

DS-GA.3001

Embodied Learning and Vision

Mengye Ren

NYU

Spring 2025

embodied-learning-vision-course.github.io



Lecture Slides for Note Taking



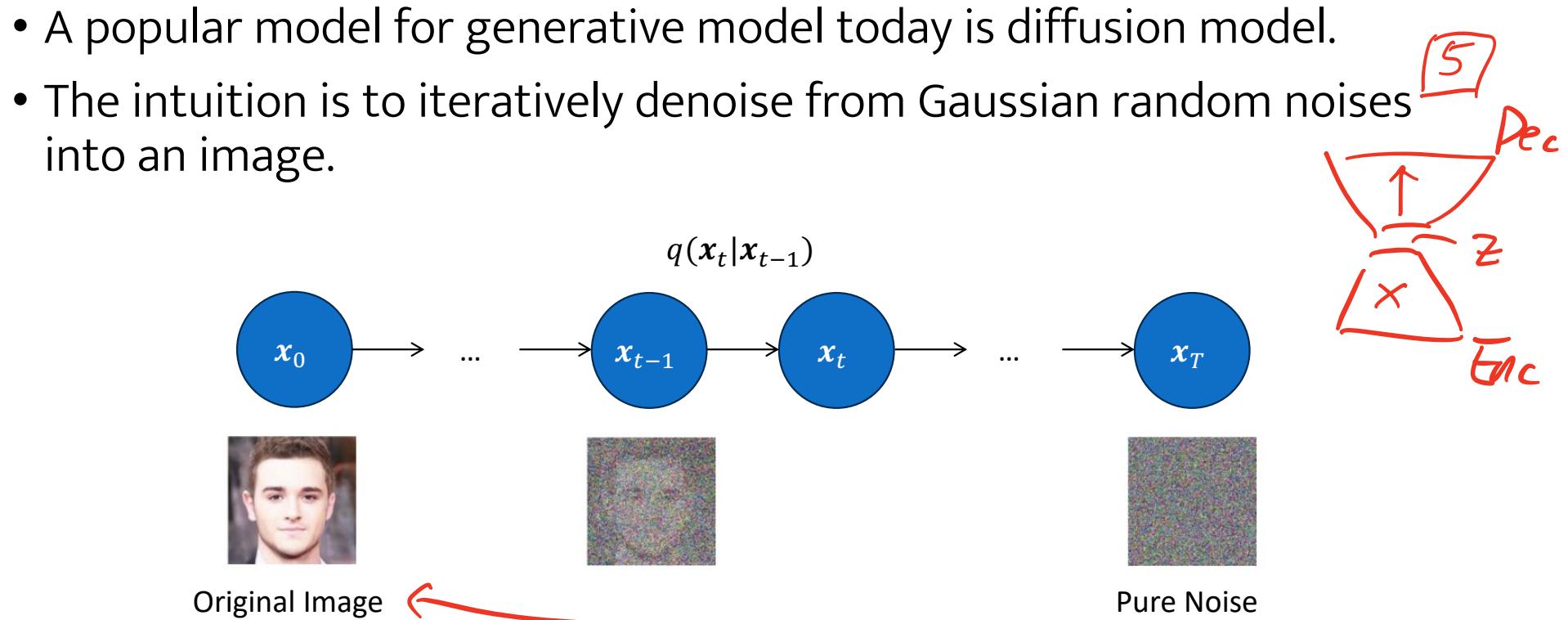
Diffusion Models

- A popular model for generative model today is diffusion model.

Diffusion Models

- A popular model for generative model today is diffusion model.
- The intuition is to iteratively denoise from Gaussian random noises into an image.

Diffusion Models



Diffusion Models

$$\sigma^2 I$$

μ



$\beta_t I$

~~prev~~

- Forward process: $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$.

\equiv

$\underbrace{\sqrt{1 - \beta_t}x_{t-1}}$

~~prev~~

$$q(x_t|x_{t-1})$$

x_0



Original Image

x_{t-1}



x_t

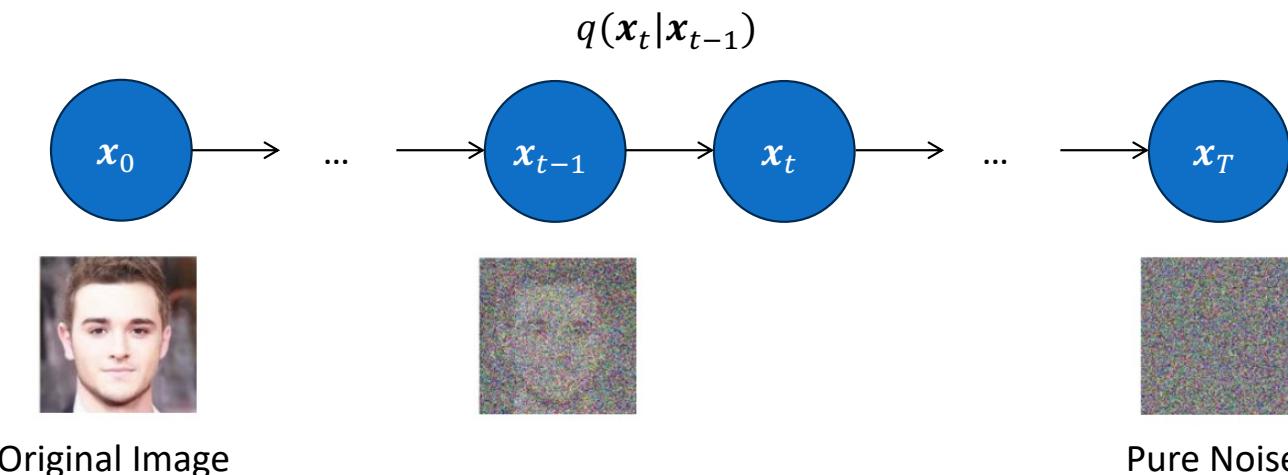
x_T



Pure Noise

Diffusion Models

- Forward process: $q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$.
 - You can also write: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t, \epsilon_t \sim \mathcal{N}(0, I)$.



Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t} \underline{x_{t-1}} + \sqrt{\beta_t} \underline{\epsilon_t}$

Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$
- Write \underline{x}_t as a function of \underline{x}_0 with larger noises:

$$q(\underline{x}_t | \underline{x}_0) = \mathcal{N}(\underline{x}_t; \sqrt{\bar{\alpha}_t}\underline{x}_0, (1 - \bar{\alpha}_t)I).$$

Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$

- Write x_t as a function of x_0 with larger noises:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, \underbrace{(1 - \bar{\alpha}_t)I}_{\text{red arrow}}).$$

- Cumulative schedule: $\alpha_t = 1 - \beta_t$. $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$.

Properties of the Forward Process

- Forward process: $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t$

- Write x_t as a function of x_0 with larger noises:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I).$$

- Cumulative schedule: $\alpha_t = 1 - \beta_t$. $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$.

- $x_t = \underbrace{\sqrt{\bar{\alpha}_t}x_0}_{\text{red underline}} + \underbrace{\sqrt{1 - \bar{\alpha}_t}\epsilon}_{\text{red underline}}$.

Cumulative Schedule

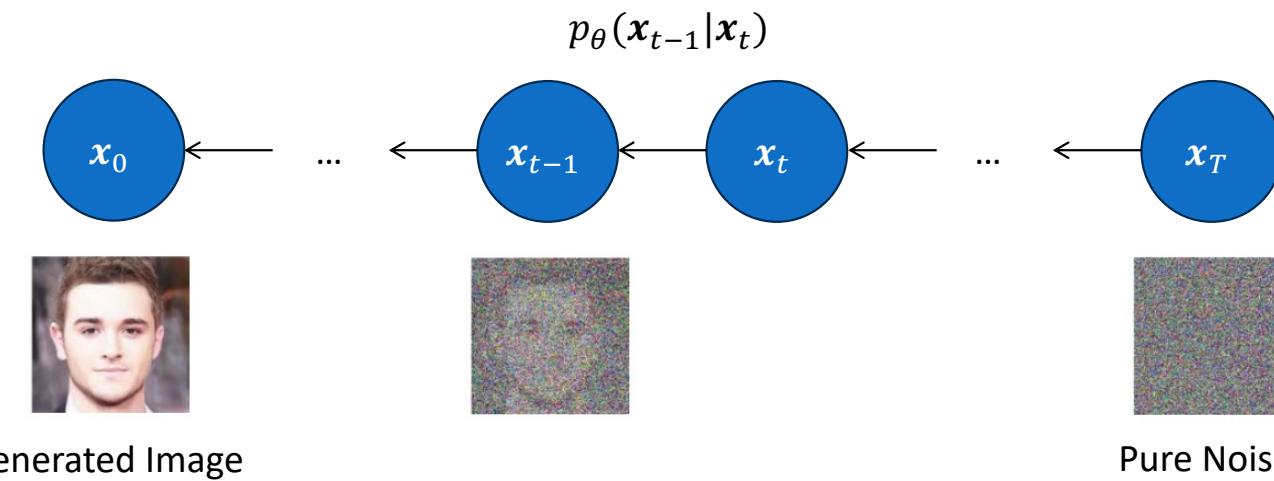
$$\alpha_t = 1 - \beta_t.$$

- Show it's true for x_2 : $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$.

$$\begin{aligned}x_2 &= \sqrt{1 - \beta_2}x_1 + \sqrt{\beta_2}\epsilon_2 = \sqrt{1 - \beta_2}\sqrt{1 - \beta_1}x_0 + \sqrt{\beta_2}\epsilon_2 + \sqrt{1 - \beta_2}\sqrt{\beta_1}\epsilon_1 \\&= \alpha_1\alpha_2x_0 + \sqrt{(1 - \beta_2)\beta_1 + \beta_2}\epsilon \\&= \bar{\alpha}_2x_0 + \sqrt{1 - (1 - \beta_1)(1 - \beta_2)}\epsilon \\&= \bar{\alpha}_2x_0 + \sqrt{1 - \bar{\alpha}_2}\epsilon.\end{aligned}$$

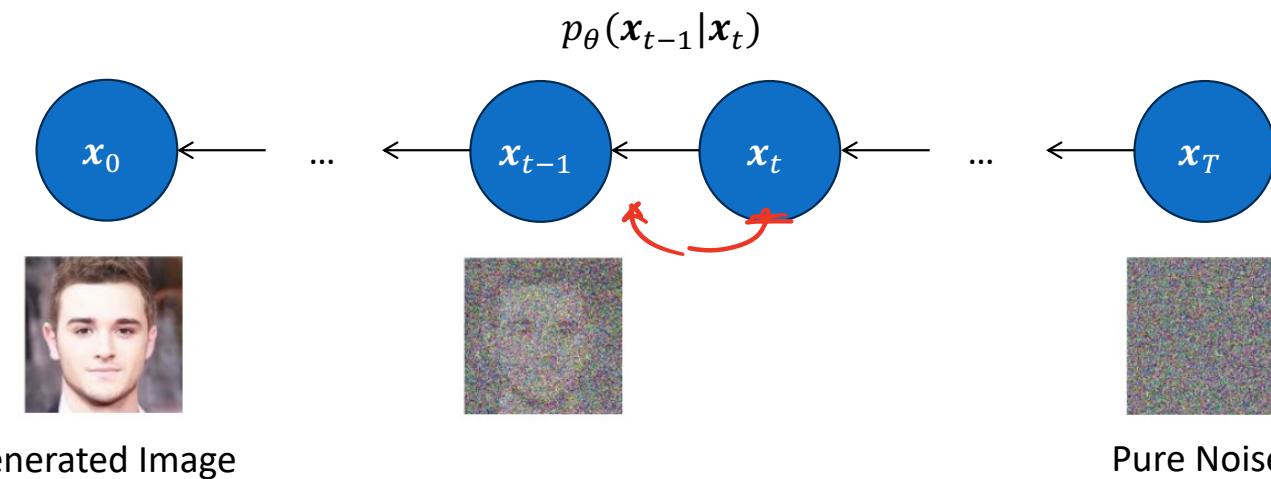
Reverse Process

- A generative model wants to predict x_0 from $\underline{x_T}$.



Reverse Process

- A generative model wants to predict x_0 from x_T .
- The reverse process transition is also Gaussian distributed. But we don't know what the transition will be like just by looking at the noisy image!

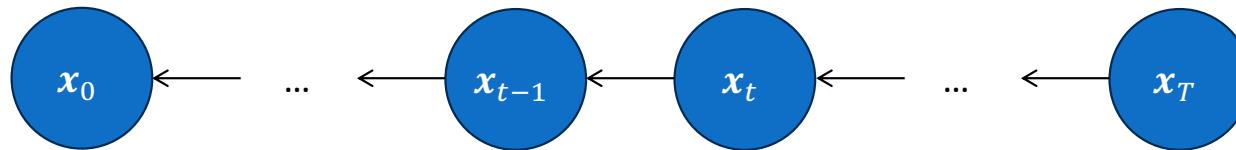


Reverse Process

- So, we need to learn a “model”:

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)).$$

$$p_{\theta}(x_{t-1}|x_t)$$



Generated Image



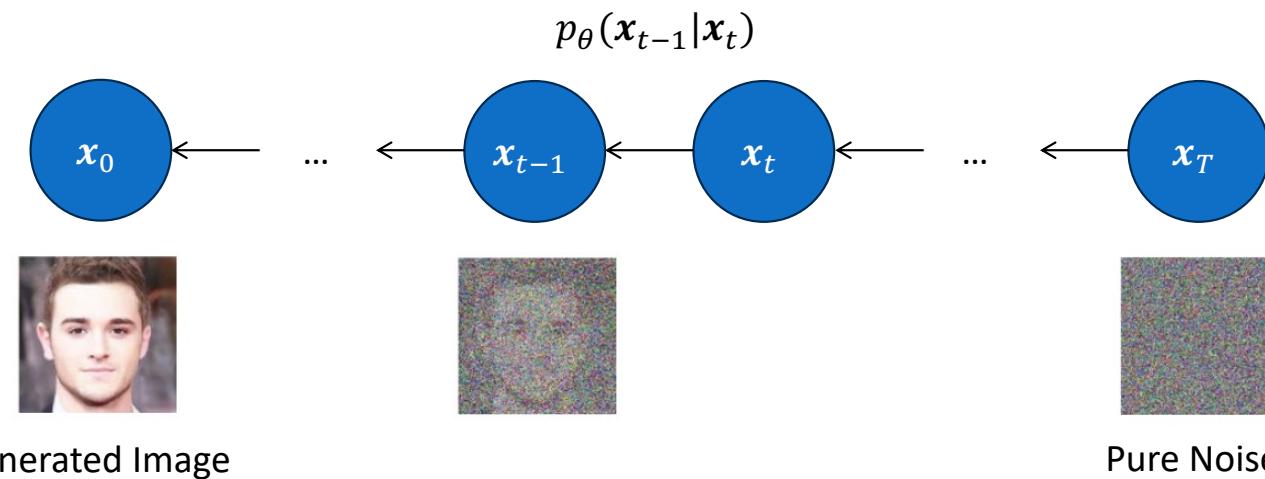
Pure Noise

Reverse Process

- So, we need to learn a “model”:

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)).$$

- μ_{θ} is the denoising vector.



Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.

Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.
- Bayes rule:
$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}.$$

Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.
- Bayes rule:
$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}.$$
- But we don't know the marginal distribution $q(x_{t-1})$. We only know q_T and $\underline{q(x_t|x_{t-1})}$.

Reverse Process

- Compute μ_θ ? Derive $p(x_{t-1}|x_t)$.
- Bayes rule:
$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}.$$
- But we don't know the marginal distribution $q(x_{t-1})$. We only know q_T and $q(x_t|x_{t-1})$.
- Solution: Condition on the original input x_0 :

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

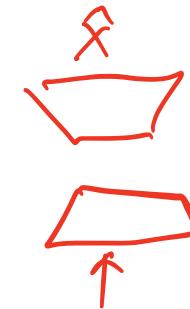
Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}I).$$



Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}. \quad X + \text{noise}$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t, \tilde{\beta} I).$$

$$\tilde{\mu}_t = \frac{\sqrt{\alpha_t} \bar{\beta}_{t-1}}{\bar{\beta}_t} x_t + \frac{\sqrt{\alpha_{t-1}} \beta_t}{\bar{\beta}_t} x_0 = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon \right).$$

training
 target

interpolation between x_t & x_0

Reverse Process

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1})q(x_{t-1}|x_0)}{q(x_t|x_0)}.$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}I).$$

$$\tilde{\mu}_t = \frac{\sqrt{\alpha_t}\bar{\beta}_{t-1}}{\bar{\beta}_t}x_t + \frac{\sqrt{\alpha_{t-1}}\beta_t}{\bar{\beta}_t}x_0 = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon\right).$$

- Want: train up a μ_θ to match with $\tilde{\mu}_t$.

 
network *target*

Training

- Sometimes it is more common to predict the denoising vector ϵ instead of μ .

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right),$$

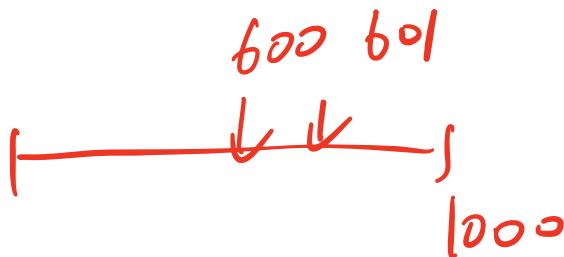
$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

Training

- Sometimes it is more common to predict the denoising vector ϵ instead of μ .
- Randomly pick at a time step and predict the difference between the noisy and the original.

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right),$$

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$



Training



- Sometimes it is more common to predict the denoising vector ϵ instead of μ .
- Randomly pick at a time step and predict the difference between the noisy and the original.

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right),$$

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

Algorithm 1 Training

- 1: **repeat**
 - 2: $\boxed{x_0} \sim q(\mathbf{x}_0)$
 - 3: $\boxed{t} \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on

$$\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$$
 - 6: **until** converged
-

network.

Sampling

- How do we sample an image?

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Sampling

- How do we sample an image?
- We know μ_θ which will help us transition from x_t to x_{t-1} .

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Sampling

- How do we sample an image?
- We know μ_θ which will help us transition from x_t to x_{t-1} .

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

- Sample from $\mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2)$. σ_t can either be β_t or $\tilde{\beta}_t$ derived from the posterior.

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

S

*sample
from posterior
Gaussian distribution*

More on Samplers

- DDPM relies on many iterations (e.g. 1000) to produce one sample.
Slower than NFs and GANs.

[Song et al. 2021]

More on Samplers

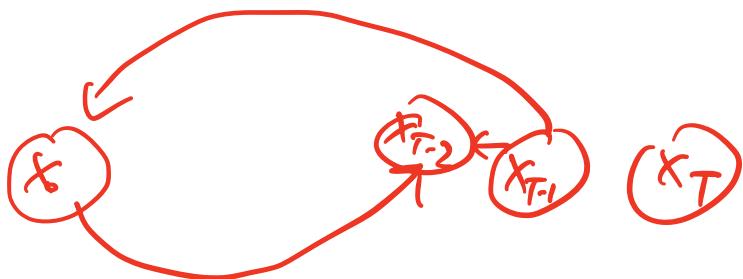
- DDPM relies on many iterations (e.g. 1000) to produce one sample.
Slower than NFs and GANs.
- In the non-Markovian model, we can first generate x_T , and based on x_T and x_0 we can generate x_{T-1} and so on.

[Song et al. 2021]

More on Samplers

- DDPM relies on many iterations (e.g. 1000) to produce one sample.
Slower than NFs and GANs.
- In the non-Markovian model, we can first generate x_T , and based on x_T and x_0 we can generate x_{T-1} and so on.
- Joint distribution:

$$q(x_{1:T}|x_0) = q(x_T|x_0) \prod_{t=2}^T q(x_{t-1}|x_t, x_0).$$



[Song et al. 2021]

More on Samplers

- DDPM relies on many iterations (e.g. 1000) to produce one sample.
Slower than NFs and GANs.
- In the non-Markovian model, we can first generate x_T , and based on x_T and x_0 we can generate x_{T-1} and so on.
- Joint distribution:

$$q(x_{1:T}|x_0) = q(x_T|x_0) \prod_{t=2}^T q(x_{t-1}|x_t, x_0).$$

- Estimate x_{t-1} based on x_0 and x_t :

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}\left(\sqrt{a_{t-1}}x_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{x_t - \sqrt{\alpha_t}x_0}{\sqrt{1-\alpha_t}}, \sigma_t^2 I\right).$$

[Song et al. 2021]

More on DDIM Samplers

- Prediction of x_0 :

$$f_\theta^{(t)}(x_t) = \frac{1}{\sqrt{\alpha_t}}(x_t - \sqrt{1 - \alpha_t} \cdot \epsilon_\theta^{(t)}(x_t)).$$

[Song et al. 2021]

More on DDIM Samplers

- Prediction of x_0 :

$$f_{\theta}^{(t)}(x_t) = \frac{1}{\sqrt{\alpha_t}}(x_t - \sqrt{1 - \alpha_t} \cdot \epsilon_{\theta}^{(t)}(x_t)).$$

- Sampling process:

$$p_{\theta}^{(t)}(x_{t-1}|x_t) = \begin{cases} \mathcal{N}(f_{\theta}^{(1)}(x_t)), \sigma_1^2 I & \text{if } t = 1 \\ q(x_{t-1}|x_t, f_{\theta}^{(t)}(x_t)) & \text{otherwise.} \end{cases}$$

\downarrow $=$ $\overbrace{x_0}$

[Song et al. 2021]

Guided Diffusion

- We can add guidance on the diffusion updates at inference time.

Classifier Guidance / External Score Model

$$x_{t-1} \leftarrow \text{sample from } \mathcal{N}(\mu + s \underbrace{\nabla_{x_t} \log p_\phi(y|x_t)}_{\hat{\epsilon}} \Sigma, \Sigma) \quad \begin{aligned} \hat{\epsilon} &\leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y|x_t) \\ x_{t-1} &\leftarrow \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon} \end{aligned}$$

ImageNet Classif.

Guided Diffusion

Inference only.

- We can add guidance on the diffusion updates at inference time.

Classifier Guidance / External Score Model

$$x_{t-1} \leftarrow \text{sample from } \mathcal{N}(\mu + s \sum \nabla_{x_t} \log p_\phi(y|x_t), \Sigma)$$

$$\hat{\epsilon} \leftarrow \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y|x_t)$$
$$x_{t-1} \leftarrow \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t} \hat{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}$$

- We also can train a conditional diffusion model.

repeat

$$(x, c) \sim p(x, c)$$

$c \leftarrow \emptyset$ with probability p_{uncond}

$$\lambda \sim p(\lambda)$$

$$\epsilon \sim \mathcal{N}(0, I)$$

$$z_\lambda = \alpha_\lambda x + \sigma_\lambda \epsilon$$

Take gradient step on $\nabla_\theta \|\epsilon_\theta(z_\lambda, c) - \epsilon\|^2$

▷ Sample data with conditioning from the dataset

▷ Randomly discard conditioning to train unconditionally

▷ Sample log SNR value

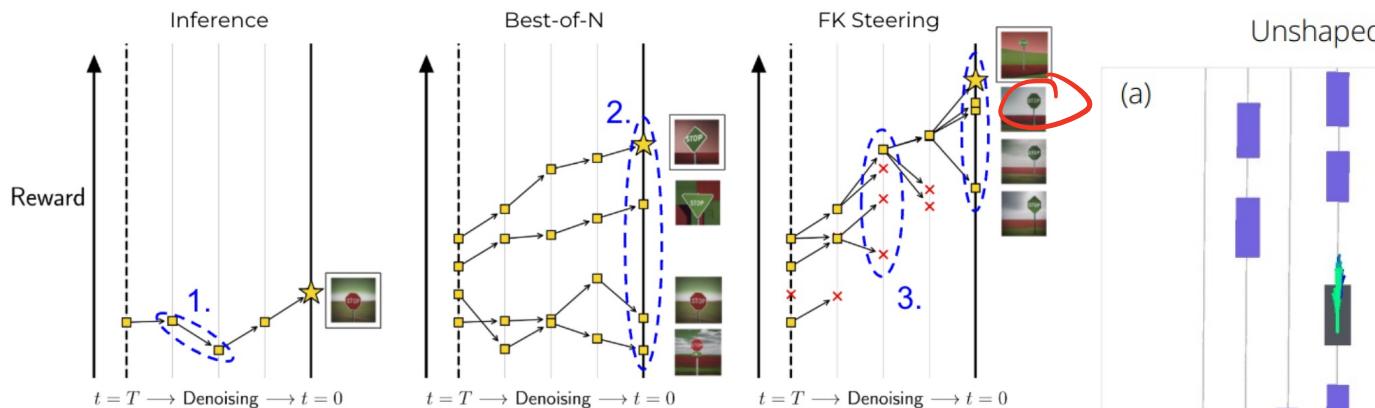
▷ Corrupt data to the sampled log SNR value

▷ Optimization of denoising model

Requires training until converged

Test-Time Adaptation

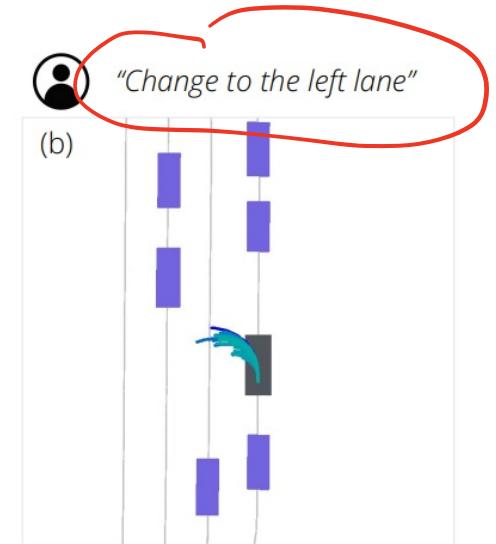
- Diffusion can be combined / guided with reward functions at test time.



Prompt: "a green stop sign in a red field"

- (1) Iteratively de-noise $x_T \rightarrow x_{T-1} \rightarrow \dots \rightarrow x_0$.
- (2) Generate multiple samples (particles).
- (3) Resample promising particles at *intermediate steps*.

[Singhal et al. 2025]



[Yang et al. 2024]

Diffusion for Detection

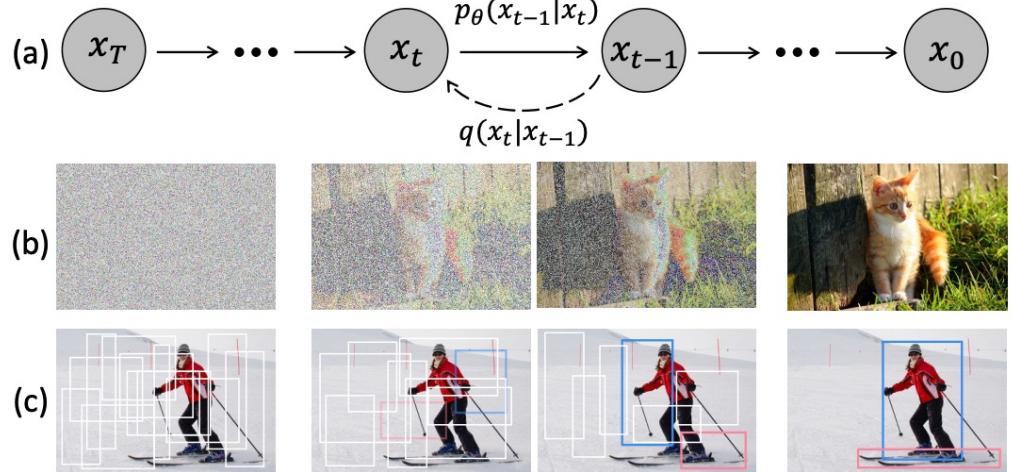
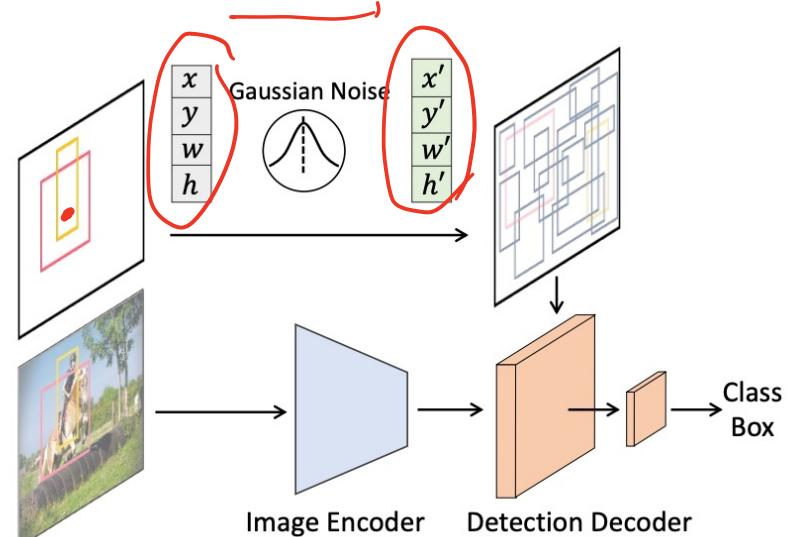
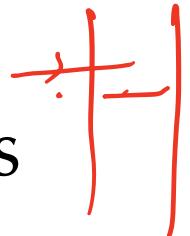


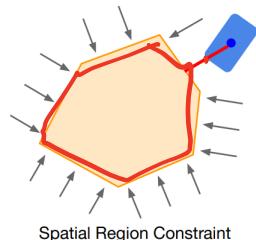
Figure 1. Diffusion model for object detection. (a) A diffusion model where q is the diffusion process and p_θ is the reverse process. (b) Diffusion model for image generation task. (c) We propose to formulate object detection as a denoising diffusion process from noisy boxes to object boxes.



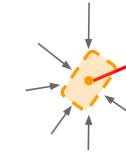
Diffusion for Generating Simulation Scenes



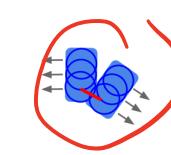
Hand Crafted



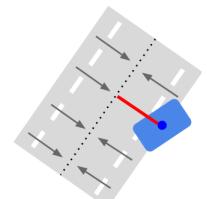
Actor Attribute Constraint



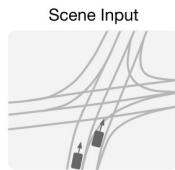
Initial Scene Constraint



Collision Constraint



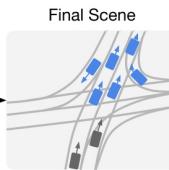
On-road Constraint



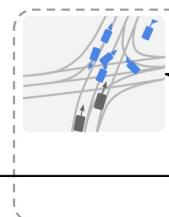
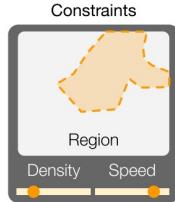
Random Initialization



T Diffusion Steps

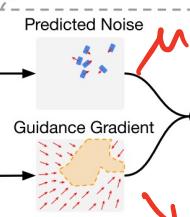


Final Scene

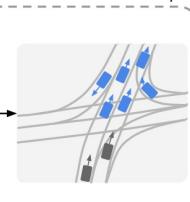


Diffusion Model

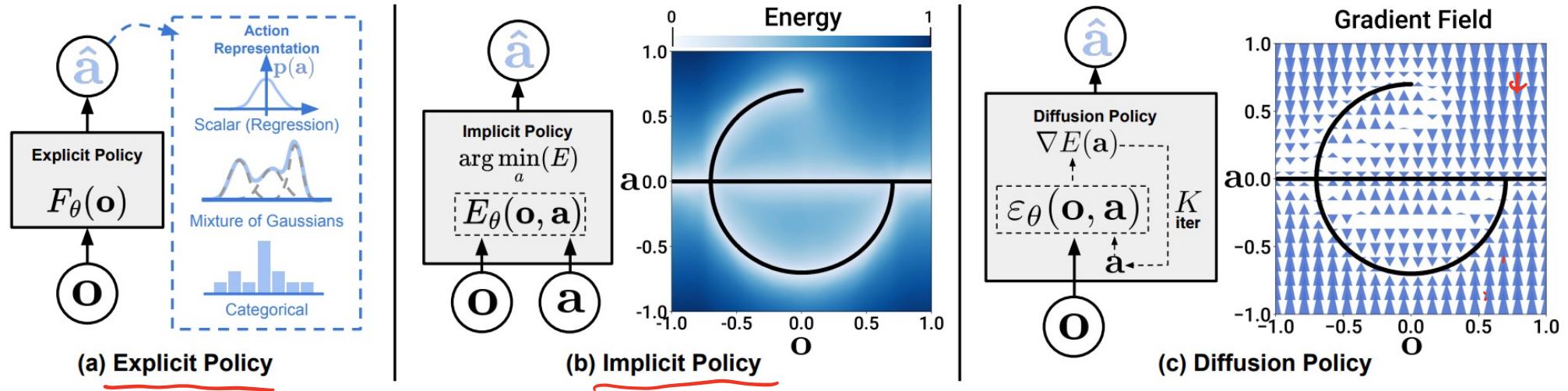
Guidance Function



guidance

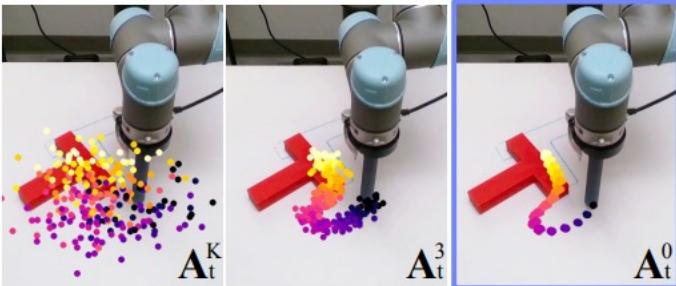
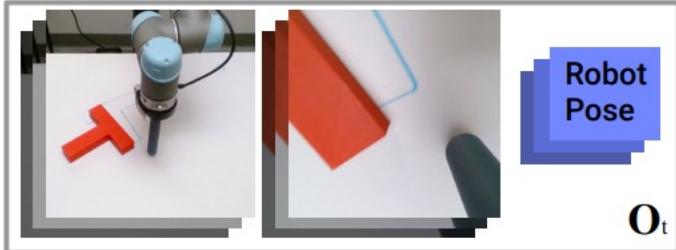


Diffusion for Planning and Control



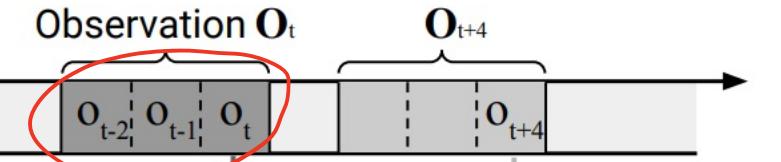
Diffusion for Planning and Control

Input: Image Observation Sequence

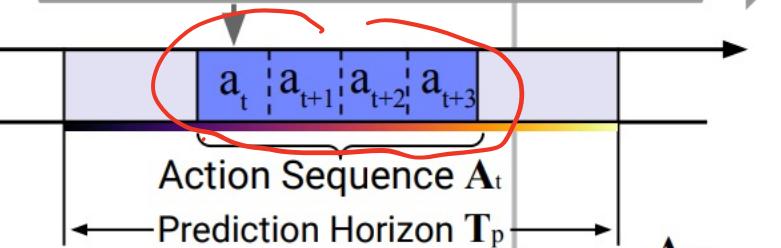


Output: Action Sequence

Observation \mathbf{O}_t

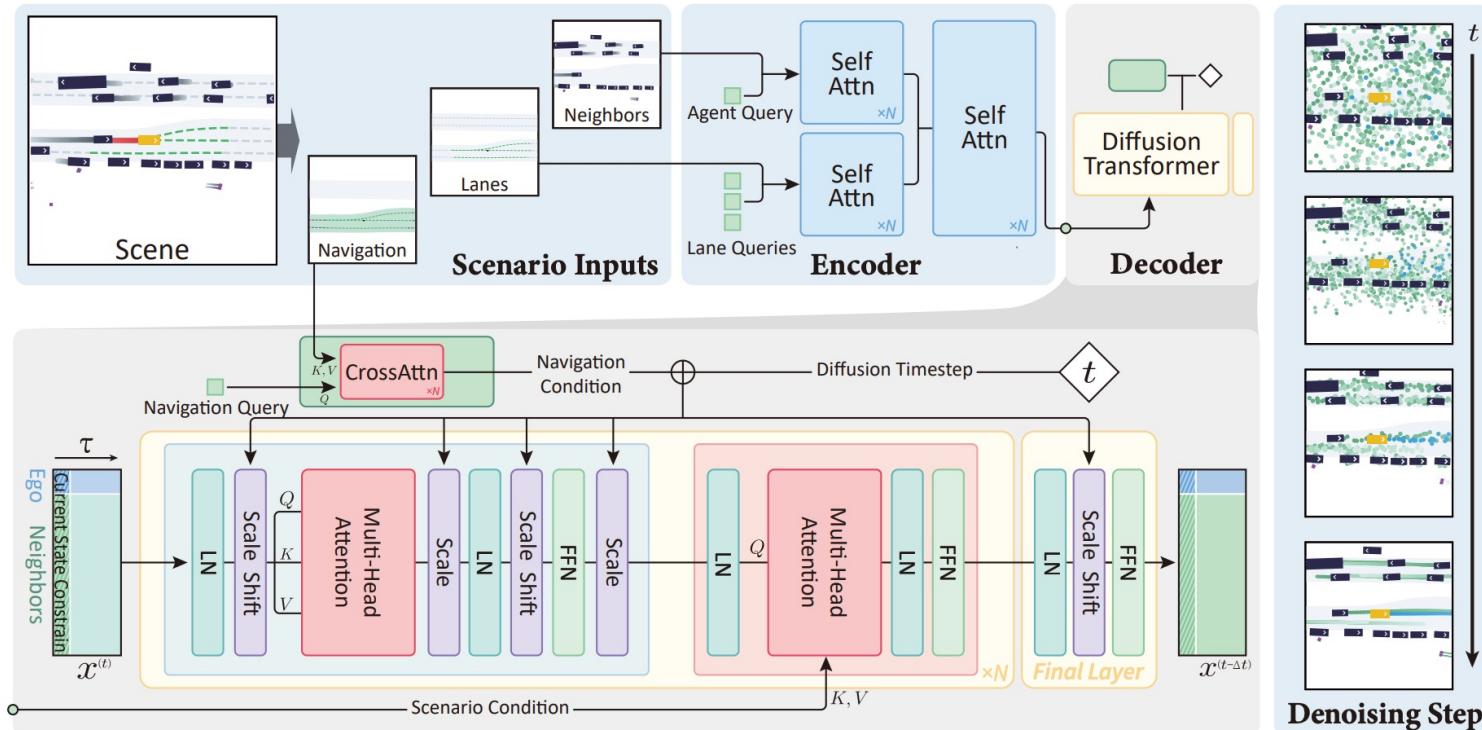


Diffusion Policy $\varepsilon_\theta(\mathbf{O}, \mathbf{A}, k)$



a) Diffusion Policy General Formulation

Diffusion Planner for Self-Driving



<https://openreview.net/forum?id=wM2sfVgMDH>

Summary: DL for Structured Outputs

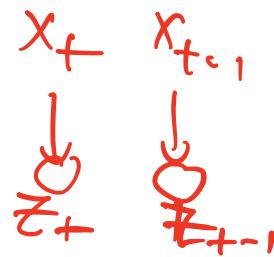
Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.

Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about generative models.

- Graphical models
- Autoregressive
- Energy-based
- Diffusion



Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about generative models.
 - Graphical models
 - Autoregressive
 - Energy-based
 - Diffusion
- Understand pros and cons. Experiment with each option.

Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about generative models.
 - Graphical models
 - Autoregressive
 - Energy-based ✓
 - Diffusion ✓
- Understand pros and cons. Experiment with each option.
- Application in embodied environments.

The World is 3D

- We have previously focused on using 2D images as input.



The World is 3D

- We have previously focused on using 2D images as input.
- But, the world is 3D. Many non-rigid in 2D becomes rigid in 3D. There are also a wide range of 3D sensors.

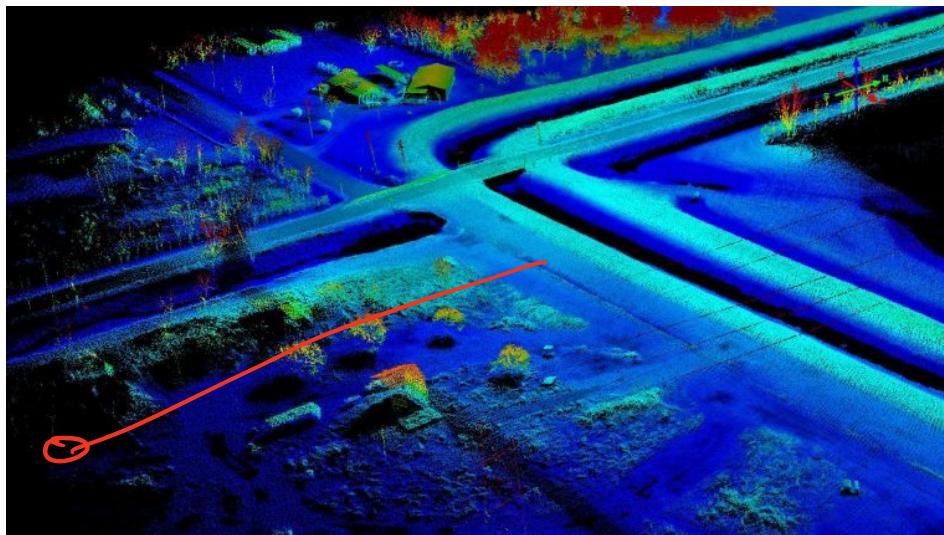


The World is 3D

- We have previously focused on using 2D images as input.
- But, the world is 3D. Many non-rigid in 2D becomes rigid in 3D. There are also a wide range of 3D sensors.
- Stereo (our binocular vision), infrared camera, LiDAR, radar, etc.

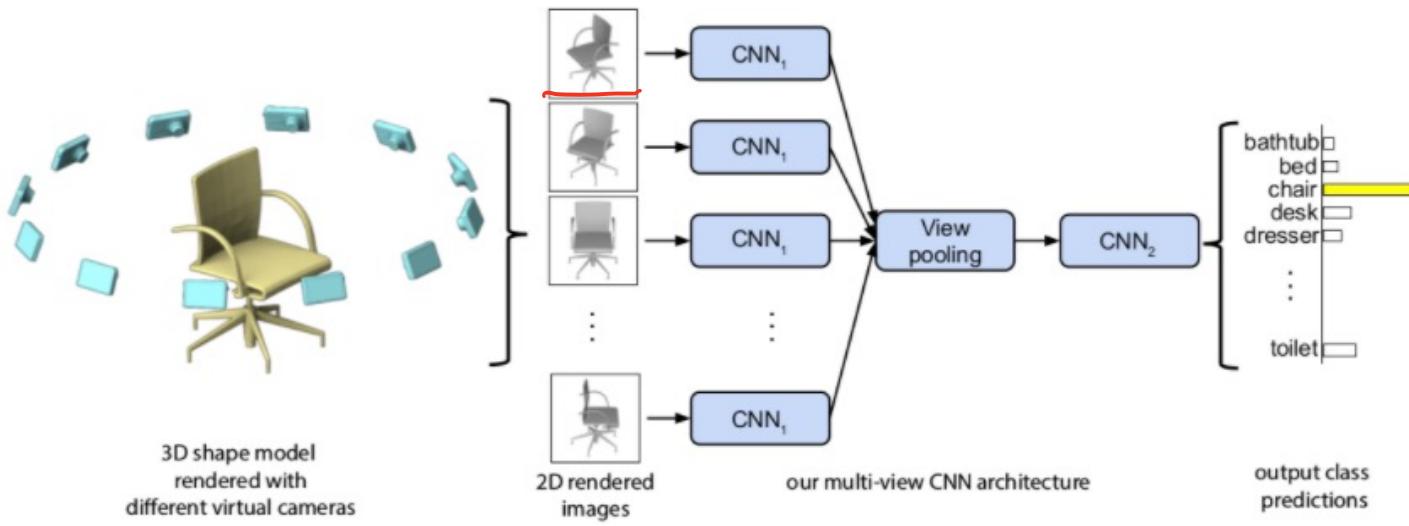


LiDAR



Multi-View CNN

- Treat it as a 2D problem.
- Aggregate the views together with a max-pooling layer.

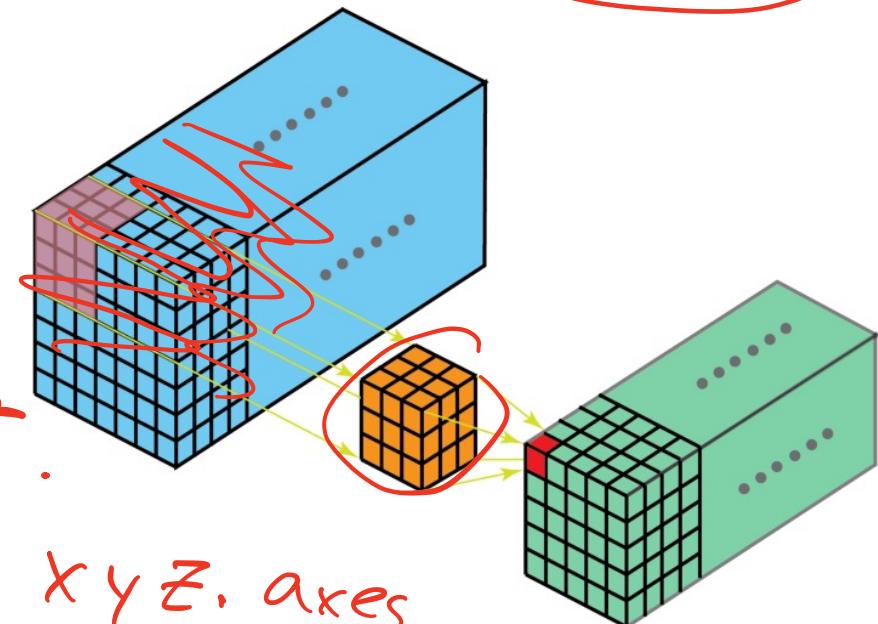
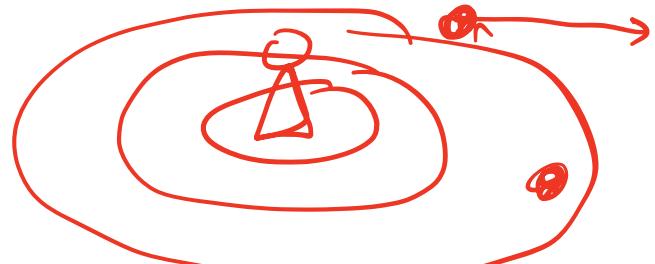


3D Convolution on Voxels

- 3D convolution on occupancy voxels.
- This can be expensive (memory + compute).

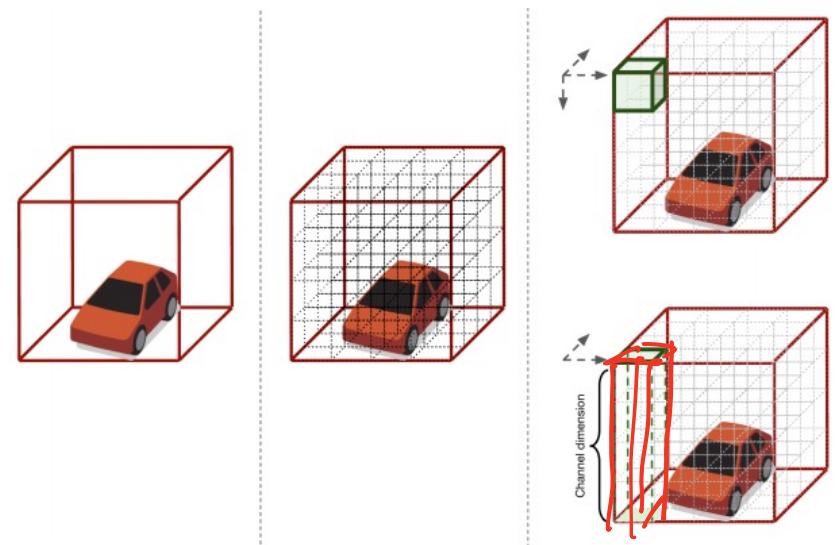
translational equivariance

$$f(T(x)) = T(f(x)).$$



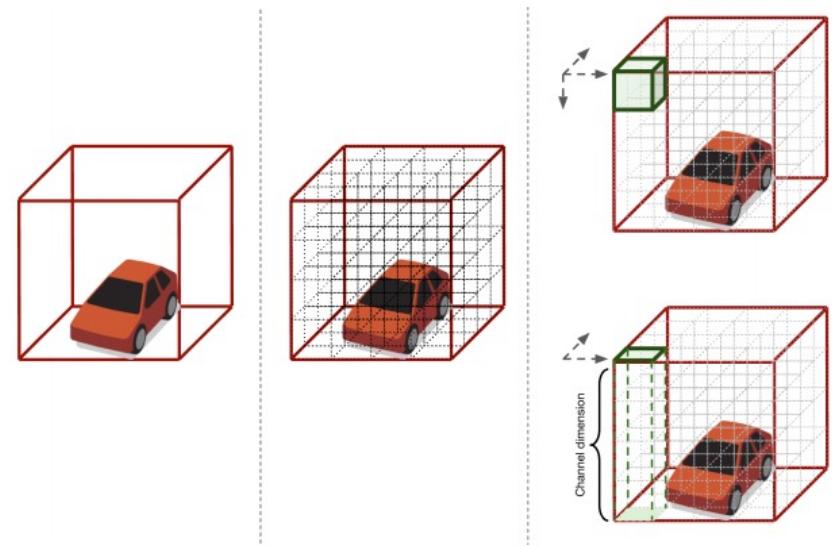
Bird's Eye View (BEV) Voxel

- Treat the z-dimension as different channels.



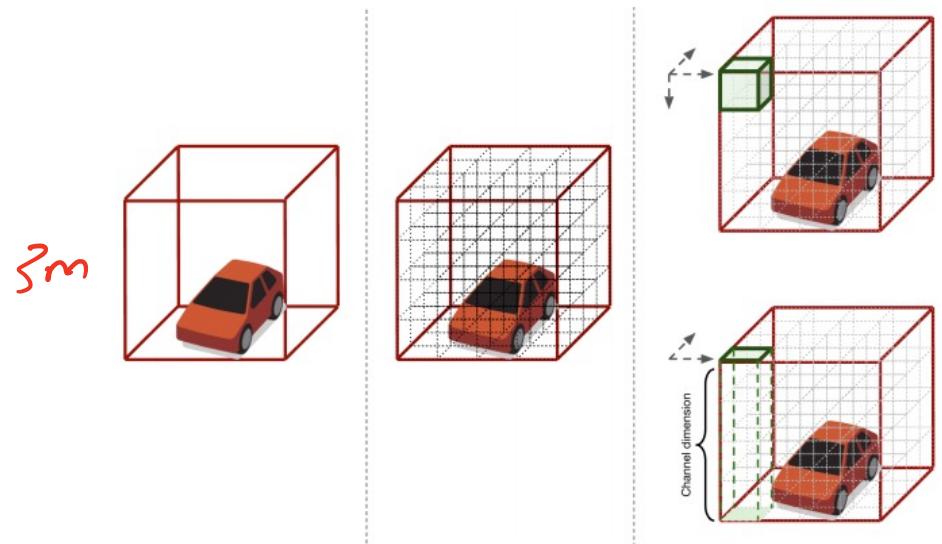
Bird's Eye View (BEV) Voxel

- Treat the z-dimension as different channels.
- 3D convolution → 2D convolution
Popular in self-driving domain, e.g.
80m x



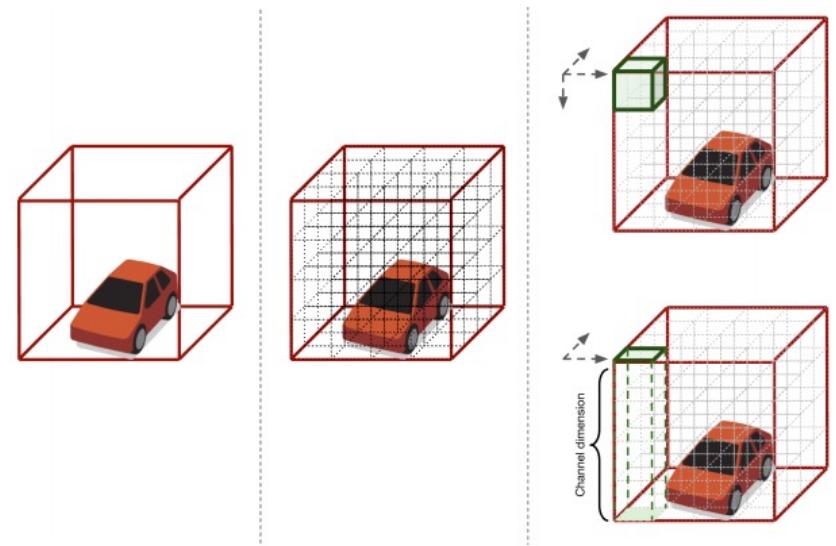
Bird's Eye View (BEV) Voxel

- Treat the z-dimension as different channels.
- 3D convolution → 2D convolution
Popular in self-driving domain, e.g.
80m x
- 140m x 3m (very thin!)
x 80m.



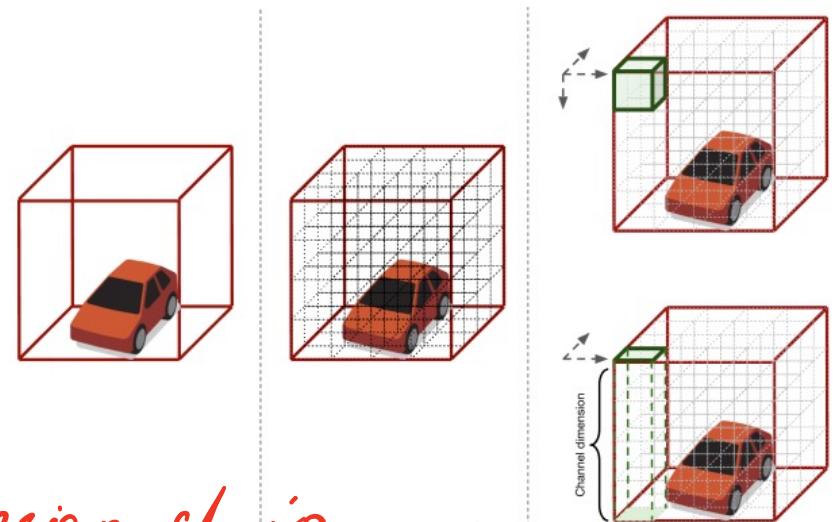
Bird's Eye View (BEV) Voxel

- Treat the z-dimension as different channels.
- 3D convolution → 2D convolution
Popular in self-driving domain, e.g.
80m x
- 140m x 3m (very thin!)
- Transformation in x-y plane is still rigid.



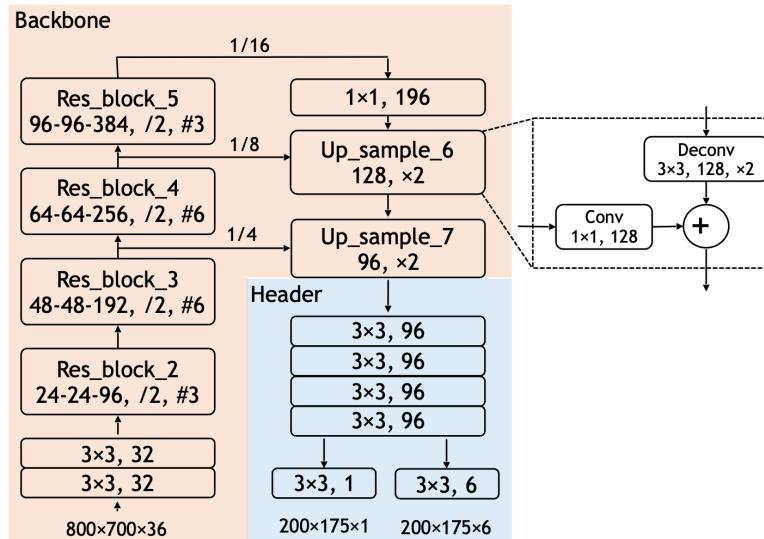
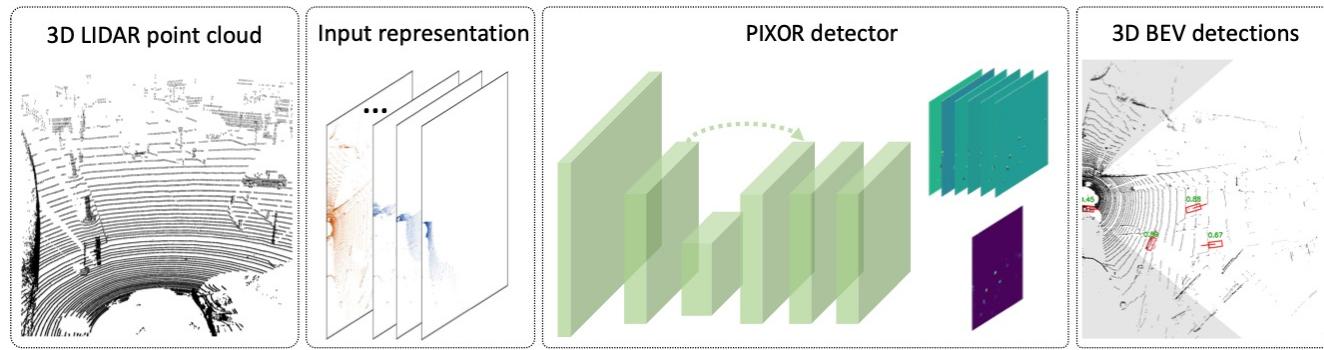
Bird's Eye View (BEV) Voxel

- Treat the z-dimension as different channels.
- 3D convolution → 2D convolution
Popular in self-driving domain, e.g.
80m x
- 140m x 3m (very thin!)
- Transformation in x-y plane is still rigid.
- Bird's eye view: Top down representation of the scene (rigid, sparse) vs. Range view (non-rigid, dense)



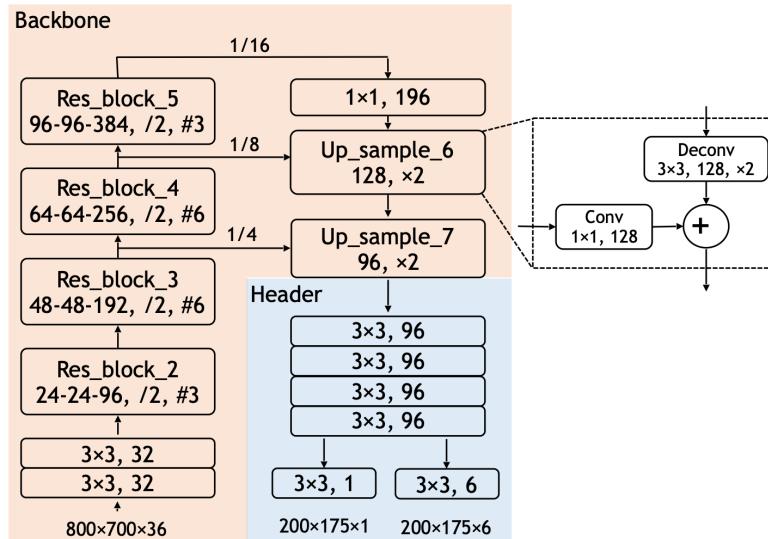
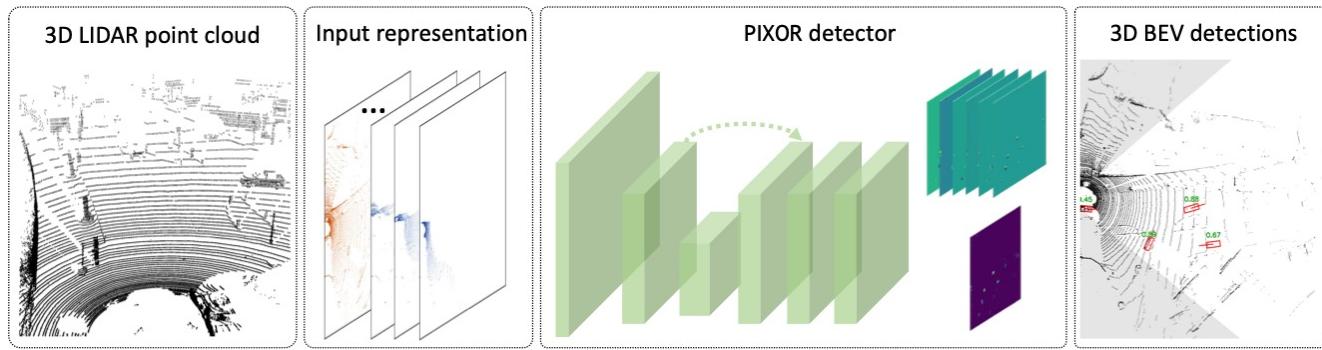
PIXOR

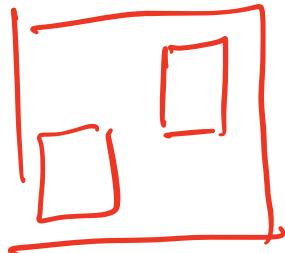
- Bird's eye view object detection.



PIXOR

- Bird's eye view object detection.
- Used the ResNet + FPN network single-stage architecture.





140m x 80m.

PIXOR

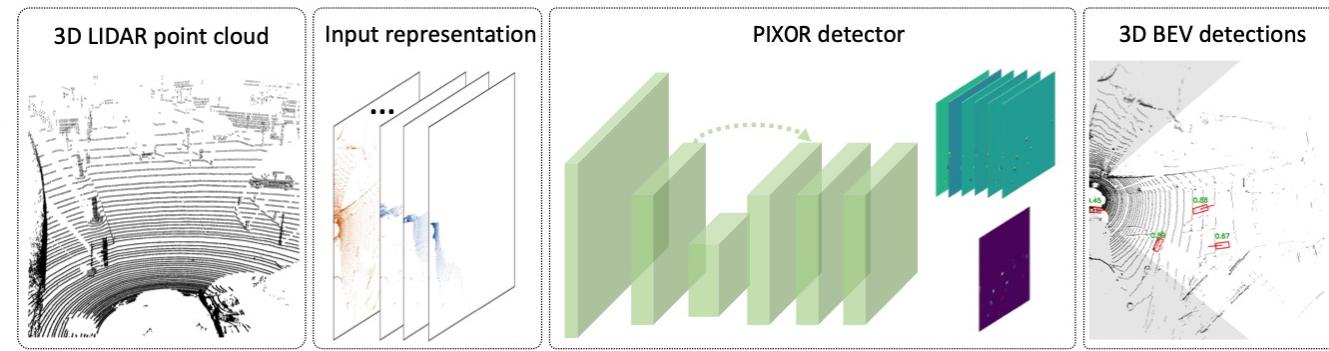
224 x 224

600 x 1000

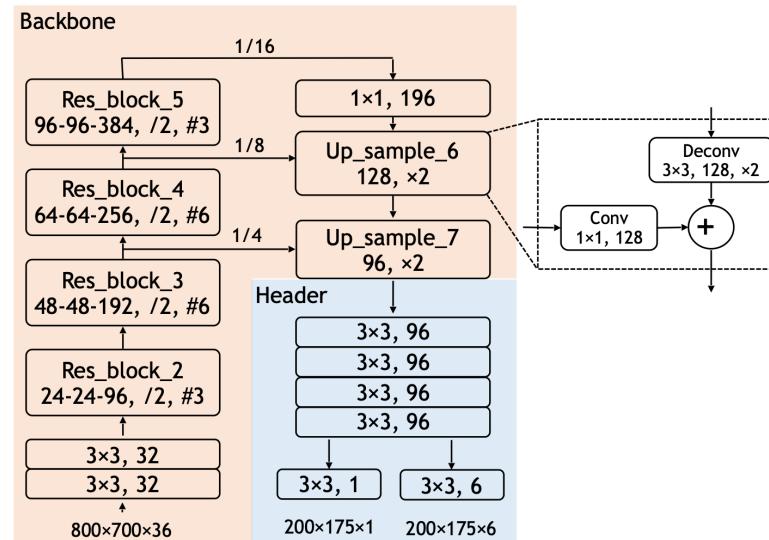
- Bird's eye view object detection.
- Used the ResNet + FPN network single-stage architecture.
- Detection:
Classification +
regression
 $\cos\theta, \sin\theta, dx, dy, w, l.$

$$O_r \xrightarrow{2\pi} r$$

$$\text{atm}\left(\frac{\sin\theta}{\cos\theta}\right).$$

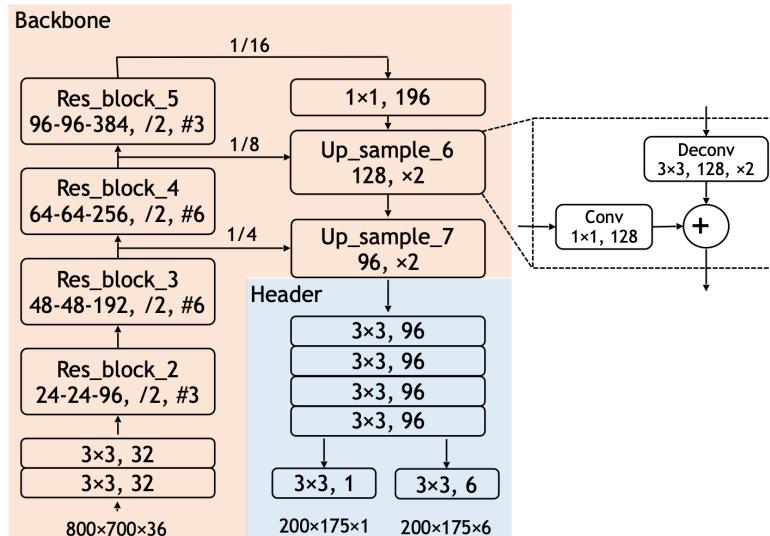
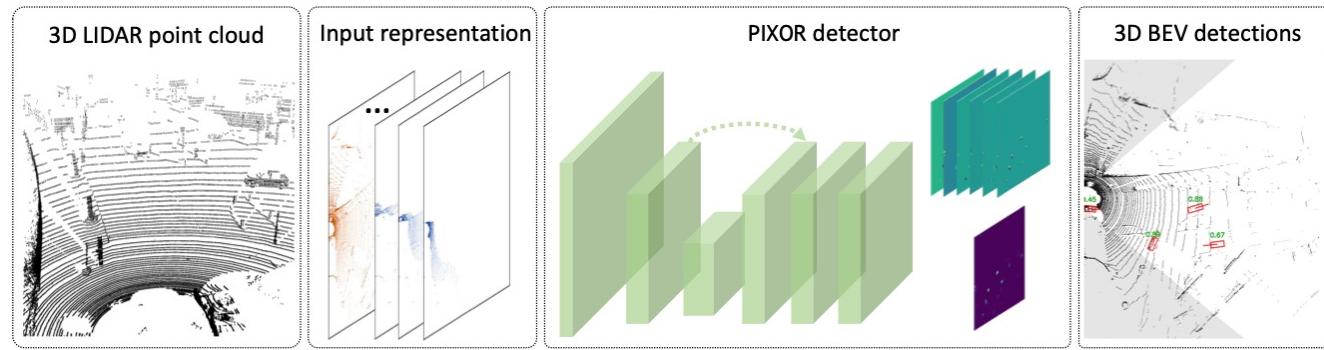


Yang et al. PIXOR: Real-time 3D Object Detection from Point Clouds. CVPR 2018.



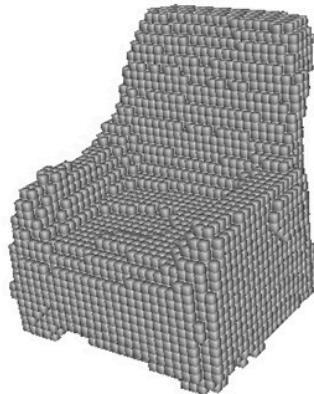
PIXOR

- Bird's eye view object detection.
- Used the ResNet + FPN network single-stage architecture.
- Detection:
Classification +
regression
 $\cos\theta, \sin\theta, dx, dy, w, l.$
- First real-time 3D detection network.



Point Cloud

- Point cloud is native for many 3D-depth sensors: RGBD sensor, LiDAR sensor, etc.



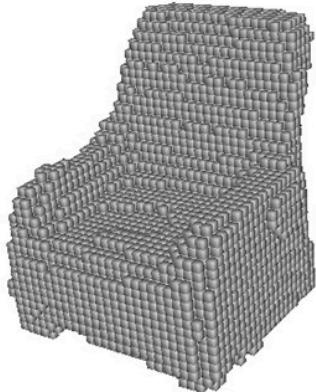
3D voxels



3D point cloud

Point Cloud

- Point cloud is native for many 3D-depth sensors: RGBD sensor, LiDAR sensor, etc.
- List of 3D points: $[(x_1, y_1, z_1), (x_2, y_2, z_2), \dots, (x_N, y_N, z_N)]$



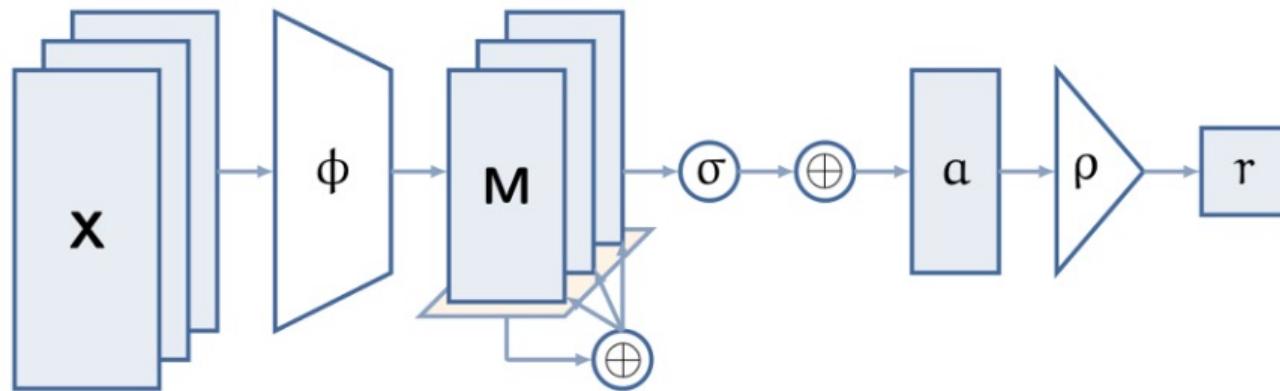
3D voxels



3D point cloud

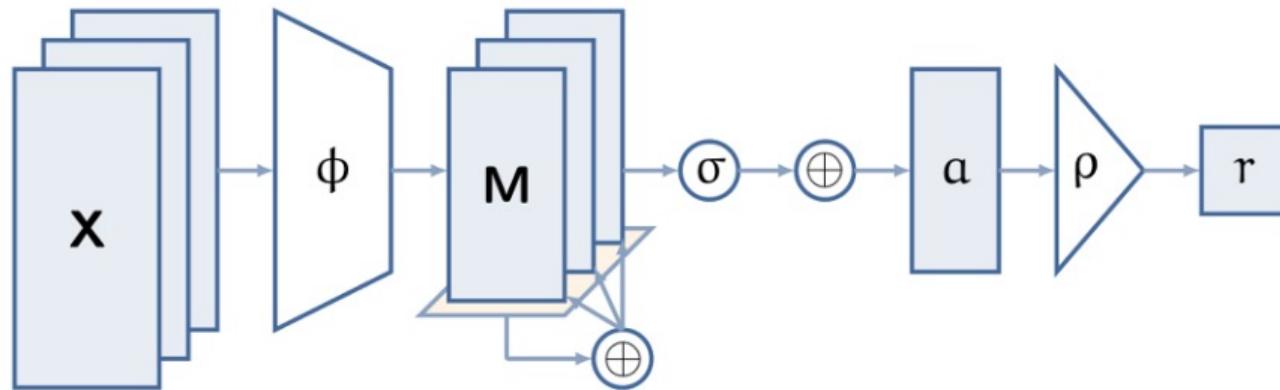
Permutation Invariance

- Point cloud is a set.



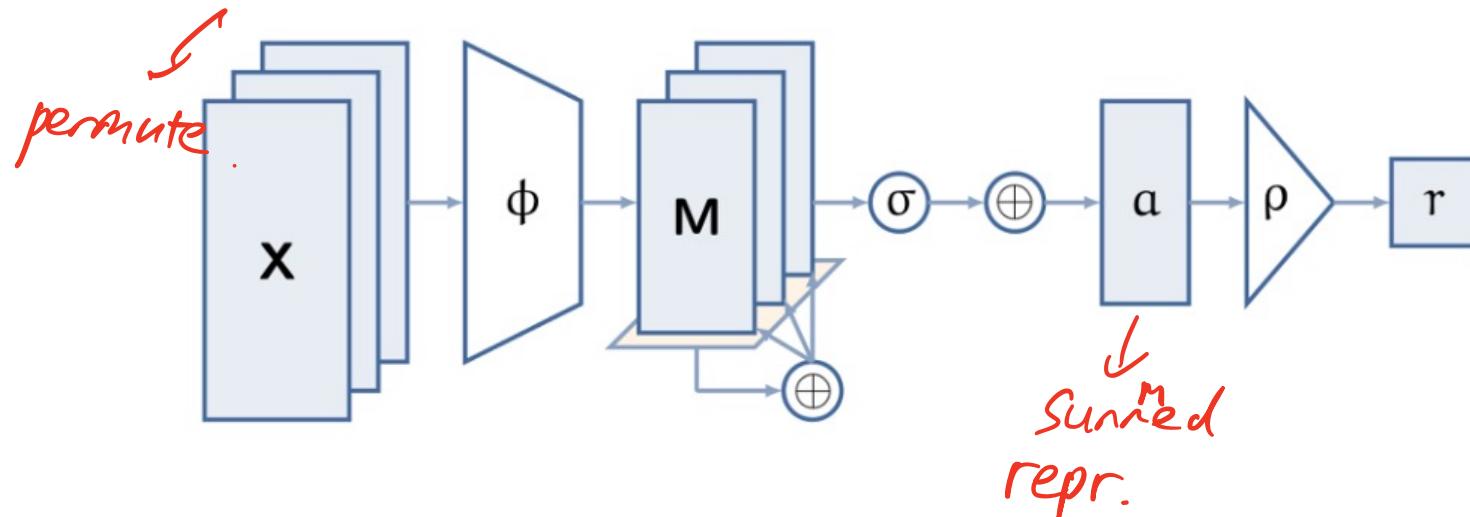
Permutation Invariance

- Point cloud is a set.
- Permutation does not affect the classification in the output.

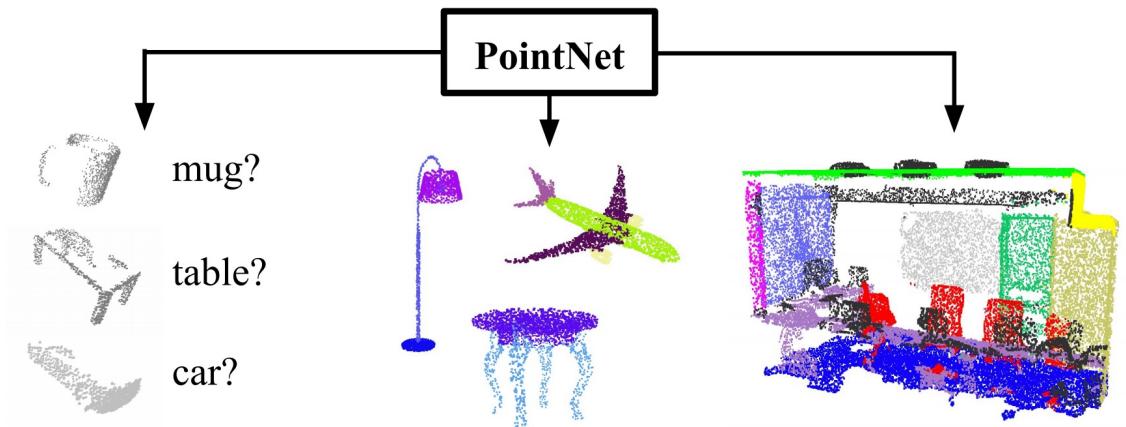


Permutation Invariance

- Point cloud is a set.
- Permutation does not affect the classification in the output.
- What operations are permutation invariant?

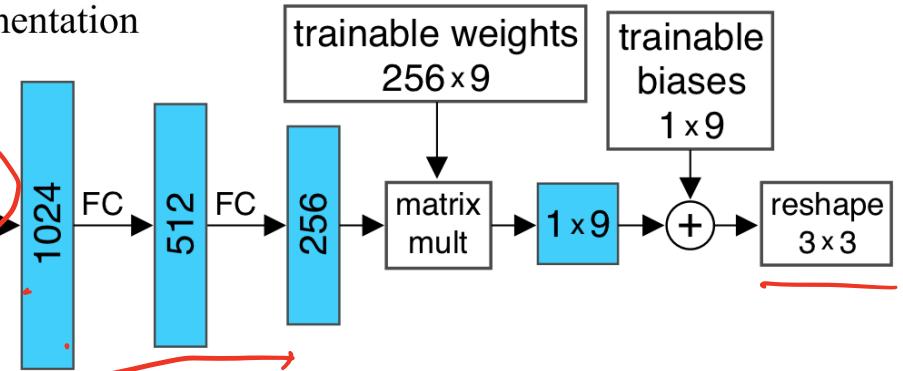
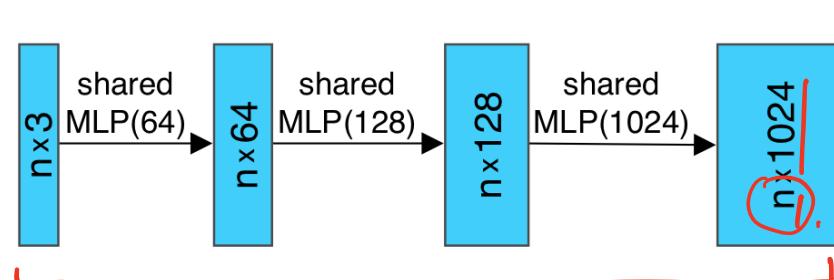


PointNet

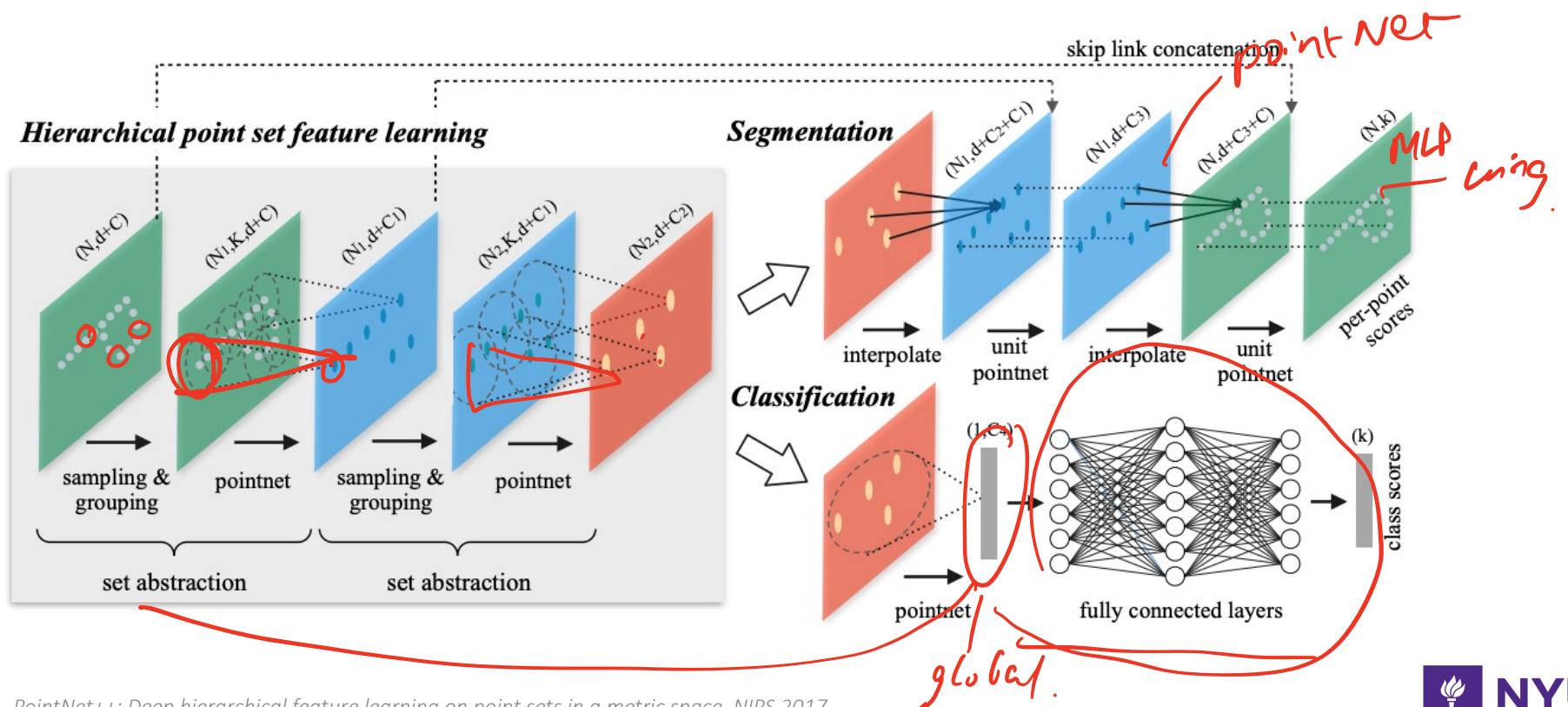


- Apply an MLP on each point.
- Max pool the features across all points.

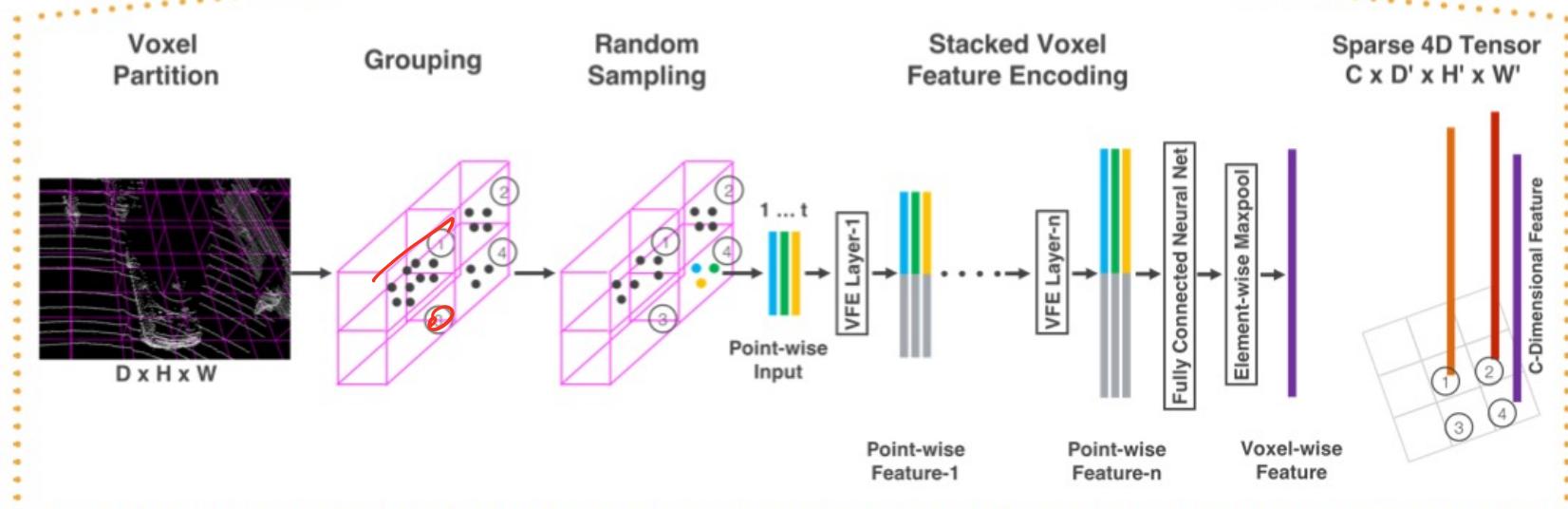
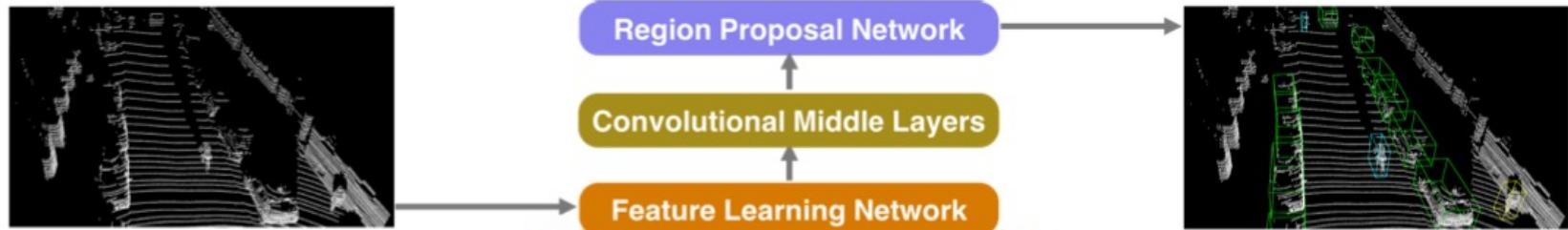
Classification Part Segmentation Semantic Segmentation



PointNet++

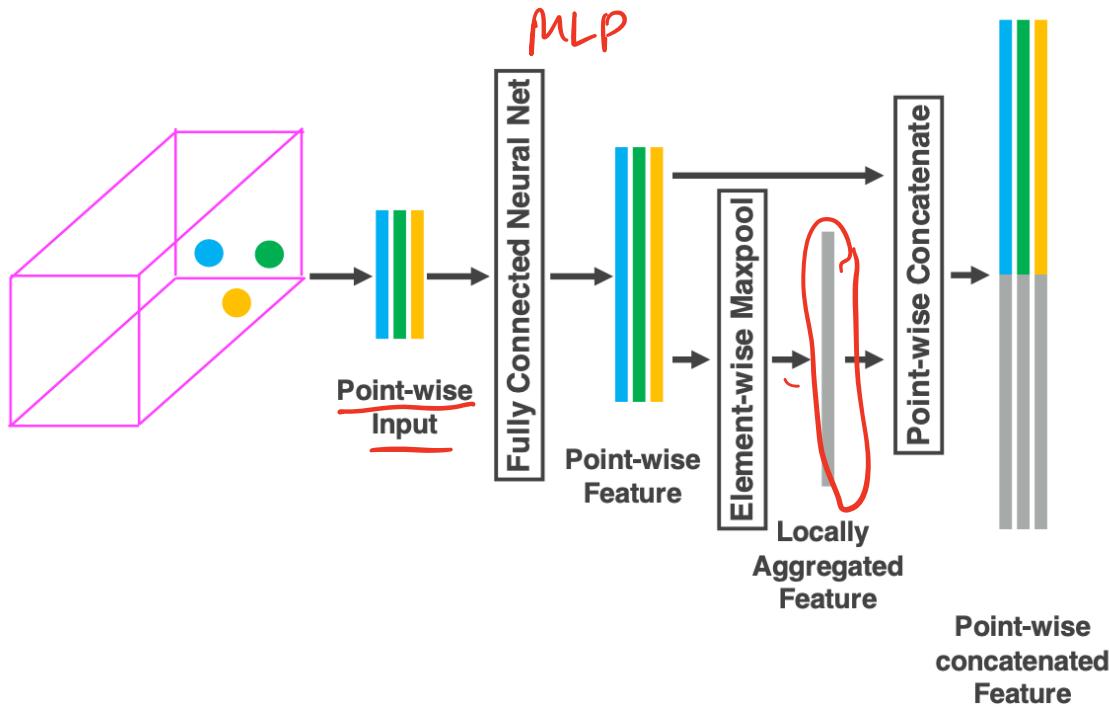


VoxelNet

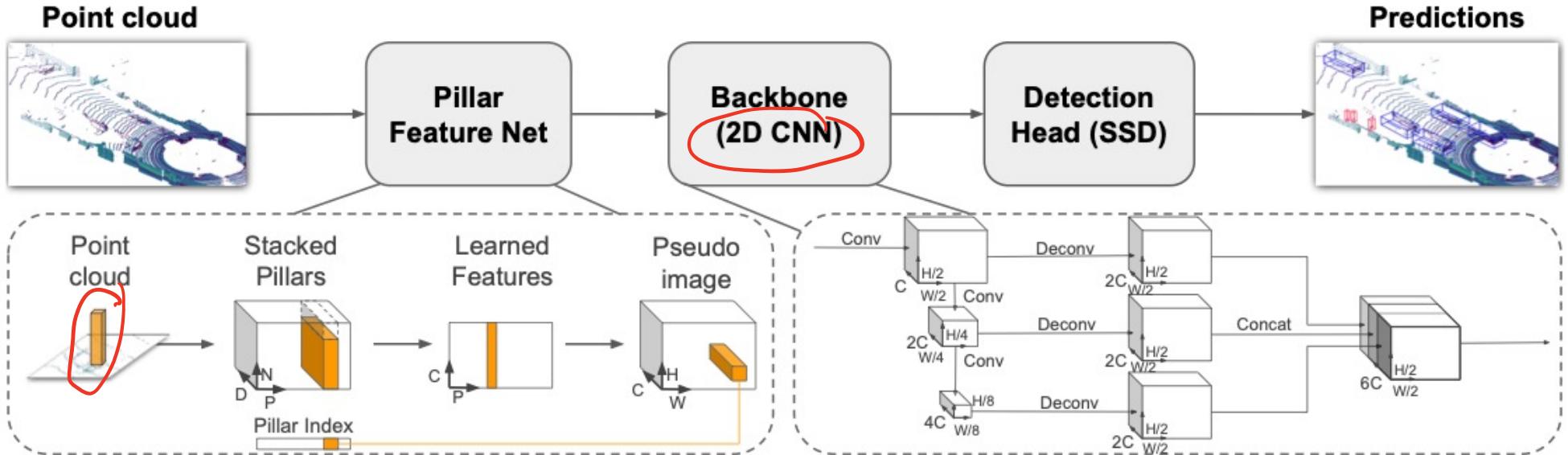


VoxelNet

- Zooming inside voxel feature encoding (VFE)

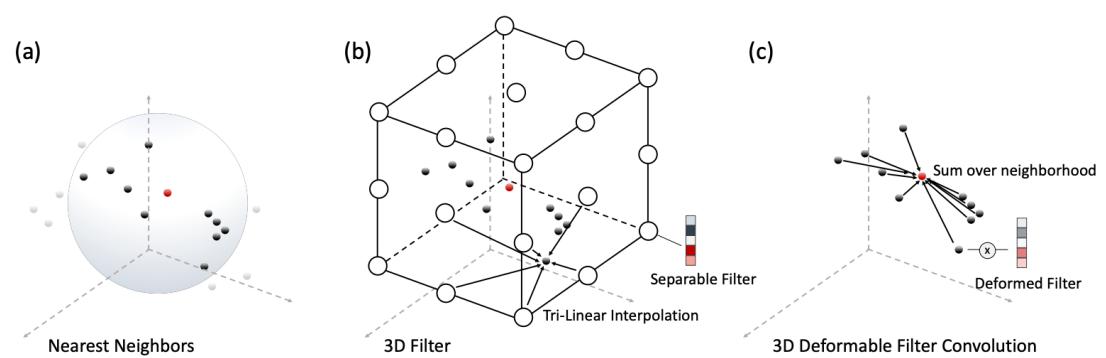
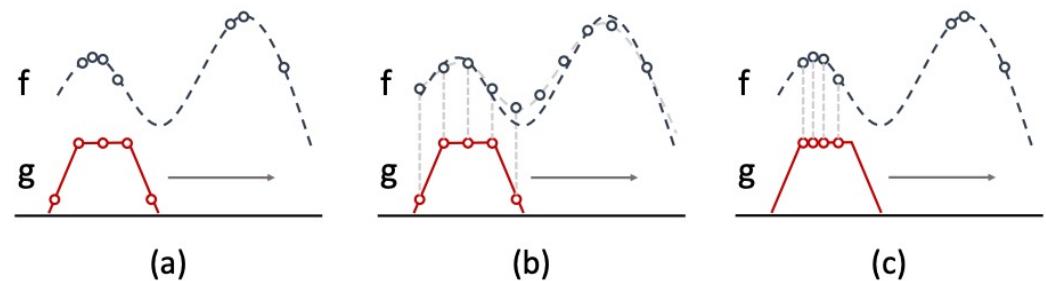


PointPillar



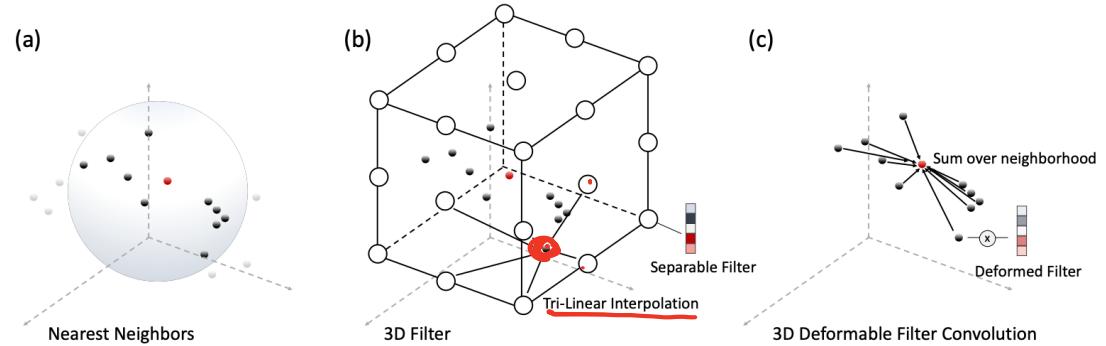
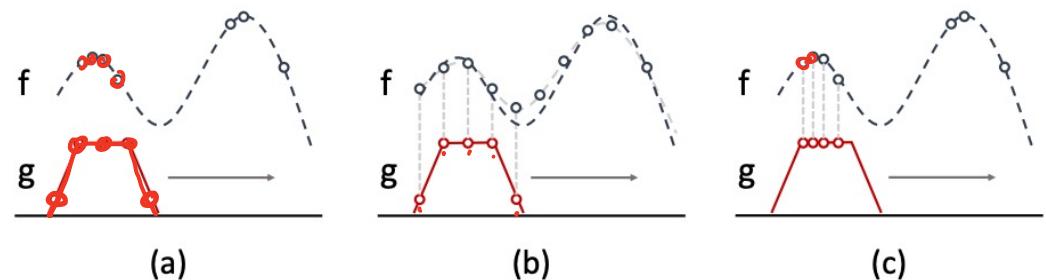
Deformable Convolution in Point Cloud

- Can we convolve a point cloud with a spatially defined kernel function?



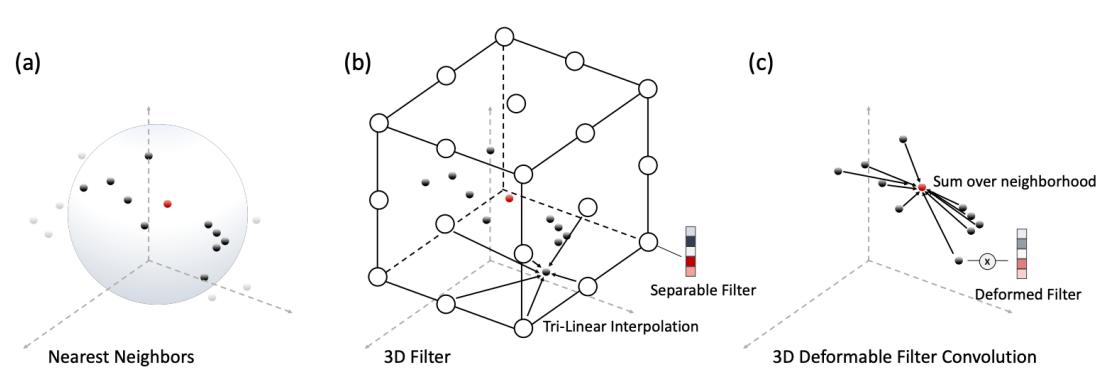
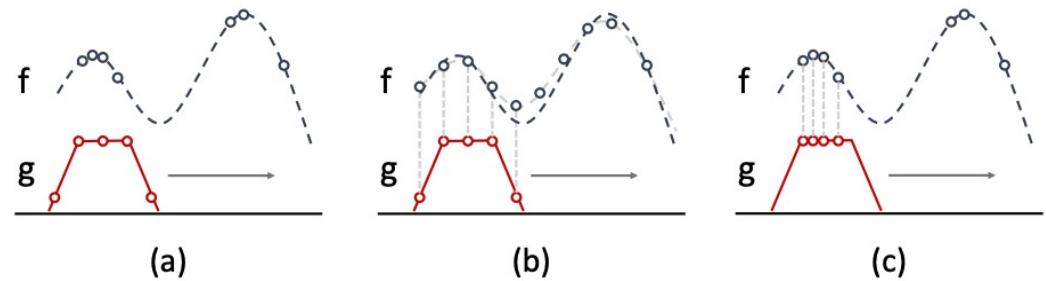
Deformable Convolution in Point Cloud

- Can we convolve a point cloud with a spatially defined kernel function?
- Resample the kernel at the point location.



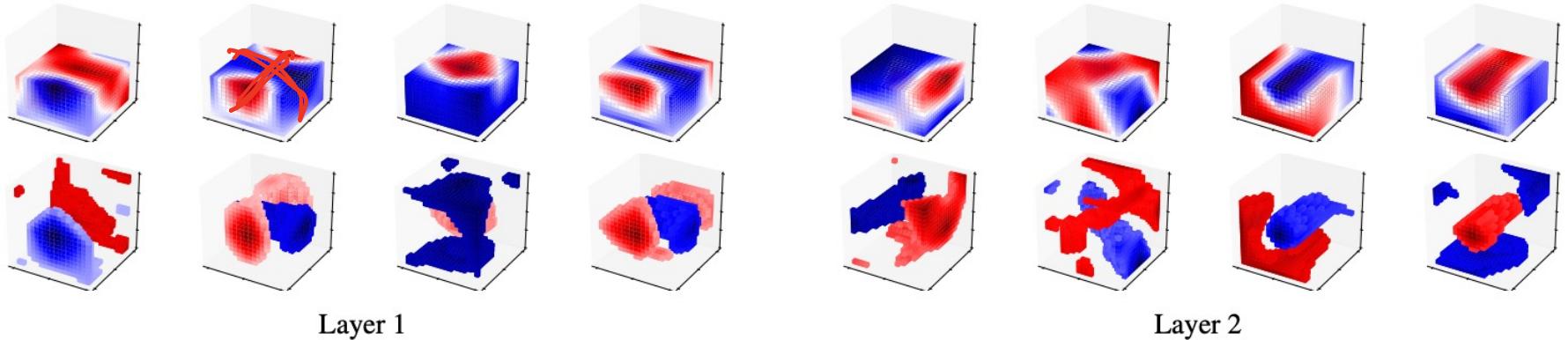
Deformable Convolution in Point Cloud

- Can we convolve a point cloud with a spatially defined kernel function?
- Resample the kernel at the point location.
- Compute the weighted sum around a neighborhood.



3D Filters

- Visualizing 3D convolution kernels.



Multi-Sensor Fusion

- LiDAR is precise in depth perception, but the point cloud format is sparse and non-uniform (dense around the ego-car and sparse in long distance.)

Multi-Sensor Fusion

- LiDAR is precise in depth perception, but the point cloud format is sparse and non-uniform (dense around the ego-car and sparse in long distance.)
- Camera provides high resolution 2D view and good for long distance but lacks 3D. Can we achieve the best of both worlds?

Multi-Sensor Fusion

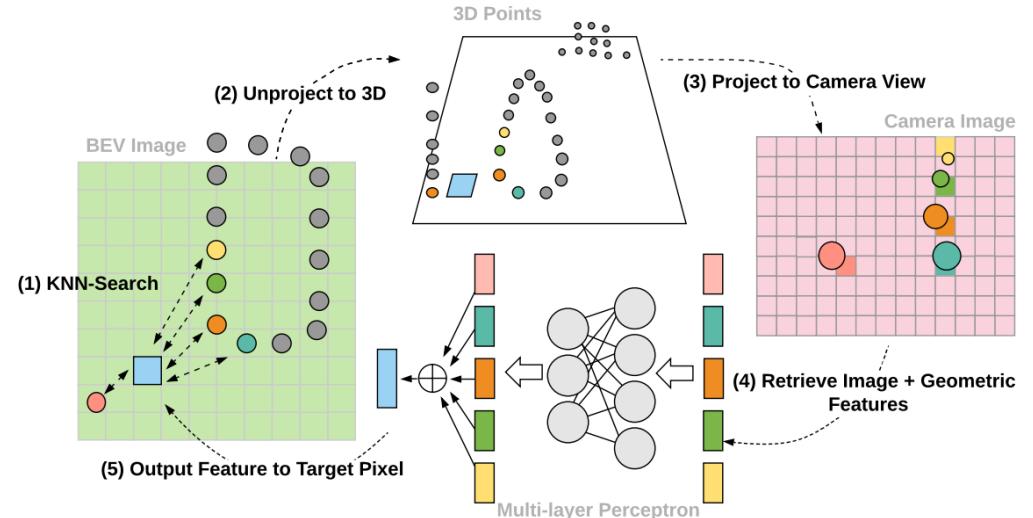
- LiDAR is precise in depth perception, but the point cloud format is sparse and non-uniform (dense around the ego-car and sparse in long distance.)
- Camera provides high resolution 2D view and good for long distance but lacks 3D. Can we achieve the best of both worlds?
- Late fusion: Generate proposals from one branch (e.g. LiDAR) and refine (e.g. using Camera).

Multi-Sensor Fusion

- LiDAR is precise in depth perception, but the point cloud format is sparse and non-uniform (dense around the ego-car and sparse in long distance.)
- Camera provides high resolution 2D view and good for long distance but lacks 3D. Can we achieve the best of both worlds?
- Late fusion: Generate proposals from one branch (e.g. LiDAR) and refine (e.g. using Camera).
- Is there a way to combine the features from both modality in lower layers?

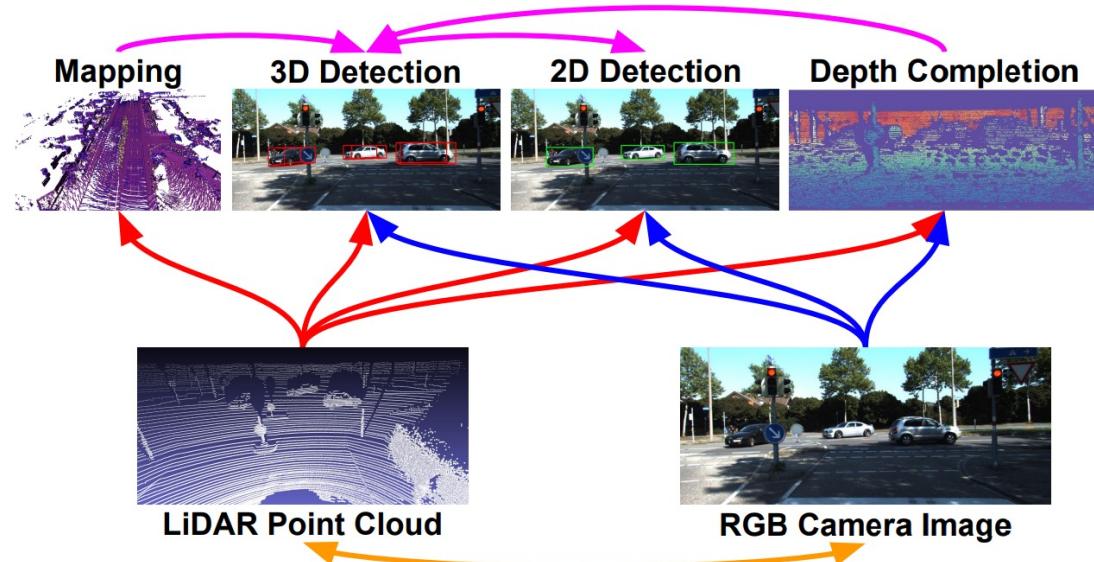
Camera-LiDAR Projection

- Unproject LiDAR points to camera view (i.e. Range View)
- Query the closest camera RGB features for each LiDAR point.
- For empty space in BEV, we can interpolate from neighboring points using kNN.
- Continuous Fusion: $h_i = \sum_j MLP([f_j, x_j - x_i])$.



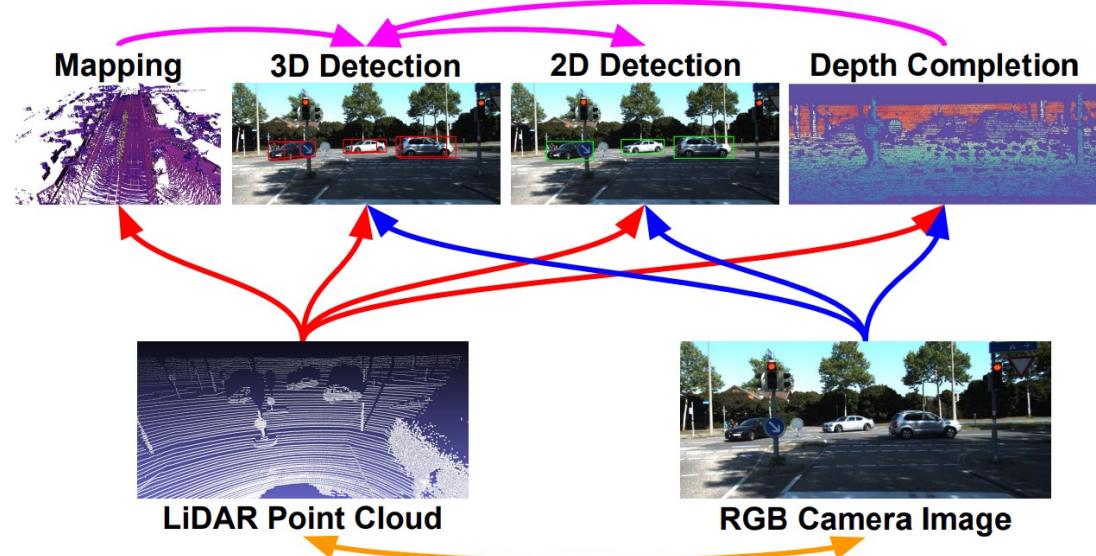
Supervised Dense Depth

- Drawback of continuous fusion: Sparse LiDAR can cause the fusion process to be less accurate. Relies on kNN.



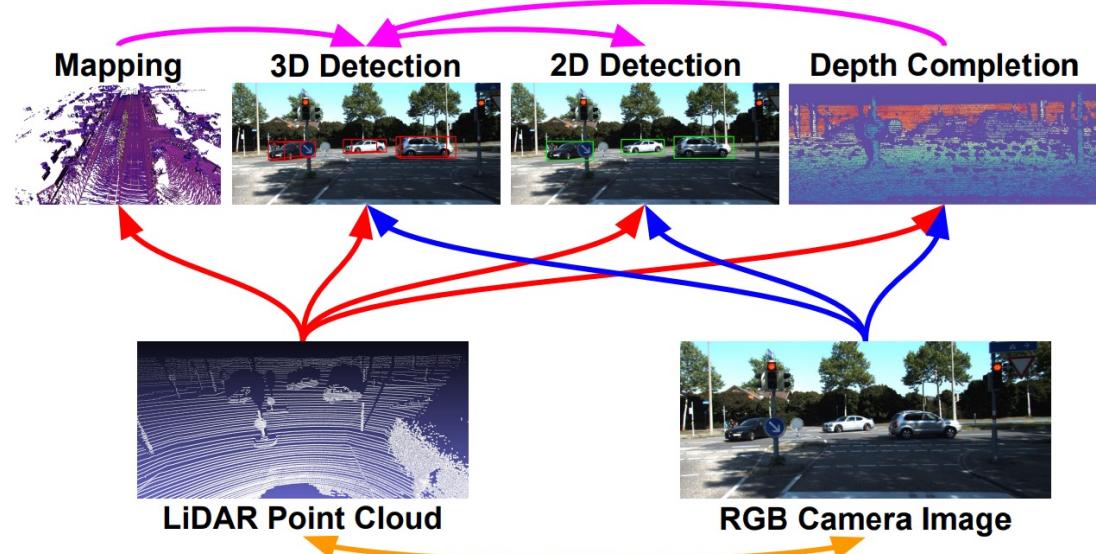
Supervised Dense Depth

- Drawback of continuous fusion: Sparse LiDAR can cause the fusion process to be less accurate. Relies on kNN.
- Why not predict a dense depth to pair with the camera image?



Supervised Dense Depth

- Drawback of continuous fusion: Sparse LiDAR can cause the fusion process to be less accurate. Relies on kNN.
- Why not predict a dense depth to pair with the camera image?
- Depth completion module is supervised by sparse LiDAR and is used for dense fusion.



3D Perception

- With the ease of use of automatic differentiation libraries, we can compose a computation graph in millions of ways.

3D Perception

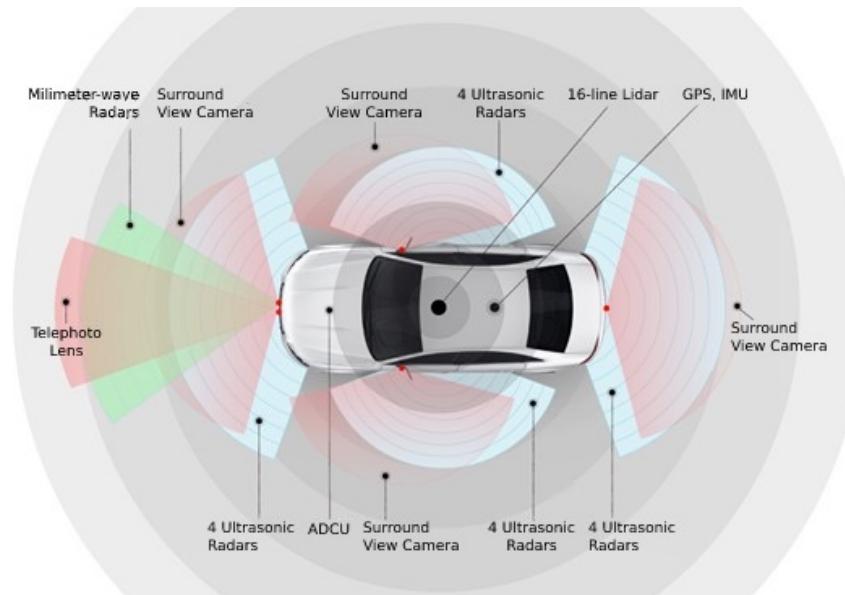
- With the ease of use of automatic differentiation libraries, we can compose a computation graph in millions of ways.
- We can design layers and operators to accomodate different types of inputs and outputs. 3D, point cloud, sparse data, etc.

3D Perception

- With the ease of use of automatic differentiation libraries, we can compose a computation graph in millions of ways.
- We can design layers and operators to accomodate different types of inputs and outputs. 3D, point cloud, sparse data, etc.
- We can fuse different modalities together too, by leveraging geometric relationships.

2D to 3D

- Not all embodied agents have the luxury to have a full set of sensors.
- Can we infer the geometric structure with 2D perception?

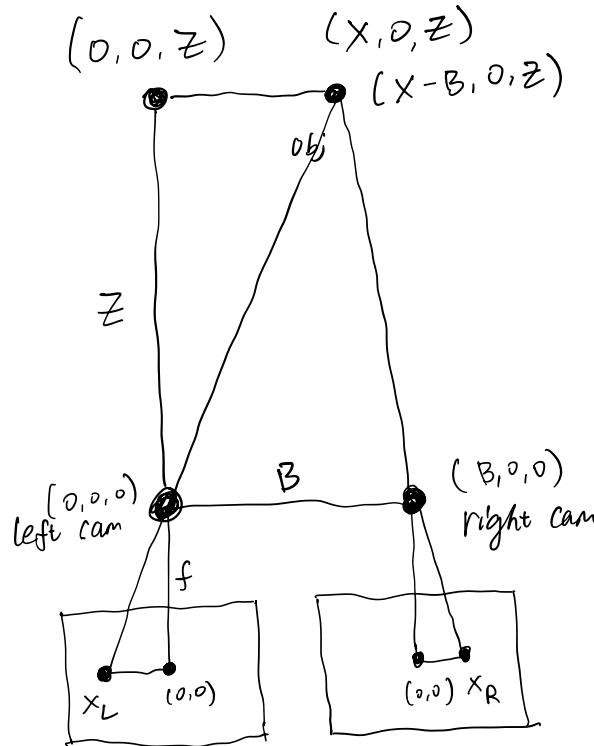


Classic Vision on Depth and Disparity

- One source of depth is from the displacement of pixels in a stereo setup.

Classic Vision on Depth and Disparity

- One source of depth is from the displacement of pixels in a stereo setup.
- But we need to estimate disparity.



$$x_L = \frac{xf}{Z}$$

$$x_R = \frac{(x-B)f}{Z}$$

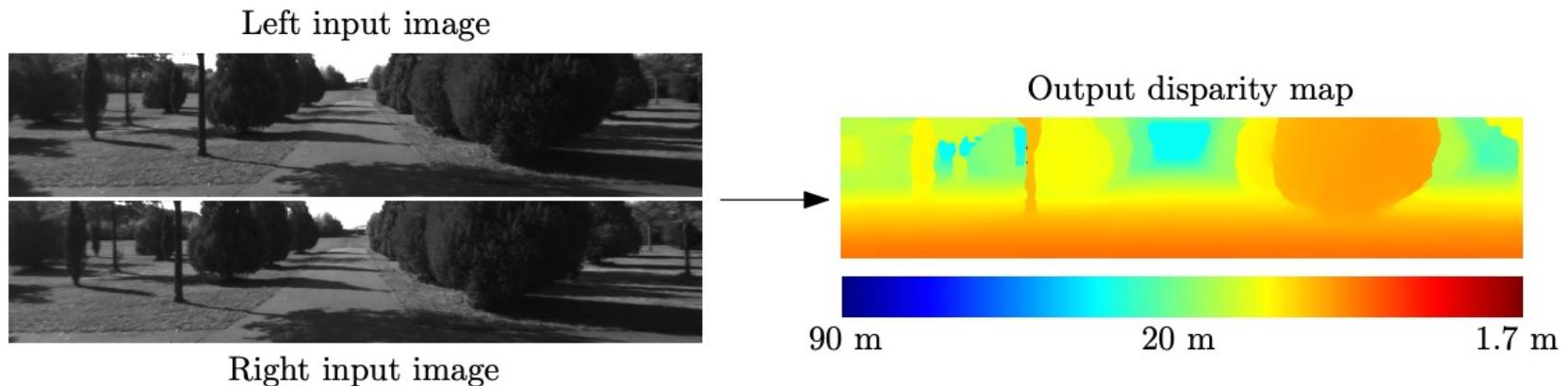
$$\begin{aligned}\delta &= x_L - x_R \\ &= \frac{[x - (x-B)]f}{Z}\end{aligned}$$

$$= \frac{Bf}{Z}$$

$$\boxed{Z = \frac{Bf}{\delta}}$$

From 2D to 3D: Depth Network

- A network that can output disparity.
- Using LiDAR or depth camera as groundtruth supervision.

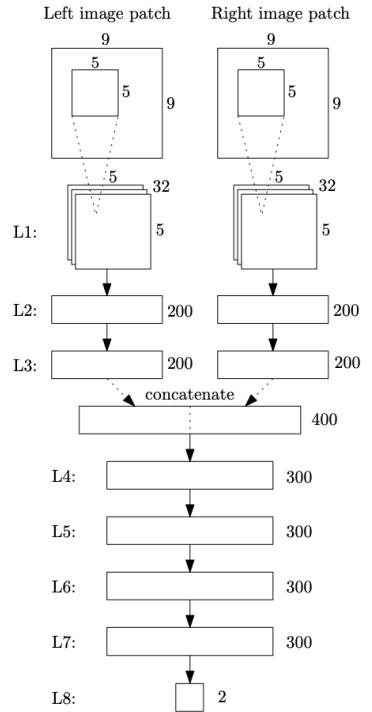


The Energy-Based Approach

- The energy penalize matching with high cost (unary), and when neighboring pixels have disparity differences greater or equal to one (pairwise).
- Cost network: Train with binary classification

$$\text{Energy} \quad E(D) = \sum_{\mathbf{p}} \left(C_{\text{CBCA}}^4(\mathbf{p}, D(\mathbf{p})) + \sum_{\mathbf{q} \in \mathcal{N}_{\mathbf{p}}} P_1 \times 1\{|D(\mathbf{p}) - D(\mathbf{q})| = 1\} + \sum_{\mathbf{q} \in \mathcal{N}_{\mathbf{p}}} P_2 \times 1\{|D(\mathbf{p}) - D(\mathbf{q})| > 1\} \right),$$

$$D(\mathbf{p}) = \operatorname{argmin}_d C(\mathbf{p}, d).$$

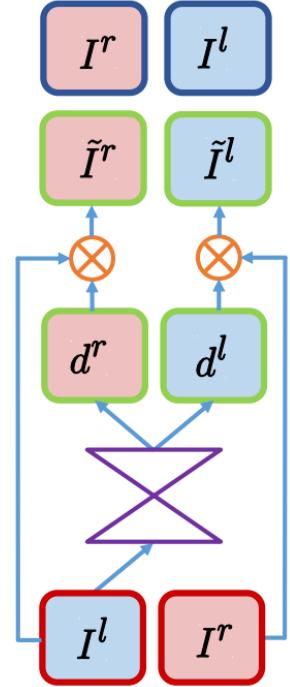
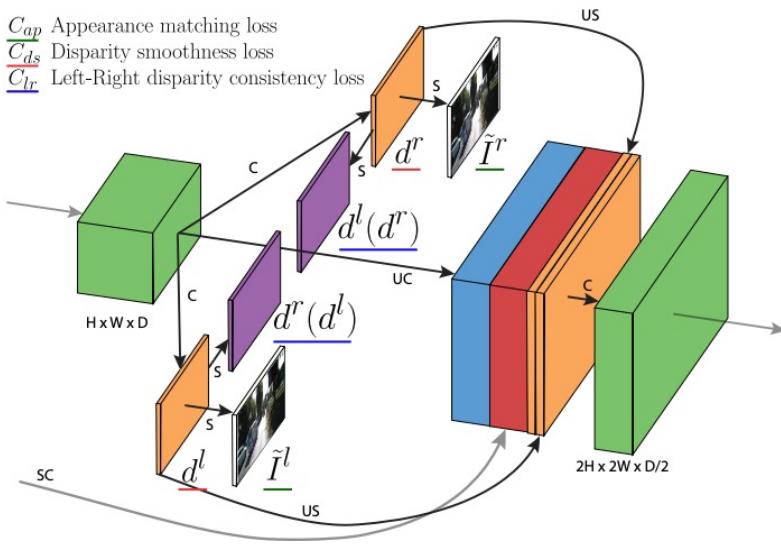


Self-Supervised Depth

- Appearance matching loss

$$C_{ap}^l = \frac{1}{N} \sum_{i,j} \alpha \frac{1 - \text{SSIM}(I_{ij}^l, \tilde{I}_{ij}^l)}{2} + (1-\alpha) \| I_{ij}^l - \tilde{I}_{ij}^l \|.$$

C_{ap} Appearance matching loss
 C_{ds} Disparity smoothness loss
 C_{lr} Left-Right disparity consistency loss



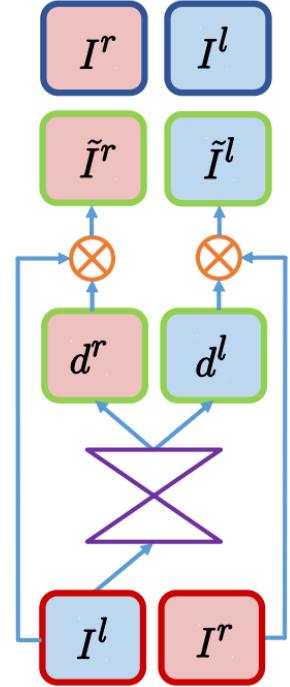
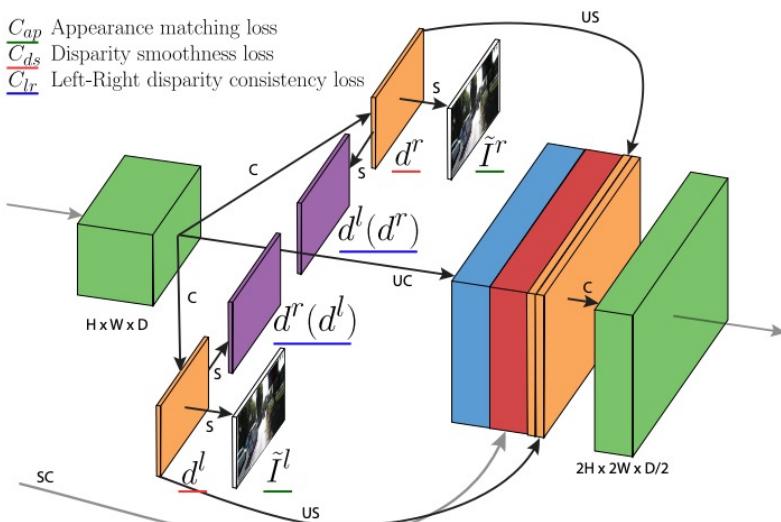
Self-Supervised Depth

- Appearance matching loss

$$C_{ap}^l = \frac{1}{N} \sum_{i,j} \alpha \frac{1 - \text{SSIM}(I_{ij}^l, \tilde{I}_{ij}^l)}{2} + (1-\alpha) \|I_{ij}^l - \tilde{I}_{ij}^l\|.$$

- Disparity smoothness loss

$$C_{ds}^l = \frac{1}{N} \sum_{i,j} |\partial_x d_{ij}^l| e^{-\|\partial_x I_{ij}^l\|} + |\partial_y d_{ij}^l| e^{-\|\partial_y I_{ij}^l\|}.$$



Self-Supervised Depth

- Appearance matching loss

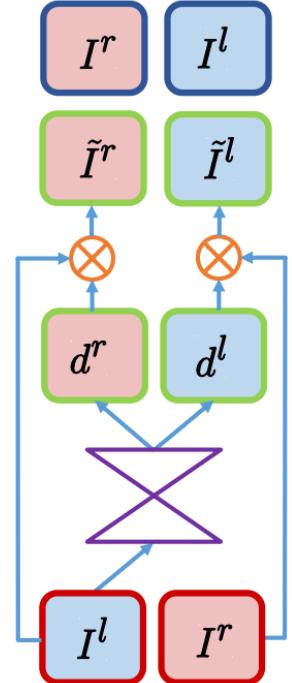
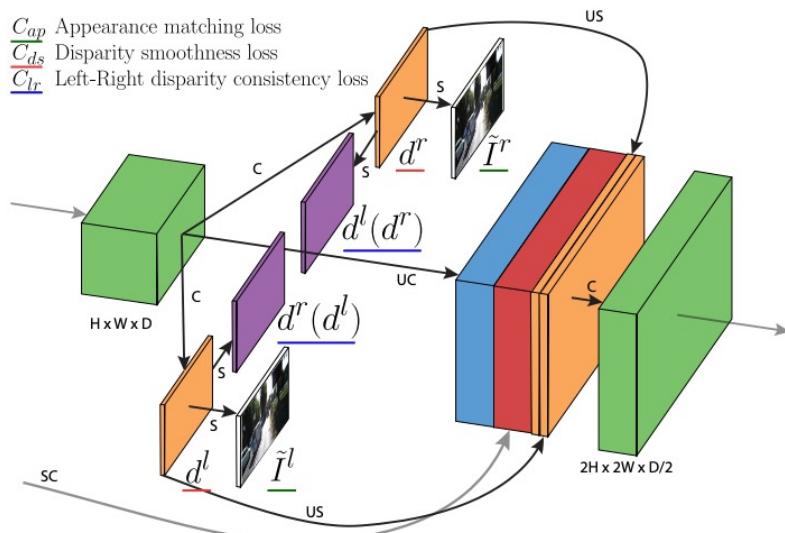
$$C_{ap}^l = \frac{1}{N} \sum_{i,j} \alpha \frac{1 - \text{SSIM}(I_{ij}^l, \tilde{I}_{ij}^l)}{2} + (1-\alpha) \|I_{ij}^l - \tilde{I}_{ij}^l\|.$$

- Disparity smoothness loss

$$C_{ds}^l = \frac{1}{N} \sum_{i,j} |\partial_x d_{ij}^l| e^{-\|\partial_x I_{ij}^l\|} + |\partial_y d_{ij}^l| e^{-\|\partial_y I_{ij}^l\|}.$$

- Left-right disparity consistency loss

$$C_{lr}^l = \frac{1}{N} \sum_{i,j} |d_{ij}^l - d_{ij}^r + d_{ij}^l|.$$



Motion, Optical Flow

- Another task that takes into a pair of image is to estimate the motion of pixels across two consecutive video frames.

Motion, Optical Flow

- Another task that takes into a pair of image is to estimate the motion of pixels across two consecutive video frames.
- Classic method uses brightness constancy assumption.

Motion, Optical Flow

- Another task that takes into a pair of image is to estimate the motion of pixels across two consecutive video frames.
- Classic method uses brightness constancy assumption.

$$I(x, y, t) = I(x + \Delta x, y + \Delta y, t + \Delta t).$$

Motion, Optical Flow

- Another task that takes into a pair of image is to estimate the motion of pixels across two consecutive video frames.
- Classic method uses brightness constancy assumption.

$$I(x, y, t) = I(x + \Delta x, y + \Delta y, t + \Delta t).$$

$$I(x + \Delta x, y + \Delta y, t + \Delta t) = I(x, y, t) + \frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t.$$

Motion, Optical Flow

- Another task that takes into a pair of image is to estimate the motion of pixels across two consecutive video frames.
- Classic method uses brightness constancy assumption.

$$I(x, y, t) = I(x + \Delta x, y + \Delta y, t + \Delta t).$$

$$I(x + \Delta x, y + \Delta y, t + \Delta t) = I(x, y, t) + \frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t.$$

$$\frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t = 0.$$

Motion, Optical Flow

- Another task that takes into a pair of image is to estimate the motion of pixels across two consecutive video frames.
- Classic method uses brightness constancy assumption.

$$I(x, y, t) = I(x + \Delta x, y + \Delta y, t + \Delta t).$$

$$I(x + \Delta x, y + \Delta y, t + \Delta t) = I(x, y, t) + \frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t.$$

$$I_x = \frac{\partial I}{\partial x} \quad \frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t = 0.$$

$$I_y = \frac{\partial I}{\partial y}$$

Motion, Optical Flow

- Another task that takes into a pair of image is to estimate the motion of pixels across two consecutive video frames.
- Classic method uses brightness constancy assumption.

$$I(x, y, t) = I(x + \Delta x, y + \Delta y, t + \Delta t).$$

$$I(x + \Delta x, y + \Delta y, t + \Delta t) = I(x, y, t) + \frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t.$$

$$I_x = \frac{\partial I}{\partial x} \quad \frac{\partial I}{\partial x} \Delta x + \frac{\partial I}{\partial y} \Delta y + \frac{\partial I}{\partial t} \Delta t = 0.$$

$$I_y = \frac{\partial I}{\partial y}$$

$$I_x u + I_y v + I_t = 0.$$

Classical Approach

- Under-constrained system

$$I_x u + I_y v + I_t = 0.$$

Classical Approach

- Under-constrained system
- Use a local patch and assume smooth motion

$$I_x u + I_y v + I_t = 0.$$

$$\mathbf{A}\mathbf{u} = \mathbf{b}$$

$$\begin{pmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{N^2}) & I_y(\mathbf{p}_{N^2}) \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix} = - \begin{pmatrix} I_t(\mathbf{p}_1) \\ \vdots \\ I_t(\mathbf{p}_{N^2}) \end{pmatrix}$$

Classical Approach

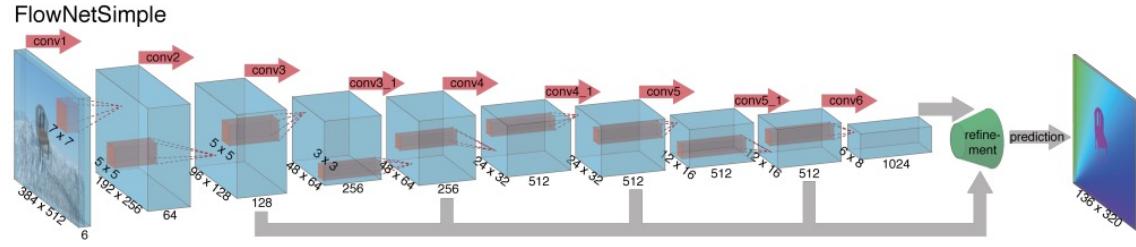
- Under-constrained system $I_x u + I_y v + I_t = 0.$
- Use a local patch and assume smooth motion
- Rigid, contains many assumptions

$$\mathbf{A}\mathbf{u} = \mathbf{b}$$

$$\begin{pmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{N^2}) & I_y(\mathbf{p}_{N^2}) \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix} = - \begin{pmatrix} I_t(\mathbf{p}_1) \\ \vdots \\ I_t(\mathbf{p}_{N^2}) \end{pmatrix}$$

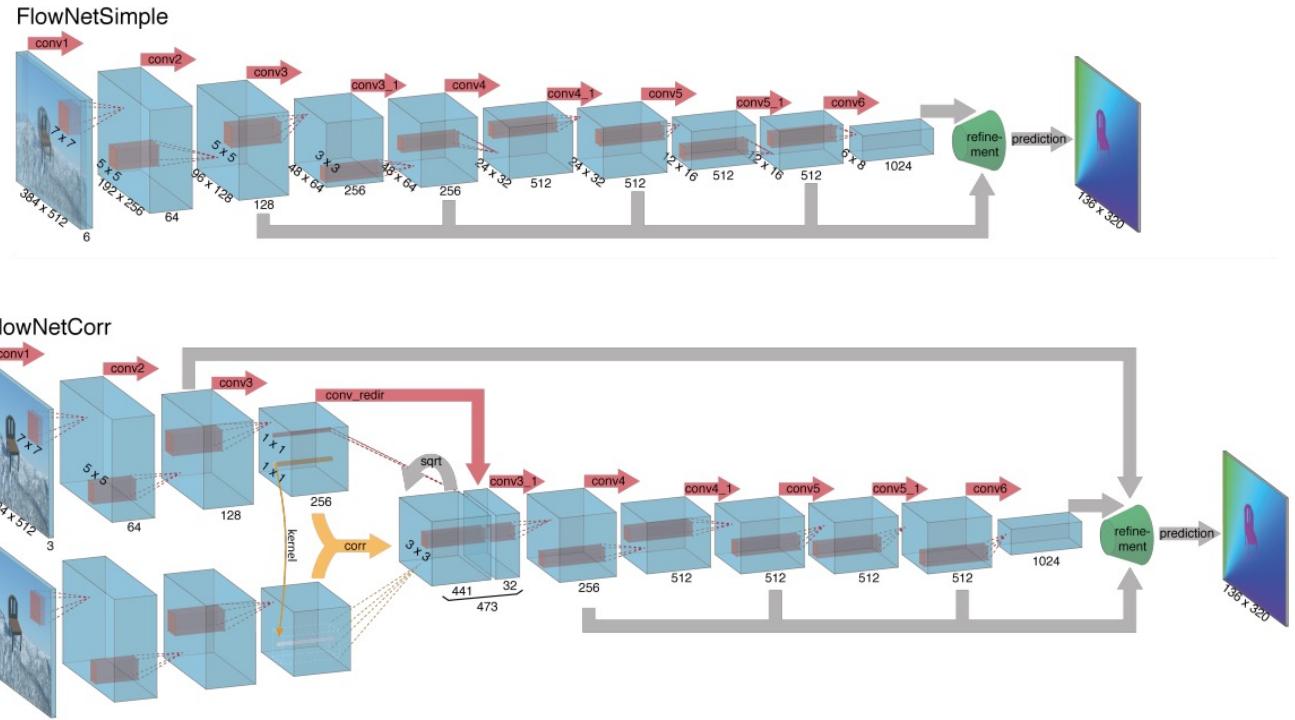
Correlation Volume Approach

- Simple Approach:
Concatenate the
two images
together.

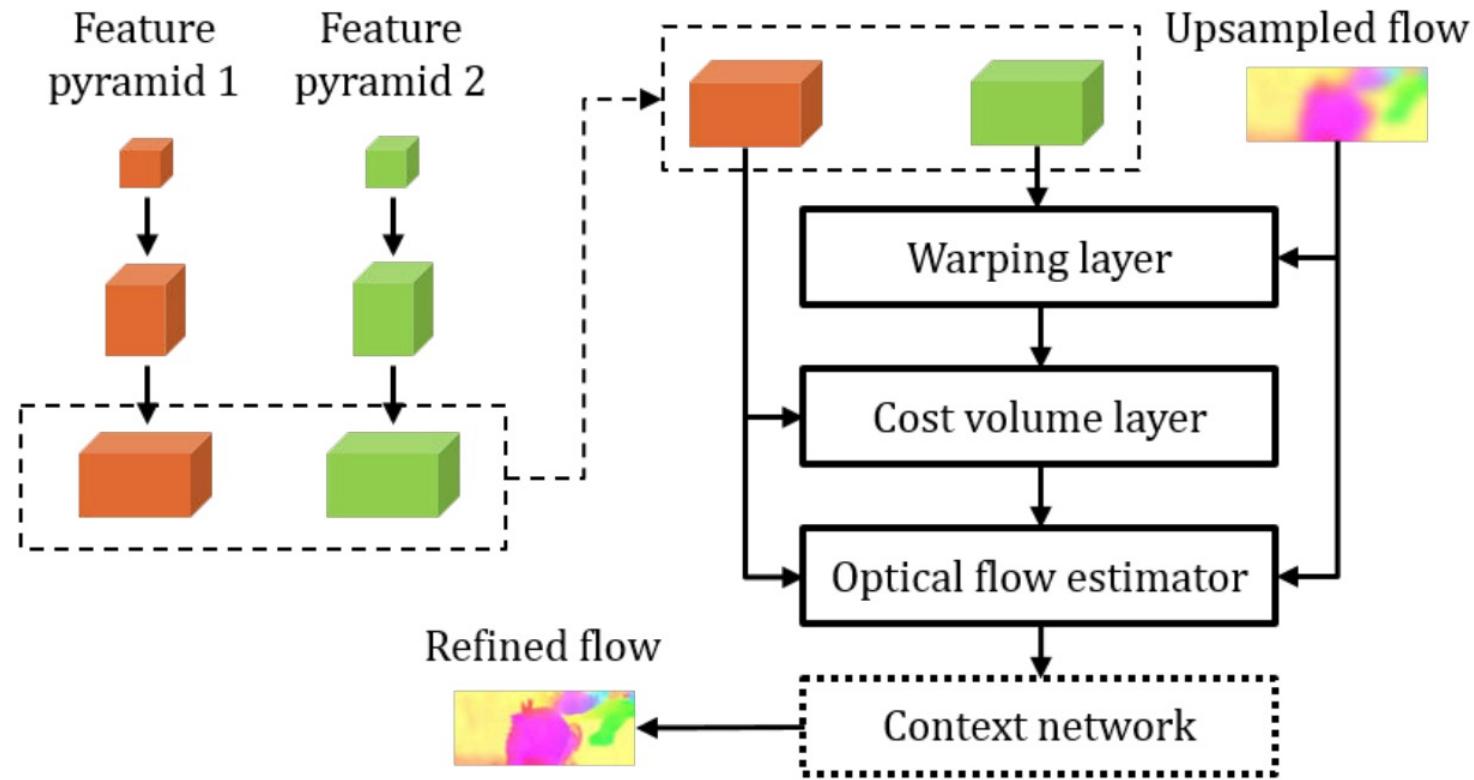


Correlation Volume Approach

- Simple Approach:
Concatenate the
two images
together.
- Correlation:
Extract some
levels of features,
and convolve one
feature on top of
another.



Iterative Refining through Feature Pyramid



Unsupervised Flow

- Photometric Consistency (Appearance)

Unsupervised Flow

- Photometric Consistency (Appearance)
- Occlusion Estimation
 - Forward-backward consistency

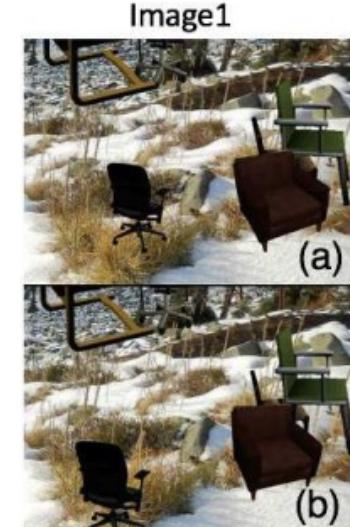


Image2

Wang et al., 2018

Unsupervised Flow

- Photometric Consistency (Appearance)
- Occlusion Estimation
 - Forward-backward consistency
- Smoothness

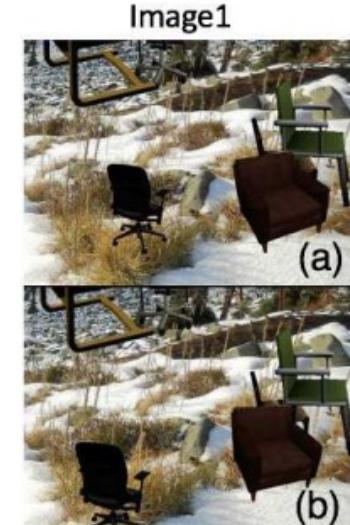


Image1

(a)

(b)

Image2

Wang et al., 2018

Unsupervised Flow

- Photometric Consistency (Appearance)
- Occlusion Estimation
 - Forward-backward consistency
- Smoothness
- Self-supervision: Ensure consistent flow at different augmentation (e.g. crops)

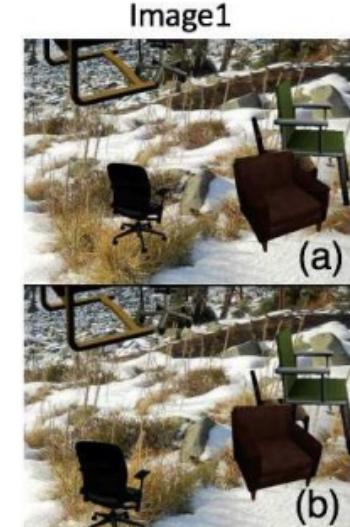


Image1

(a)

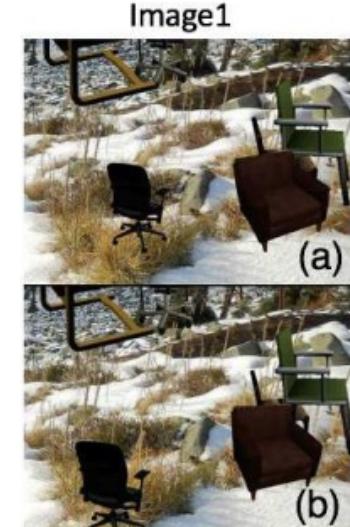
(b)

Image2

Wang et al., 2018

Unsupervised Flow

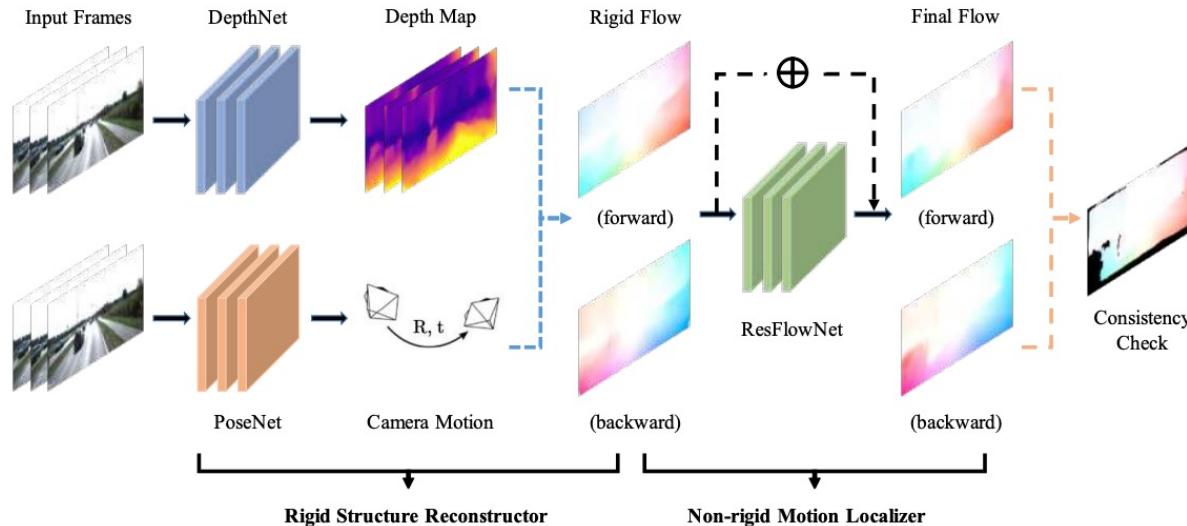
- Photometric Consistency (Appearance)
- Occlusion Estimation
 - Forward-backward consistency
- Smoothness
- Self-supervision: Ensure consistent flow at different augmentation (e.g. crops)
- Can 3D information help us reason about motion?



Wang et al., 2018

Depth, Flow, and Pose Movement

- The static objects follow rigid flow: determined by camera motion and depth.
$$f_{t \mapsto s}^{rig}(p_t) = K T_{t \mapsto s} D_t(p_t) K^{-1} p_t - p_t.$$



Training Losses

- Appearance Loss (Warping):

$$\mathcal{L}_{rw} = \alpha \frac{1 - SSIM(I_t, \tilde{I}_s^{rig})}{2} + (1 - \alpha) \|I_t - \tilde{I}_s^{rig}\|_1.$$

$$\mathcal{L}_{fw} = \alpha \frac{1 - SSIM(I_t, \tilde{I}_s^{full})}{2} + (1 - \alpha) \|I_t - \tilde{I}_s^{full}\|_1.$$

Training Losses

- Appearance Loss (Warping):

$$\mathcal{L}_{rw} = \alpha \frac{1 - SSIM(I_t, \tilde{I}_s^{rig})}{2} + (1 - \alpha) \|I_t - \tilde{I}_s^{rig}\|_1.$$
$$\mathcal{L}_{fw} = \alpha \frac{1 - SSIM(I_t, \tilde{I}_s^{full})}{2} + (1 - \alpha) \|I_t - \tilde{I}_s^{full}\|_1.$$

- Smoothness Loss:

$$\mathcal{L} = \sum_{p_t} |\nabla D(p_t)| \cdot (\exp(-|\nabla I(p(t))|))^T.$$

Training Losses

- Appearance Loss (Warping):

$$\mathcal{L}_{rw} = \alpha \frac{1 - SSIM(I_t, \tilde{I}_s^{rig})}{2} + (1 - \alpha) \|I_t - \tilde{I}_s^{rig}\|_1.$$
$$\mathcal{L}_{fw} = \alpha \frac{1 - SSIM(I_t, \tilde{I}_s^{full})}{2} + (1 - \alpha) \|I_t - \tilde{I}_s^{full}\|_1.$$

- Smoothness Loss:

$$\mathcal{L} = \sum_{p_t} |\nabla D(p_t)| \cdot (\exp(-|\nabla I(p(t))|))^T.$$

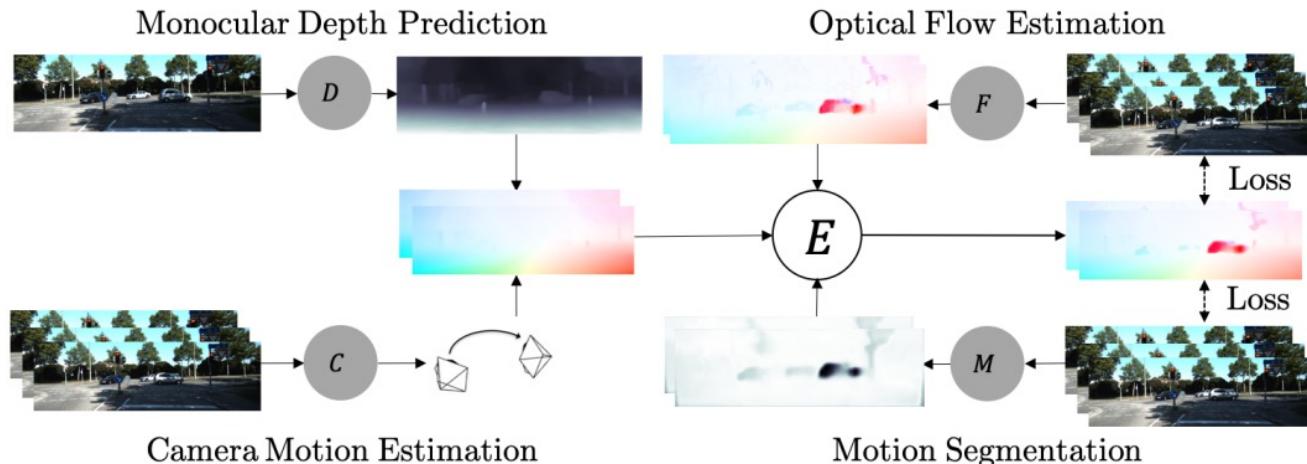
- Forward-Backward Consistency:

$$\mathcal{L} = \sum_{p_t} [\delta(p_t)] \cdot \|\Delta f_{t \rightarrow s}^{full}(p_t)\|_1.$$

$$\delta(p_t) = \|f_{t \rightarrow s}^{full}(p_t)\|_2 \max\{\alpha, \beta \|f_{t \rightarrow s}^{full}(p_t)\|_2\}.$$

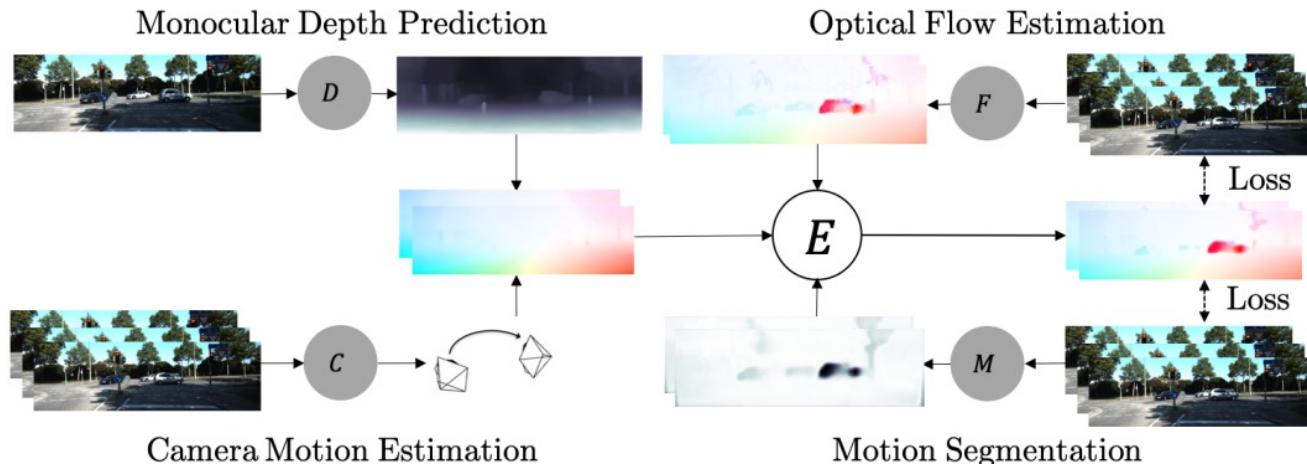
Summary

- Leverage cross correlation structure for spatial similarity matching.



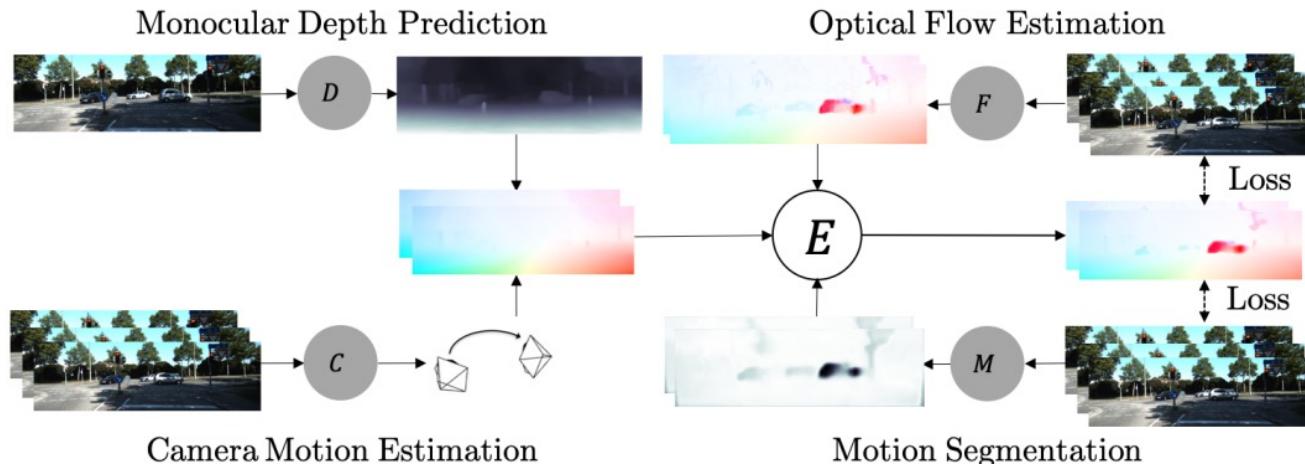
Summary

- Leverage cross correlation structure for spatial similarity matching.
- Can be used towards: depth, flow, and pose prediction.



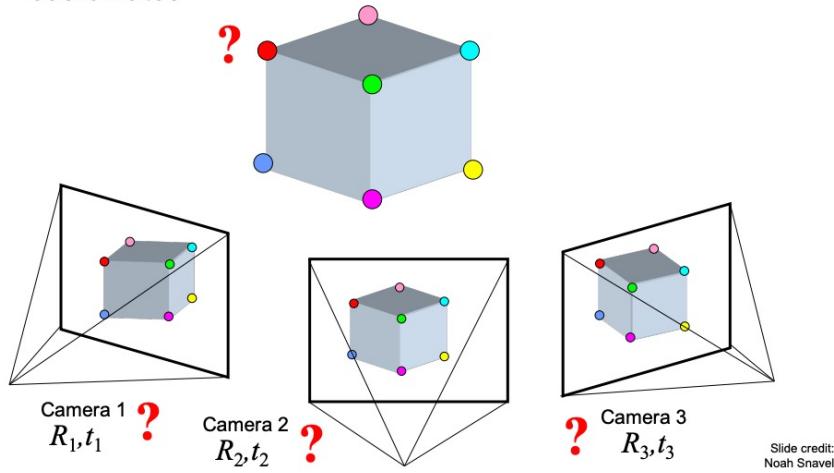
Summary

- Leverage cross correlation structure for spatial similarity matching.
- Can be used towards: depth, flow, and pose prediction.
- Can form triangulation for self-supervision.

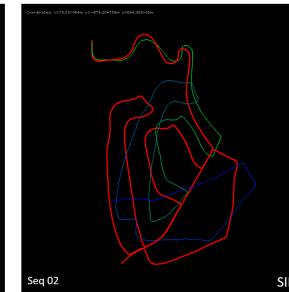
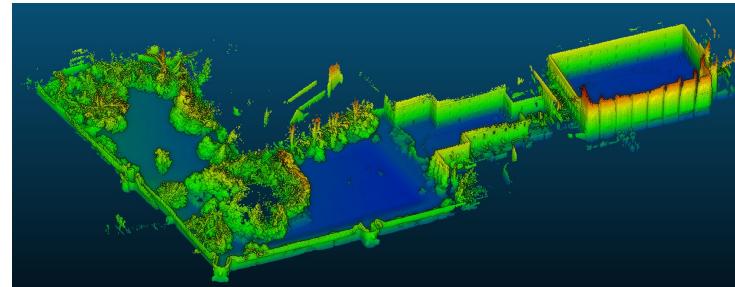


Classical Mapping

- Estimating 3D structure and location from 2D observations.
 - Simultaneous Localization and Mapping.
 - Common Techniques: Extended Kalman Filter, GraphSLAM
- Given a set of corresponding points in two or more images, compute the camera parameters and the 3D point coordinates



Slide credit:
Noah Snavely



Common Drawbacks

- Probabilistic inference can take long to compute, and mapping takes a large memory storage.

Common Drawbacks

- Probabilistic inference can take long to compute, and mapping takes a large memory storage.
- Great for 3D reconstruction but downstream tasks may not need a full precision explicit map.

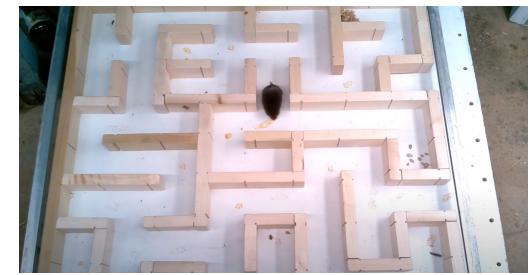
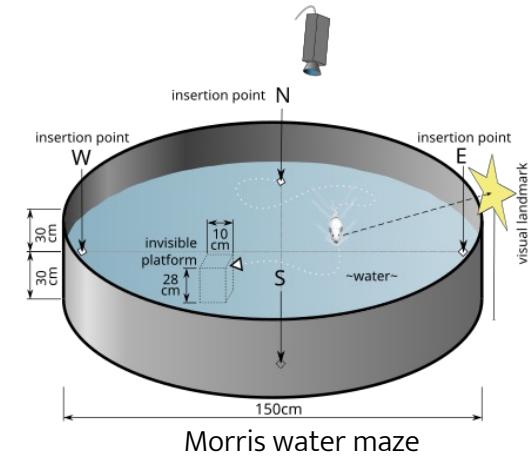
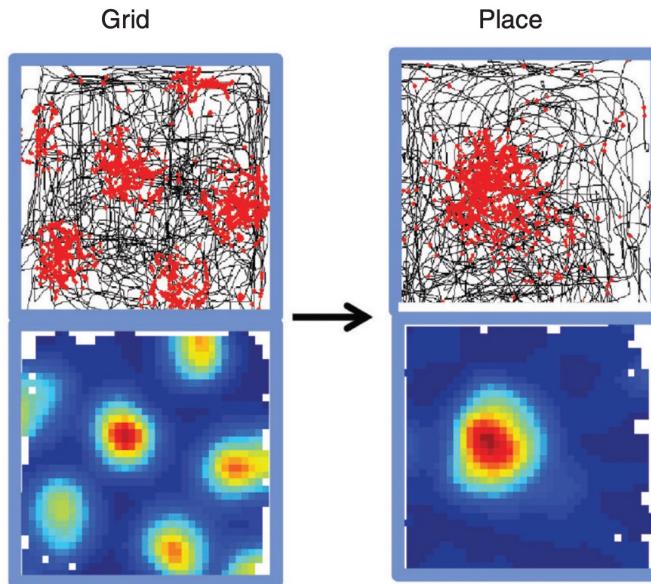
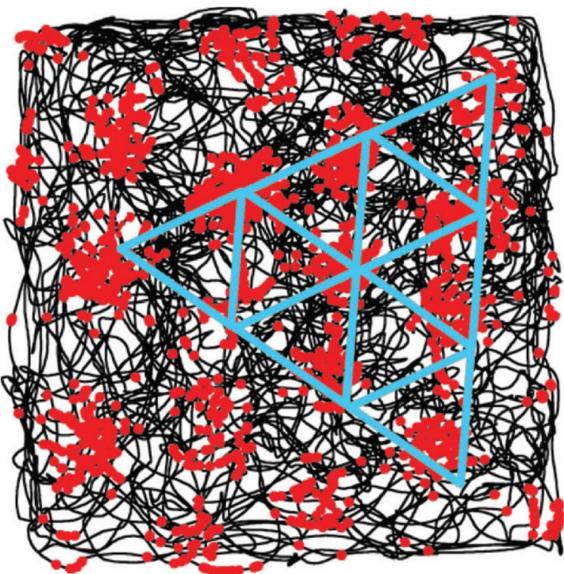
Common Drawbacks

- Probabilistic inference can take long to compute, and mapping takes a large memory storage.
- Great for 3D reconstruction but downstream tasks may not need a full precision explicit map.
- May not fully understand dynamic objects (averaging across multiple scans).

Common Drawbacks

- Probabilistic inference can take long to compute, and mapping takes a large memory storage.
- Great for 3D reconstruction but downstream tasks may not need a full precision explicit map.
- May not fully understand dynamic objects (averaging across multiple scans).
- Is there a more end-to-end version?

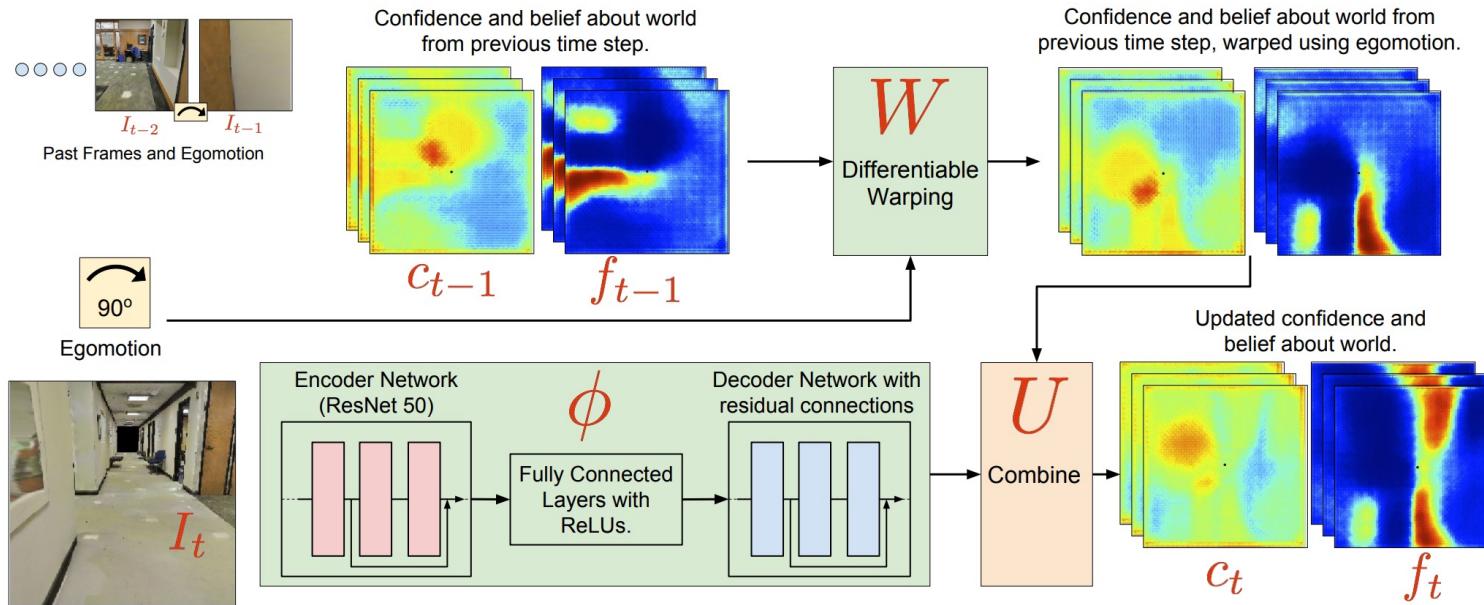
Mapping in the Brain: Grid and Place Cells



Matthias Wandel, 2018

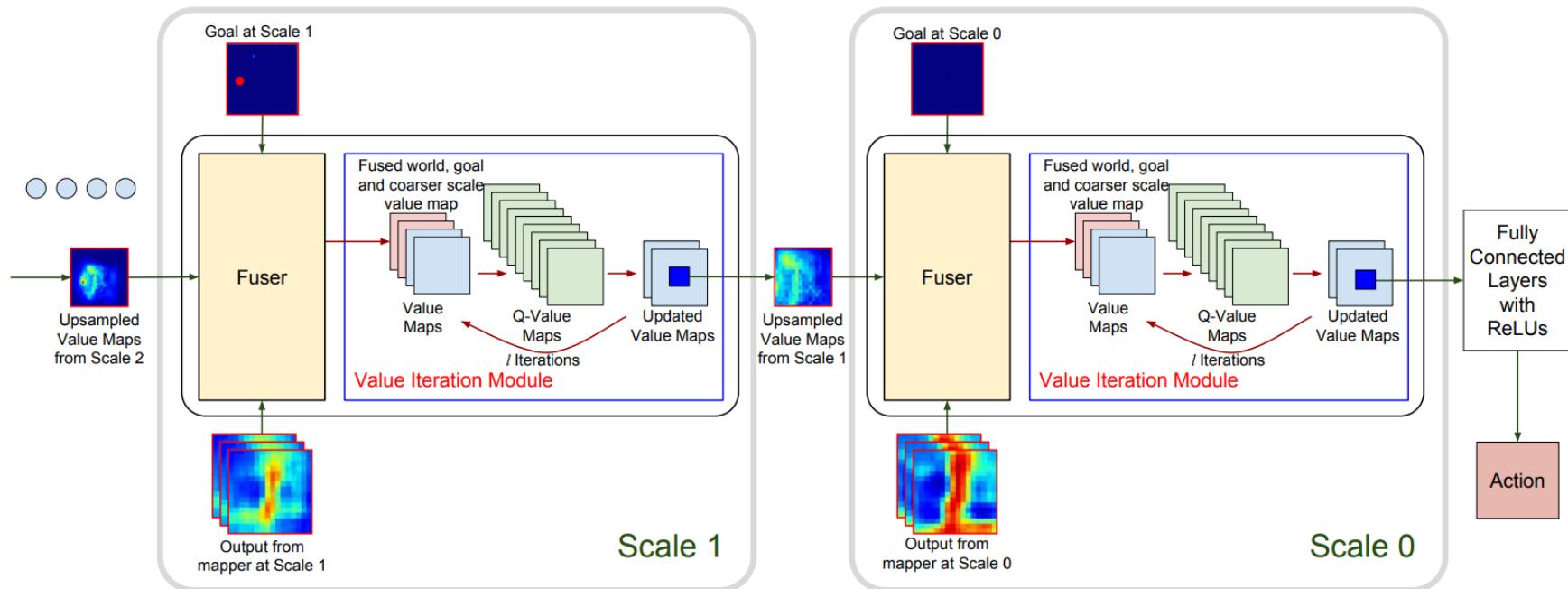
Neural Mapping

- Can we learn a mapping representation?
- Metric space, top-down warping (known egomotion).

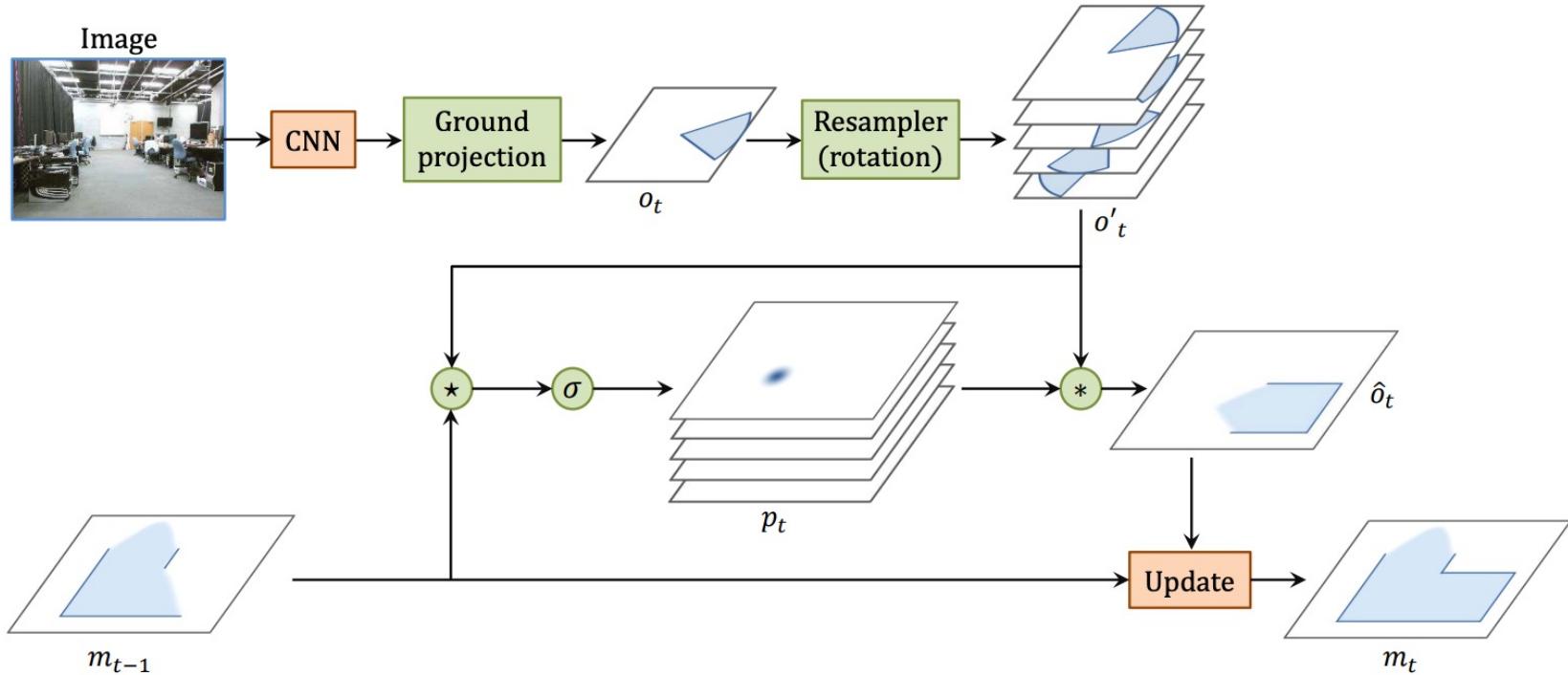


Hierarchical Planning

- How do we use the learned map (allocentric) feature of the world?



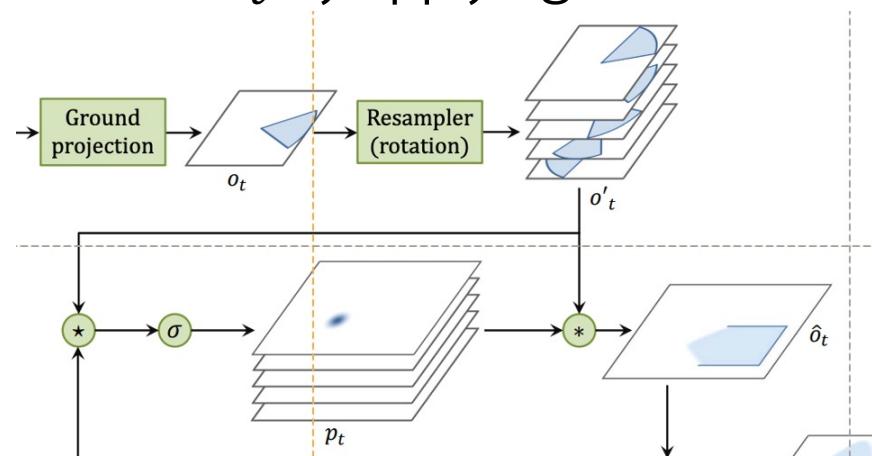
Simultaneous Localization and Registration



Simultaneous Localization and Registration

- The observations o_t are transformed into a stack o'_t by applying a rotation resampler.

$$o'_{ijkl} = [R(o, 2\pi l/r)]_{ijk}.$$



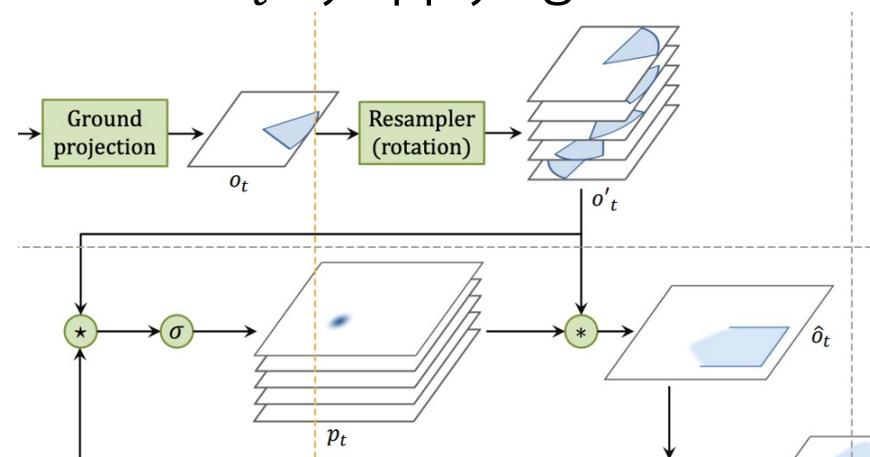
Simultaneous Localization and Registration

- The observations o_t are transformed into a stack o'_t by applying a rotation resampler.

$$o'_{ijkl} = [R(o, 2\pi l/r)]_{ijk}.$$

- o'_t convolve with the base feature.

$$p_t = \text{Softmax}(m_{t-1} * o'_t).$$



Simultaneous Localization and Registration

- The observations o_t are transformed into a stack o'_t by applying a rotation resampler.

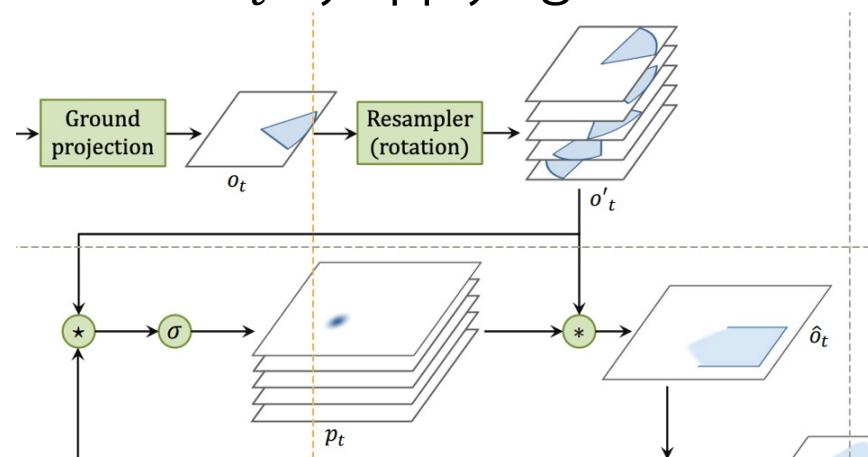
$$o'_{ijkl} = [R(o, 2\pi l/r)]_{ijk}.$$

- o'_t convolve with the base feature.

$$p_t = \text{Softmax}(m_{t-1} * o'_t).$$

- Transform observations into allocentric

$$\hat{o}_t = \sum_{uvw} p_{uvw} T(o|u, v, w).$$



Simultaneous Localization and Registration

- The observations o_t are transformed into a stack o'_t by applying a rotation resampler.

$$o'_{ijkl} = [R(o, 2\pi l/r)]_{ijk}.$$

- o'_t convolve with the base feature.

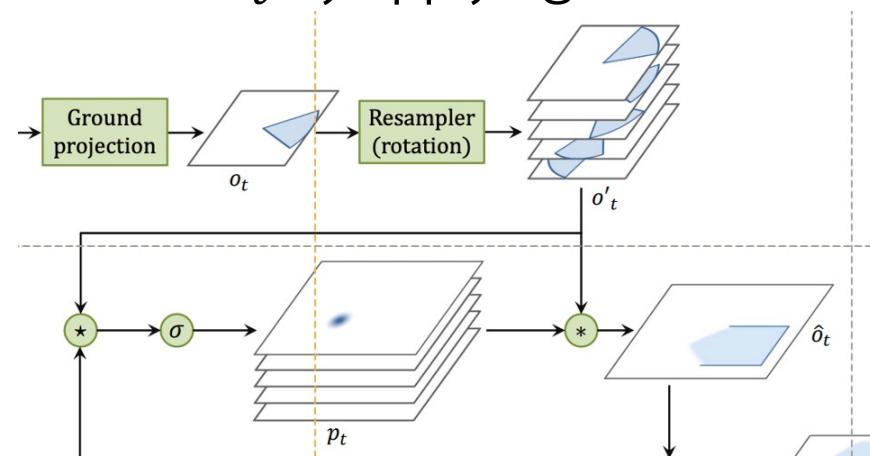
$$p_t = \text{Softmax}(m_{t-1} * o'_t).$$

- Transform observations into allocentric

$$\hat{o}_t = \sum_{uvw} p_{uvw} T(o|u, v, w).$$

- Update belief:

$$m_{i,j,t+1} = \text{LSTM}(m_{i,j,t}, \hat{o}_{i,j,t}).$$



Simultaneous Localization and Registration

- The observations o_t are transformed into a stack o'_t by applying a rotation resampler.

$$o'_{ijkl} = [R(o, 2\pi l/r)]_{ijk}.$$

- o'_t convolve with the base feature.

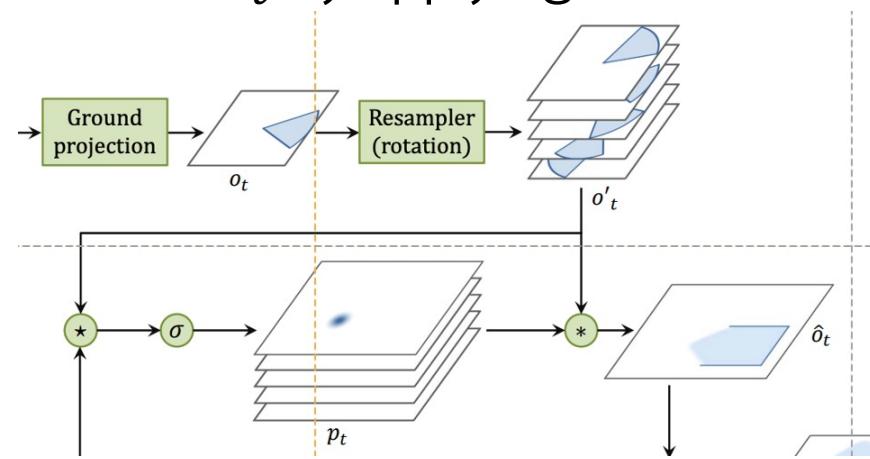
$$p_t = \text{Softmax}(m_{t-1} * o'_t).$$

- Transform observations into allocentric

$$\hat{o}_t = \sum_{uvw} p_{uvw} T(o|u, v, w).$$

- Update belief:

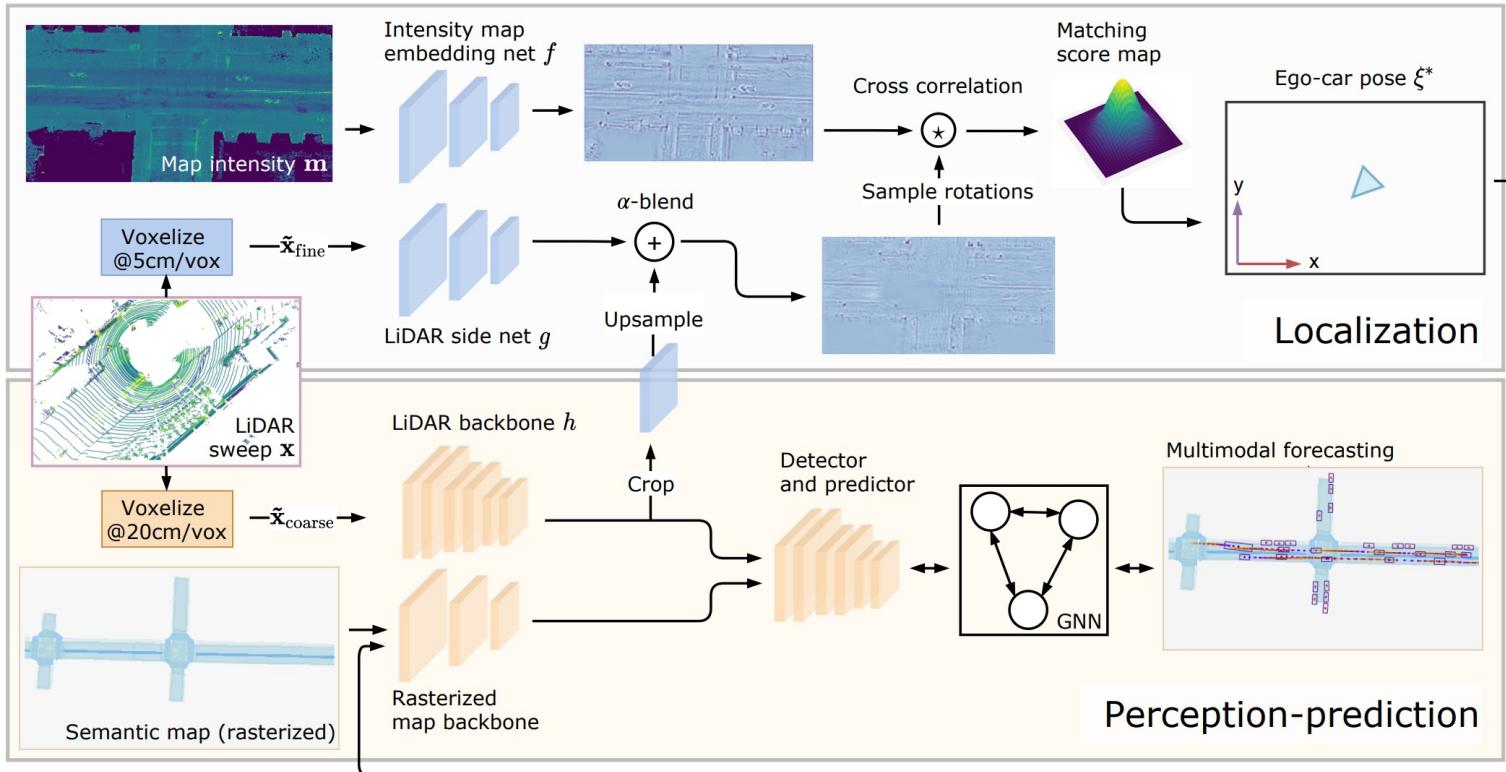
$$m_{i,j,t+1} = \text{LSTM}(m_{i,j,t}, \hat{o}_{i,j,t}).$$



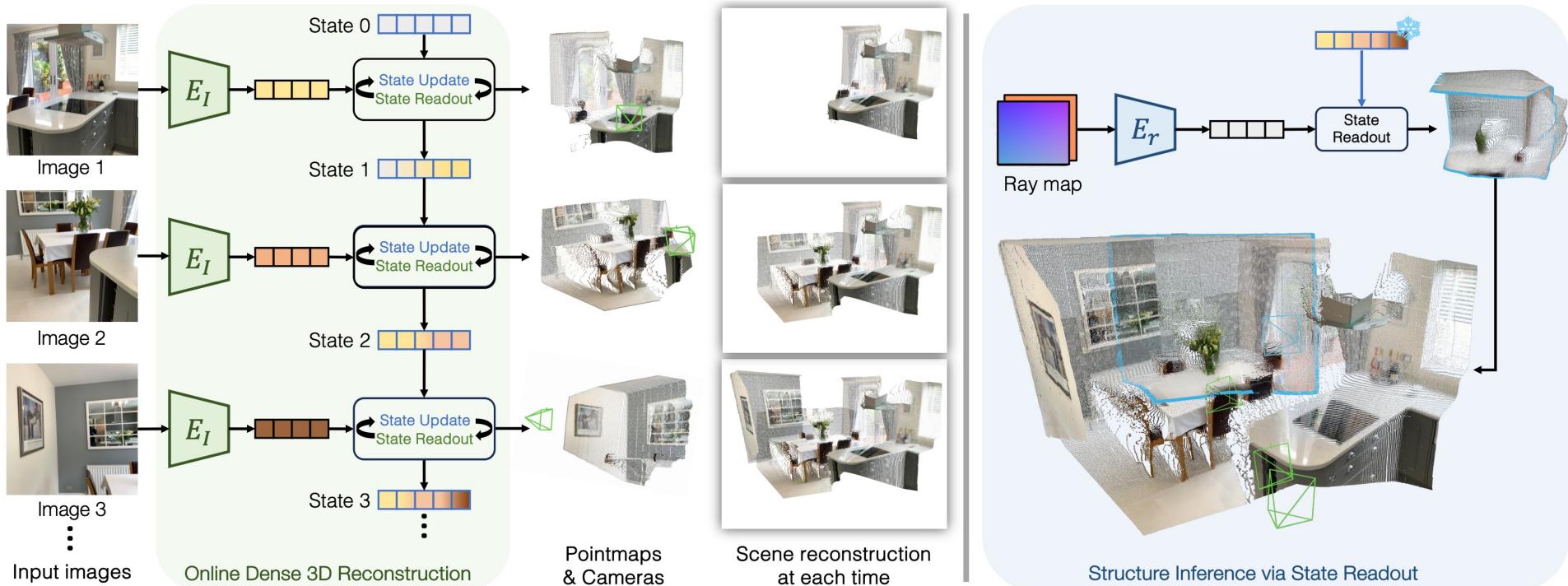
Loss:

$$\mathcal{L}(p) = -\log \sum_t p_{H_t W_t R_t t}.$$

Joint Localization, Perception and Prediction

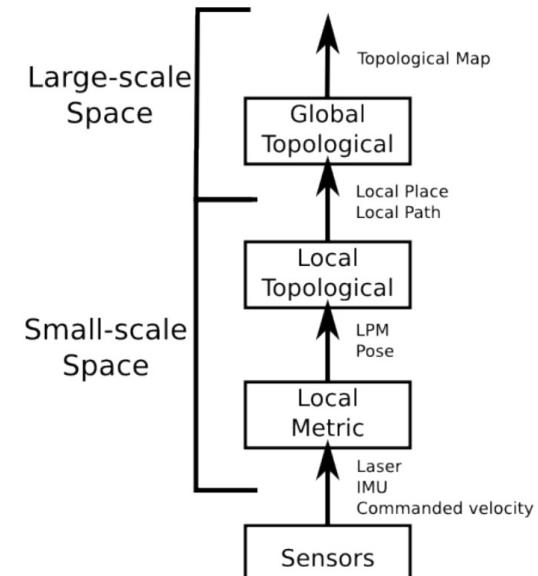
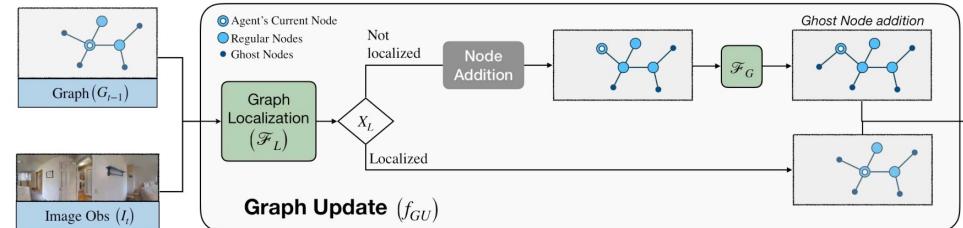
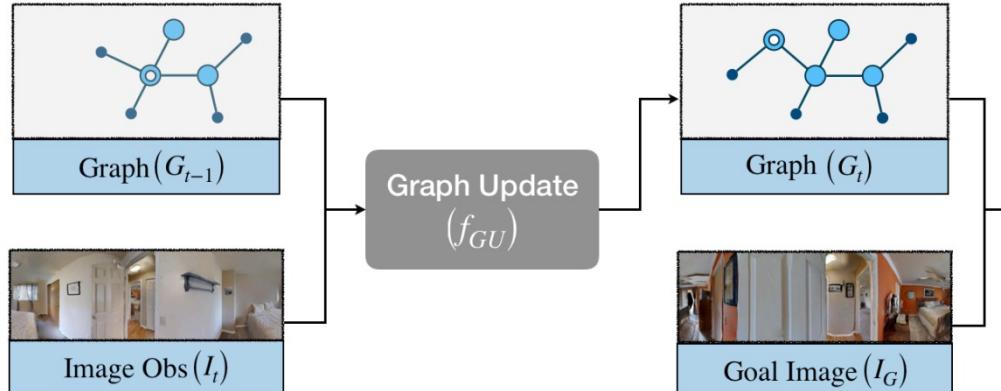


Continuous 3D Perception and Mapping



Topological Mapping

- High-level graph representation
- Each node contains more summarized information
- Enables global planning



Summary

Summary

- Covers 3D, motion, depth, and mapping.

Summary

- Covers 3D, motion, depth, and mapping.
- Still needs high-level features (recognizing the object and semantics):
Spatial pyramid.

Summary

- Covers 3D, motion, depth, and mapping.
- Still needs high-level features (recognizing the object and semantics):
Spatial pyramid.
- Can be made unsupervised

Summary

- Covers 3D, motion, depth, and mapping.
- Still needs high-level features (recognizing the object and semantics):
Spatial pyramid.
- Can be made unsupervised
- Design end-to-end modules that contain rich features.

Summary

- Covers 3D, motion, depth, and mapping.
- Still needs high-level features (recognizing the object and semantics): Spatial pyramid.
- Can be made unsupervised
- Design end-to-end modules that contain rich features.
- Design joint learning frameworks.

Summary

- Covers 3D, motion, depth, and mapping.
- Still needs high-level features (recognizing the object and semantics):
Spatial pyramid.
- Can be made unsupervised
- Design end-to-end modules that contain rich features.
- Design joint learning frameworks.
- Using geometric transformation to ground representations.

Summary

- Covers 3D, motion, depth, and mapping.
- Still needs high-level features (recognizing the object and semantics): Spatial pyramid.
- Can be made unsupervised
- Design end-to-end modules that contain rich features.
- Design joint learning frameworks.
- Using geometric transformation to ground representations.
- Useful for planning (a few weeks from now).