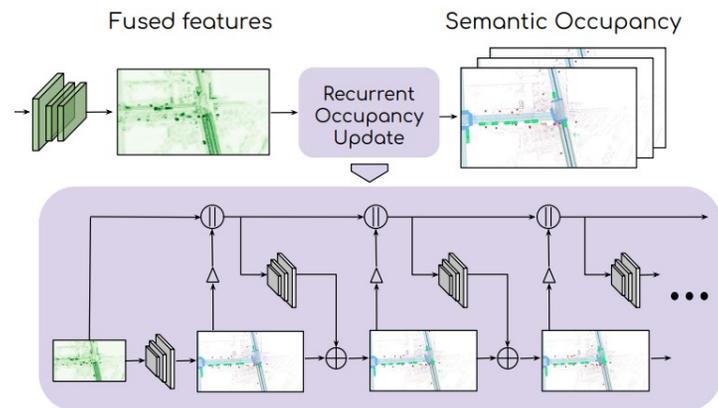


Semantic Occupancy Volume Prediction

- 4-D volume: H, W, T, C (class) $o_{i,j}^{t,c}$
- Doesn't grow with the increasing number of actors.
- Recurrent occupancy updates for further into the future.

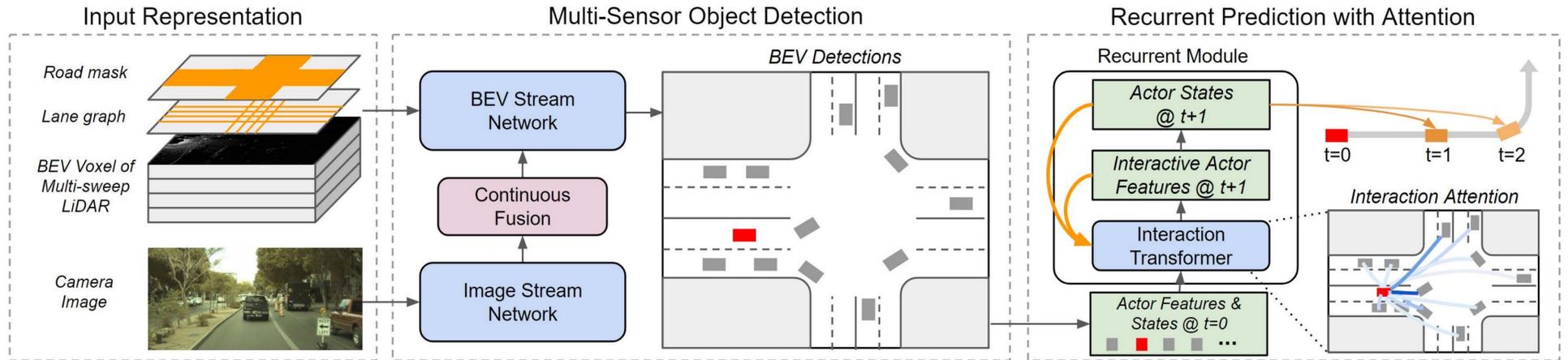
$$l^{t,c} = l^{t-1,c} + \mathcal{U}_{\theta}^t(f_{\text{occ}}, \underline{l^{0:t-1,c}})$$

↳ render
cost volume.



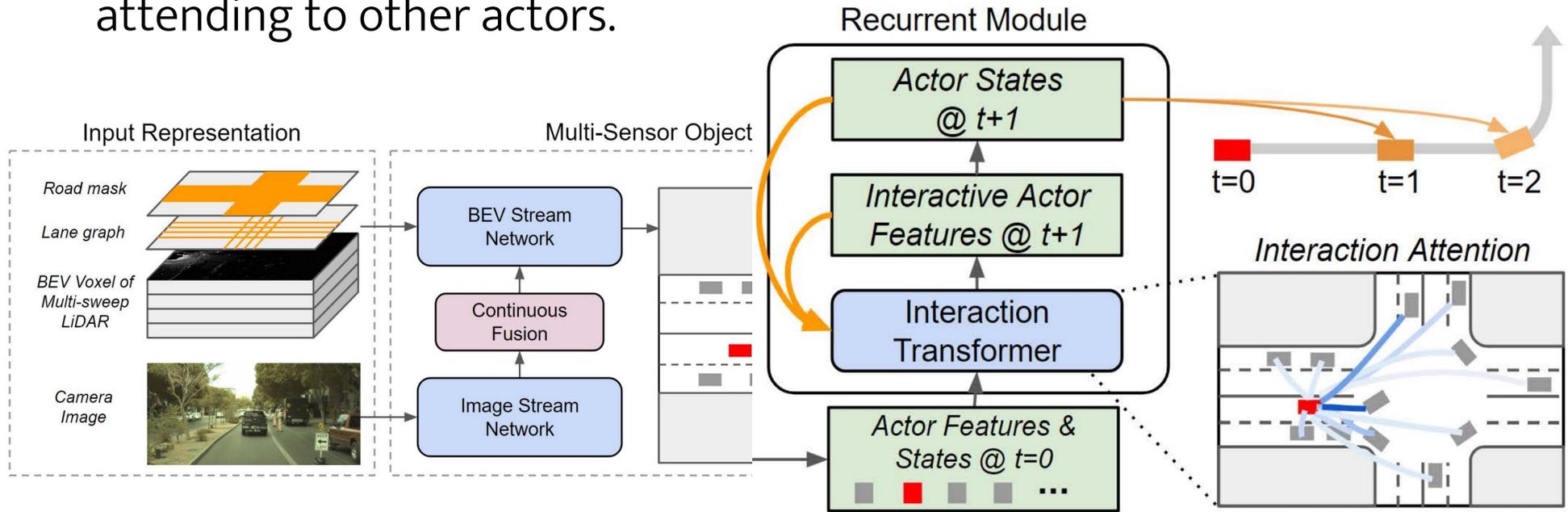
Multi-Agents Joint Prediction

- Joint predict future trajectories by attending to other actors.



Multi-Agents Joint Prediction

- Joint predict future trajectories by attending to other actors.



Generating Raw Pixels

<https://deepmind.google/blog/genie-3-a-new-frontier-for-world-models/>

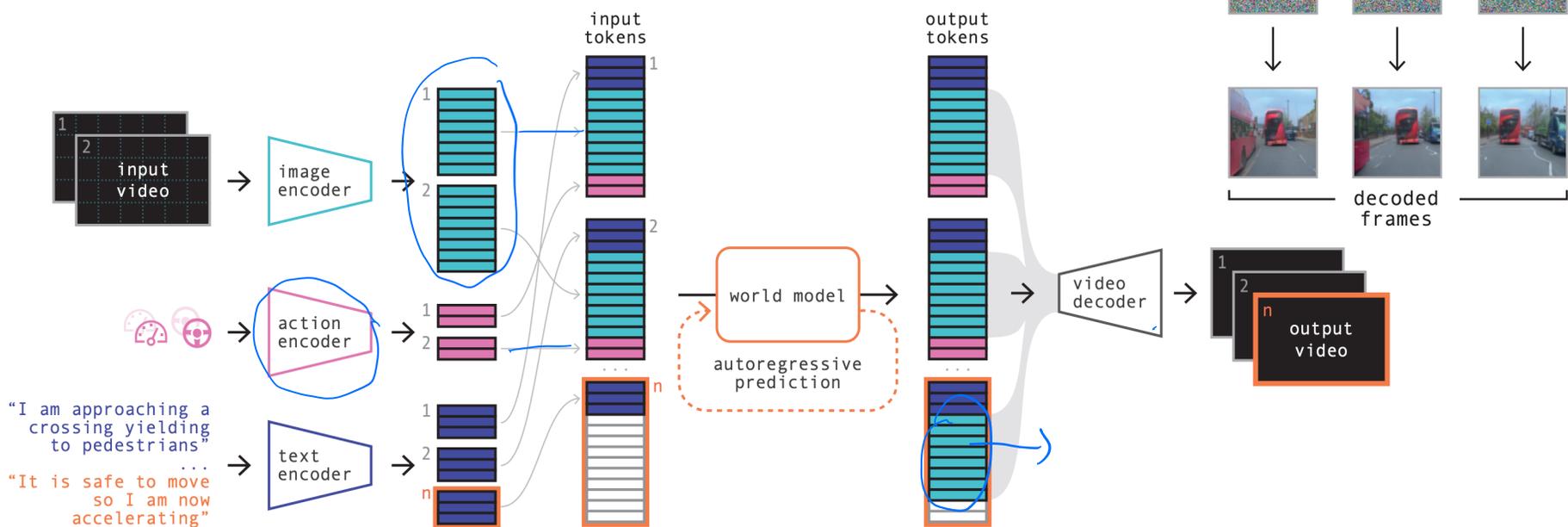
Generating Raw Pixels



<https://deepmind.google/blog/genie-3-a-new-frontier-for-world-models/>

Action-Conditioned Video Generation

- Autoregressive backbone with a diffusion decoder.



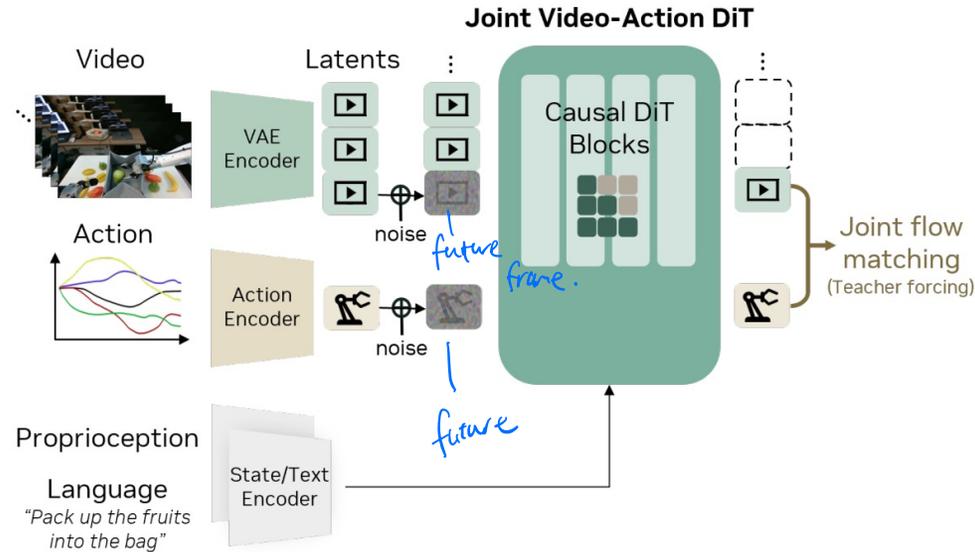
<https://www.youtube.com/watch?v=5Jx2QgEUZUI>



<https://www.youtube.com/watch?v=5Jx2QgEUZUI>

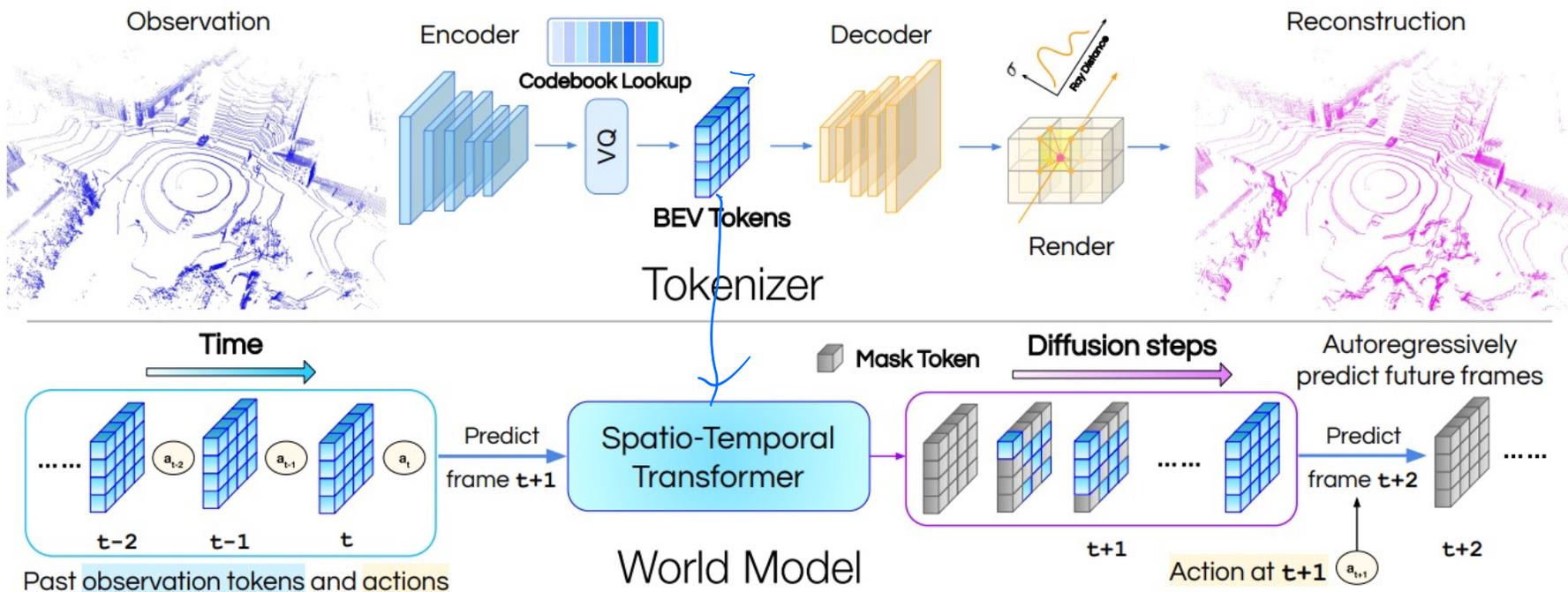
Action-Conditioned Video Generation

- DiT backbone with teacher forcing flow matching objective.
- Single actor robotic videos.



World Model in 3D volume prediction

- Autoregressively predict future 3D point clouds.

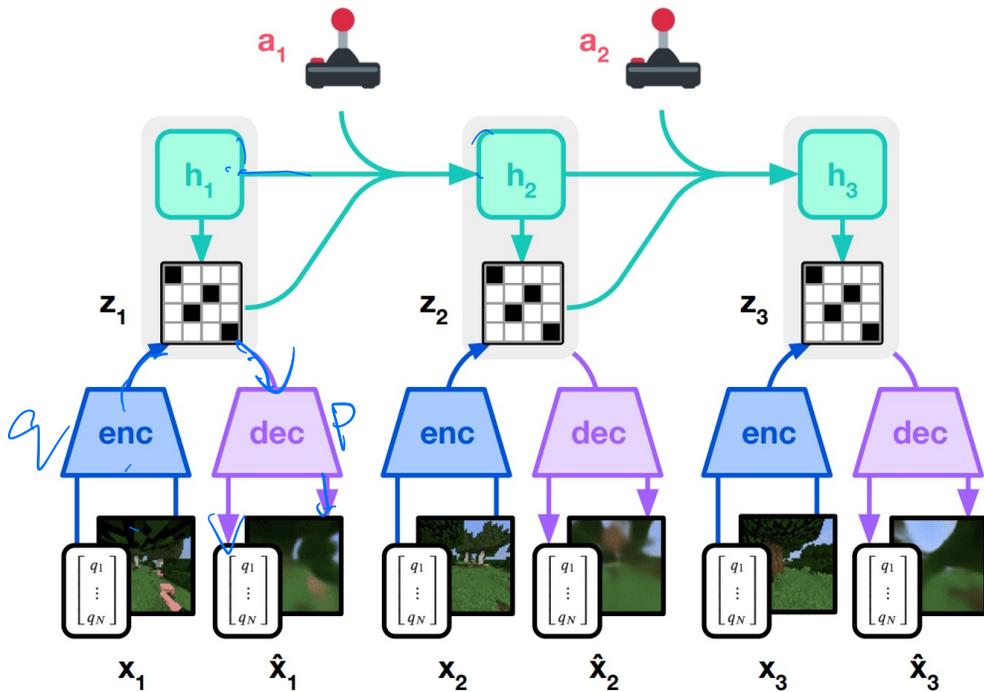


Latent Sequence World Model

- Autoencoder to ensure the latent representations are meaningful.

$$z_t \sim q_\phi(z_t | h_t, x_t)$$

$$\hat{x}_t \sim p_\phi(\hat{x}_t | h_t, z_t)$$



J&PA

Latent Sequence World Model

- Autoencoder to ensure the latent representations are meaningful.

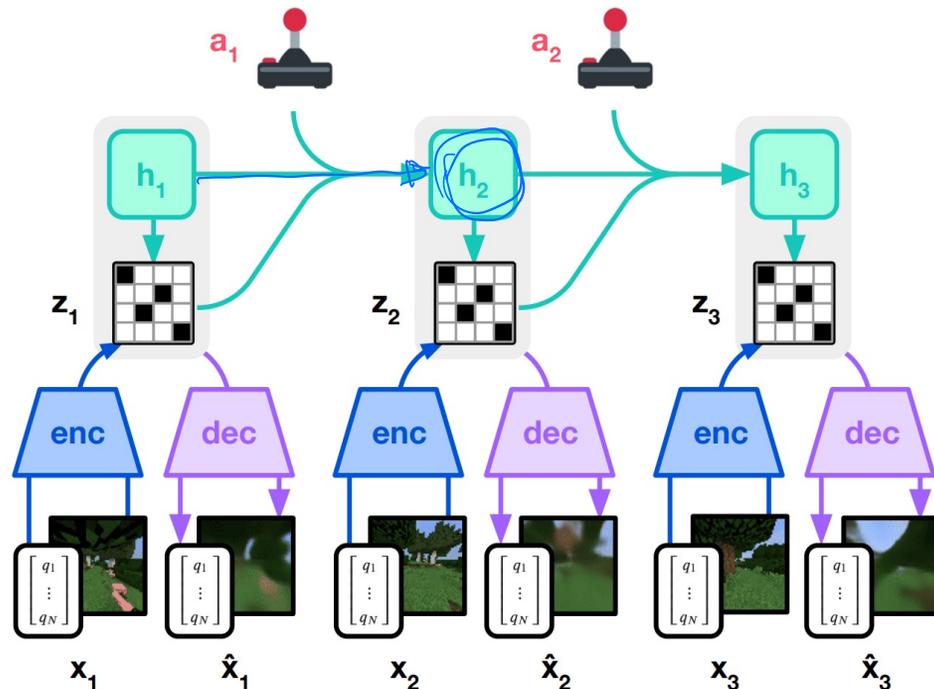
$$z_t \sim q_\phi(z_t | h_t, x_t)$$

$$\hat{x}_t \sim p_\phi(\hat{x}_t | h_t, z_t)$$

- Learn a sequence model to predict the latent conditioned on previous action.

$$h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1})$$

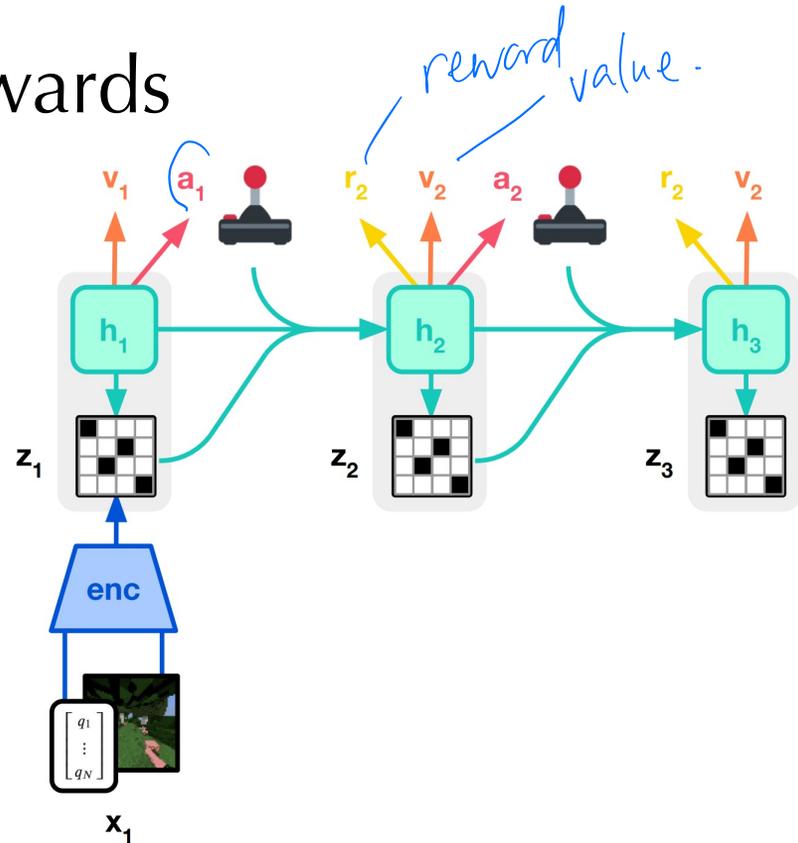
$$\hat{z}_t \sim p_\phi(\hat{z}_t | h_t)$$



Incorporating Rewards

- Predicting reward

$$\hat{r}_t \sim p_\phi(\hat{r}_t | h_t, z_t)$$



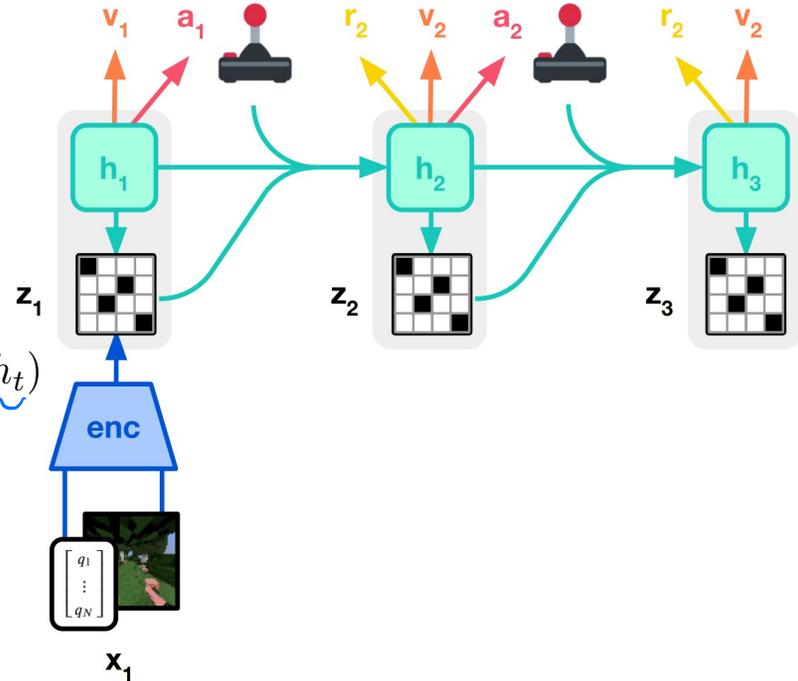
Incorporating Rewards

- Predicting reward

$$\hat{r}_t \sim p_\phi(\hat{r}_t | h_t, z_t)$$

- Reconstruction, reward, continue

$$\mathcal{L}_{\text{pred}}(\phi) = \underbrace{-\log p_\phi(x_t | z_t, h_t)}_{\text{reconstruction}} - \log p_\phi(r_t | z_t, h_t) - \log p_\phi(c_t | z_t, h_t)$$



Incorporating Rewards

- Predicting reward

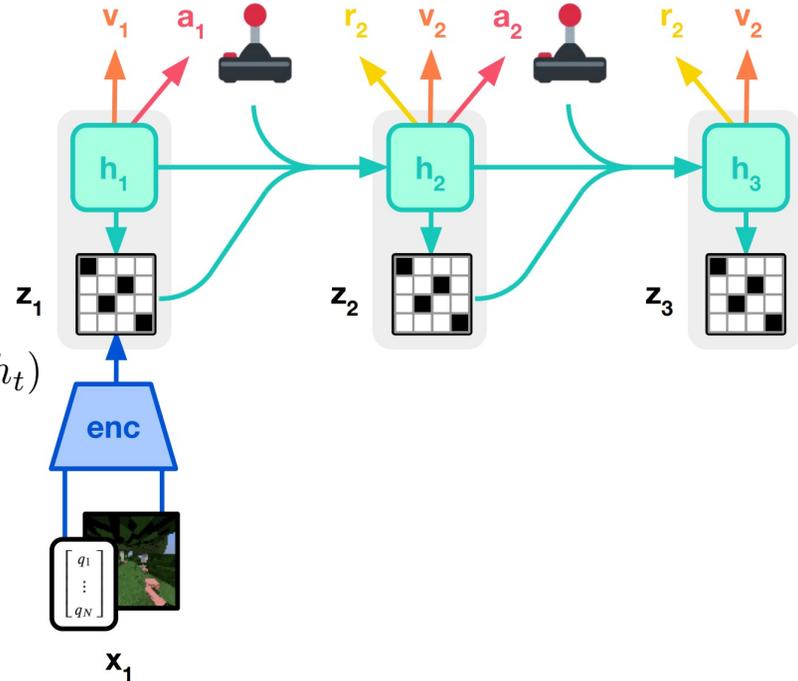
$$\hat{r}_t \sim p_\phi(\hat{r}_t | h_t, z_t)$$

- Reconstruction, reward, continue

$$\mathcal{L}_{\text{pred}}(\phi) = -\log p_\phi(x_t | z_t, h_t) - \log p_\phi(r_t | z_t, h_t) - \log p_\phi(c_t | z_t, h_t)$$

- Dynamics: Predicting future z

$$\mathcal{L}_{\text{dyn}}(\phi) = \max(1, \text{KL}[\text{sg}(q_\phi(z_t | h_t, x_t)) \| p_\phi(z_t | h_t)])$$



Incorporating Rewards

- Predicting reward

$$\hat{r}_t \sim p_\phi(\hat{r}_t | h_t, z_t)$$

- Reconstruction, reward, continue

$$\mathcal{L}_{\text{pred}}(\phi) = -\log p_\phi(x_t | z_t, h_t) - \log p_\phi(r_t | z_t, h_t) - \log p_\phi(c_t | z_t, h_t)$$

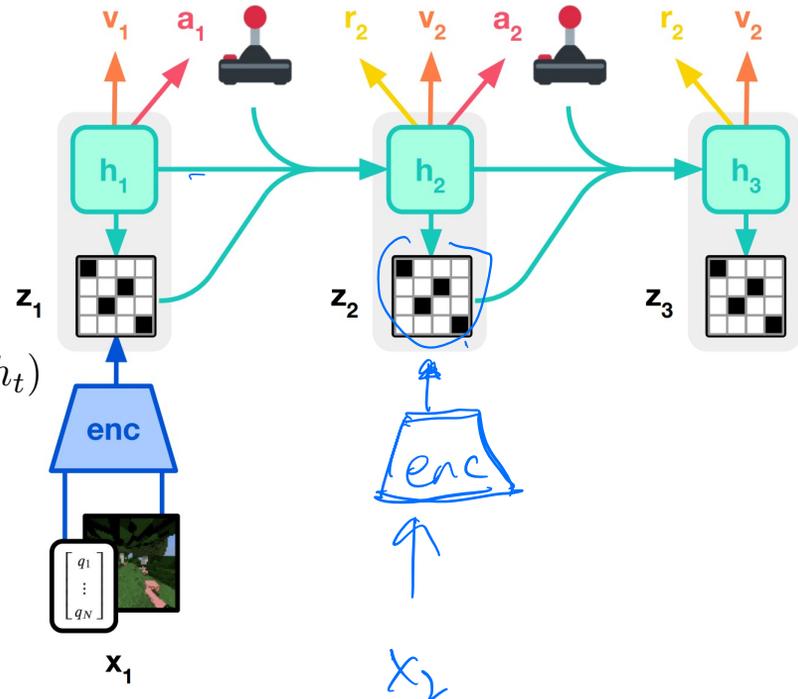
- Dynamics: Predicting future z

$$\mathcal{L}_{\text{dyn}}(\phi) = \max(1, \text{KL}[\text{sg}(q_\phi(z_t | h_t, x_t)) \| p_\phi(z_t | h_t)])$$

- Align representation

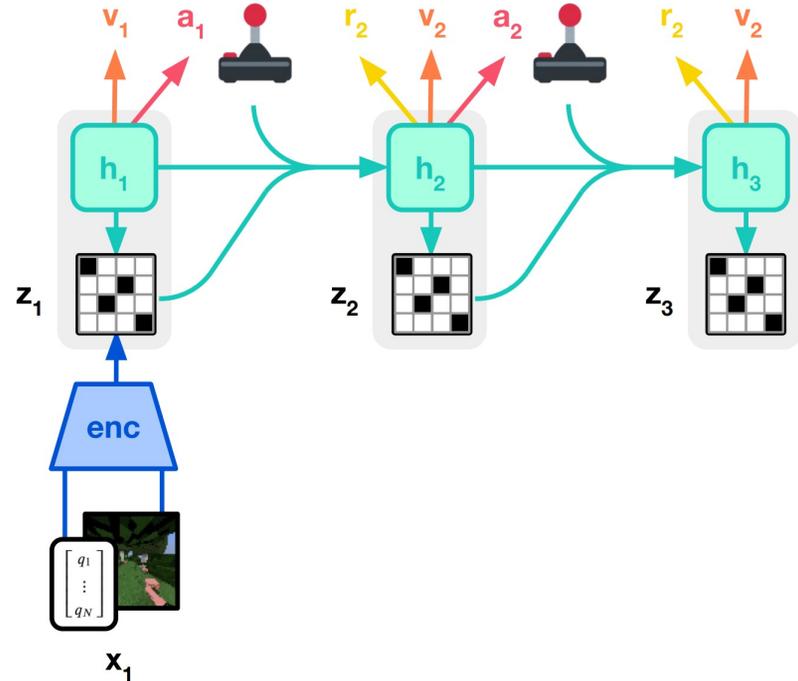
$$\mathcal{L}_{\text{rep}} = \max(1, \text{KL}[q_\phi(z_t | h_t, x_t) \| \text{sg}(p_\phi(z_t | h_t))])$$

JEPA



Incorporating Rewards

- How to use WM in planning?
Predicting value by simulate a batch of trajectories.



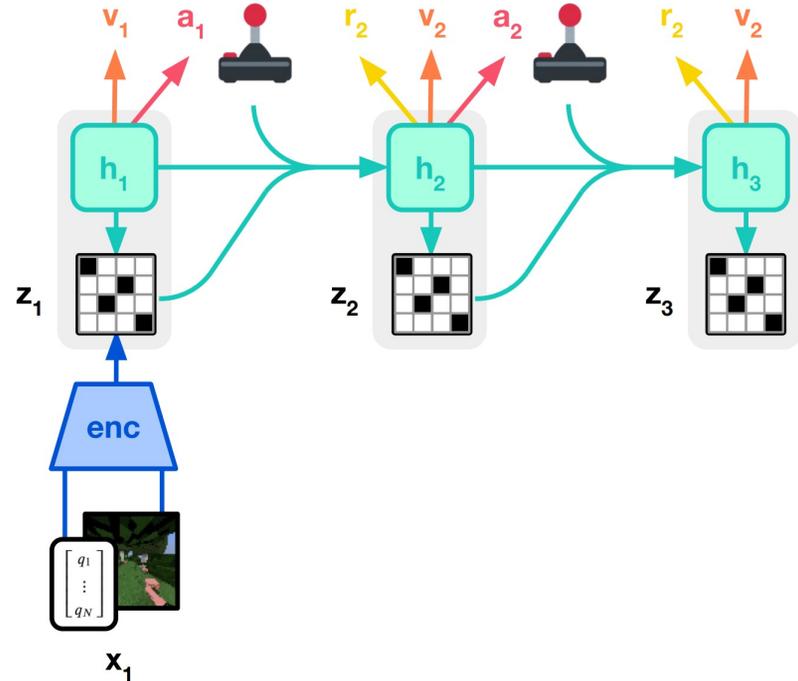
Incorporating Rewards

- How to use WM in planning?
Predicting value by simulate a batch of trajectories.

- Actor-Critic RL:

$$a_t \sim \pi_{\theta}(a_t | s_t)$$

$$v_{\psi}(R_t | s_t)$$



Incorporating Rewards

- How to use WM in planning?
Predicting value by simulate a batch of trajectories.

- Actor-Critic RL:

$$a_t \sim \pi_{\theta}(a_t | s_t) \quad v_{\psi}(R_t | s_t)$$

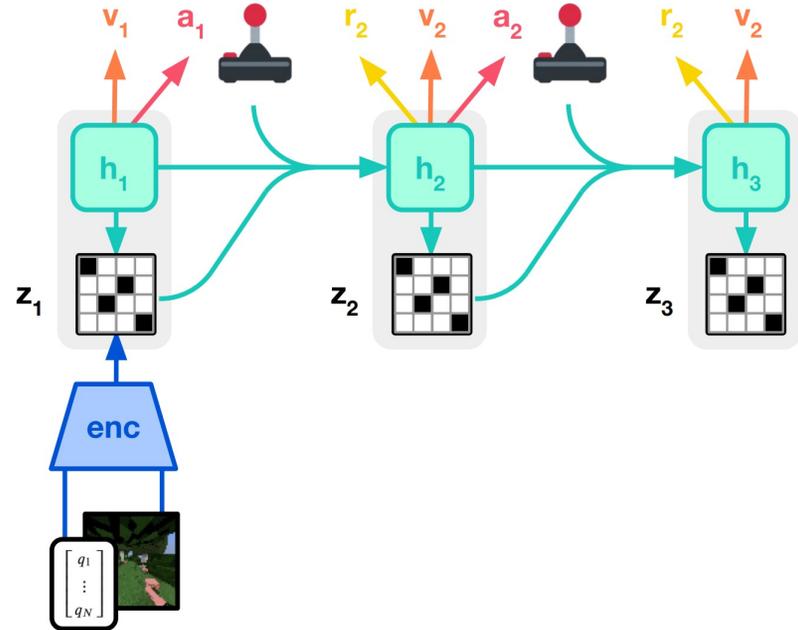
- Learning a critic:

Categorical distribution

$$\mathcal{L}_{\phi} = - \sum_{t=1}^T \log p_{\phi}(R_t^{\lambda} | s_t) \quad s_t = \{h_t, z_t\} \quad x_1$$

Sum of discounted future rewards

$$R_t^{\lambda} = r_t + \gamma c_t [(1 - \lambda)v_t + \lambda R_{t+1}^{\lambda}] \quad R_T^{\lambda} = v_T$$



Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$

Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$
- REINFORCE algorithm [Williams 1992]

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \underbrace{\text{sg}((R_t^\lambda - v_\phi(s_t)) / \text{max}(1, S))}_{\text{critic.}} \log \pi_\theta(a_t | s_t) + \eta H[\pi_\theta(a_t | s_t)]$$

Normalization Entropy Regularizer

Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$
- REINFORCE algorithm [Williams 1992]

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \text{sg}((R_t^\lambda - v_\phi(s_t)) / \overset{\text{Normalization}}{\max(1, S)}) \log \pi_\theta(a_t | s_t) + \overset{\text{Entropy Regularizer}}{\eta H[\pi_\theta(a_t | s_t)]}$$

- Notes on policy gradient:

Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$

- REINFORCE algorithm [Williams 1992]

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \text{sg}((R_t^\lambda - v_\phi(s_t)) / \overset{\text{Normalization}}{\max(1, S)}) \log \pi_\theta(a_t | s_t) + \overset{\text{Entropy Regularizer}}{\eta H[\pi_\theta(a_t | s_t)]}$$

- Notes on policy gradient:

$$\pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) = \pi_\theta(\tau) \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} = \nabla_\theta \pi_\theta(\tau).$$

Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$
- REINFORCE algorithm [Williams 1992]

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \text{sg}((R_t^\lambda - v_\phi(s_t)) / \overset{\text{Normalization}}{\max(1, S)}) \log \pi_\theta(a_t | s_t) + \overset{\text{Entropy Regularizer}}{\eta H[\pi_\theta(a_t | s_t)]}$$

- Notes on policy gradient:

$$\pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) = \pi_\theta(\tau) \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} = \nabla_\theta \pi_\theta(\tau).$$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t).$$

Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$
- REINFORCE algorithm [Williams 1992]

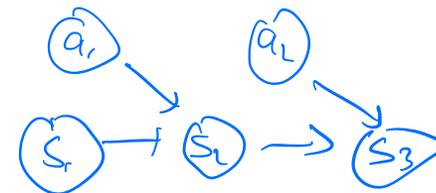
$$\mathcal{L}(\theta) = - \sum_{t=1}^T \text{sg}((R_t^\lambda - v_\phi(s_t)) / \overset{\text{Normalization}}{\max(1, S)}) \log \pi_\theta(a_t | s_t) + \overset{\text{Entropy Regularizer}}{\eta H[\pi_\theta(a_t | s_t)]}$$

- Notes on policy gradient:

$$\pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) = \pi_\theta(\tau) \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} = \nabla_\theta \pi_\theta(\tau).$$

$$\pi_\theta(\tau) = \pi_\theta(\underline{s_1}, \underline{a_1}, \dots, \underline{s_T}, \underline{a_T}) = p(s_1) \prod_{t=1}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t).$$

$$\log \pi_\theta(\tau) = \log p(s_1) + \sum_{t=1}^T \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t).$$



Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$

- REINFORCE algorithm [Williams 1992]

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \text{sg}((R_t^\lambda - v_\phi(s_t)) / \overset{\text{Normalization}}{\max(1, S)}) \log \pi_\theta(a_t | s_t) + \overset{\text{Entropy Regularizer}}{\eta H[\pi_\theta(a_t | s_t)]}$$

- Notes on policy gradient:

$$\pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) = \pi_\theta(\tau) \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} = \nabla_\theta \pi_\theta(\tau).$$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t).$$

$$\log \pi_\theta(\tau) = \log p(s_1) + \sum_{t=1}^T \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t).$$

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)} [r(\tau)] = \int \pi_\theta r(\tau) d\tau.$$

Actor Learning

- Learn a policy network $a_t \sim \pi_\theta(a_t | s_t)$

- REINFORCE algorithm [Williams 1992]

$$\mathcal{L}(\theta) = - \sum_{t=1}^T \underbrace{\text{sg}((R_t^\lambda - v_\phi(s_t)) / \text{max}(1, S))}_{\text{Normalization}} \underbrace{\log \pi_\theta(a_t | s_t)}_{\text{Entropy Regularizer}} + \eta H[\pi_\theta(a_t | s_t)]$$

- Notes on policy gradient:

$$\pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau) = \pi_\theta(\tau) \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} = \nabla_\theta \pi_\theta(\tau).$$

$$\pi_\theta(\tau) = \pi_\theta(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t).$$

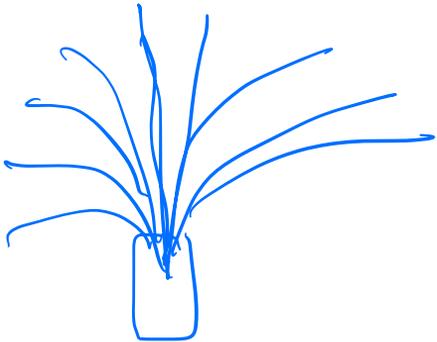
$$\log \pi_\theta(\tau) = \log p(s_1) + \sum_{t=1}^T \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t).$$

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)}[r(\tau)] = \int \pi_\theta r(\tau) d\tau.$$

$$\nabla \mathcal{L}(\theta) = \int \nabla \pi_\theta r(\tau) d\tau = \int \pi_\theta(\tau) \nabla \log \pi_\theta r(\tau) d\tau = \mathbb{E}_{\tau \sim \pi}[\nabla \log \pi_\theta r(\tau)].$$

When Do We Need A Learned Actor?

- For low dimensional or discrete problems, we can directly take the argmax of the value function.
- For problems with a good model, we can roll out and sample many future trajectories. Evaluation can be done in real time with GPU.



When Do We Need A Learned Actor?

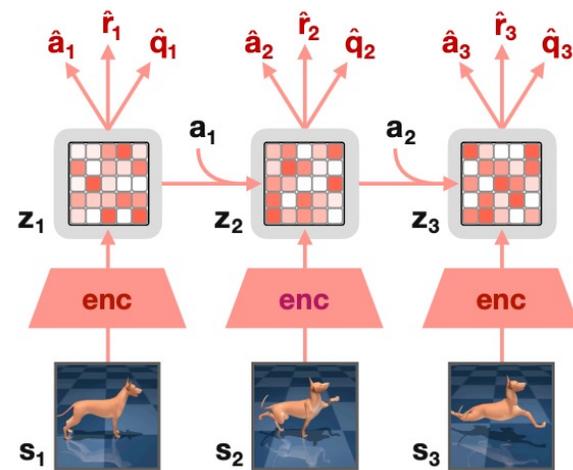
- For low dimensional or discrete problems, we can directly take the argmax of the value function.
- For problems with a good model, we can roll out and sample many future trajectories. Evaluation can be done in real time with GPU.
- It compresses planning into a reactive policy.

When Do We Need A Learned Actor?

- For low dimensional or discrete problems, we can directly take the argmax of the value function.
- For problems with a good model, we can roll out and sample many future trajectories. Evaluation can be done in real time with GPU.
- It compresses planning into a reactive policy.
- For general control problems, learning a separate actor can be a general solution without invoking domain knowledge.

Temporal Difference MPC

- Instead of learning an actor, we can directly update the actor network at test-time by optimizing the learned Q function from TD.

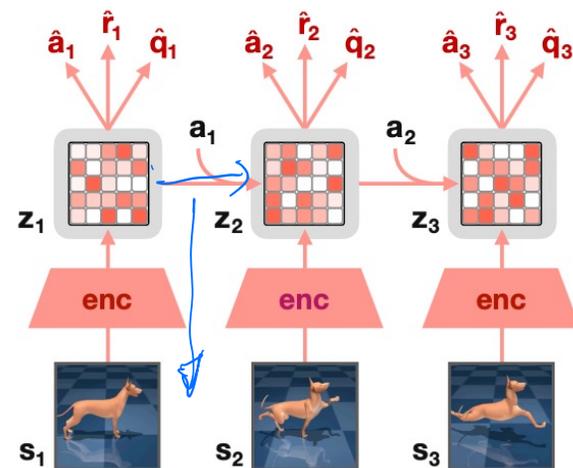


No reconstruction. JEPA.

Temporal Difference MPC

- Instead of learning an actor, we can directly update the actor network at test-time by optimizing the learned Q function from TD.
- Dynamics Learning:

$$\begin{aligned} z_t &= h_{\theta}(x_t) \\ \hat{z}_{t+k} &= d_{\theta}(\hat{z}_{t+k-1}, a_{t+k-1}) \end{aligned}$$



No reconstruction. JEPa.

Temporal Difference MPC

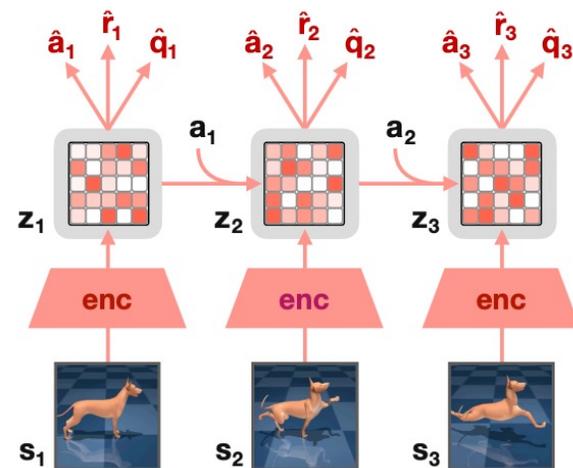
- Instead of learning an actor, we can directly update the actor network at test-time by optimizing the learned Q function from TD.

- Dynamics Learning:

$$z_t = h_{\theta}(x_t)$$

$$\hat{z}_{t+k} = d_{\theta}(\hat{z}_{t+k-1}, a_{t+k-1})$$

- TD: $\hat{Q}(z_t, a_t) \approx r_t + \gamma \bar{Q}(z_{t+1}, a_{t+1})$



No reconstruction. JEPA.

Temporal Difference MPC

- Instead of learning an actor, we can directly update the actor network at test-time by optimizing the learned Q function from TD.
- Dynamics Learning:

$$z_t = h_{\theta}(x_t)$$

$$\hat{z}_{t+k} = d_{\theta}(\hat{z}_{t+k-1}, a_{t+k-1})$$

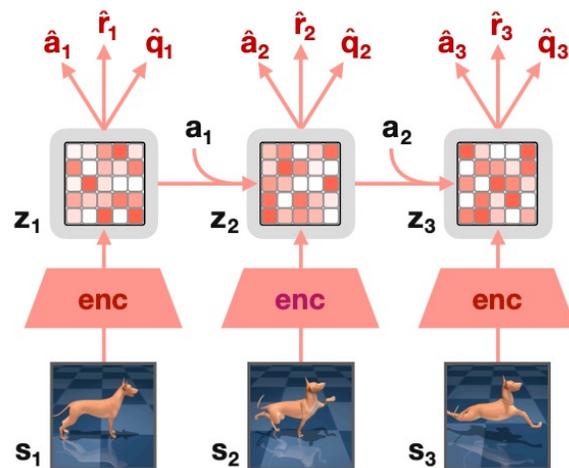
- TD: $\hat{Q}(z_t, a_t) \approx r_t + \gamma \bar{Q}(z_{t+1}, a_{t+1})$

- Planning:

$$J(A) = \sum_{k=0}^{H-1} \gamma^k \hat{r}(z_{t+k}, a_{t+k}) + \gamma^H \hat{Q}(z_{t+H}, a_{t+H})$$

estimate reward within horizon.

$$\mu^* = \arg \min_{\mu} (\mathbb{E}_{A \sim \mathcal{N}(\mu, \Sigma)} [-J(A)] + \lambda D_{\text{KL}}(\mathcal{N}(\mu, \Sigma) \parallel \mathbb{P}_{\text{prior}}))$$



No reconstruction. JEPA.

Sample Diverse Trajectories From Learned Prior

- Initialize and sample from policy prior

$$a_{t,i} \sim \mathcal{N}(\mu_{\theta}(z_t), \sigma_{\theta}(z_t)^2)$$

$$\mathcal{L}_{\text{prior}} = -\mathbb{E}_{a \sim \pi_{\theta}(\cdot|z)} \left[\hat{Q}_{\theta}(z, a) - \alpha \log \pi_{\theta}(a|z) \right]$$

Max-Entropy

Sample Diverse Trajectories From Learned Prior

- Initialize and sample from policy prior

$$a_{t,i} \sim \mathcal{N}(\mu_\theta(z_t), \sigma_\theta(z_t)^2)$$

$$\mathcal{L}_{\text{prior}} = -\mathbb{E}_{a \sim \pi_\theta(\cdot|z)} \left[\hat{Q}_\theta(z, a) - \alpha \log \pi_\theta(a|z) \right]$$

Max-Entropy

- Roll out

$$R(\tau_i) = \sum_{k=0}^{H-1} \gamma^k \hat{r}_{t+k,i} + \gamma^H \hat{Q}(\hat{z}_{t+H,i}, \underline{a}_{t+H,i})$$

Sample Diverse Trajectories From Learned Prior

- Initialize and sample from policy prior

$$a_{t,i} \sim \mathcal{N}(\mu_\theta(z_t), \sigma_\theta(z_t)^2) \quad \mathcal{L}_{\text{prior}} = -\mathbb{E}_{a \sim \pi_\theta(\cdot|z)} \left[\hat{Q}_\theta(z, a) - \alpha \log \pi_\theta(a|z) \right] \quad \text{Max-Entropy}$$

- Roll out

$$R(\tau_i) = \sum_{k=0}^{H-1} \gamma^k \hat{r}_{t+k,i} + \gamma^H \hat{Q}(\hat{z}_{t+H,i}, a_{t+H,i})$$

gradient-free.

- Weighted average smoothing and adaptation

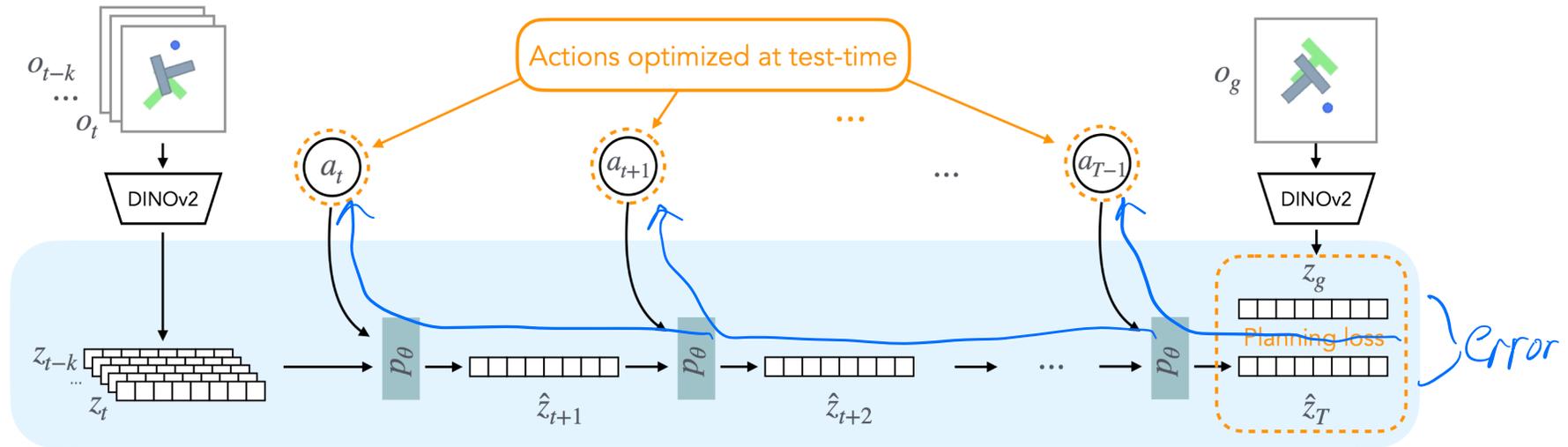
$$w_i = \frac{\exp\left(\frac{1}{\lambda} R(\tau_i)\right)}{\sum_{j=1}^N \exp\left(\frac{1}{\lambda} R(\tau_j)\right)} \quad a_t^{\text{new}} = \sum_{i=1}^N w_i a_{t,i}$$

Latent MPC Using Gradient Descent

- MPC is a differentiable computation graph.
- We can directly optimize action variables.

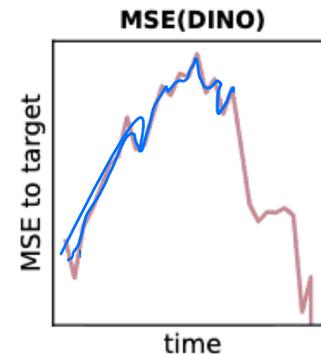
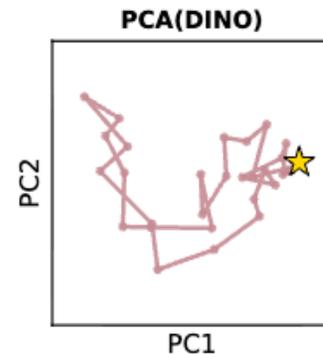
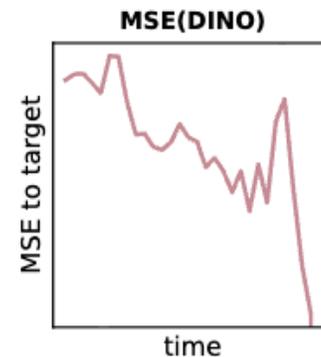
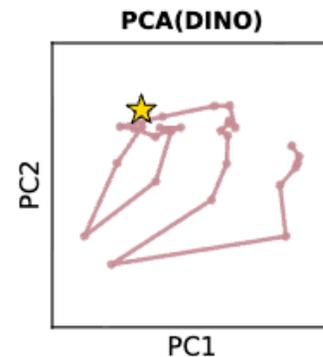
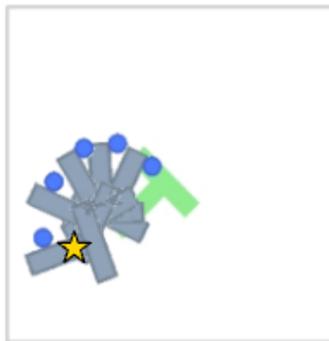
Latent MPC Using Gradient Descent

- MPC is a differentiable computation graph.
- We can directly optimize action variables.

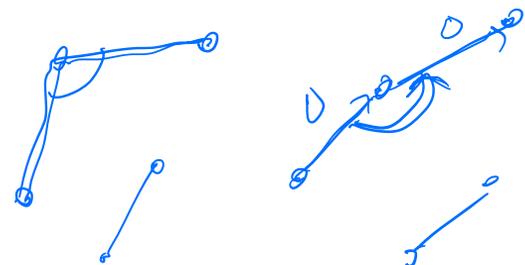


Learn Better Latent WM Representations

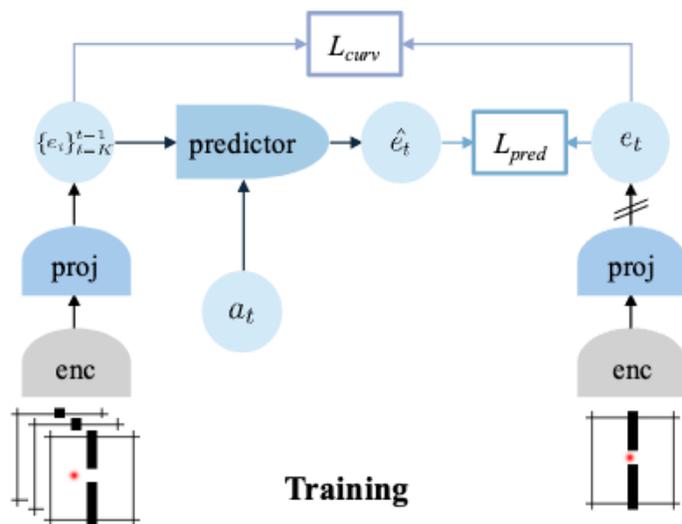
- The representation space is often jagged.
- Hard to optimize action sequences at test-time with gradient descent.



Straightening Objective



- In addition to predictors, maximizing cosine similarity across adjacent velocities (minimizing curvature).



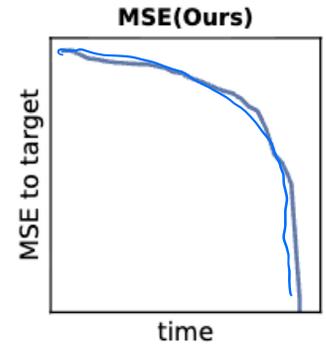
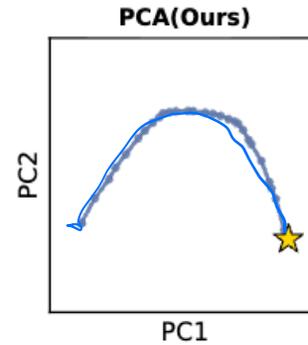
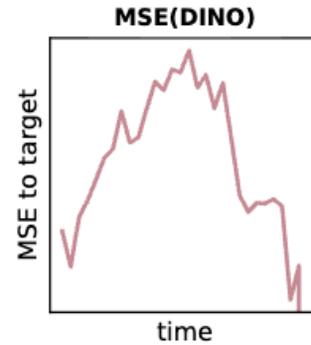
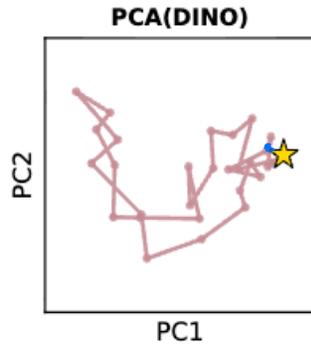
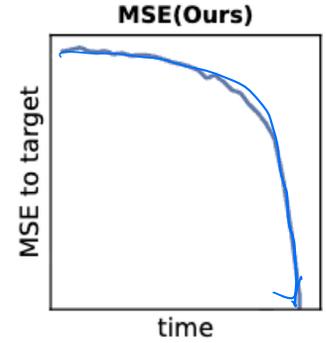
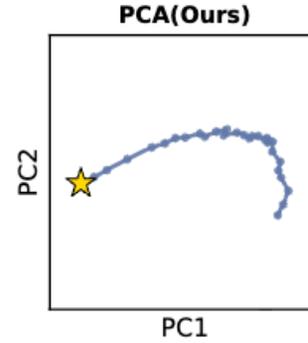
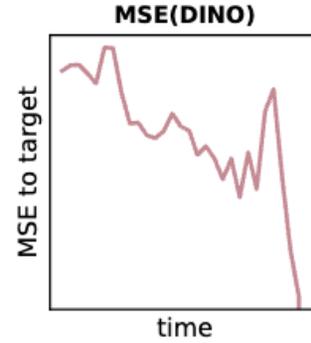
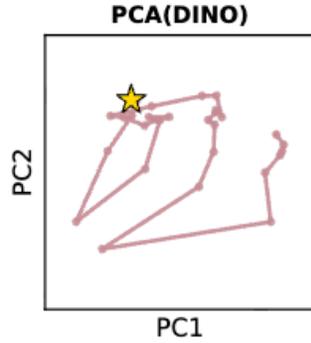
JEPA

$$\mathcal{L}_{pred} = \|\hat{z}_{t+1} - \text{sg}(z_{t+1})\|_2^2$$

$$v_t = z_{t+1} - z_t, \quad v_{t+1} = z_{t+2} - z_{t+1}$$

$$\mathcal{L}_{curve} = 1 - \frac{v_t \cdot v_{t+1}}{\|v_t\|_2 \cdot \|v_{t+1}\|_2}$$

Straightened Representations



Why Straightening Helps?

- Encourages the predictor to be linear – more generalizable.

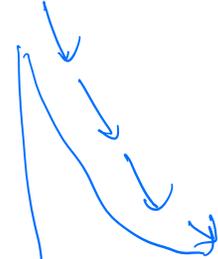
Why Straightening Helps?



- Encourages the predictor to be linear – more generalizable.
- Shortest path in the state space (with obstacle) – shortest path in the latent space (without obstacle).



Why Straightening Helps?



- Encourages the predictor to be linear – more generalizable.
- Shortest path in the state space (with obstacle) – shortest path in the latent space (without obstacle).
- Convexification effect:

$$\nabla_{\mathbf{a}} C(z_T(\mathbf{a})) = \left(\frac{\partial z_T}{\partial \mathbf{a}} \right)^T \nabla_z C(z_T)$$
$$H = \underbrace{\left(\frac{\partial z_T}{\partial \mathbf{a}} \right)^T \nabla_z^2 C(z_T) \left(\frac{\partial z_T}{\partial \mathbf{a}} \right)}_{\mathcal{T}_1: \text{Cost Convexity}} + \underbrace{\nabla_z C(z_T) \cdot \frac{\partial^2 z_T}{\partial \mathbf{a}^2}}_{\mathcal{T}_2: \text{Dynamics Curvature}}$$

↓
Closer zero.

Why Straightening Helps?

- Encourages the predictor to be linear – more generalizable.
- Shortest path in the state space (with obstacle) – shortest path in the latent space (without obstacle).

- Convexification effect: $\nabla_{\mathbf{a}} C(z_T(\mathbf{a})) = \left(\frac{\partial z_T}{\partial \mathbf{a}} \right)^T \nabla_z C(z_T)$

$$H = \underbrace{\left(\frac{\partial z_T}{\partial \mathbf{a}} \right)^T \nabla_z^2 C(z_T) \left(\frac{\partial z_T}{\partial \mathbf{a}} \right)}_{\mathcal{T}_1: \text{Cost Convexity}} + \underbrace{\nabla_z C(z_T) \cdot \frac{\partial^2 z_T}{\partial \mathbf{a}^2}}_{\mathcal{T}_2: \text{Dynamics Curvature}}$$

O(1) wrt time.

- Conditioning effect: Convergence of convex gradient optimization depends on κ (condition number). For straight trajectories, condition number is constant.

$$\|\mathbf{a}^{(k)} - \mathbf{a}^*\| \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^k \|\mathbf{a}^{(0)} - \mathbf{a}^*\|$$

Summary: World Models

- Explicit object representation
 - Traditional, light weight, instance-specific, hard to learn jointly
- WM for external actors.
Low-dim spaces
- WM for simple, single actor environment, visualization
- WM for general planning and test-time MPC

Summary: World Models



- Explicit object representation
 - Traditional, light weight, instance-specific, hard to learn jointly
- Differentiable occupancy, motion field
 - Relatively heavy, spatially grounded, end-to-end learnable

WM for external actors.
Low-dim spaces

WM for simple, single
actor environment,
visualization

WM for general
planning and
test-time MPC

Summary: World Models

- Explicit object representation
 - Traditional, light weight, instance-specific, hard to learn jointly
 - Differentiable occupancy, motion field
 - Relatively heavy, spatially grounded, end-to-end learnable
 - Raw video/3D prediction
 - Expensive, good for simulation and safety for short horizons
 - Good for single actor environments
- WM for external actors.
Low-dim spaces
- WM for simple, single actor environment, visualization
- WM for general planning and test-time MPC

Summary: World Models

- Explicit object representation
 - Traditional, light weight, instance-specific, hard to learn jointly
- Differentiable occupancy, motion field
 - Relatively heavy, spatially grounded, end-to-end learnable
- Raw video/3D prediction
 - Expensive, good for simulation and safety for short horizons
 - Good for single actor environments
- Global latent, RNNs
 - General-purpose, flexible
 - Planning-aware representations

* Collect TD reward and train value function / critic.
* Model predictive control.

WM for external actors.
Low-dim spaces

WM for simple, single actor environment, visualization

WM for general planning and test-time MPC

Learnable Planning and Cost

Imitation Learning

- The explicit policy model, supervised learning (behavior cloning)

$$\hat{a} = f_{\theta}(x) \quad \mathcal{L} = \min_i \|a_i - \hat{a}\|_2^2 \quad \mathcal{L} = -\log \hat{a}_j$$

Imitation Learning

- The explicit policy model, supervised learning (behavior cloning)

$$\hat{a} = f_{\theta}(x) \quad \mathcal{L} = \min_i \|a_i - \hat{a}\|_2^2 \quad \mathcal{L} = -\log \hat{a}_j$$

- Energy-based (cost-based) approach

$$\tau^* = \operatorname{argmin}_{\tau} E(x, \tau)$$
$$p(\tau | x) = \frac{\exp(E(x, \tau))}{\int_{\tau} \exp(E(x, \tau))}$$

Imitation Learning

- The explicit policy model, supervised learning (behavior cloning)

$$\hat{a} = f_{\theta}(x) \quad \mathcal{L} = \min_i \|a_i - \hat{a}\|_2^2 \quad \mathcal{L} = -\log \hat{a}_j$$

- Energy-based (cost-based) approach

$$\tau^* = \operatorname{argmin}_{\tau} E(x, \tau)$$
$$p(\tau | x) = \frac{\exp(E(x, \tau))}{\int_{\tau} \exp(E(x, \tau))}$$

- Dataset Aggregation (DAGger)
 - Learned policy may deviate from experts
 - Need to collect more groundtruths

```
Initialize  $\mathcal{D} \leftarrow \emptyset$ .
Initialize  $\hat{\pi}_1$  to any policy in  $\Pi$ .
for  $i = 1$  to  $N$  do
  Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}_i$ .
  Sample  $T$ -step trajectories using  $\pi_i$ .
  Get dataset  $\mathcal{D}_i = \{(s, \pi^*(s))\}$  of visited states by  $\pi_i$ 
  and actions given by expert.
  Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .
  Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .
end for
Return best  $\hat{\pi}_i$  on validation.
```

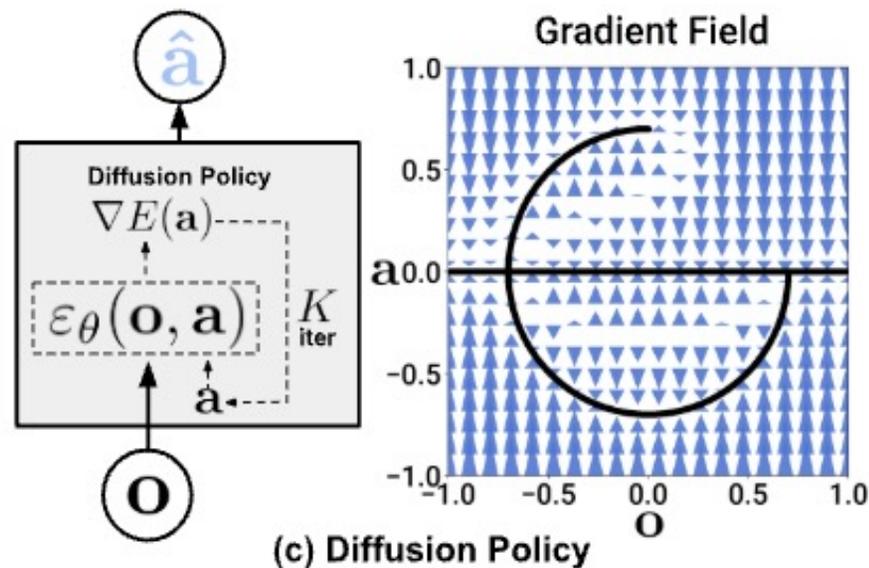
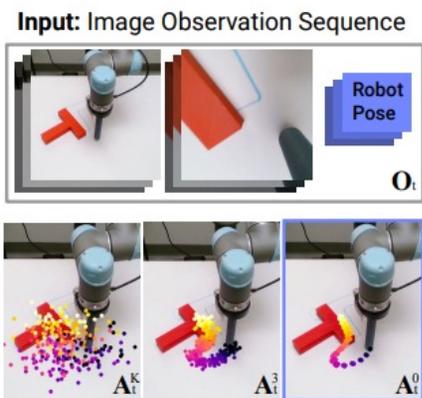
Algorithm 3.1: DAGGER Algorithm.

Direct Policy Learning from Diffusion

- Error prediction network is conditioned on observation features.

$$A_t^{k-1} = \alpha(A_t^k - \gamma \epsilon_\theta(O_t, A_t^k, k) + \mathcal{N}(0, \sigma^2 I)).$$

$$\mathcal{L} = \text{MSE}(\epsilon^k, \epsilon_\theta(O_t, A_t + \epsilon^k, k)).$$



Learning Cost

- Typically, we have some rough ideas on what the cost should look like
 - E.g. Avoid obstacle with squared distance barrier.

Learning Cost

- Typically, we have some rough ideas on what the cost should look like
 - E.g. Avoid obstacle with squared distance barrier.
- But there are costs that are implicitly defined.
 - E.g. Human comfort

Learning Cost

- Typically, we have some rough ideas on what the cost should look like
 - E.g. Avoid obstacle with squared distance barrier.
- But there are costs that are implicitly defined.
 - E.g. Human comfort
- Can be learned from human demonstrations (also called IRL)

Max-Margin Planning with Explicit Cost Volume

- If we have an explicit cost volume, the cost of a trajectory can be directly queried.

Max-Margin Planning with Explicit Cost Volume

- If we have an explicit cost volume, the cost of a trajectory can be directly queried.
- We can use the max-margin objective to make the groundtruth trajectory have lower costs.

$$\operatorname{argmin}_{\theta} \sum_{\{(\hat{x}_i^t, \hat{y}_i^t)\}_{i=1 \dots N}} \max_i \sum_{t=1}^T C_{\theta}^t[x_t, y_t] - C_{\theta}^t[\hat{x}_i^t, \hat{y}_i^t] + d_i^t$$

Max-Margin Planning with Explicit Cost Volume

- If we have an explicit cost volume, the cost of a trajectory can be directly queried.
- We can use the max-margin objective to make the groundtruth trajectory have lower costs.
- Find the lowest cost trajectory among a batch of samples.

$$\operatorname{argmin}_{\theta} \sum_{\{(\hat{x}_i^t, \hat{y}_i^t)\}_{i=1 \dots N}} \max_i \sum_{t=1}^T C_{\theta}^t[x_t, y_t] - C_{\theta}^t[\hat{x}_i^t, \hat{y}_i^t] + d_i^t$$

Max-Margin Planning with Explicit Cost Volume

- If we have an explicit cost volume, the cost of a trajectory can be directly queried.
- We can use the max-margin objective to make the groundtruth trajectory have lower costs.
- Find the lowest cost trajectory among a batch of samples.
- Low-dimensional/known dynamics problems: External samplers

$$\operatorname{argmin}_{\theta} \sum_{\{(\hat{x}_i^t, \hat{y}_i^t)\}_{i=1 \dots N}} \max_i \sum_{t=1}^T C_{\theta}^t[x_t, y_t] - C_{\theta}^t[\hat{x}_i^t, \hat{y}_i^t] + d_i^t$$

Max-Margin Planning with Explicit Cost Volume

- If we have an explicit cost volume, the cost of a trajectory can be directly queried.
- We can use the max-margin objective to make the groundtruth trajectory have lower costs.
- Find the lowest cost trajectory among a batch of samples.
- Low-dimensional/known dynamics problems: External samplers
- In general, needs to perform optimization (e.g. DP)

$$\operatorname{argmin}_{\theta} \sum_{\{(\hat{x}_i^t, \hat{y}_i^t)\}_{i=1 \dots N}} \max_i \sum_{t=1}^T C_{\theta}^t[x_t, y_t] - C_{\theta}^t[\hat{x}_i^t, \hat{y}_i^t] + d_i^t$$

EBM Planning

- Energy-based framework also needs negative samples.

EBM Planning

- Energy-based framework also needs negative samples.
- “Pick” the groundtruth sample among others.

EBM Planning

- Energy-based framework also needs negative samples.
- “Pick” the groundtruth sample among others.
- If there isn’t an external sampler, we can either use autoregressive energy of sampling one dimension at a time, or gradient-based Langevin MCMC.

EBM Planning

- Energy-based framework also needs negative samples.
- “Pick” the groundtruth sample among others.
- If there isn’t an external sampler, we can either use autoregressive energy of sampling one dimension at a time, or gradient-based Langevin MCMC.

Langevin MCMC: $\tilde{\mathbf{y}}_i^k = \tilde{\mathbf{y}}_i^{k-1} - \lambda \left(\frac{1}{2} \nabla_{\mathbf{y}} E_{\theta}(\mathbf{x}_i, \mathbf{y}_i^{k-1}) + \omega^k \right), \omega^k \sim \mathcal{N}(0, \sigma).$

EBM Planning

- Energy-based framework also needs negative samples.
- “Pick” the groundtruth sample among others.
- If there isn’t an external sampler, we can either use autoregressive energy of sampling one dimension at a time, or gradient-based Langevin MCMC.

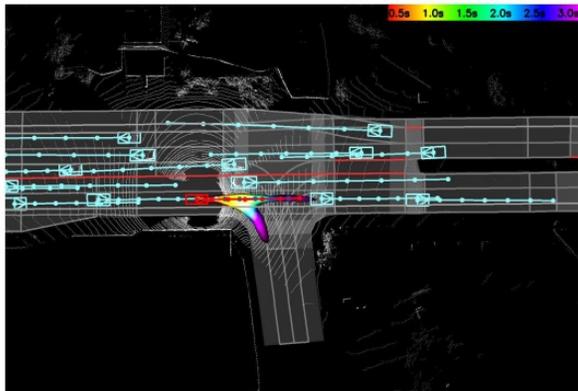
Langevin MCMC: $\tilde{\mathbf{y}}_i^k = \tilde{\mathbf{y}}_i^{k-1} - \lambda \left(\frac{1}{2} \nabla_{\mathbf{y}} E_{\theta}(\mathbf{x}_i, \mathbf{y}_i^{k-1}) + \omega^k \right), \omega^k \sim \mathcal{N}(0, \sigma).$

$$\text{Loss: } \mathcal{L} = \sum_i -\log(p_{\theta}(\mathbf{y}_i | \mathbf{x}, \{\tilde{\mathbf{y}}_i\}_j)) \quad p_{\theta}(\mathbf{y}_i | \mathbf{x}, \{\tilde{\mathbf{y}}_i\}_j) = \frac{e^{-E_{\theta}(\mathbf{x}_i, \mathbf{y}_i)}}{e^{-E_{\theta}(\mathbf{x}_i, \mathbf{y}_i)} + \sum_j e^{-E_{\theta}(\mathbf{x}_i, \tilde{\mathbf{y}}_{i,j})}}$$

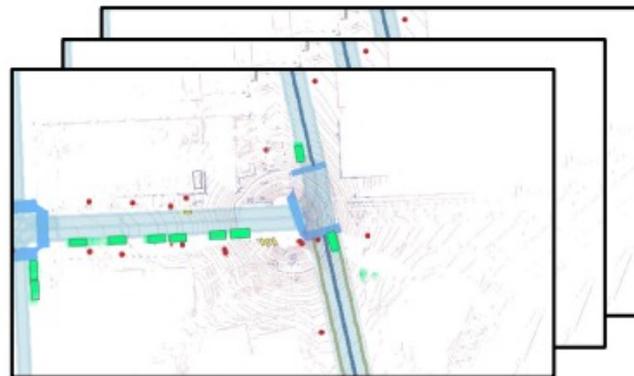
Non-Parametric Cost Volume for 2D Planning

- Interpretability (both costs and planner inputs)

Rasterization



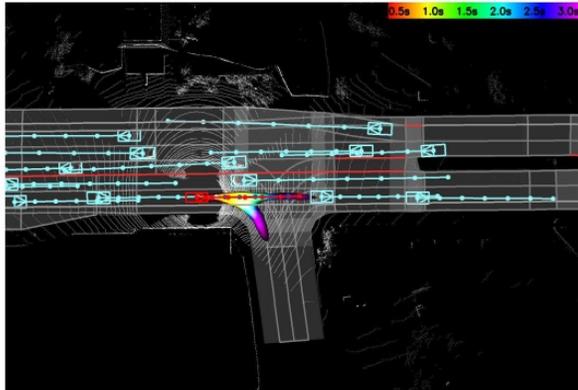
Semantic Occupancy



Non-Parametric Cost Volume for 2D Planning

- Interpretability (both costs and planner inputs)
- Use spatial geometry to form cost from explicit objects

Rasterization



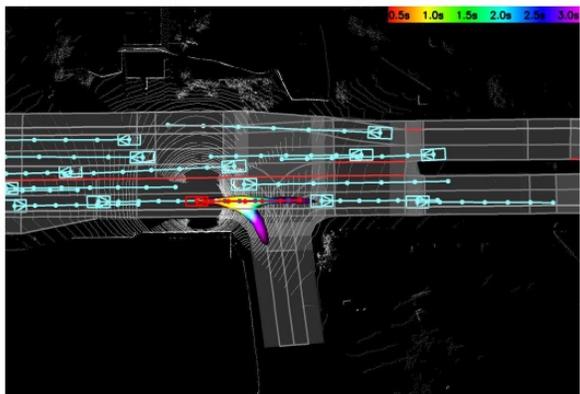
Semantic Occupancy



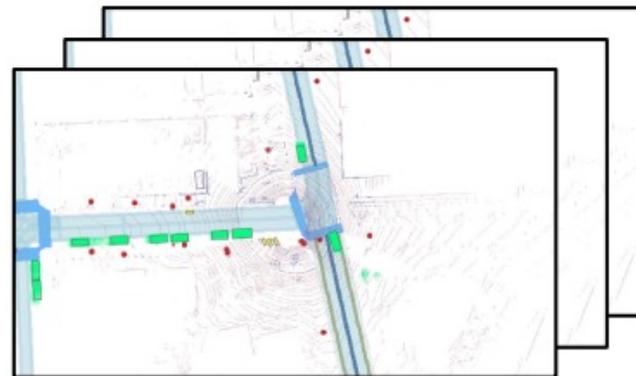
Non-Parametric Cost Volume for 2D Planning

- Interpretability (both costs and planner inputs)
- Use spatial geometry to form cost from explicit objects
- Predict spatial cost volume
 - Rasterize the scene for spatial inputs
 - Predict soft occupancy volumes (present and future)

Rasterization

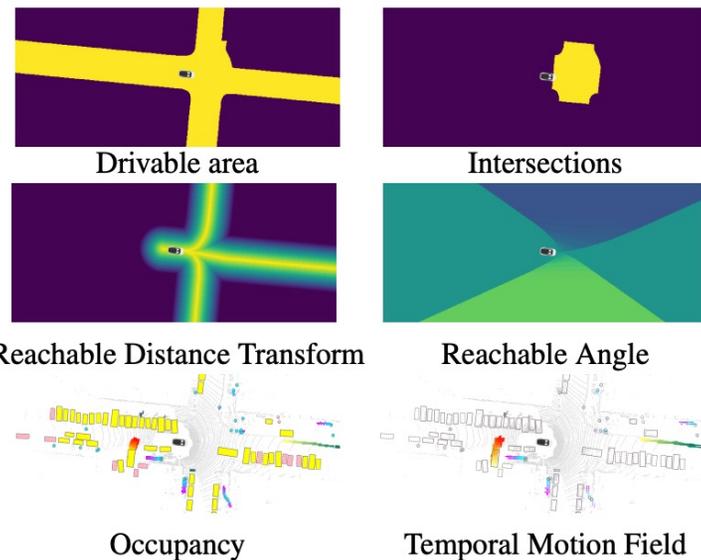


Semantic Occupancy



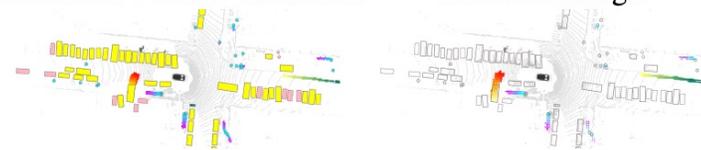
Learning Through Interpretable Predictions

- Semantic occupancy, motion field, mapping, etc. as intermediate predictions.



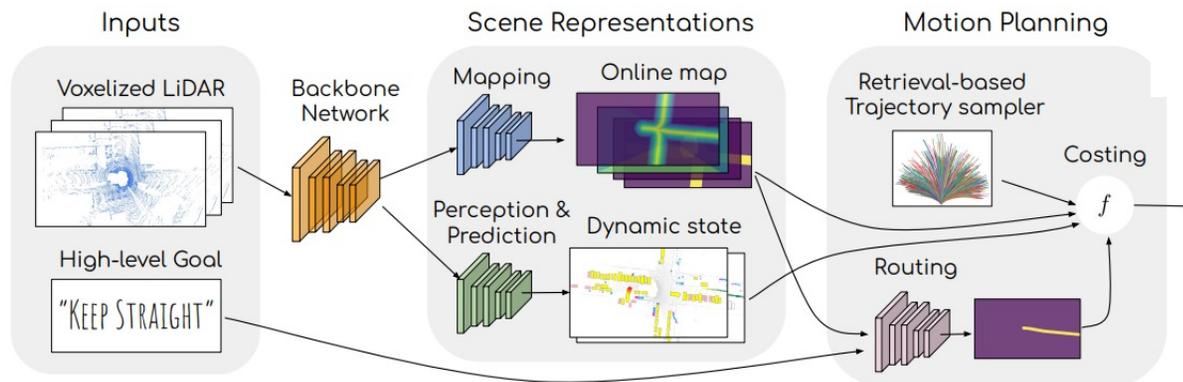
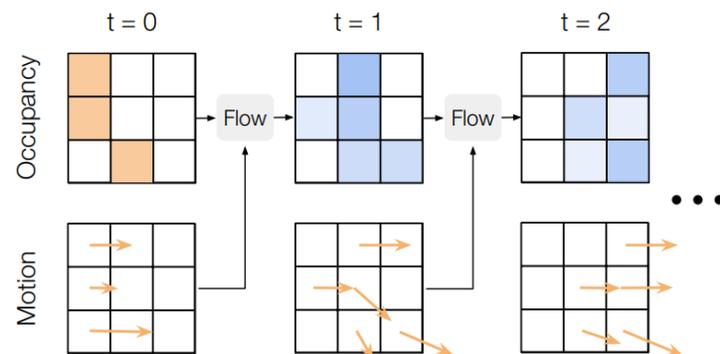
Reachable Distance Transform

Reachable Angle



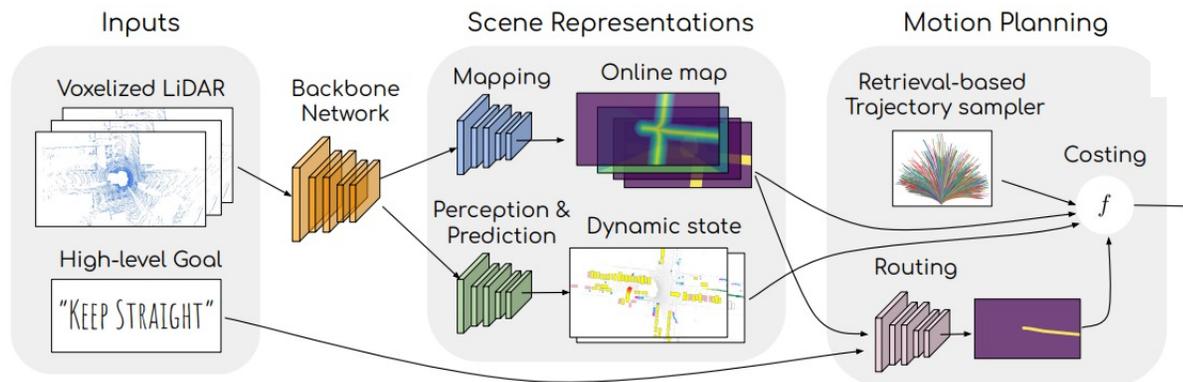
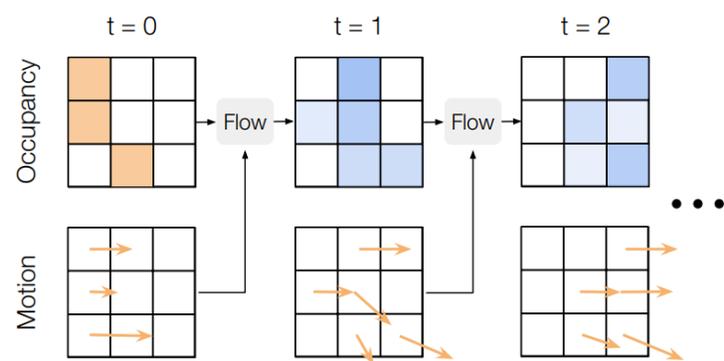
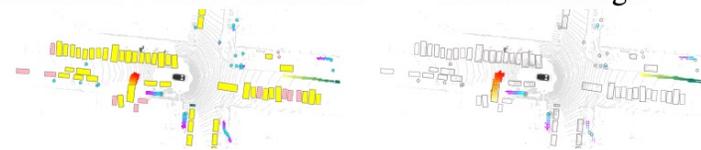
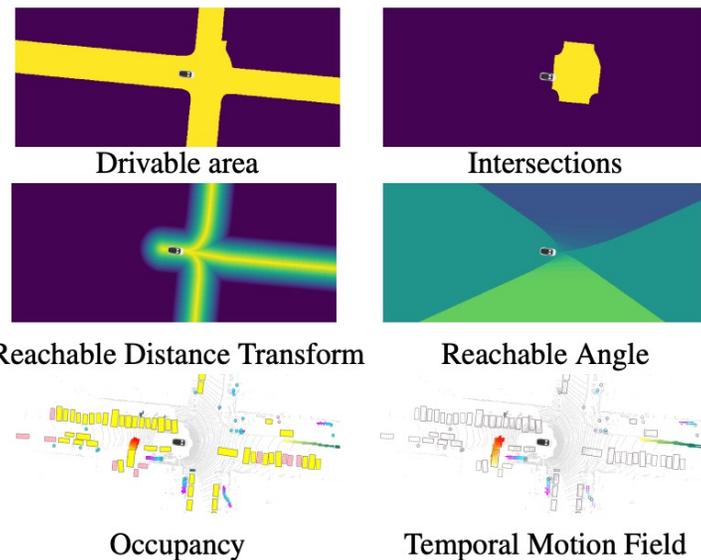
Occupancy

Temporal Motion Field



Learning Through Interpretable Predictions

- Semantic occupancy, motion field, mapping, etc. as intermediate predictions.
- Differentiable, supports end-to-end interpretable learning from perception to planning.



Value Iteration Networks

- A network design for predicting cost volumes that are grounded from the classic value iteration algorithm.

Value Iteration Networks

- A network design for predicting cost volumes that are grounded from the classic value iteration algorithm.
- Classic VI algorithm:

Value Iteration Networks

- A network design for predicting cost volumes that are grounded from the classic value iteration algorithm.
- Classic VI algorithm:

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

Value Iteration Networks

- A network design for predicting cost volumes that are grounded from the classic value iteration algorithm.
- Classic VI algorithm:

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$V^{\pi}(s) = \mathbb{E}^{\pi} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

Value Iteration Networks

- A network design for predicting cost volumes that are grounded from the classic value iteration algorithm.
- Classic VI algorithm:

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$V^{\pi}(s) = \mathbb{E}^{\pi} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

$$Q_n(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) V_n(s')$$

Value Iteration Networks

- A network design for predicting cost volumes that are grounded from the classic value iteration algorithm.
- Classic VI algorithm:

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$V^{\pi}(s) = \mathbb{E}^{\pi} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

$$Q_n(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) V_n(s')$$

$$V_{n+1}(s) = \max_a Q_n(s, a)$$

Value Iteration Networks

- A network design for predicting cost volumes that are grounded from the classic value iteration algorithm.

- Classic VI algorithm:

$$V^*(s) = \max_{\pi} V^{\pi}(s)$$

$$V^{\pi}(s) = \mathbb{E}^{\pi} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

$$Q_n(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) V_n(s')$$

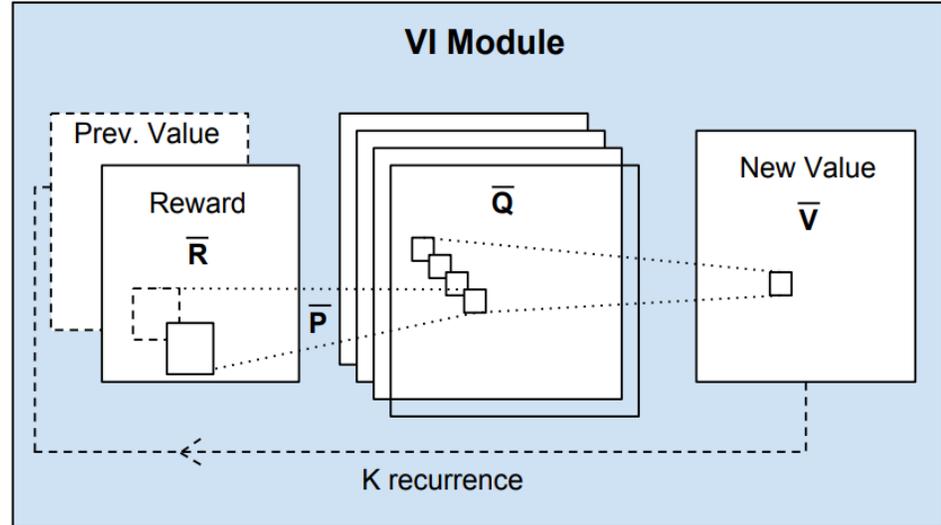
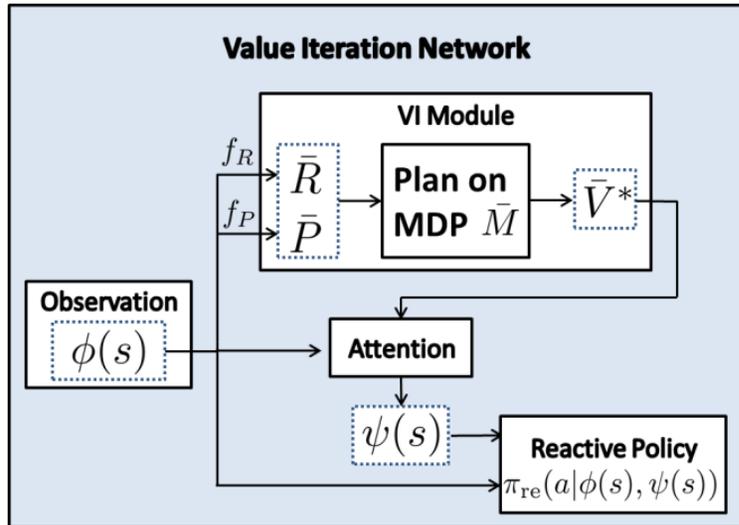
$$V_{n+1}(s) = \max_a Q_n(s, a)$$

$$\pi^*(s) = \operatorname{argmax}_a Q_{\infty}(s, a)$$

Value Iteration Networks

- Reward and previous value are fed into a network to generate Q of A channels. Transition matrix is convolutional kernel. Then Max-Pooling.

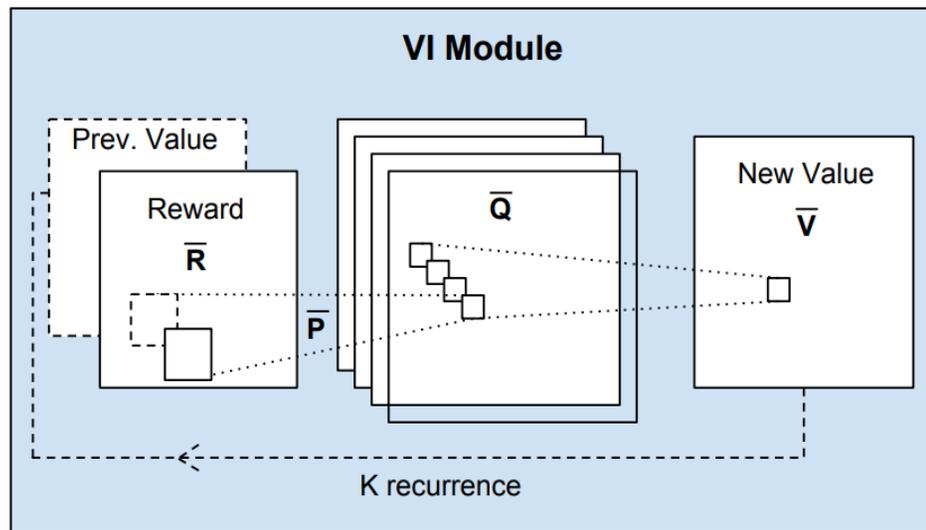
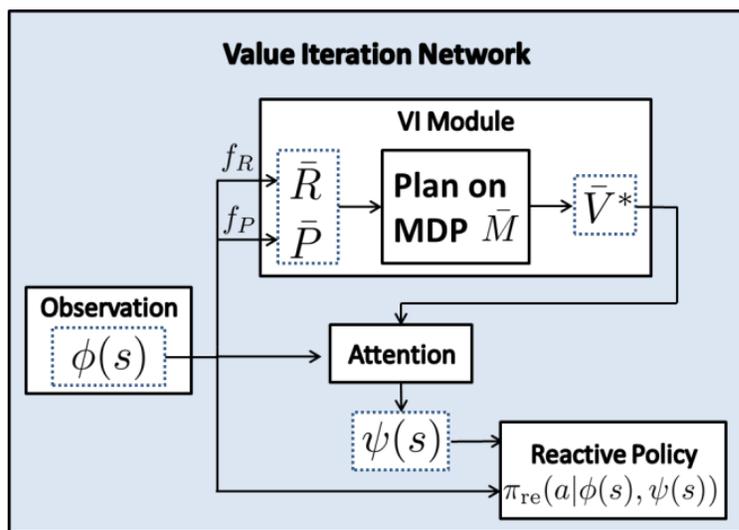
$$Q_n(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) V_n(s') \quad V_{n+1}(s) = \max_a Q_n(s, a)$$



Value Iteration Networks

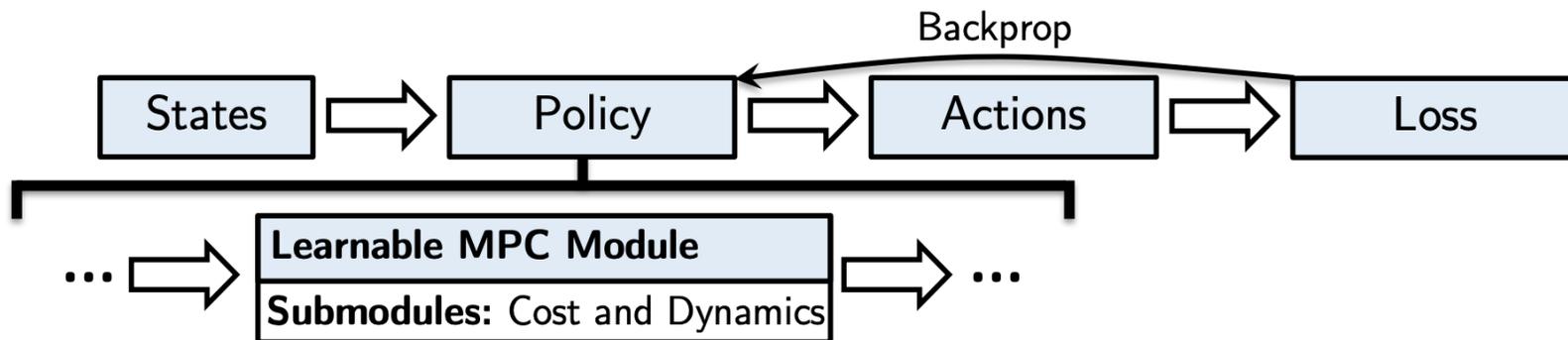
- Select the current state and choose an action from softmax.

$$\hat{a} \sim \text{softmax}_a(Q(s, a))$$



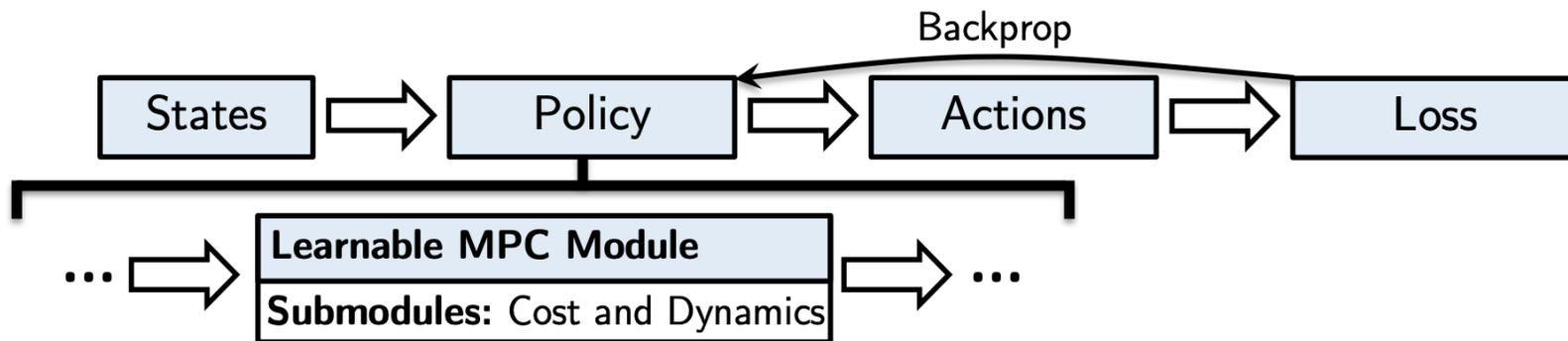
Backprop through Planning

- Treat planning as an end-to-end layer. Can be used for RL/Imitation.



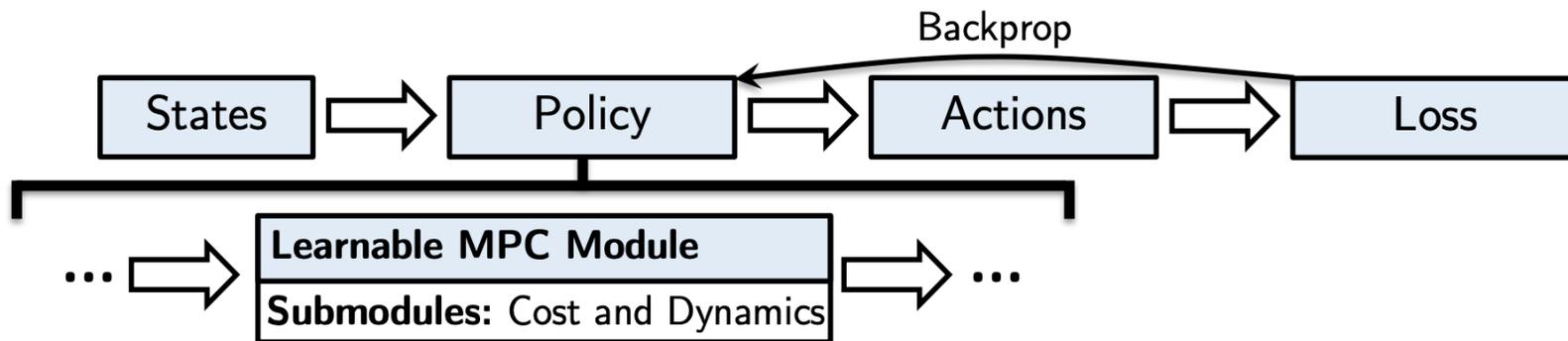
Backprop through Planning

- Treat planning as an end-to-end layer. Can be used for RL/Imitation.
- Option 1: Unrolling a finite number of steps



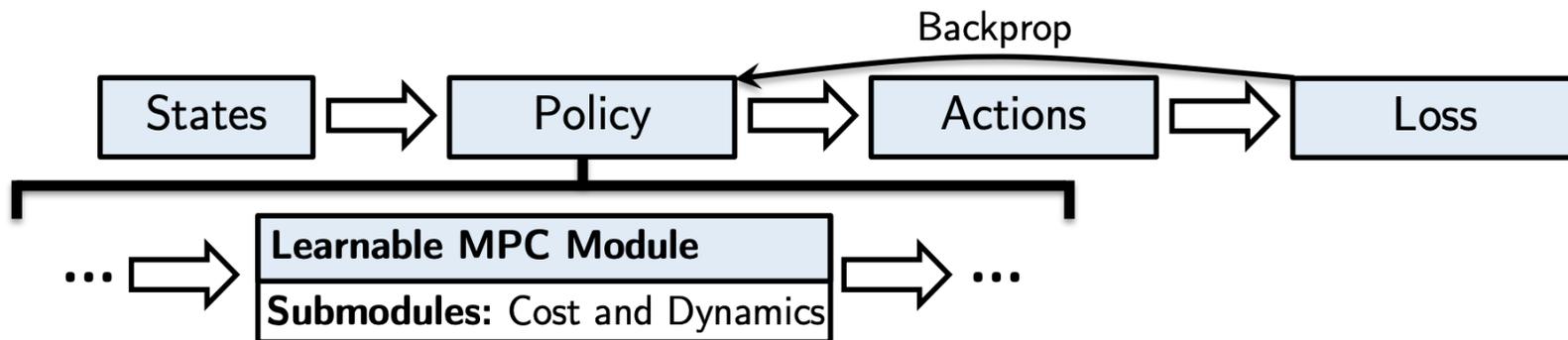
Backprop through Planning

- Treat planning as an end-to-end layer. Can be used for RL/Imitation.
- Option 1: Unrolling a finite number of steps
- Option 2: Solve till convergence, backprop for a finite step



Backprop through Planning

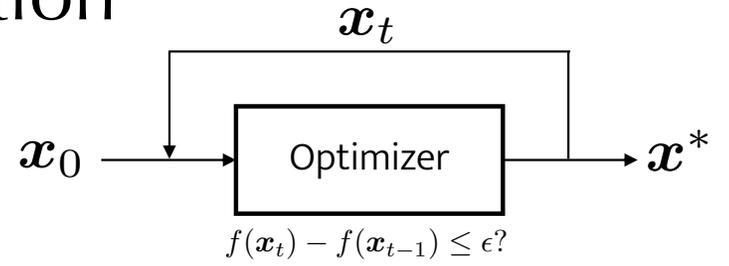
- Treat planning as an end-to-end layer. Can be used for RL/Imitation.
- Option 1: Unrolling a finite number of steps
- Option 2: Solve till convergence, backprop for a finite step
- Option 3: Converged at fixed point: Implicit differentiation



Implicit Differentiation

- Unconstrained case

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}; \boldsymbol{\theta}).$$

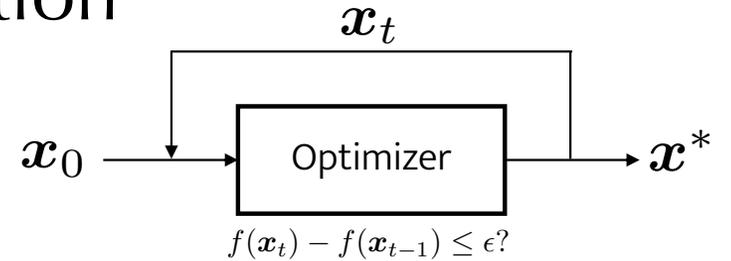


Implicit Differentiation

- Unconstrained case

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta}).$$

$$\mathbf{0} = \frac{d}{d\boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$



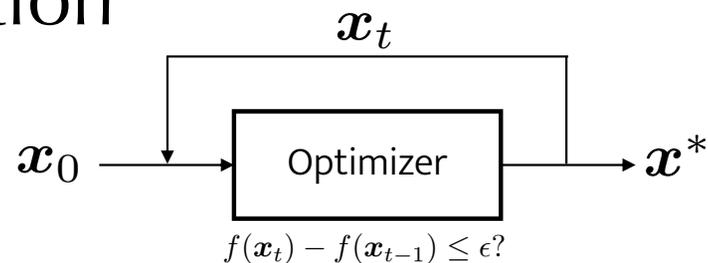
Implicit Differentiation

- Unconstrained case

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta}).$$

$$\mathbf{0} = \frac{d}{d\boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$

$$\mathbf{0} = \frac{\partial}{\partial \mathbf{x}^*} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta}) \frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\theta}} + \frac{\partial}{\partial \boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$



Implicit Differentiation

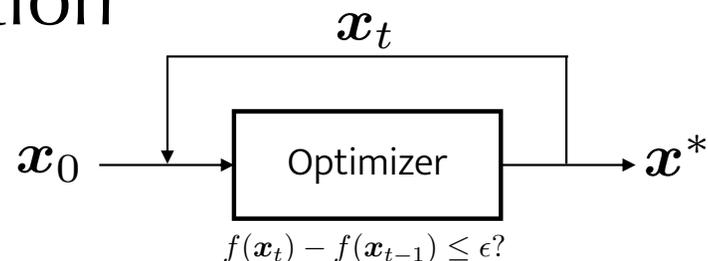
- Unconstrained case

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}; \boldsymbol{\theta}).$$

$$\mathbf{0} = \frac{d}{d\boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$

$$\mathbf{0} = \frac{\partial}{\partial \mathbf{x}^*} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta}) \frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\theta}} + \frac{\partial}{\partial \boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$

$$\mathbf{0} = H_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta}) \frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\theta}} + \frac{\partial}{\partial \boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$



Implicit Differentiation

- Unconstrained case

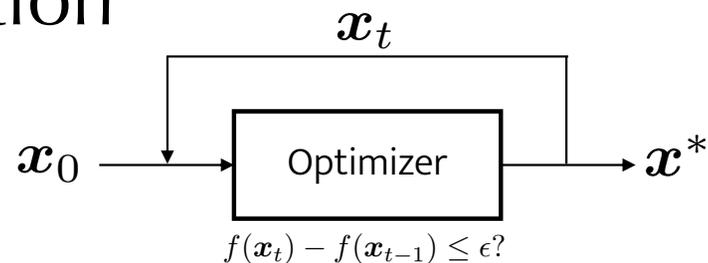
$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}; \boldsymbol{\theta}).$$

$$\mathbf{0} = \frac{d}{d\boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$

$$\mathbf{0} = \frac{\partial}{\partial \mathbf{x}^*} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta}) \frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\theta}} + \frac{\partial}{\partial \boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$

$$\mathbf{0} = H_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta}) \frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\theta}} + \frac{\partial}{\partial \boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})$$

$$\frac{\partial \mathbf{x}^*}{\partial \boldsymbol{\theta}} = H_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta})^{-1} \frac{\partial}{\partial \boldsymbol{\theta}} J_{f, \mathbf{x}^*}(\mathbf{x}^*; \boldsymbol{\theta}).$$



Implicit Differentiation

- How to compute Hessian inverse vector product?

Implicit Differentiation

- How to compute Hessian inverse vector product?
- Conjugate gradient, solve $Ax = b$

Conjugate Gradient Method

$$\mathbf{r}_0 := \mathbf{b} - \mathbf{A}\mathbf{x}_0$$

if \mathbf{r}_0 is sufficiently small, then return \mathbf{x}_0 as the result

$$\mathbf{p}_0 := \mathbf{r}_0$$

$$k := 0$$

repeat

$$\alpha_k := \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k}$$

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

$$\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$$

if \mathbf{r}_{k+1} is sufficiently small, then exit loop

$$\beta_k := \frac{\mathbf{r}_{k+1}^\top \mathbf{r}_{k+1}}{\mathbf{r}_k^\top \mathbf{r}_k}$$

$$\mathbf{p}_{k+1} := \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k$$

$$k := k + 1$$

end repeat

return \mathbf{x}_{k+1} as the result

Implicit Differentiation

- How to compute Hessian inverse vector product?
- Conjugate gradient, solve $A\mathbf{x} = \mathbf{b}$
- Neumann series (finite truncation)

$$(I - A)^{-1} = \sum_{k=0}^{\infty} A^k.$$

- Same as backprop the last K steps (Option 2).
- Memory savings.

Conjugate Gradient Method

$$\mathbf{r}_0 := \mathbf{b} - \mathbf{A}\mathbf{x}_0$$

if \mathbf{r}_0 is sufficiently small, then return \mathbf{x}_0 as the result

$$\mathbf{p}_0 := \mathbf{r}_0$$

$$k := 0$$

repeat

$$\alpha_k := \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{p}_k^\top \mathbf{A} \mathbf{p}_k}$$

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

$$\mathbf{r}_{k+1} := \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$$

if \mathbf{r}_{k+1} is sufficiently small, then exit loop

$$\beta_k := \frac{\mathbf{r}_{k+1}^\top \mathbf{r}_{k+1}}{\mathbf{r}_k^\top \mathbf{r}_k}$$

$$\mathbf{p}_{k+1} := \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k$$

$$k := k + 1$$

end repeat

return \mathbf{x}_{k+1} as the result

Differentiable LQR

- Now add linear equality constraints on the dynamics and initialization.

$$\tau_{1:T} = \{x_t, u_t\}_{1:T}$$

$$\operatorname{argmin}_{\tau_{1:T}} \sum_{t=1}^T \frac{1}{2} \tau_t^\top C_t \tau_t + c_t^\top \tau_t$$

$$\text{subject to } x_{t+1} = F_t \tau_t + f_t, x_1 = x_{\text{init}}.$$

Differentiable LQR

- Now add linear equality constraints on the dynamics and initialization.

$$\tau_{1:T} = \{x_t, u_t\}_{1:T}$$

- Chain rule: $\frac{\partial \ell}{\partial \theta} = \frac{\partial \ell}{\partial \tau_{1:T}^*} \frac{\partial \tau_{1:T}^*}{\partial \theta}$.

$$\operatorname{argmin}_{\tau_{1:T}} \sum_{t=1}^T \frac{1}{2} \tau_t^\top C_t \tau_t + c_t^\top \tau_t$$

$$\text{subject to } x_{t+1} = F_t \tau_t + f_t, x_1 = x_{\text{init}}.$$

Differentiable LQR

- Now add linear equality constraints on the dynamics and initialization.

$$\tau_{1:T} = \{x_t, u_t\}_{1:T}$$

- Chain rule: $\frac{\partial \ell}{\partial \theta} = \frac{\partial \ell}{\partial \tau_{1:T}^*} \frac{\partial \tau_{1:T}^*}{\partial \theta}$.

$$\operatorname{argmin}_{\tau_{1:T}} \sum_{t=1}^T \frac{1}{2} \tau_t^\top C_t \tau_t + c_t^\top \tau_t$$

subject to $x_{t+1} = F_t \tau_t + f_t, x_1 = x_{\text{init}}$.

General QP:

$$\mathbf{x}^* = \operatorname{argmin} \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \mathbf{c}^\top \mathbf{x}$$

subject to $A\mathbf{x} = \mathbf{b}$.

Differentiable LQR

- Now add linear equality constraints on the dynamics and initialization.

$$\tau_{1:T} = \{x_t, u_t\}_{1:T}$$

- Chain rule: $\frac{\partial \ell}{\partial \theta} = \frac{\partial \ell}{\partial \tau_{1:T}^*} \frac{\partial \tau_{1:T}^*}{\partial \theta}$.
- KKT:

$$\operatorname{argmin}_{\tau_{1:T}} \sum_{t=1}^T \frac{1}{2} \tau_t^\top C_t \tau_t + c_t^\top \tau_t$$

subject to $x_{t+1} = F_t \tau_t + f_t, x_1 = x_{\text{init}}$.

General QP:

$$\mathbf{x}^* = \operatorname{argmin} \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \mathbf{c}^\top \mathbf{x}$$

subject to $A \mathbf{x} = \mathbf{b}$.

Differentiable LQR

- Now add linear equality constraints on the dynamics and initialization.

$$\tau_{1:T} = \{x_t, u_t\}_{1:T}$$

- Chain rule: $\frac{\partial \ell}{\partial \theta} = \frac{\partial \ell}{\partial \tau_{1:T}^*} \frac{\partial \tau_{1:T}^*}{\partial \theta}$.
- KKT:

$$\begin{bmatrix} Q & A^\top \\ A & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}^* \\ \boldsymbol{\lambda}^* \end{bmatrix} = \begin{bmatrix} -\mathbf{c} \\ \mathbf{b} \end{bmatrix}$$

$$K \begin{bmatrix} \mathbf{x}^* \\ \boldsymbol{\lambda}^* \end{bmatrix} = \mathbf{v}.$$

$$\operatorname{argmin}_{\tau_{1:T}} \sum_{t=1}^T \frac{1}{2} \tau_t^\top C_t \tau_t + c_t^\top \tau_t$$

subject to $x_{t+1} = F_t \tau_t + f_t, x_1 = x_{\text{init}}$.

General QP:

$$\mathbf{x}^* = \operatorname{argmin} \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \mathbf{c}^\top \mathbf{x}$$

subject to $A\mathbf{x} = \mathbf{b}$.

In classic LQR solver, the Riccati recursion solves this linear system.

Differentiable LQR

- Apply differentiation $\frac{d}{d\theta} \left(K \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} \right) = \frac{dv}{d\theta}. \quad \frac{dK}{d\theta} \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} + K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\lambda^*}{d\theta} \end{bmatrix} = \frac{dv}{d\theta}$

Differentiable LQR

• Apply differentiation $\frac{d}{d\theta} \left(K \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} \right) = \frac{dv}{d\theta} \cdot \frac{dK}{d\theta} \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} + K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\lambda^*}{d\theta} \end{bmatrix} = \frac{dv}{d\theta}$

$$K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\lambda^*}{d\theta} \end{bmatrix} = K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\mathbf{c}} & \frac{d\mathbf{x}^*}{d\mathbf{b}} & \frac{d\mathbf{x}^*}{dQ} & \frac{d\mathbf{x}^*}{dA} \\ \frac{d\lambda^*}{d\mathbf{c}} & \frac{d\lambda^*}{d\mathbf{b}} & \frac{d\lambda^*}{dQ} & \frac{d\lambda^*}{dA} \end{bmatrix} = \begin{bmatrix} -I & 0 & -x^* & -\lambda^* \\ 0 & I & 0 & -x^* \end{bmatrix}$$

Differentiable LQR

• Apply differentiation $\frac{d}{d\theta} \left(K \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} \right) = \frac{dv}{d\theta} \cdot \frac{dK}{d\theta} \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} + K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\lambda^*}{d\theta} \end{bmatrix} = \frac{dv}{d\theta}$

$$K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\lambda^*}{d\theta} \end{bmatrix} = K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\mathbf{c}} & \frac{d\mathbf{x}^*}{d\mathbf{b}} & \frac{d\mathbf{x}^*}{dQ} & \frac{d\mathbf{x}^*}{dA} \\ \frac{d\lambda^*}{d\mathbf{c}} & \frac{d\lambda^*}{d\mathbf{b}} & \frac{d\lambda^*}{dQ} & \frac{d\lambda^*}{dA} \end{bmatrix} = \begin{bmatrix} -I & 0 & -\mathbf{x}^* & -\lambda^* \\ 0 & I & 0 & -\mathbf{x}^* \end{bmatrix}$$

$$K \frac{\partial \ell}{\partial \mathbf{z}^*} \begin{bmatrix} \frac{d\mathbf{x}^*}{d\mathbf{c}} & \frac{d\mathbf{x}^*}{d\mathbf{b}} \\ \frac{d\lambda^*}{d\mathbf{c}} & \frac{d\lambda^*}{d\mathbf{b}} \end{bmatrix} = \begin{bmatrix} -\frac{\partial \ell}{\partial \mathbf{x}^*} \\ \mathbf{0} \end{bmatrix} \quad K \mathbf{d}^* = \begin{bmatrix} -\frac{\partial \ell}{\partial \mathbf{x}^*} \\ \mathbf{0} \end{bmatrix}$$

Differentiable LQR

- Apply differentiation $\frac{d}{d\theta} \left(K \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} \right) = \frac{dv}{d\theta}$. $\frac{dK}{d\theta} \begin{bmatrix} \mathbf{x}^* \\ \lambda^* \end{bmatrix} + K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\lambda^*}{d\theta} \end{bmatrix} = \frac{dv}{d\theta}$

$$K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\lambda^*}{d\theta} \end{bmatrix} = K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\mathbf{c}} & \frac{d\mathbf{x}^*}{d\mathbf{b}} & \frac{d\mathbf{x}^*}{dQ} & \frac{d\mathbf{x}^*}{dA} \\ \frac{d\lambda^*}{d\mathbf{c}} & \frac{d\lambda^*}{d\mathbf{b}} & \frac{d\lambda^*}{dQ} & \frac{d\lambda^*}{dA} \end{bmatrix} = \begin{bmatrix} -I & 0 & -\mathbf{x}^* & -\lambda^* \\ 0 & I & 0 & -\mathbf{x}^* \end{bmatrix}$$

$$K \frac{\partial \ell}{\partial \mathbf{z}^*} \begin{bmatrix} \frac{d\mathbf{x}^*}{d\mathbf{c}} & \frac{d\mathbf{x}^*}{d\mathbf{b}} \\ \frac{d\lambda^*}{d\mathbf{c}} & \frac{d\lambda^*}{d\mathbf{b}} \end{bmatrix} = \begin{bmatrix} -\frac{\partial \ell}{\partial \mathbf{x}^*} \\ \mathbf{0} \end{bmatrix} \quad K \mathbf{d}^* = \begin{bmatrix} -\frac{\partial \ell}{\partial \mathbf{x}^*} \\ \mathbf{0} \end{bmatrix}$$

- Equivalent QP:

$$\mathbf{d}^* = \underset{\mathbf{d}}{\operatorname{argmin}} \frac{1}{2} \mathbf{d}^\top Q \mathbf{d} + \frac{\partial \ell}{\partial \mathbf{x}^*}^\top \mathbf{d},$$

subject to $A\mathbf{d} = \mathbf{0}$.

Differentiable LQR

- Apply differentiation $\frac{d}{d\theta} \left(K \begin{bmatrix} \mathbf{x}^* \\ \boldsymbol{\lambda}^* \end{bmatrix} \right) = \frac{dv}{d\theta}$. $\frac{dK}{d\theta} \begin{bmatrix} \mathbf{x}^* \\ \boldsymbol{\lambda}^* \end{bmatrix} + K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\boldsymbol{\lambda}^*}{d\theta} \end{bmatrix} = \frac{dv}{d\theta}$

$$K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\theta} \\ \frac{d\boldsymbol{\lambda}^*}{d\theta} \end{bmatrix} = K \begin{bmatrix} \frac{d\mathbf{x}^*}{d\mathbf{c}} & \frac{d\mathbf{x}^*}{d\mathbf{b}} & \frac{d\mathbf{x}^*}{dQ} & \frac{d\mathbf{x}^*}{dA} \\ \frac{d\boldsymbol{\lambda}^*}{d\mathbf{c}} & \frac{d\boldsymbol{\lambda}^*}{d\mathbf{b}} & \frac{d\boldsymbol{\lambda}^*}{dQ} & \frac{d\boldsymbol{\lambda}^*}{dA} \end{bmatrix} = \begin{bmatrix} -I & 0 & -\mathbf{x}^* & -\boldsymbol{\lambda}^* \\ 0 & I & 0 & -\mathbf{x}^* \end{bmatrix}$$

$$K \frac{\partial \ell}{\partial \mathbf{z}^*} \begin{bmatrix} \frac{d\mathbf{x}^*}{d\mathbf{c}} & \frac{d\mathbf{x}^*}{d\mathbf{b}} \\ \frac{d\boldsymbol{\lambda}^*}{d\mathbf{c}} & \frac{d\boldsymbol{\lambda}^*}{d\mathbf{b}} \end{bmatrix} = \begin{bmatrix} -\frac{\partial \ell}{\partial \mathbf{x}^*} \\ \mathbf{0} \end{bmatrix} \quad K \mathbf{d}^* = \begin{bmatrix} -\frac{\partial \ell}{\partial \mathbf{x}^*} \\ \mathbf{0} \end{bmatrix}$$

- Equivalent QP:

$$\mathbf{d}^* = \underset{\mathbf{d}}{\operatorname{argmin}} \frac{1}{2} \mathbf{d}^\top Q \mathbf{d} + \frac{\partial \ell}{\partial \mathbf{x}^*}{}^\top \mathbf{d}, \quad \frac{\partial \ell}{\partial Q} = \frac{1}{2} (\mathbf{d}_x^* \otimes \mathbf{x}^* + \mathbf{x}^* \otimes \mathbf{d}_x^*)$$

subject to $A\mathbf{d} = \mathbf{0}$.

$$\frac{\partial \ell}{\partial A} = \mathbf{d}_\lambda^* \otimes \mathbf{x}^* + \boldsymbol{\lambda}^* \otimes \mathbf{d}_x^*.$$

Differentiable LQR

- The backward pass can also be formulated as a LQR problem.
- Swap c to $\nabla_{\tau^*} \ell$ and f to 0.

Module 1 Differentiable LQR

(The LQR algorithm is defined in [appendix A](#))

Input: Initial state x_{init}

Parameters: $\theta = \{C, c, F, f\}$

Forward Pass:

- 1: $\tau_{1:T}^* = \text{LQR}_T(x_{\text{init}}; C, c, F, f)$
- 2: Compute $\lambda_{1:T}^*$ with (7)

▷ Solve (2)

Backward Pass:

- 1: $d_{\tau_{1:T}}^* = \text{LQR}_T(0; C, \nabla_{\tau^*} \ell, F, 0)$ ▷ Solve (9), ideally reusing the factorizations from the forward pass
- 2: Compute $d_{\lambda_{1:T}}^*$ with (7)
- 3: Compute the derivatives of ℓ with respect to C, c, F, f , and x_{init} with (8)

Differentiable MPC

- What about general MPC?

$$\operatorname{argmin}_{x_{1:T} \in \mathcal{X}, u_{1:T} \in \mathcal{U}} \sum_{t=1}^T C_t(x_t, u_t)$$

subject to $x_{t+1} = f(x_t, u_t), x_1 = x_{\text{init}}$.

Differentiable MPC

- What about general MPC?

$$\operatorname{argmin}_{x_{1:T} \in \mathcal{X}, u_{1:T} \in \mathcal{U}} \sum_{t=1}^T C_t(x_t, u_t)$$

subject to $x_{t+1} = f(x_t, u_t), x_1 = x_{\text{init}}$.

- Use Taylor expansion to approximate.

$$\tilde{C}_{\theta,t}^i = C_{\theta,t}(\tau_t^i) + p_t^{i\top} (\tau_t - \tau_t^i) + \frac{1}{2} (\tau_t - \tau_t^i)^\top H_t^i (\tau_t - \tau_t^i).$$

Differentiable MPC

- What about general MPC?

$$\operatorname{argmin}_{x_{1:T} \in \mathcal{X}, u_{1:T} \in \mathcal{U}} \sum_{t=1}^T C_t(x_t, u_t)$$

subject to $x_{t+1} = f(x_t, u_t), x_1 = x_{\text{init}}$.

- Use Taylor expansion to approximate.

$$\tilde{C}_{\theta,t}^i = C_{\theta,t}(\tau_t^i) + p_t^{i\top} (\tau_t - \tau_t^i) + \frac{1}{2} (\tau_t - \tau_t^i)^\top H_t^i (\tau_t - \tau_t^i).$$

- Fixed point iteration

$$\tau^{i+1} = \operatorname{argmin}_{\tau} \sum_t^T \tilde{C}_t^i(\tau_t^i).$$

Differentiable MPC

- What about general MPC?

$$\operatorname{argmin}_{x_{1:T} \in \mathcal{X}, u_{1:T} \in \mathcal{U}} \sum_{t=1}^T C_t(x_t, u_t)$$

subject to $x_{t+1} = f(x_t, u_t), x_1 = x_{\text{init}}$.

- Use Taylor expansion to approximate.

$$\tilde{C}_{\theta,t}^i = C_{\theta,t}(\tau_t^i) + p_t^{i\top} (\tau_t - \tau_t^i) + \frac{1}{2} (\tau_t - \tau_t^i)^\top H_t^i (\tau_t - \tau_t^i).$$

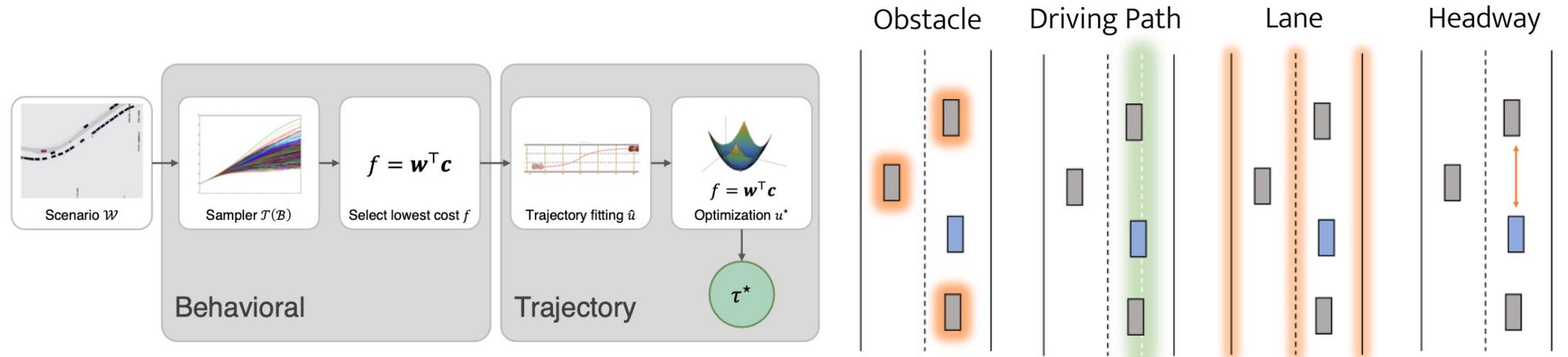
- Fixed point iteration

$$\tau^{i+1} = \operatorname{argmin}_{\tau} \sum_t^T \tilde{C}_t^i(\tau_t^i).$$

- Backward only depends on final quadratic approximation.

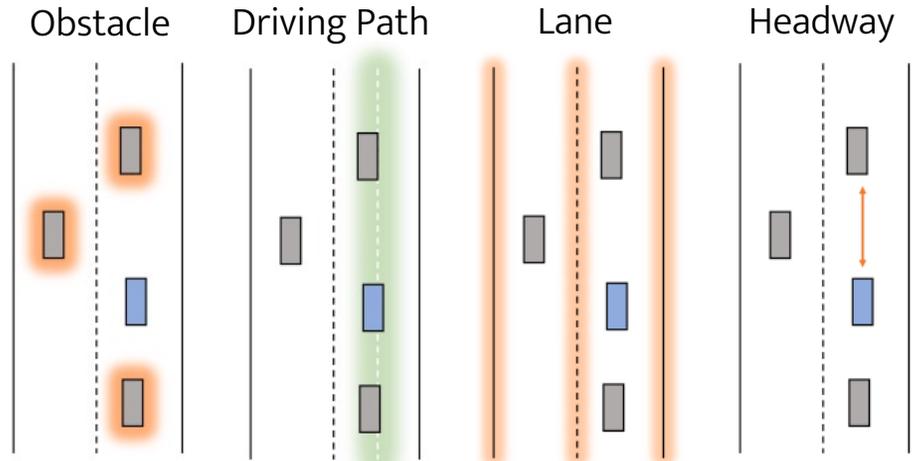
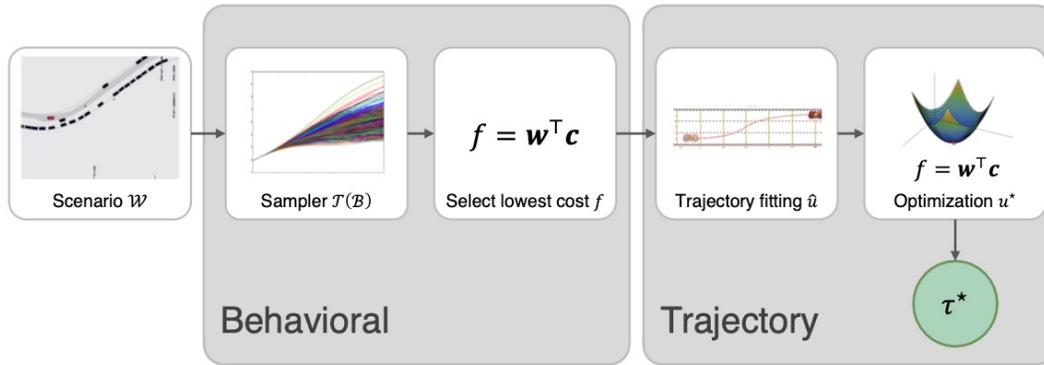
Behavioral vs. Trajectory Planning

- Gradient-based optimization provides a locally optimized trajectory.



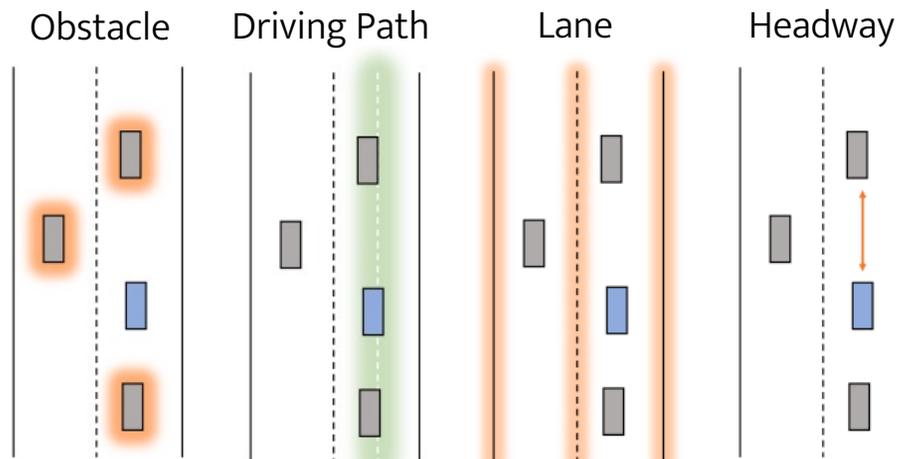
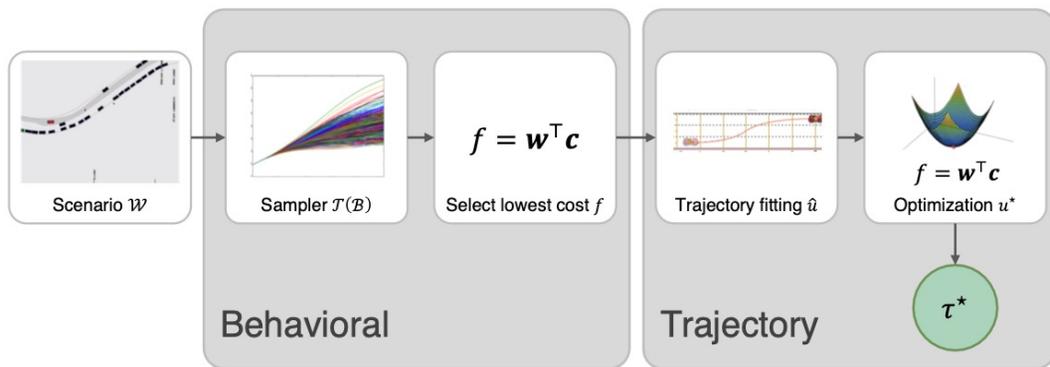
Behavioral vs. Trajectory Planning

- Gradient-based optimization provides a locally optimized trajectory.
- Samples may be needed for reasoning global structure.



Behavioral vs. Trajectory Planning

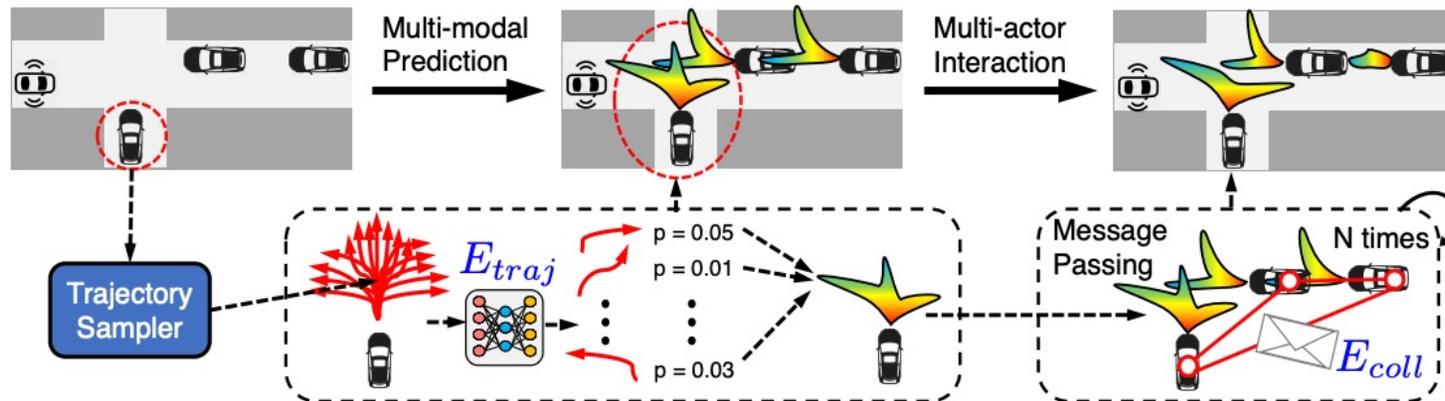
- Gradient-based optimization provides a locally optimized trajectory.
- Samples may be needed for reasoning global structure.
- Can learn together using the same learned costs.



Planning with Social Reasoning

- Jointly reason the future trajectories of multiple agents as an energy-based graphical model.

$$p(\mathbf{s}_1, \dots, \mathbf{s}_N \mid \mathbf{X}) = \frac{1}{Z} \exp(-E_{\theta}(\mathbf{s}_1, \dots, \mathbf{s}_N) \mid \mathbf{X})$$



Planning with Social Reasoning

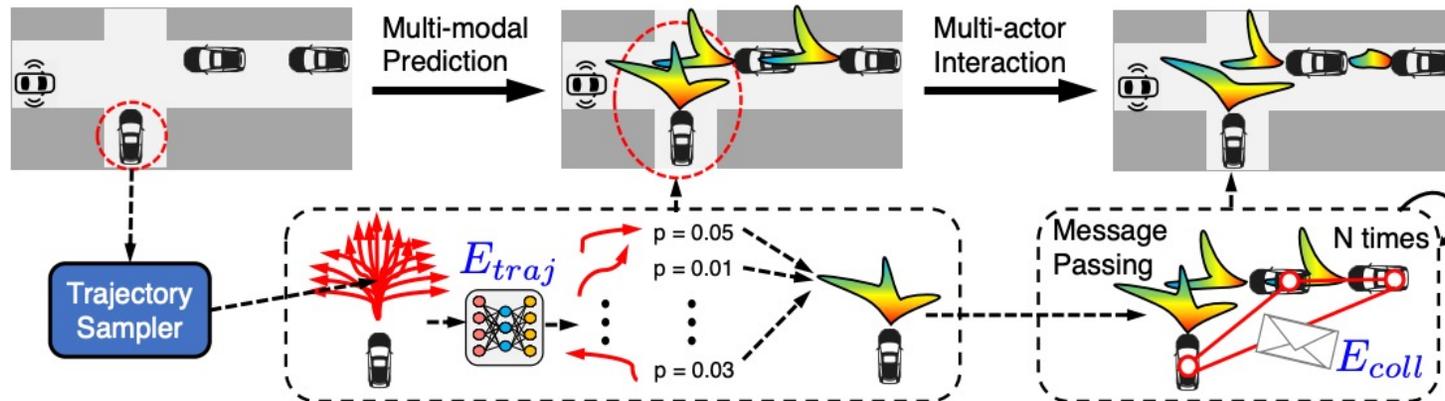
- Jointly reason the future trajectories of multiple agents as an energy-based graphical model.

$$p(\mathbf{s}_1, \dots, \mathbf{s}_N | \mathbf{X}) = \frac{1}{Z} \exp(-E_\theta(\mathbf{s}_1, \dots, \mathbf{s}_N) | \mathbf{X})$$

- Trajectory Goodness + Collision.

$$\sum_i E_\theta(s_i | X) + \sum_{i \neq j} E(s_i, s_j)$$

$$E(\mathbf{s}_i, \mathbf{s}_j) = \gamma \quad \text{if } \mathbf{s}_i \text{ collides } \mathbf{s}_j$$



Planning with Social Reasoning

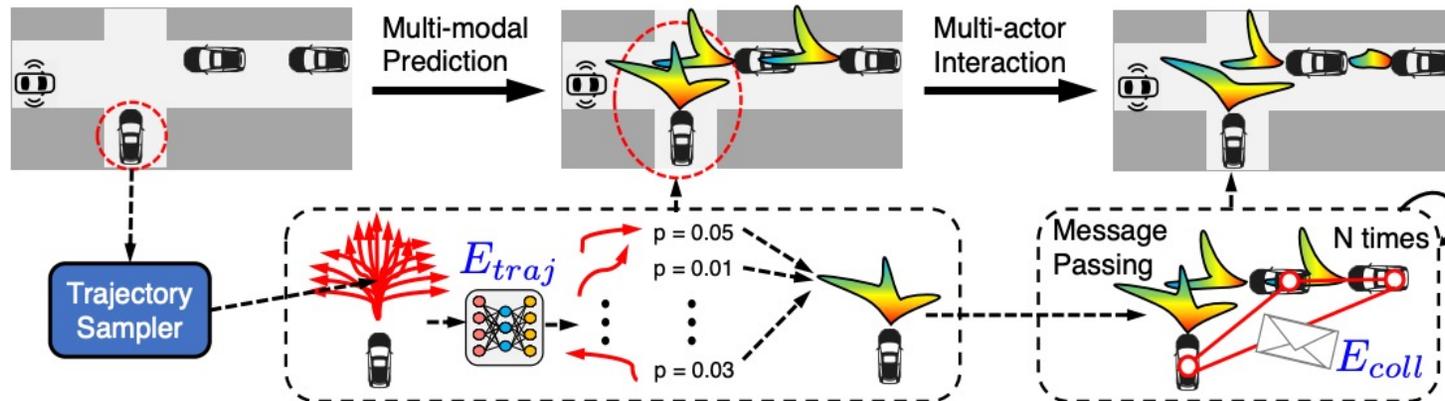
- Jointly reason the future trajectories of multiple agents as an energy-based graphical model.

$$p(\mathbf{s}_1, \dots, \mathbf{s}_N \mid \mathbf{X}) = \frac{1}{Z} \exp(-E_\theta(\mathbf{s}_1, \dots, \mathbf{s}_N) \mid \mathbf{X})$$

- Trajectory Goodness + Collision.
- Batch of trajectory samples.

$$\sum_i E_\theta(s_i \mid X) + \sum_{i \neq j} E(s_i, s_j)$$

$$E(\mathbf{s}_i, \mathbf{s}_j) = \gamma \quad \text{if } \mathbf{s}_i \text{ collides } \mathbf{s}_j$$



Planning with Social Reasoning

- Jointly reason the future trajectories of multiple agents as an energy-based graphical model.

$$p(\mathbf{s}_1, \dots, \mathbf{s}_N | \mathbf{X}) = \frac{1}{Z} \exp(-E_\theta(\mathbf{s}_1, \dots, \mathbf{s}_N) | \mathbf{X})$$

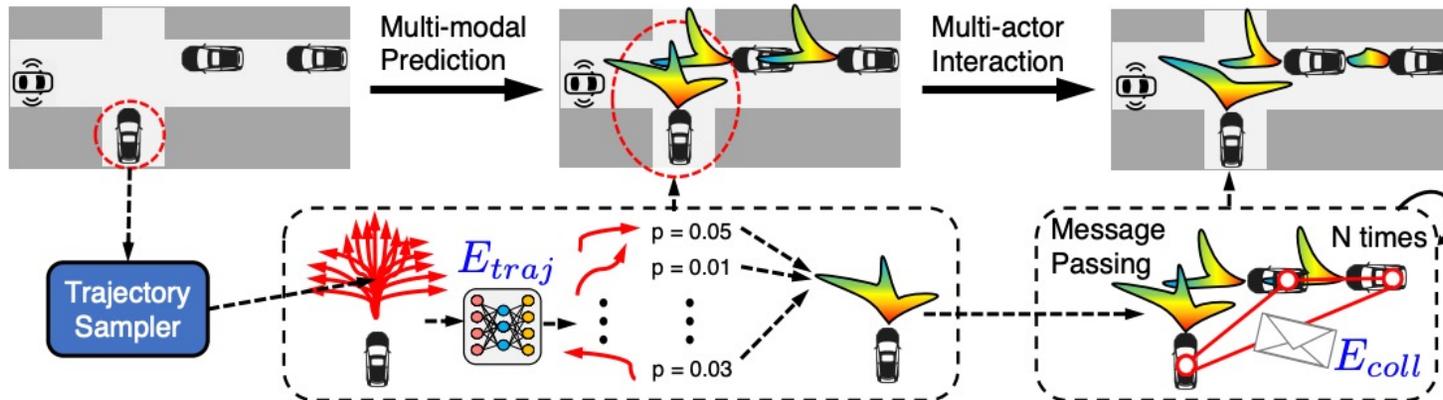
- Trajectory Goodness + Collision.

$$\sum_i E_\theta(s_i | X) + \sum_{i \neq j} E(s_i, s_j)$$

- Batch of trajectory samples.

$$E(\mathbf{s}_i, \mathbf{s}_j) = \gamma \text{ if } \mathbf{s}_i \text{ collides } \mathbf{s}_j$$

- Classification of groundtruth trajectory.



Summary: End-to-End Planning and Control

- Direct Policy Prediction
 - Condition perception features into the model
 - Use of diffusion models

Summary: End-to-End Planning and Control

- Direct Policy Prediction
 - Condition perception features into the model
 - Use of diffusion models
- Cost Learning (IRL) from Experts
 - Max-margin, max-entropy/EBM
 - Need negative samples
 - Can be combined with efficient external samplers
 - Cost volume prediction: parametric + non-parametric

Summary: End-to-End Planning and Control

- Direct Policy Prediction
 - Condition perception features into the model
 - Use of diffusion models
- Cost Learning (IRL) from Experts
 - Max-margin, max-entropy/EBM
 - Need negative samples
 - Can be combined with efficient external samplers
 - Cost volume prediction: parametric + non-parametric
- Differentiable Planner
 - Backprop through local optimization, learnable cost
 - Can be memory efficient, implicit differentiation