

Summary: Normalizing Flow

- Requires invertibility and tractability on Jacobian determinant
- Constraint on the functional form
- Or estimation of Jacobian / iterative inverse (numerical approximation).



Inference as Integration

- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

Inference as Integration

- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

- Can be written in continuous time (continuous normalizing flow).

$$x_t = \phi_t(x_0) = x_0 + \int_0^t \delta u_s(x_s) ds.$$

Inference as Integration

- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

- Can be written in continuous time (continuous normalizing flow).

$$x_t = \phi_t(x_0) = x_0 + \int_0^t \delta u_s(x_s) ds.$$

- Integration: Euler, or directly run an ODE solver

Inference as Integration

- Residual Flow: $\phi(x) = x + \delta u(x)$

$$\frac{\phi(x) - x}{\delta} = u(x)$$

- Can be written in continuous time (continuous normalizing flow).

$$x_t = \phi_t(x_0) = x_0 + \int_0^t \delta u_s(x_s) ds.$$

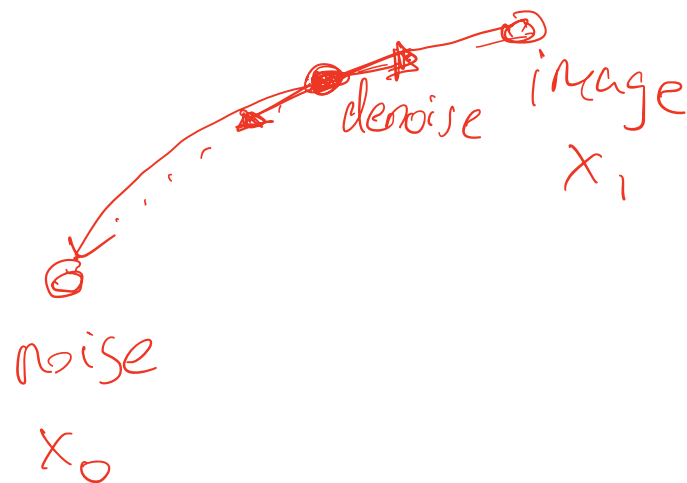
- Integration: Euler, or directly run an ODE solver
- Backprop solving an ODE without storing intermediate states

Flow Matching

- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.

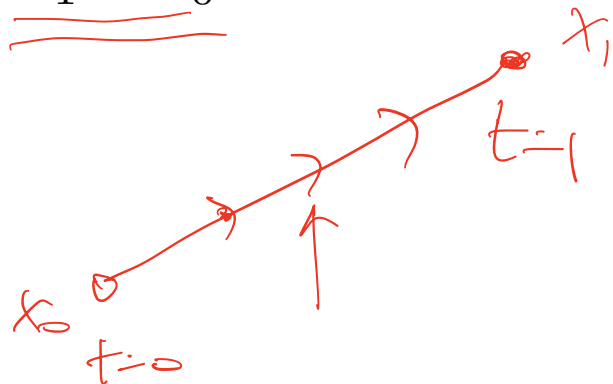
Flow Matching

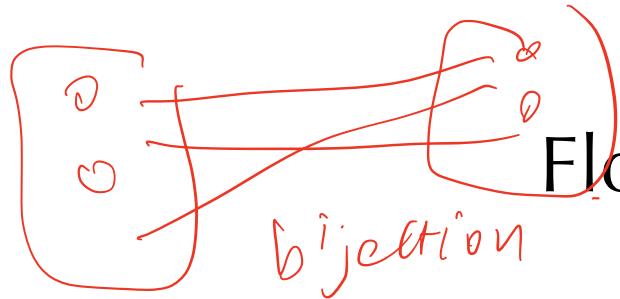
- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.
- Given an original input x_1 and a noise x_0 , $x_t = (1 - t)x_0 + tx_1$



Flow Matching

- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.
- Given an original input x_1 and a noise x_0 , $x_t = (1 - t)x_0 + tx_1$
- Want to fit a velocity $v_t(x)$ to match $u_t = x_1 - x_0$





Flow Matching

deterministic
 ○ noise.
 ○

- Taking integrals during training is cumbersome, and the learning signal from the top latent is weak.
- Given an original input x_1 and a noise x_0 , $x_t = (1 - t)x_0 + tx_1$
- Want to fit a velocity $v_t(x)$ to match $u_t = x_1 - x_0$
- Flow Matching learning objective:

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t, x_0, x_1} [\| \underbrace{v(x_t, t)}_{\text{targets}} - \underbrace{u_t(x_0, x_1)}_{\text{ideal velocity}} \|^2]$$

targets
 ideal velocity.
 Euler integration →

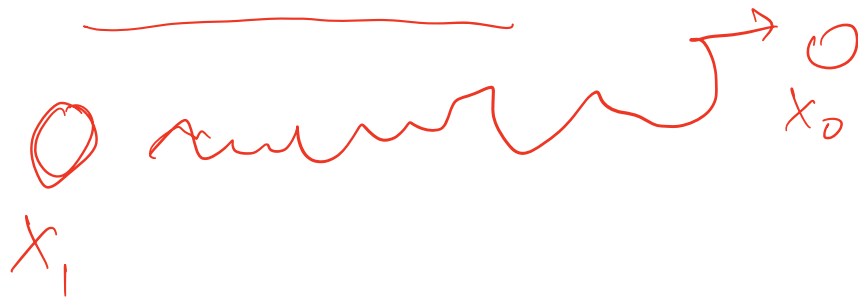
efficient.
 10-20 steps.

Diffusion Models

- A popular generative model formulation.

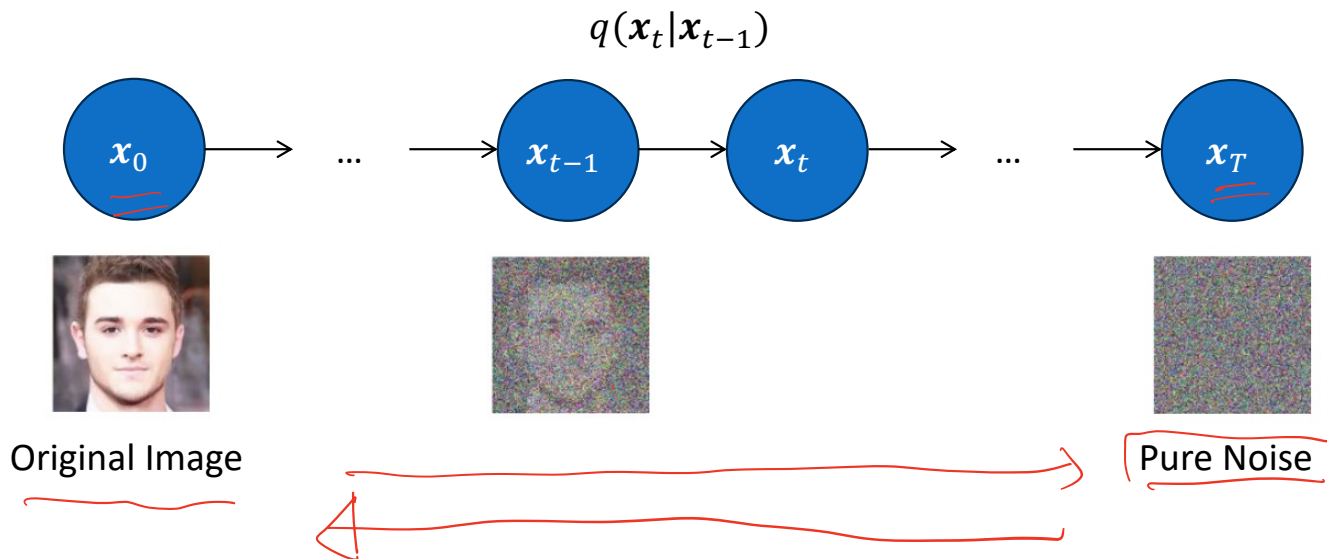
Diffusion Models

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



Diffusion Models

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



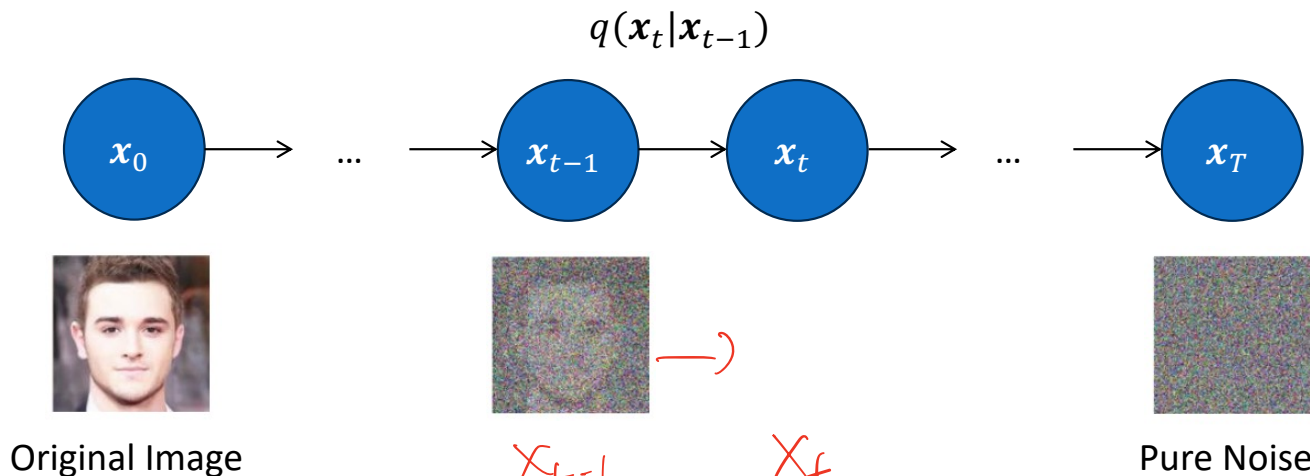
Diffusion Models

volume preserving
variance previous

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

$\sqrt{\beta_t}$

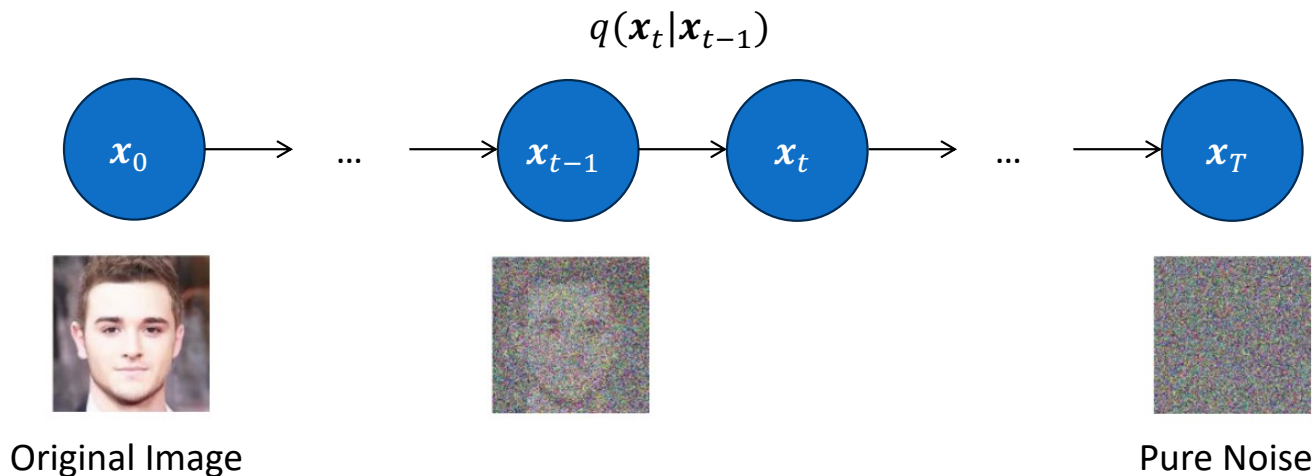
Mean. Co-Variance



x_t
noisier.

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

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:

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

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

*amount
of original.*

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:

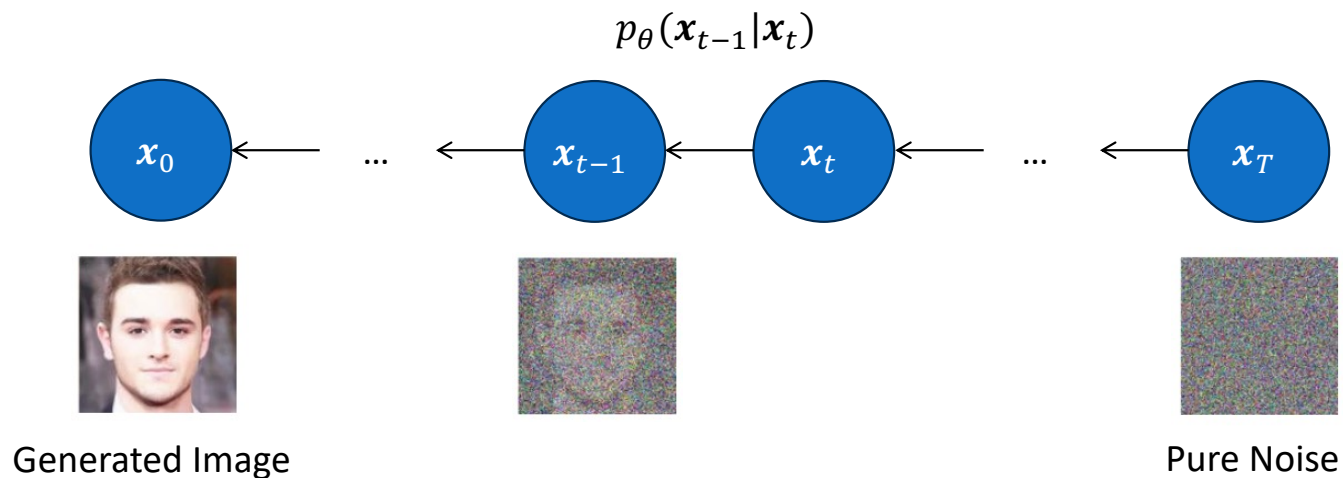
$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I).$$

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

- $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$.

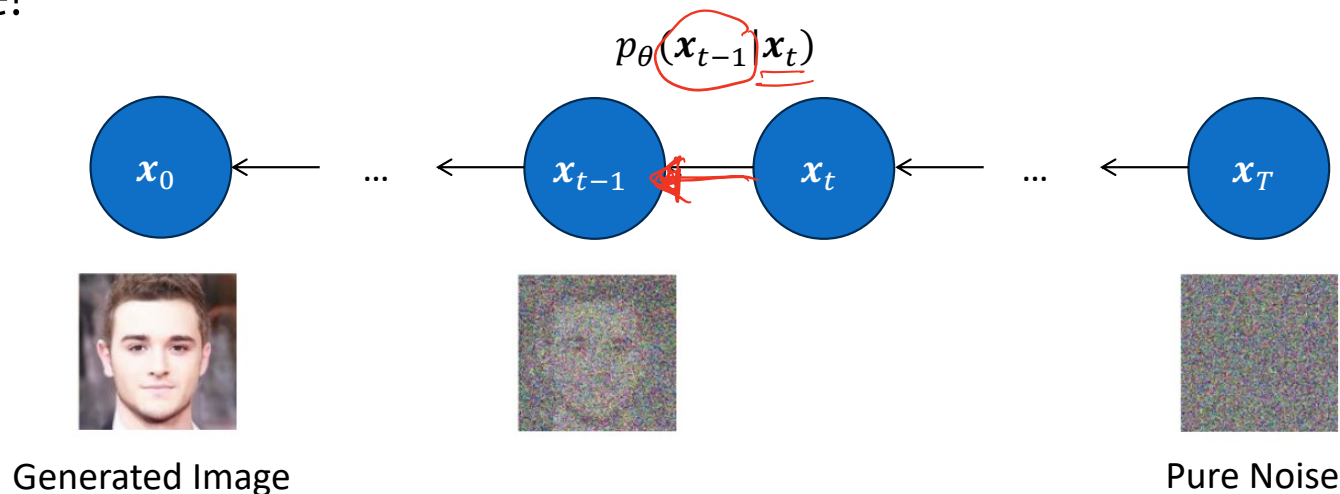
Reverse Process

- A generative model wants to predict x_0 from x_T .



Reverse Process

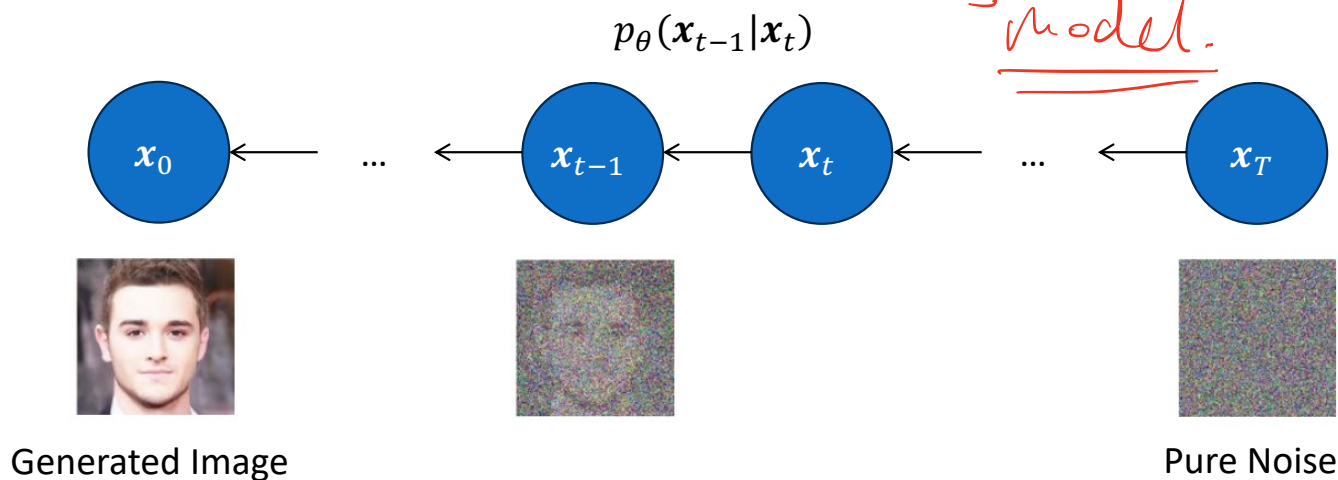
- A generative model wants to predict x_0 from x_T .
- The reverse process transition is also Gaussian distributed. But we don't know what the transition will be like just by looking at the noisy image!



Reverse Process

- So, we need to learn a “model”:

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



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:

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

forward process,

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

learning approx. $p(x_{t-1}|x_t) \rightarrow \mu_\theta$

↳ training data $\rightarrow x_0$'s. learn.

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 .

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

— during training

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

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

small jump.

big jump.

big jump.

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

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

noisy

orig. img.

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 $\tilde{\mu}_t$.

Training

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

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

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

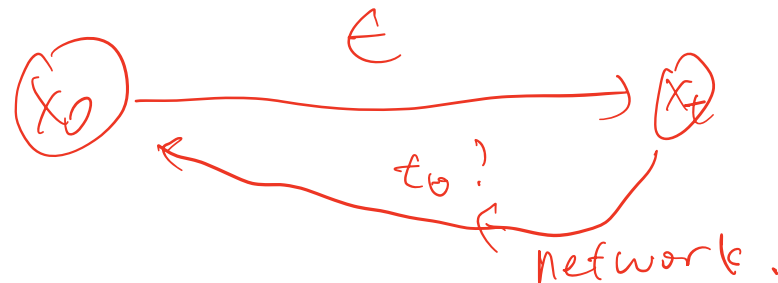
Training

- 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.

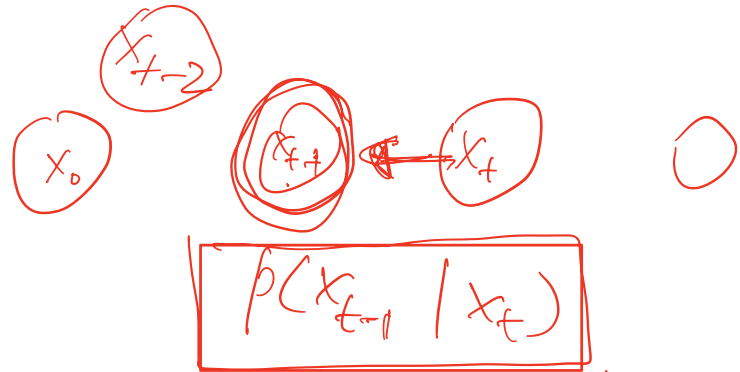
Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 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
-

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

free noise network noisy ver.

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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return** \mathbf{x}_0



denoising

Schedule.

mean

add noise

Sampler

Current image

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} .

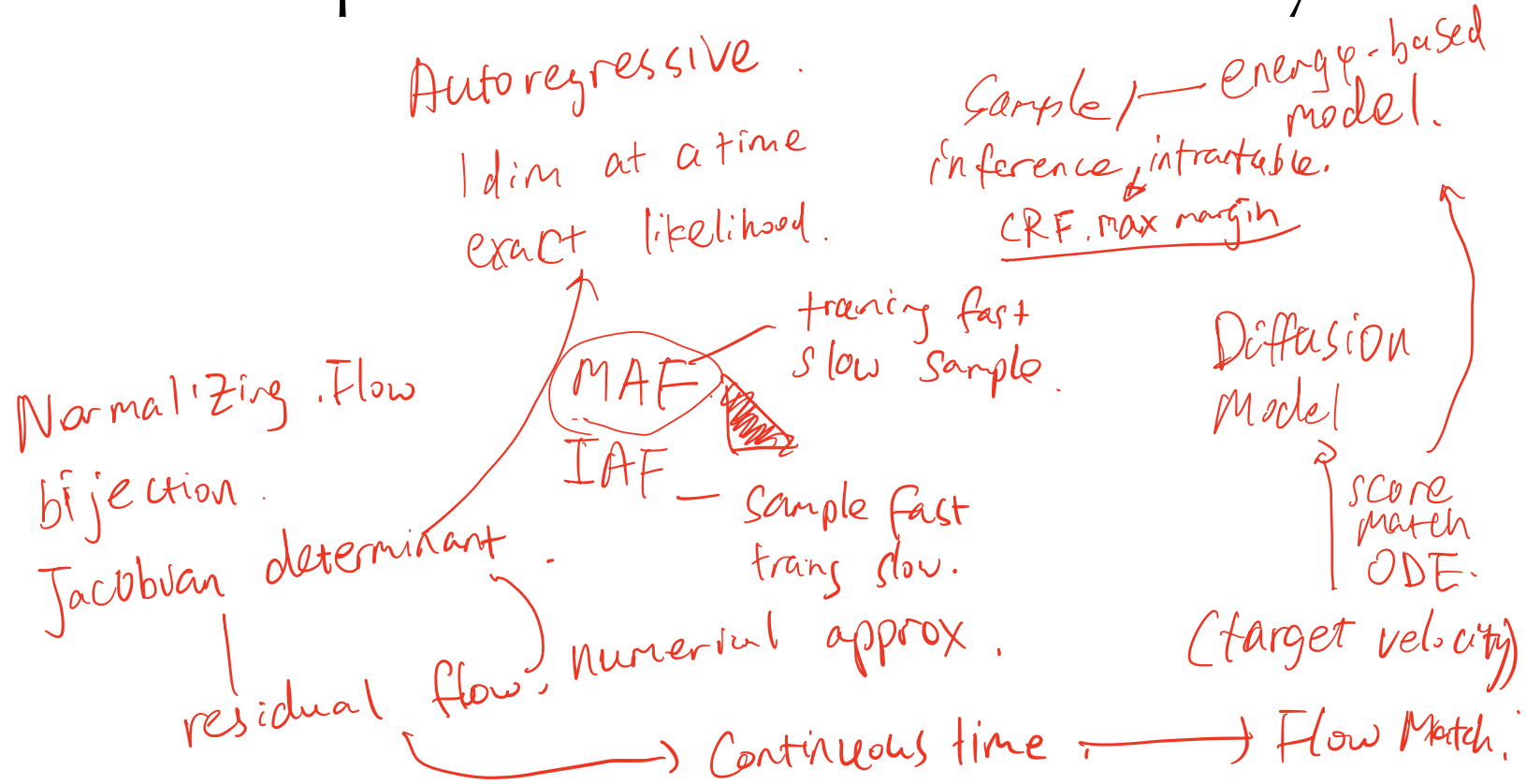
$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t) \right).$$

- 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
-

Deep Generative Models Summary



Guided Diffusion

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

Classifier Guidance / External Score Model

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

gradient from another criterion

Guided Diffusion

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

Classifier Guidance / External Score Model

$$\hat{\epsilon} \leftarrow \epsilon_{\theta}(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_{\phi}(y|x_t)$$

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

- We also can train a conditional diffusion model.

repeat

$$(\mathbf{x}, \mathbf{c}) \sim p(\mathbf{x}, \mathbf{c})$$

▷ Sample data with conditioning from the dataset

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

▷ Randomly discard conditioning to train unconditionally

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

▷ Sample log SNR value

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

$$\mathbf{z}_{\lambda} = \alpha_{\lambda} \mathbf{x} + \sigma_{\lambda} \epsilon$$

▷ Corrupt data to the sampled log SNR value

$$\text{Take gradient step on } \nabla_{\theta} \|\epsilon_{\theta}(\mathbf{z}_{\lambda}, \mathbf{c}) - \epsilon\|^2$$

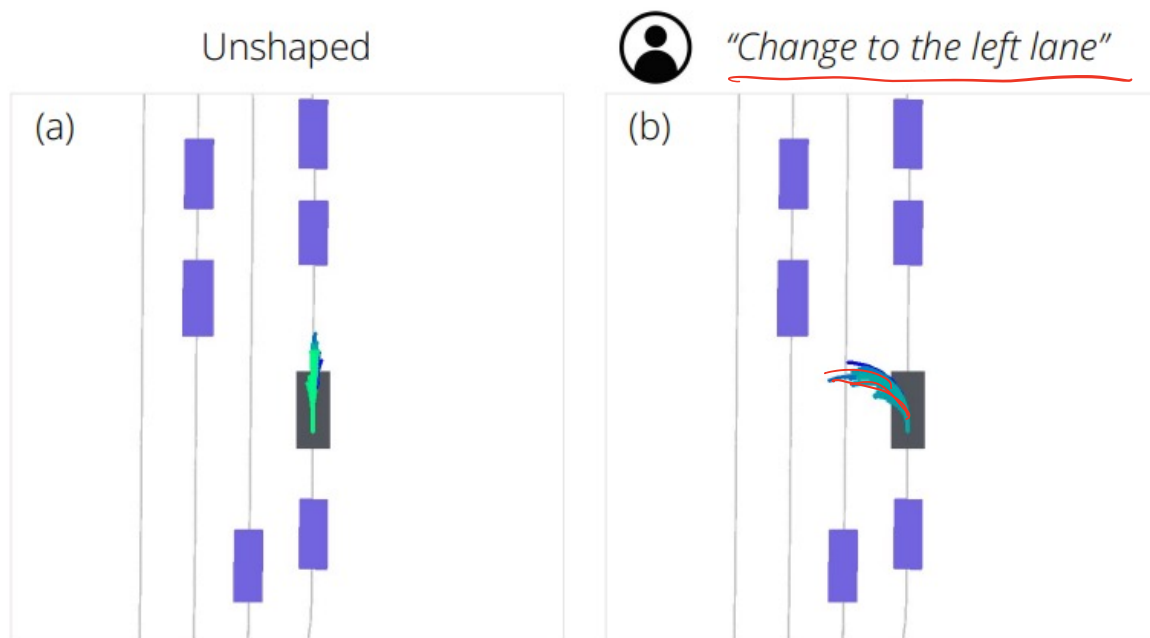
▷ Optimization of denoising model

until converged

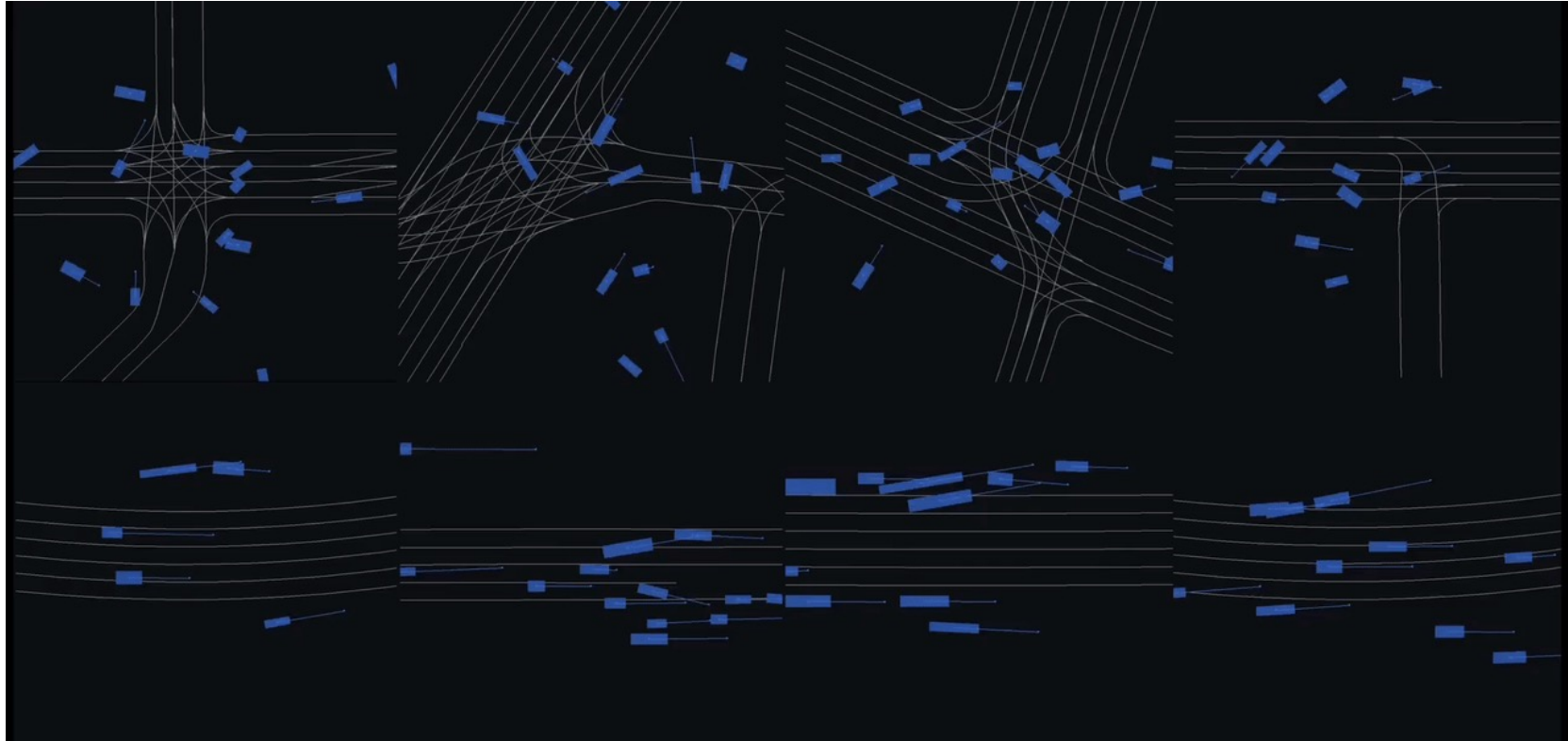
Condition

Test-Time Guidance / Adaptation

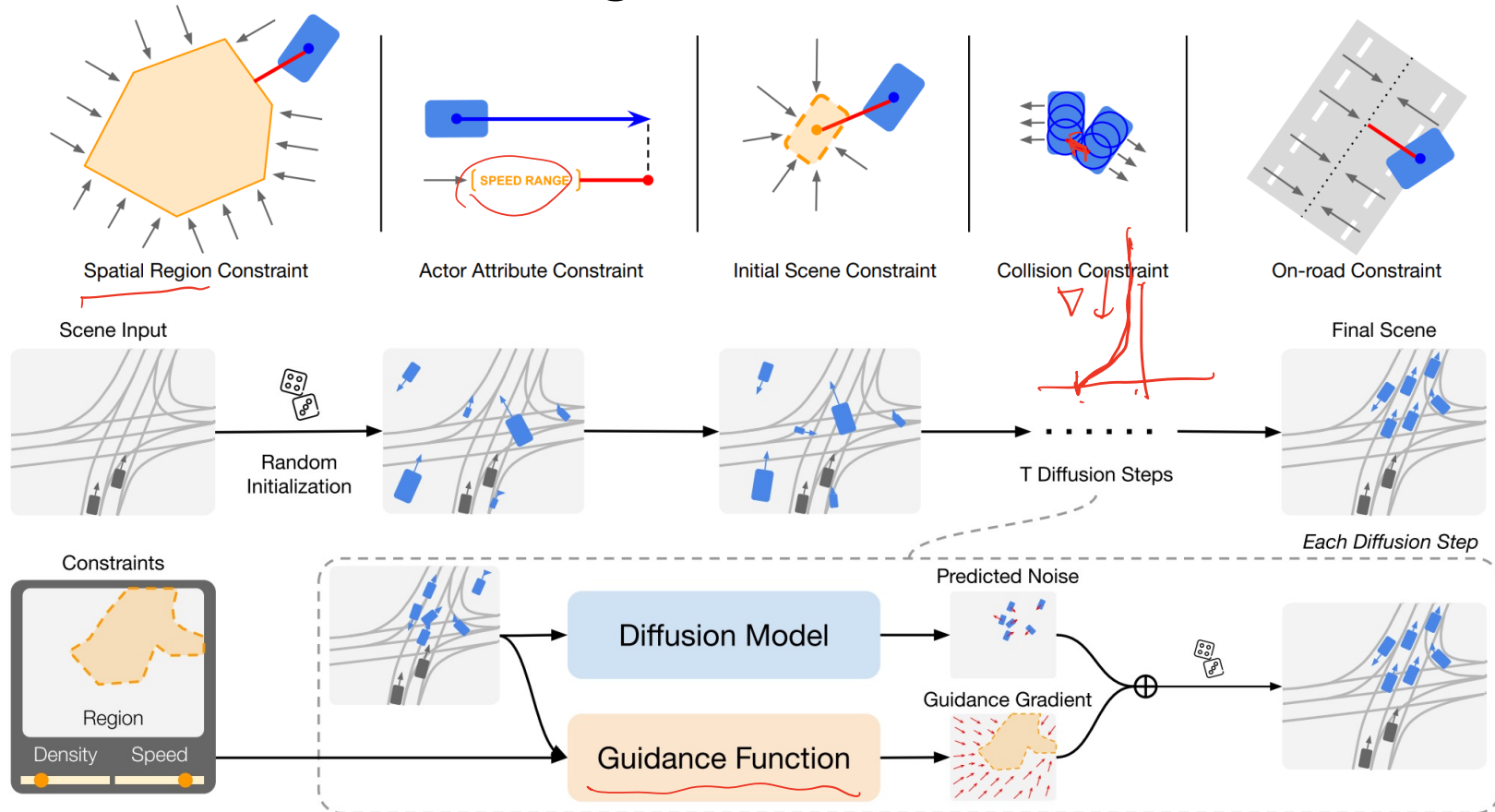
- Diffusion can be combined / guided at test time.



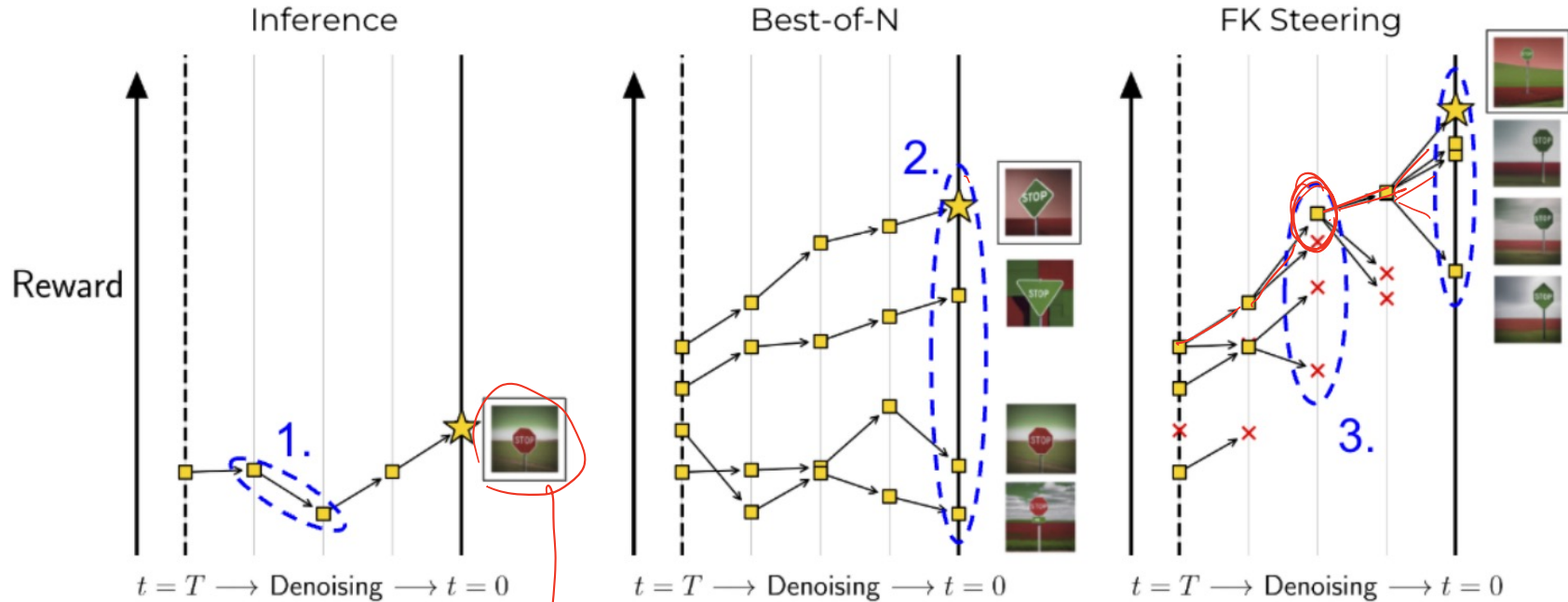
Generating Simulation Scenes



Generating Simulation Scenes



Non-Differentiable Rewards?



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.

Diffusion Models for Detection

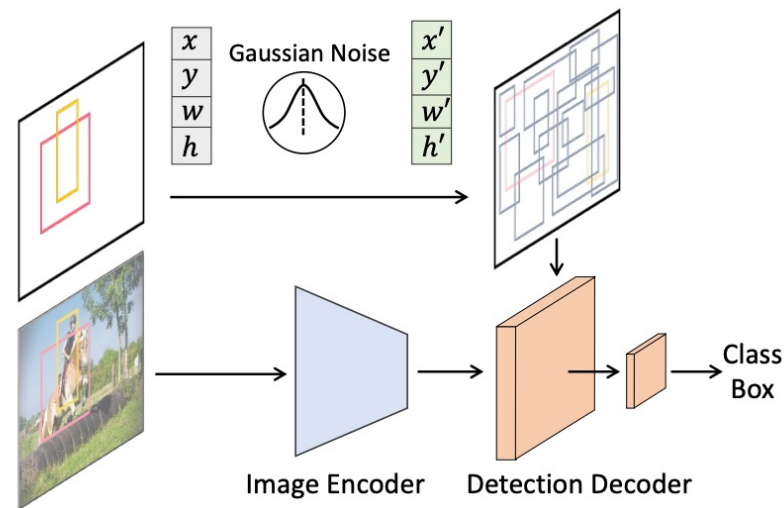
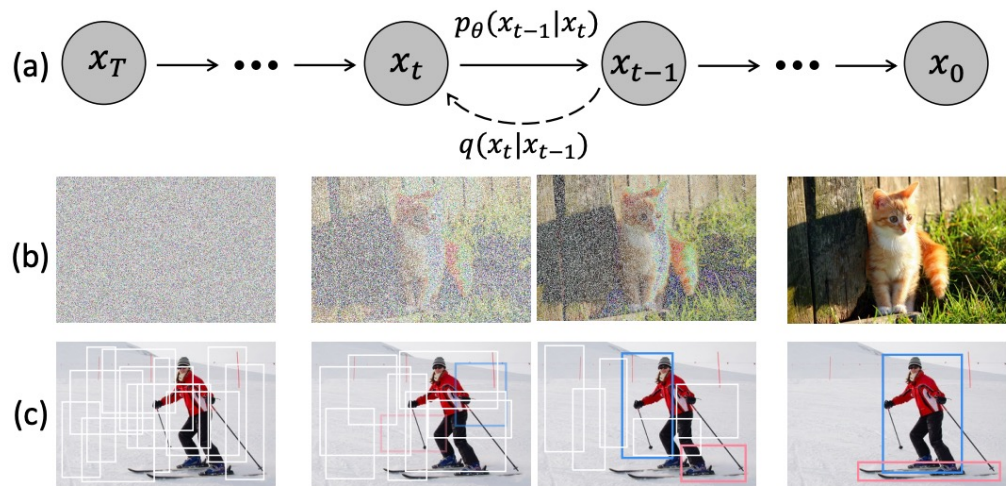
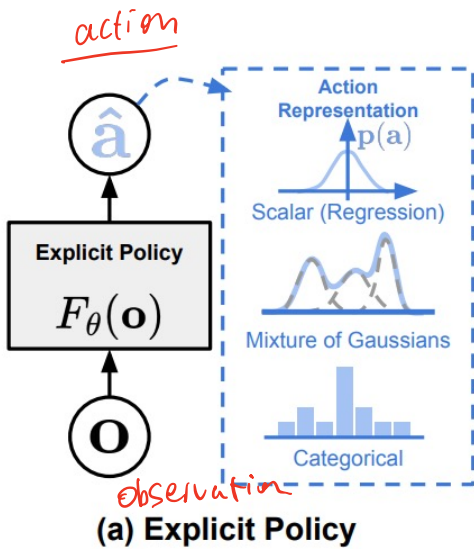
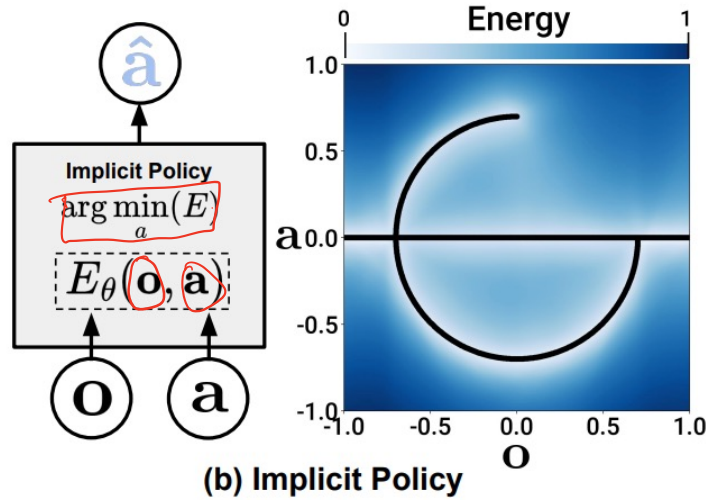
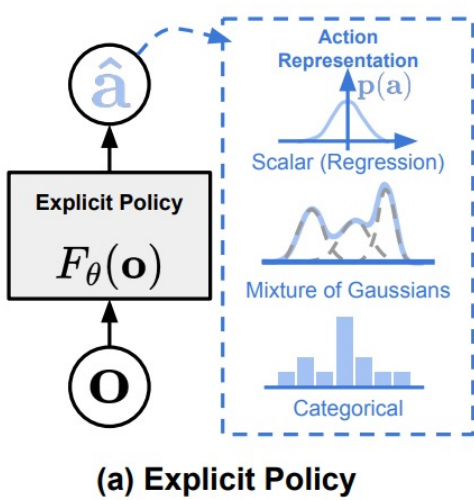


Figure 1. **Diffusion model for object detection.** (a) A diffusion model where q is the diffusion process and p_θ is the reverse process. (b) Diffusion model for image generation task. (c) We propose to formulate object detection as a denoising diffusion process from noisy boxes to object boxes.

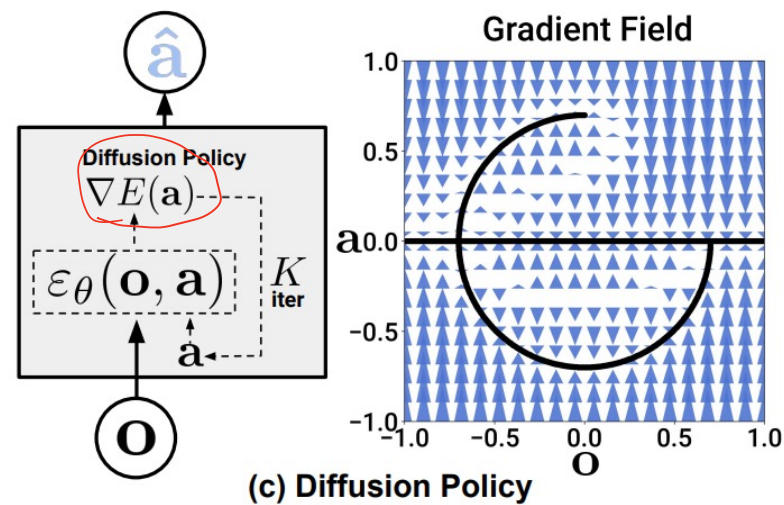
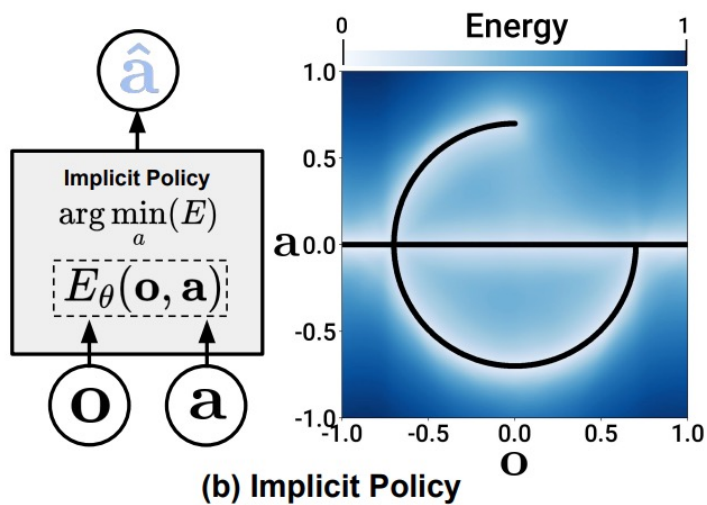
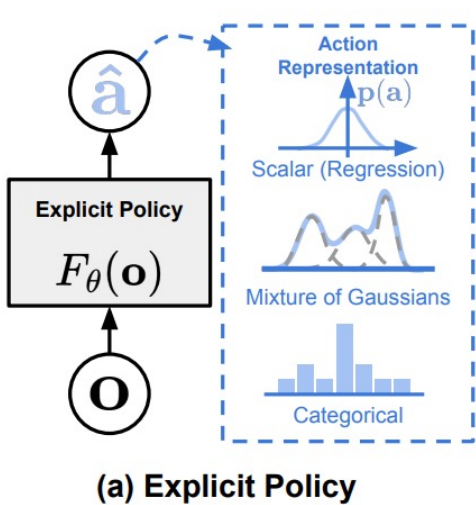
Planning and Control



Planning and Control



Planning and Control



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 deep generative models.
 - Graphical models
 - Autoregressive
 - Energy-based
 - Normalizing flow
 - Diffusion

Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about deep generative models.
 - Graphical models
 - Autoregressive
 - Energy-based
 - Normalizing flow
 - Diffusion
- Understand relations, pros and cons.

Summary: DL for Structured Outputs

- Expanding the output dimension has limitations.
- Requires us thinking about deep generative models.
 - Graphical models
 - Autoregressive
 - Energy-based
 - Normalizing flow
 - Diffusion
- Understand relations, pros and cons.
- Application in embodied environments.

Module 2: 3D Vision and Mapping

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.



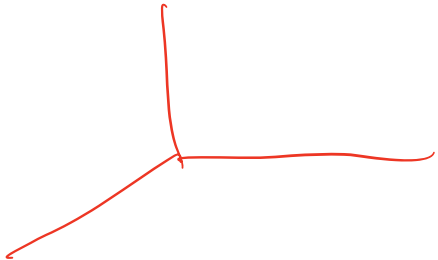
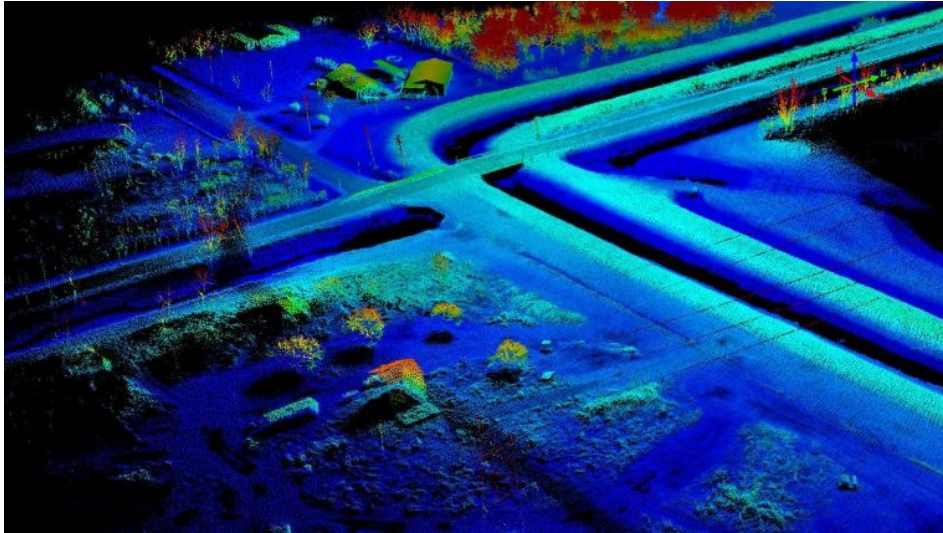
The World is 3D

- 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.



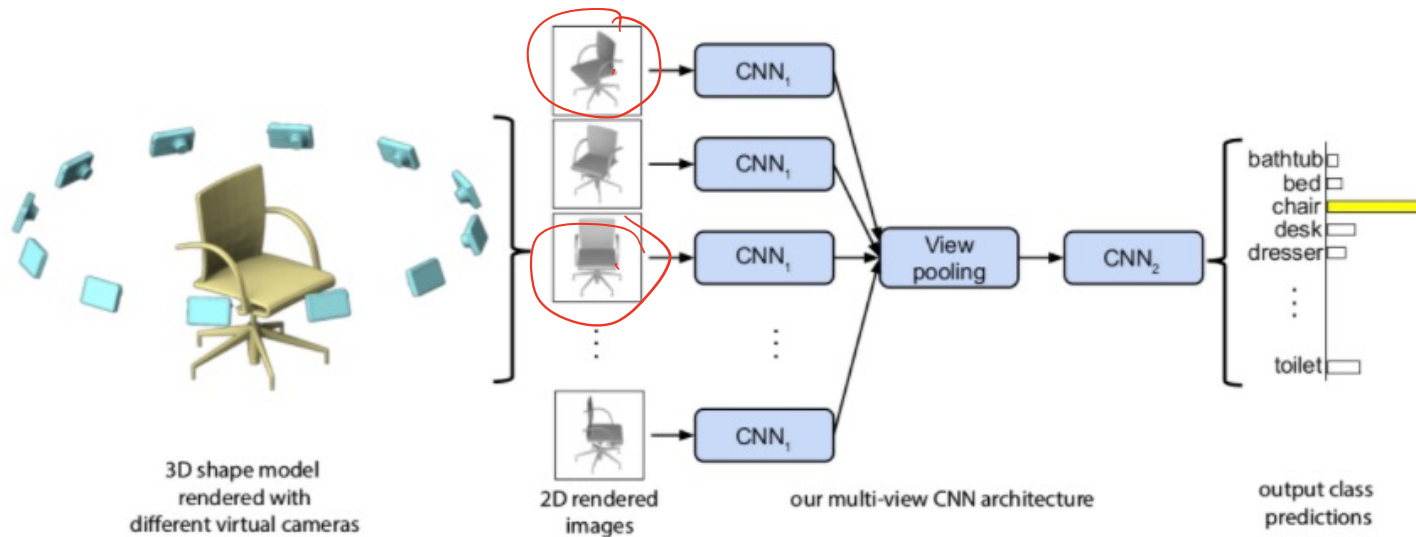
LiDAR

10k. (x,y,z) intensity.



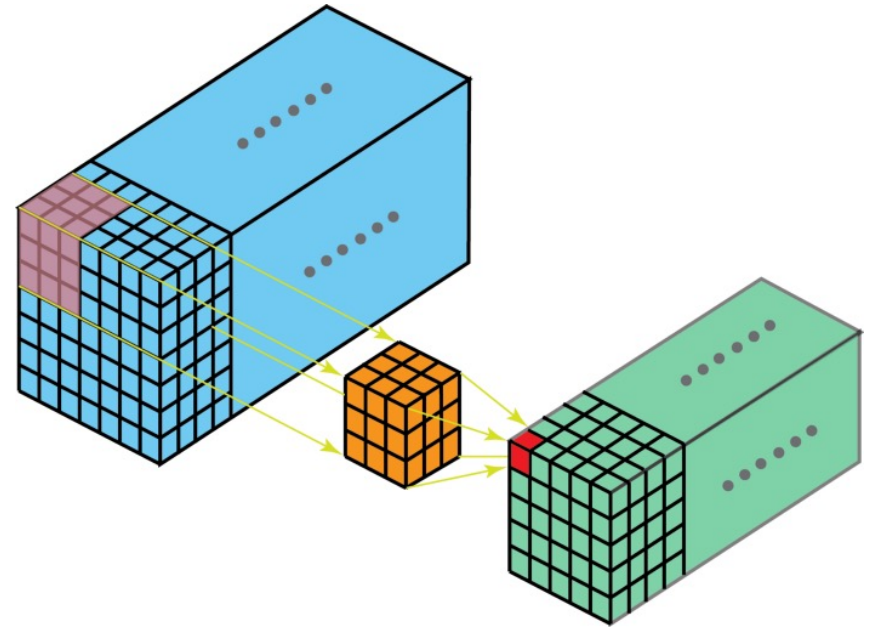
Multi-View CNN

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



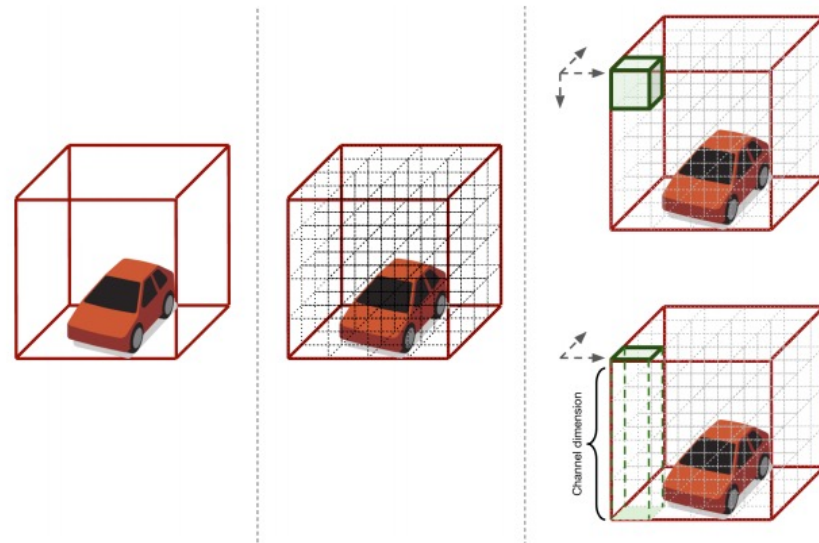
3D Convolution on Voxels

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



Bird's Eye View (BEV) Voxel

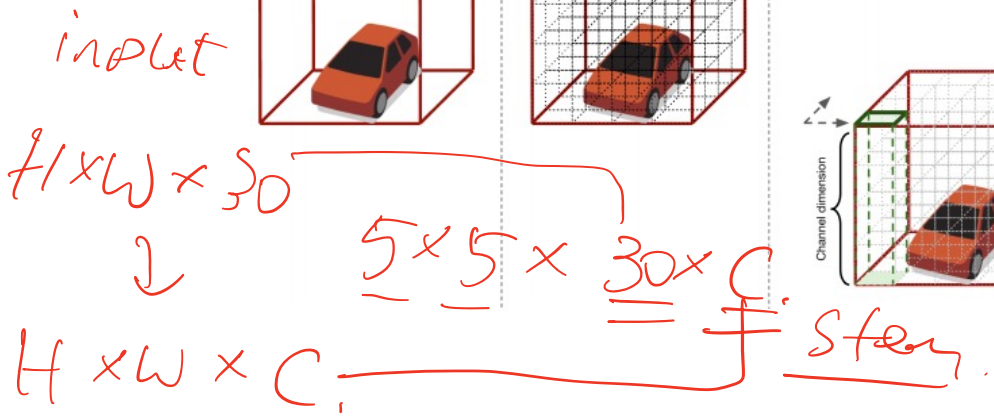
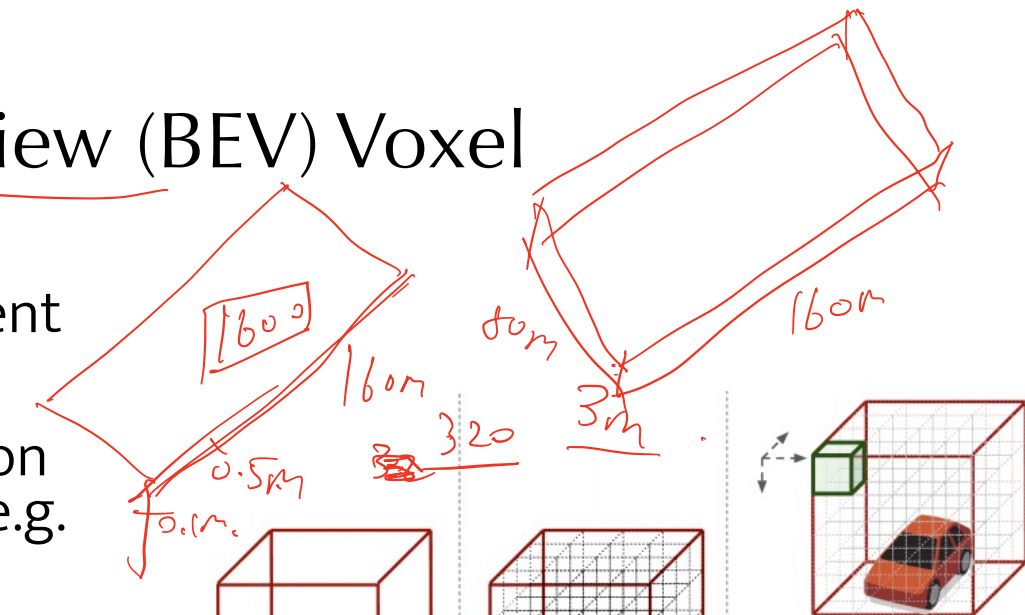
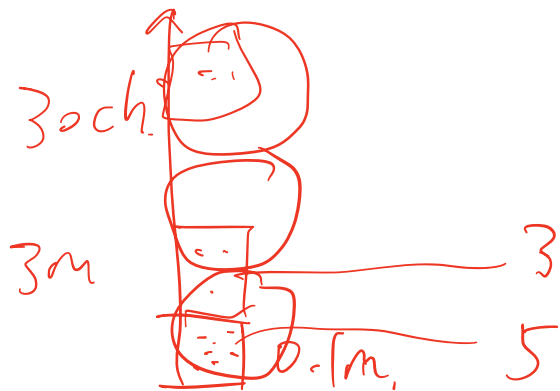
- Treat the z-dimension as different channels.



20000

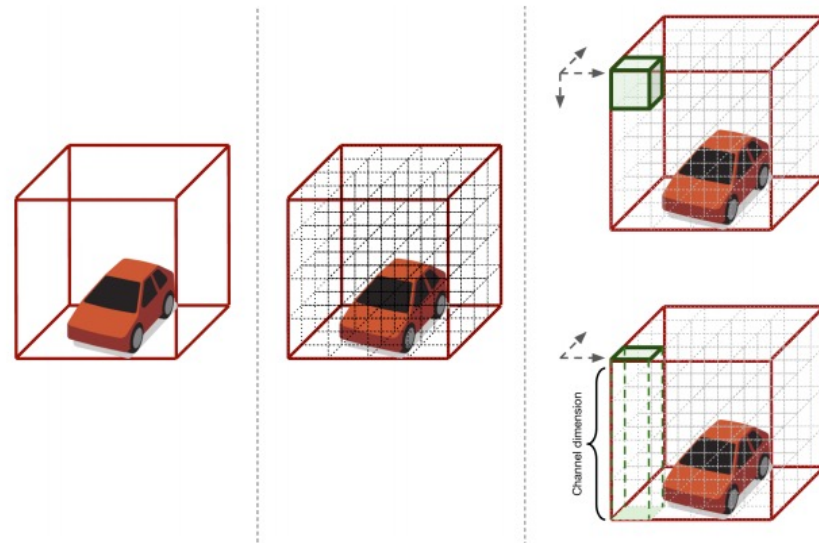
Bird's Eye View (BEV) Voxel

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



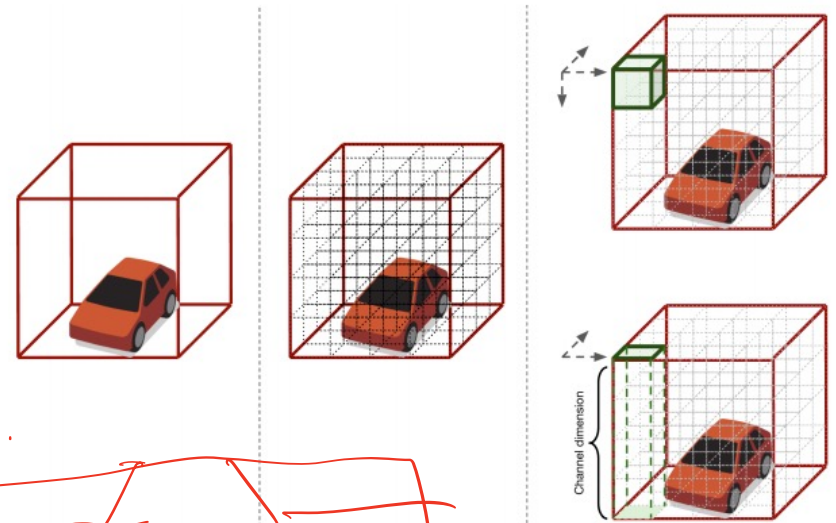
Bird's Eye View (BEV) Voxel

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

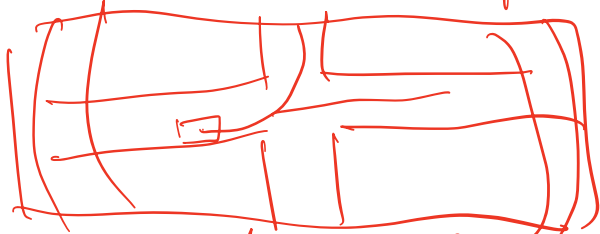


Bird's Eye View (BEV) Voxel

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



BEV.



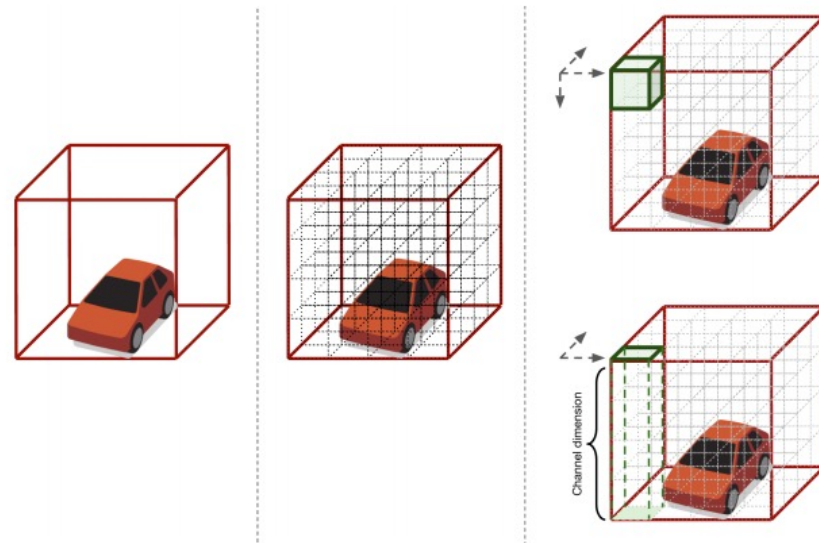
sparse (outer ring)



dense input

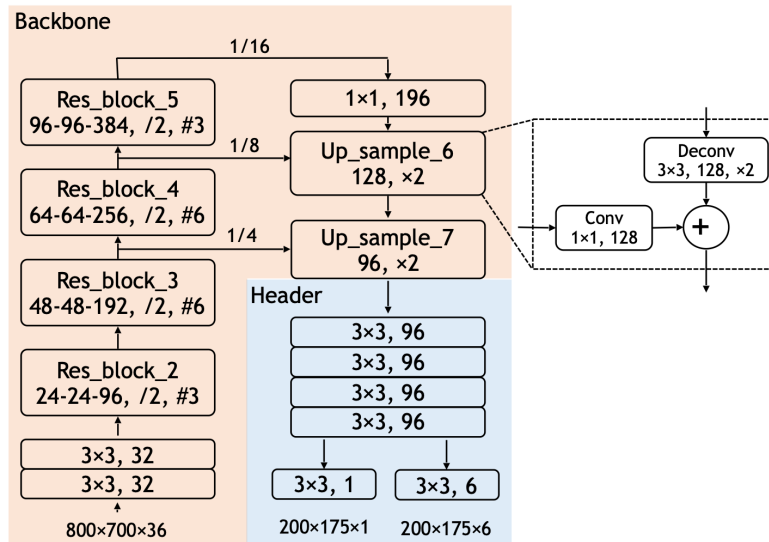
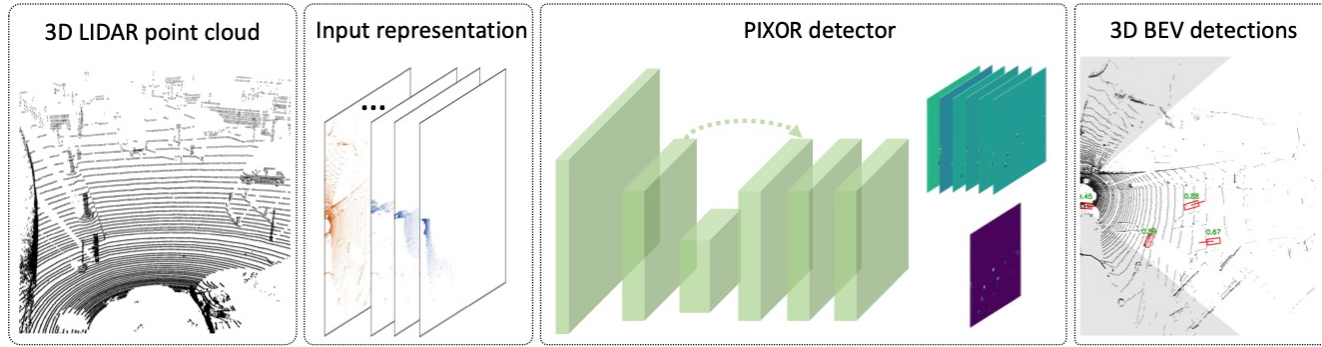
Bird's Eye View (BEV) Voxel

- Treat the z-dimension as different channels.
- 3D convolution \rightarrow 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)



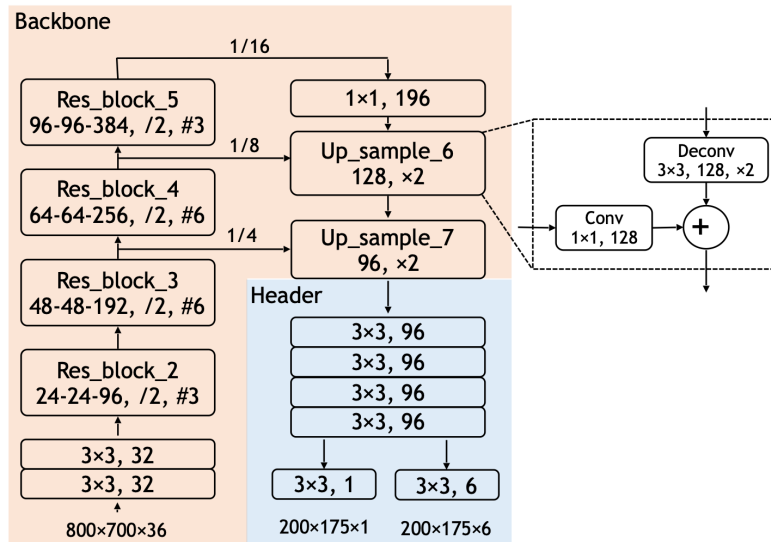
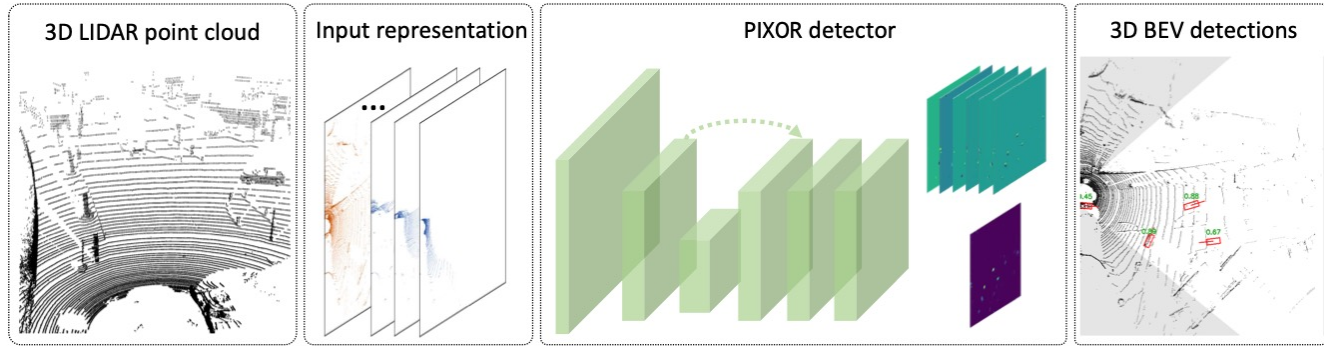
PIXOR

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



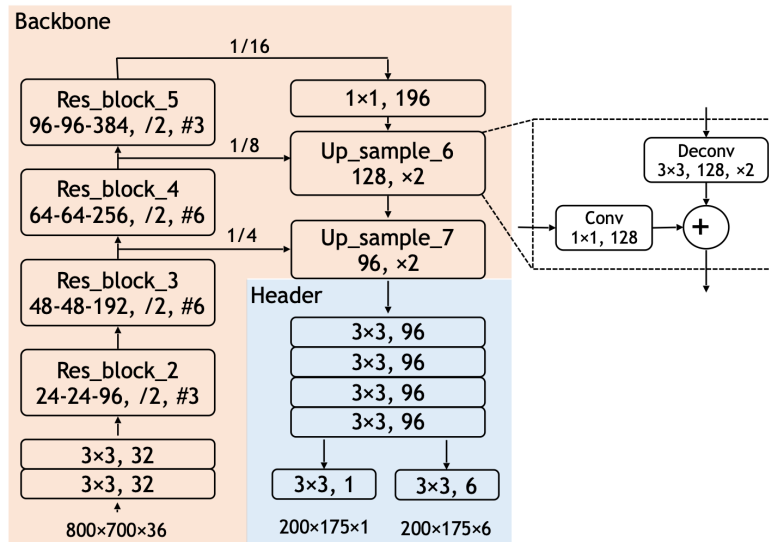
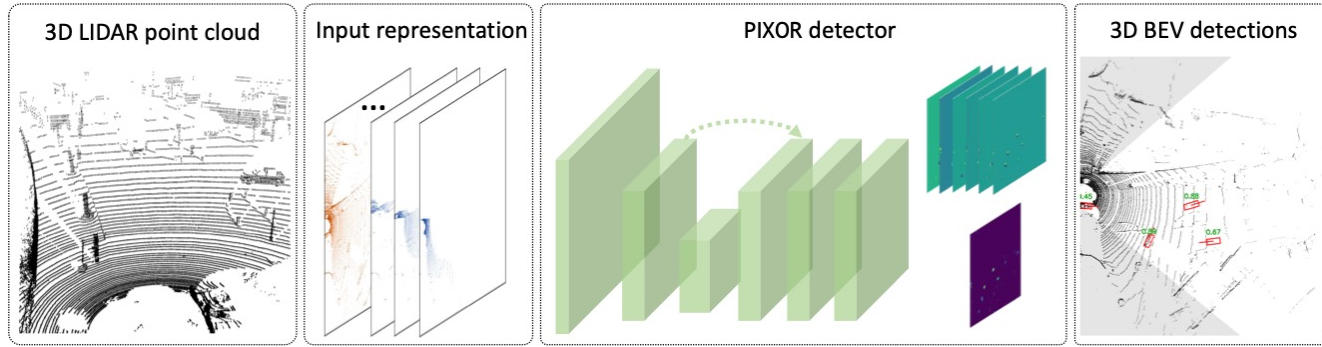
PIXOR

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



PIXOR

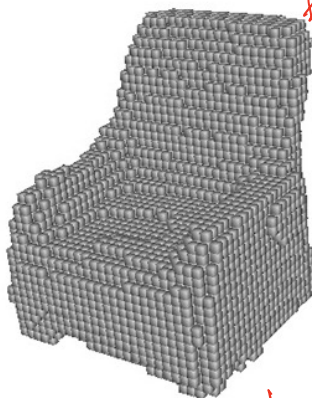
- Bird's eye view object detection.
- Used the ResNet + FPN network single-stage architecture.
- Detection: Classification + regression $\cos\theta, \sin\theta, dx, dy, w, l$.
- First real-time 3D detection network.



Point Cloud

— samples of
3D world

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



3D voxels

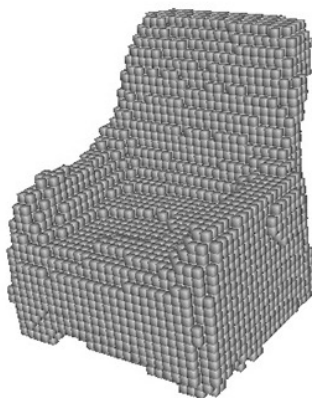


3D point cloud

network

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

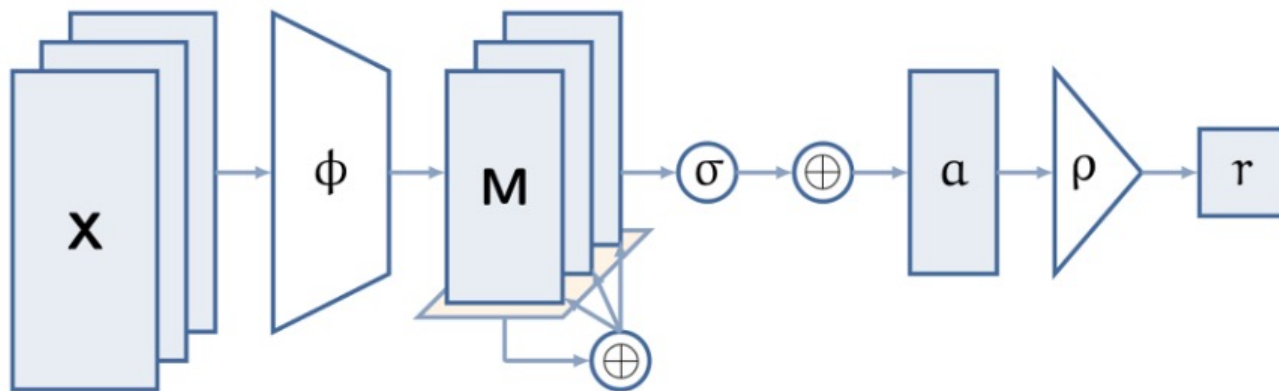


3D point cloud

↓
2000, ..., or more.

Permutation Invariance

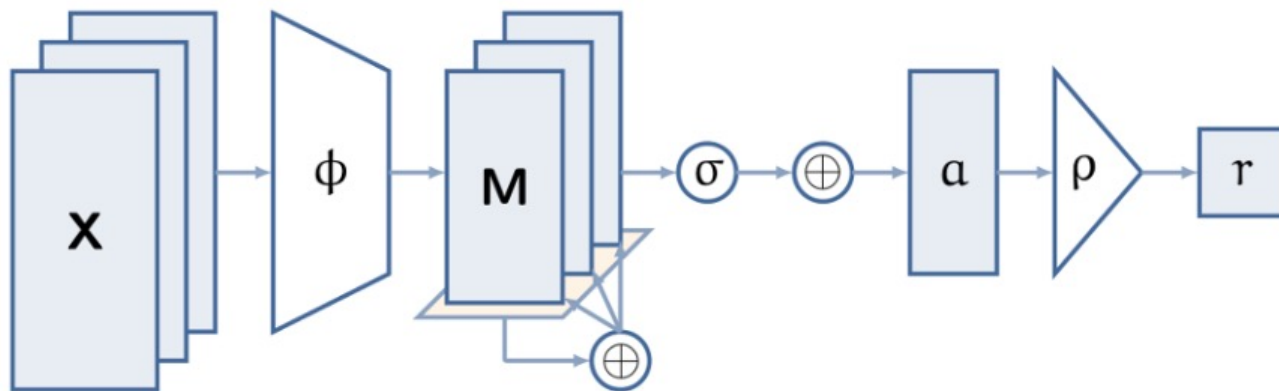
- Point cloud is a set.



Permutation Invariance

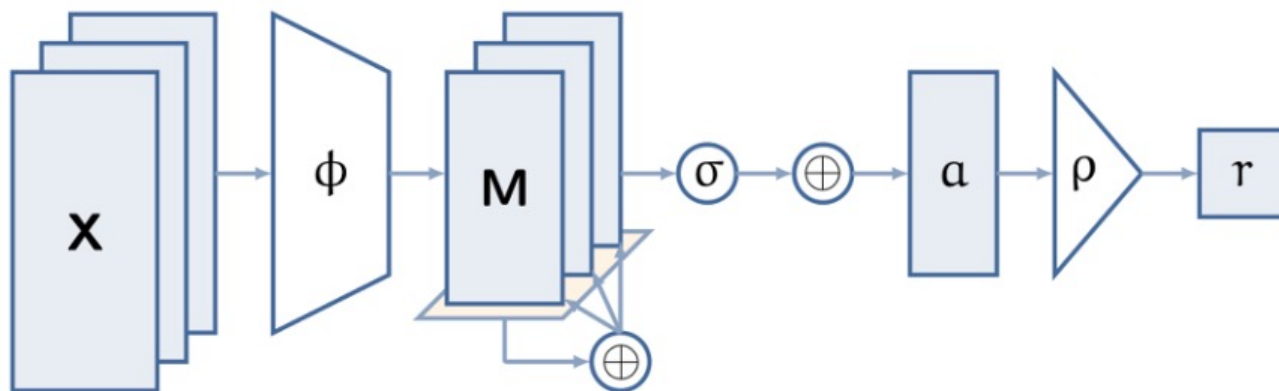
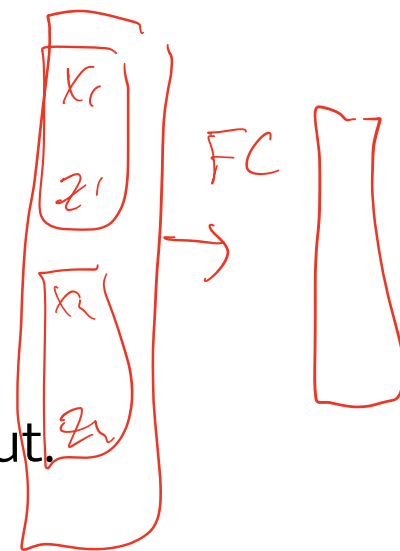
- Point cloud is a set.
- Permutation does not affect the classification in the output.

(x_1, y_1, z_1) (x_2, y_2, z_2)
↔

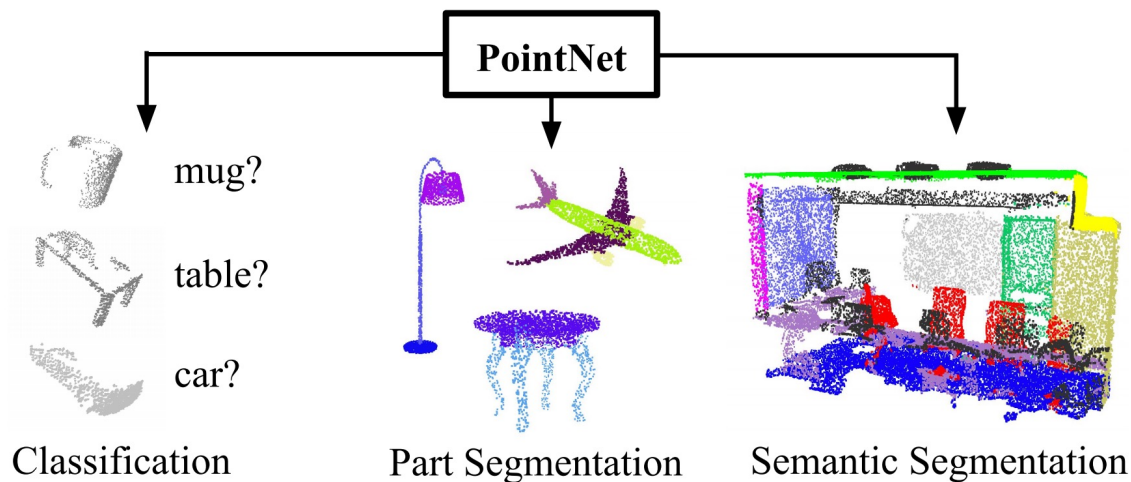


Permutation Invariance

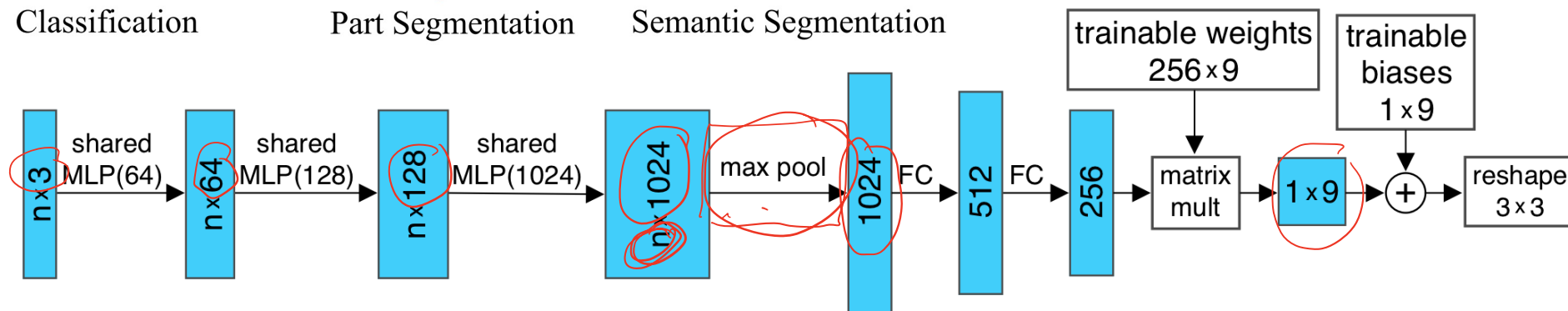
- Point cloud is a set.
- Permutation does not affect the classification in the output.
- What operations are permutation invariant?

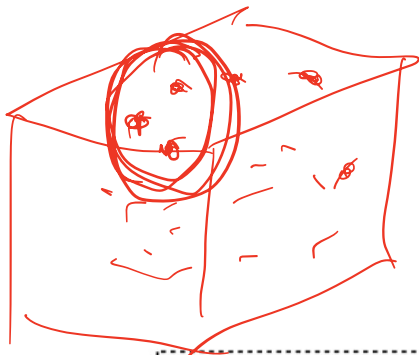


PointNet

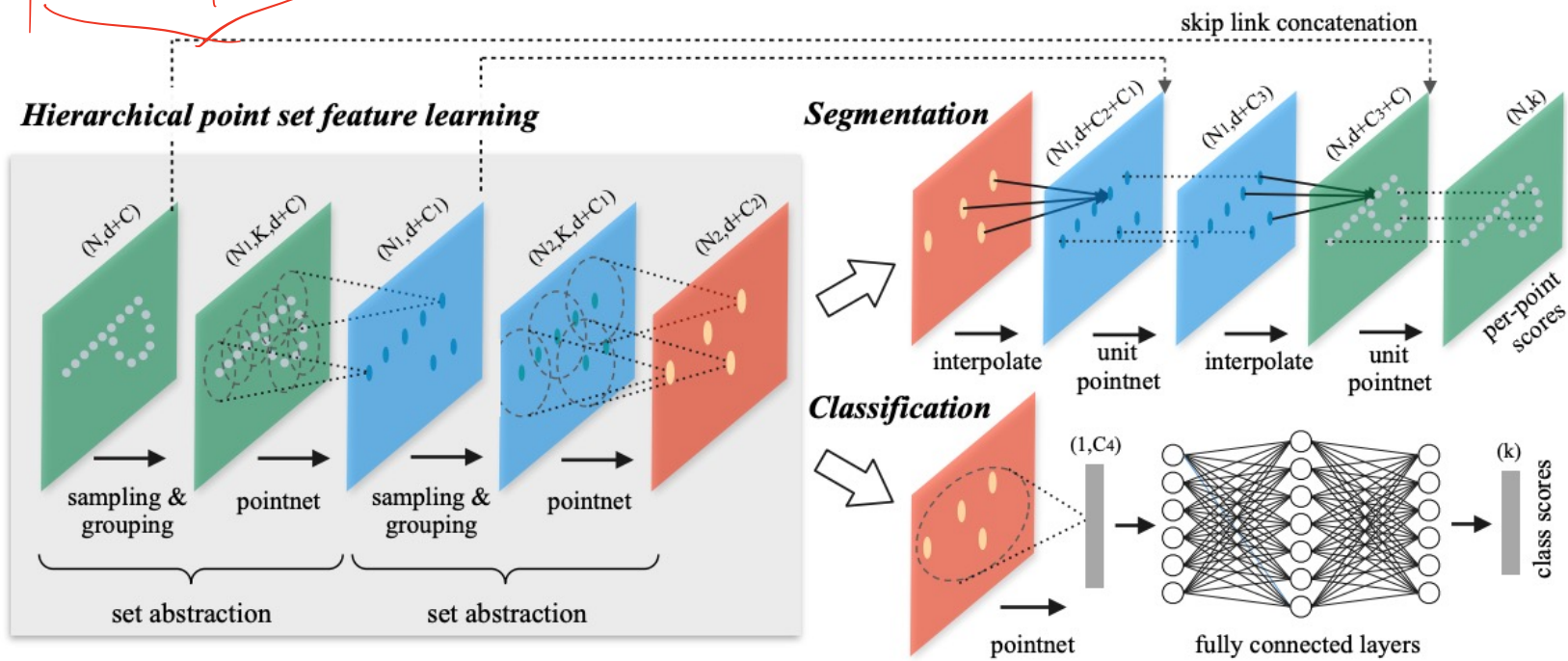


- Apply an MLP on each point.
- Max pool the features across all points.

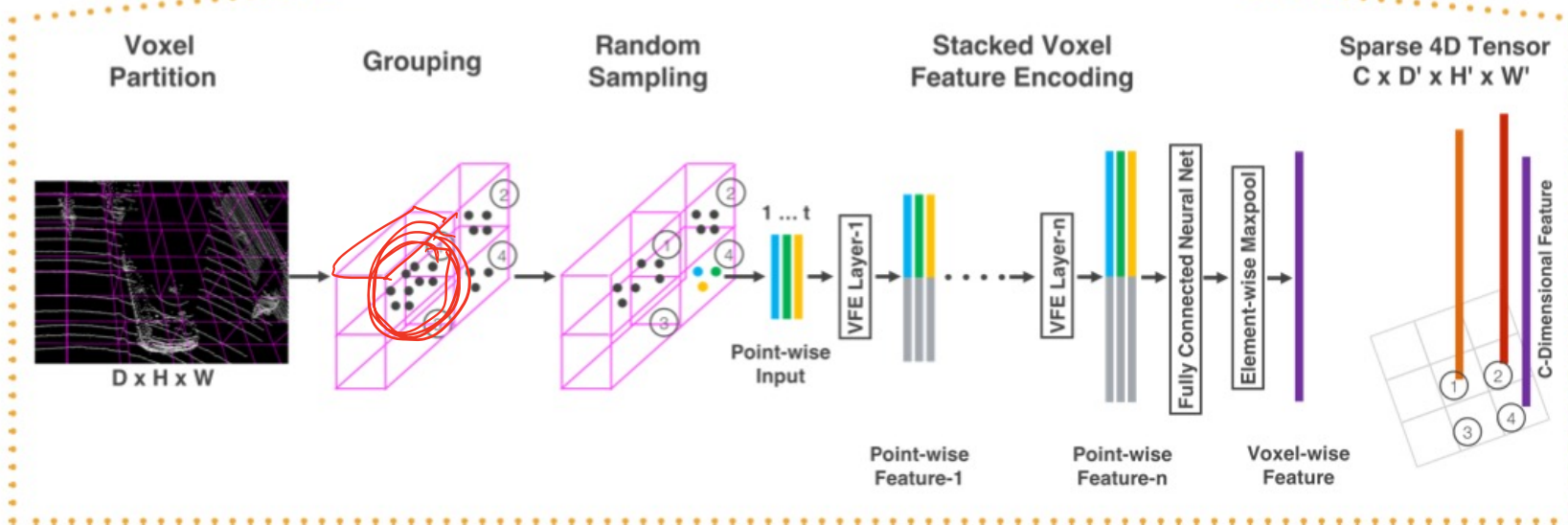
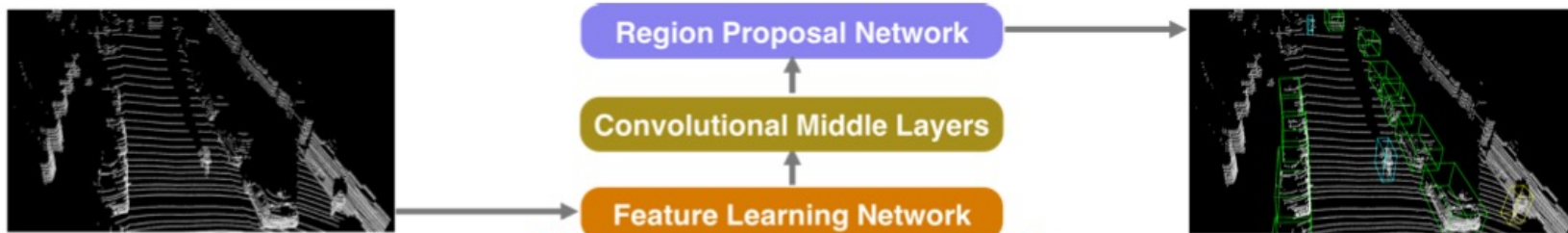




PointNet++

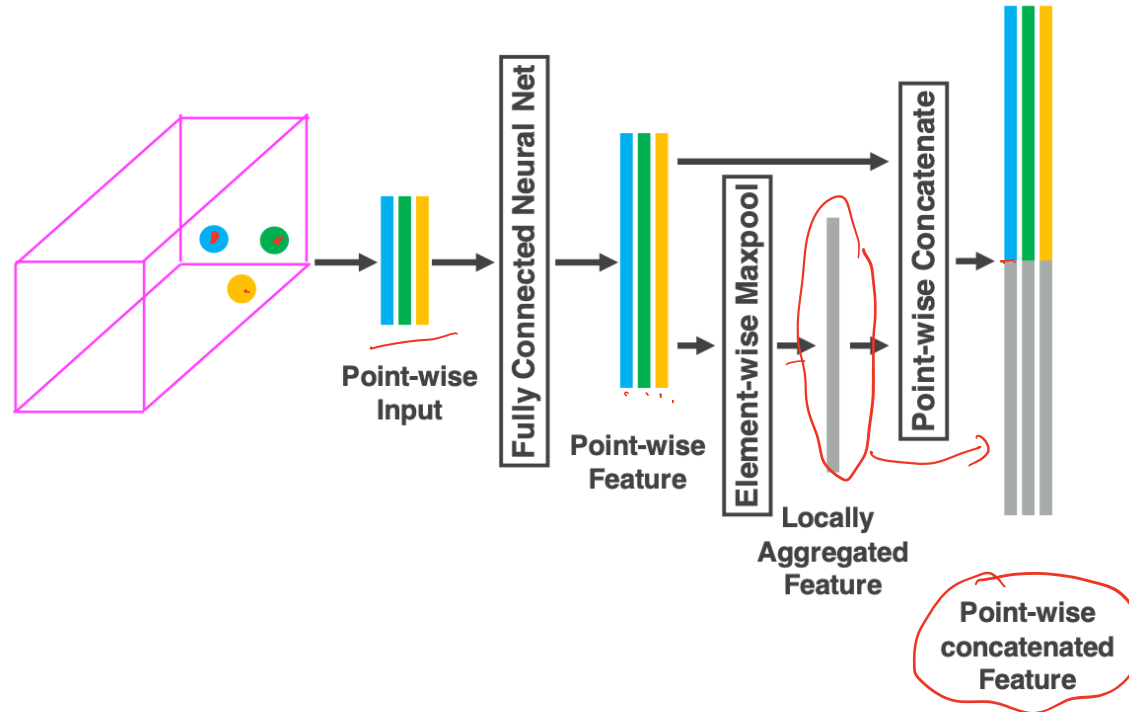


VoxelNet



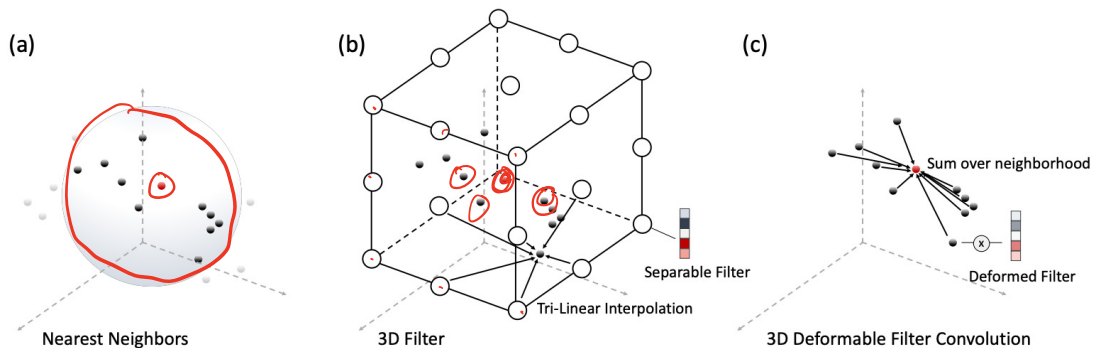
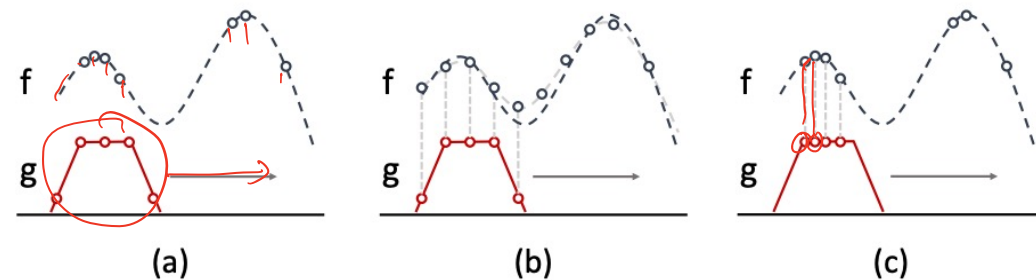
VoxelNet

- Zooming inside voxel feature encoding (VFE)



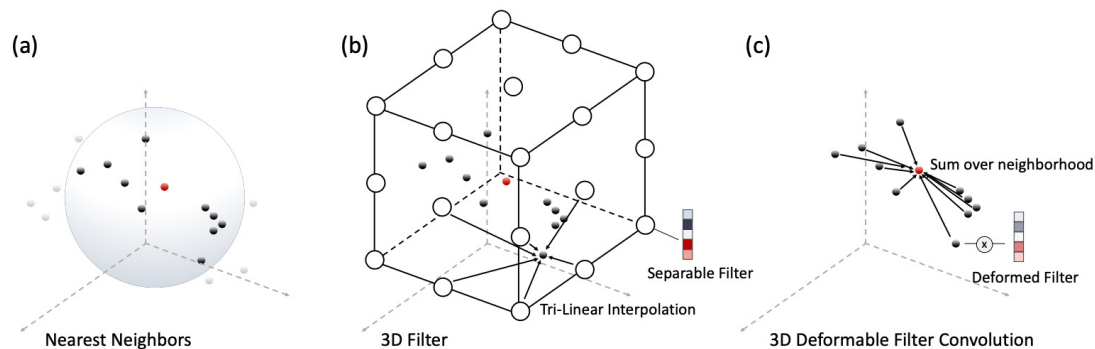
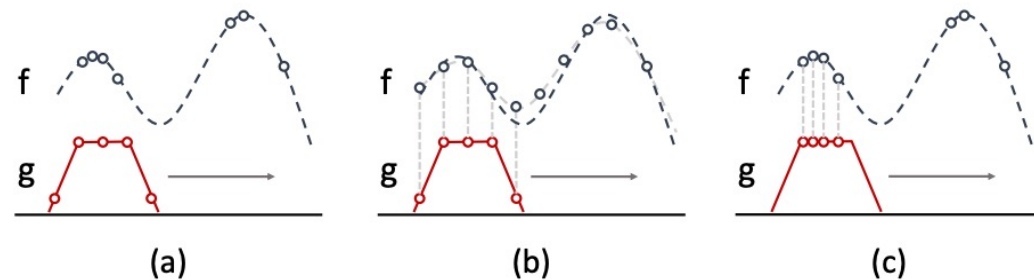
Deformable Convolution in Point Cloud

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



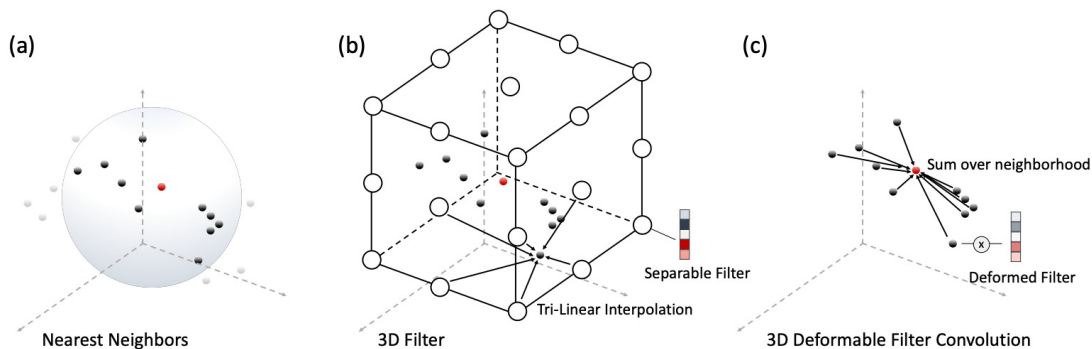
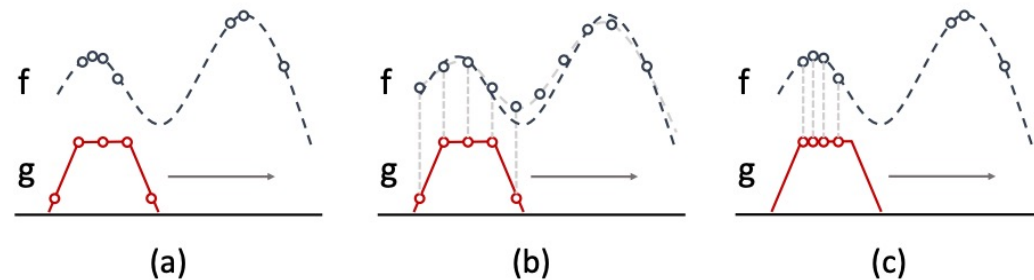
Deformable Convolution in Point Cloud

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



Deformable Convolution in Point Cloud

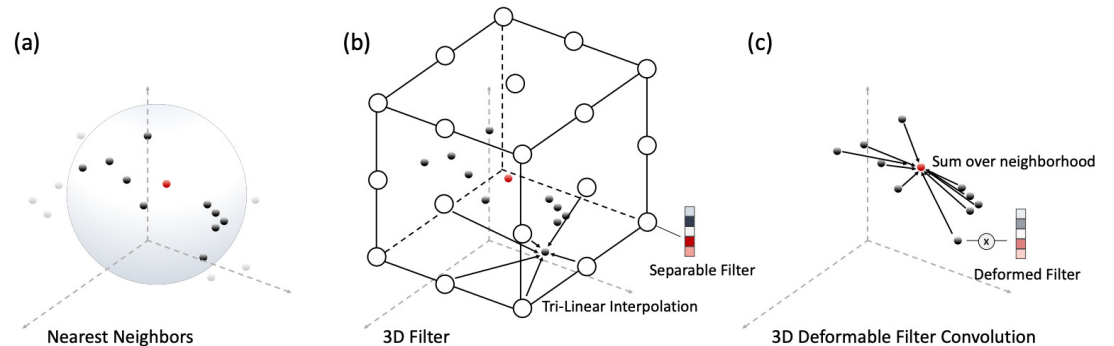
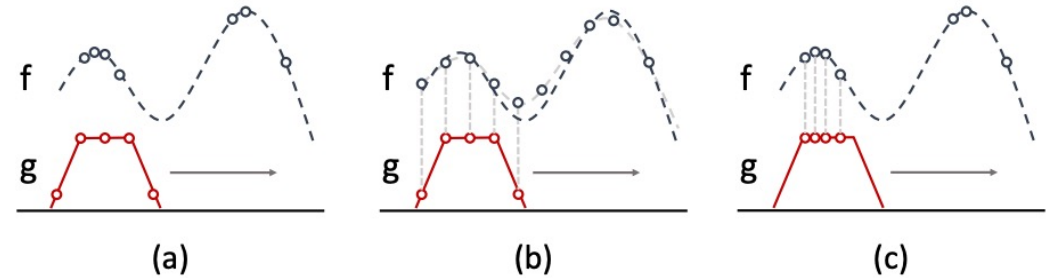
- Can we convolve a point cloud with a spatially defined kernel function?
- Resample the kernel at the point location.
- Compute the weighted sum around a neighborhood.



(x, y, z)

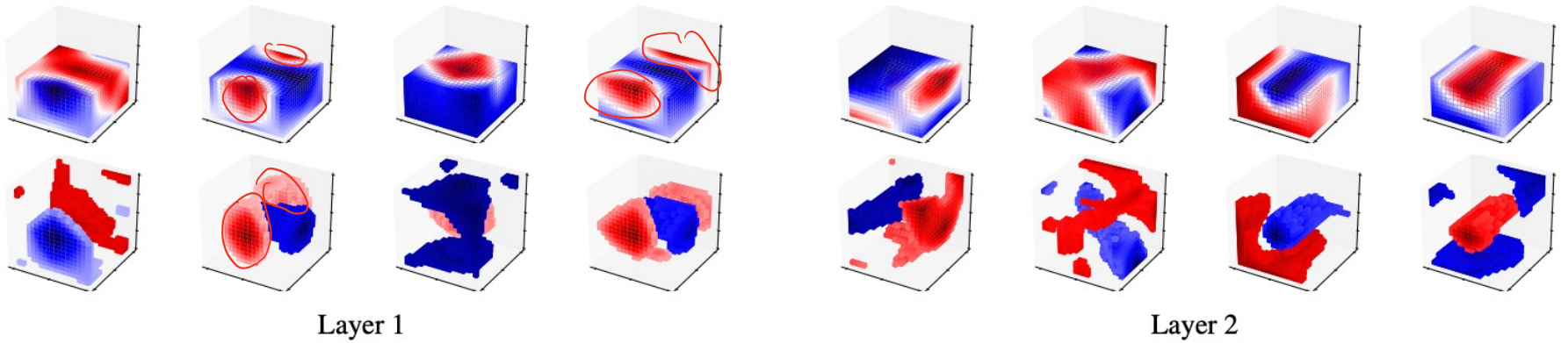
Deformable Convolution in Point Cloud

- 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.
- Translational equivariance.

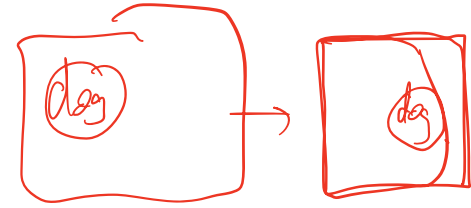


3D Filters

- Visualizing 3D convolution kernels.



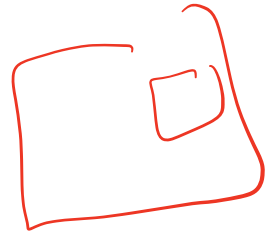
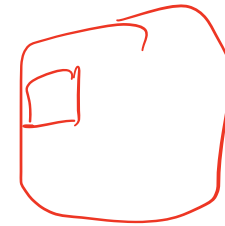
More on Equivariance



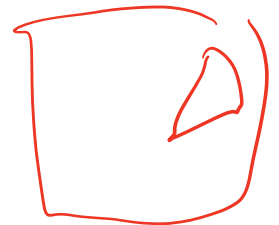
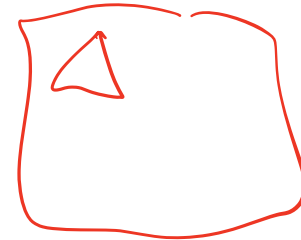
- Equivariance is a symmetry property on functions and transformations.

$$f(\mathcal{T}x) = \mathcal{T}'f(x)$$

translation



rotation



$H \times W \times C$.

rotate.

new feature map.

More on Equivariance

- Equivariance is a symmetry property on functions and transformations.

$$f(\mathcal{T}x) = \mathcal{T}'f(x)$$

- Useful property to constrain the hidden representation space.

data efficiency, \leftrightarrow Compute efficiency
(not always):

More on Equivariance

- Equivariance is a symmetry property on functions and transformations.

$$f(\mathcal{T}x) = \mathcal{T}' f(x)$$

- Useful property to constrain the hidden representation space.
- Convolution has the property of translational equivariance. Translated input = translated output. (up to discretization)

More on Equivariance

- Equivariance is a symmetry property on functions and transformations.

$$f(\mathcal{T}x) = \mathcal{T}'f(x)$$

- Useful property to constrain the hidden representation space.
- Convolution has the property of translational equivariance. Translated input = translated output. (up to discretization)
- 3D convolution handles translational equivariance.

More on Equivariance

- Equivariance is a symmetry property on functions and transformations.

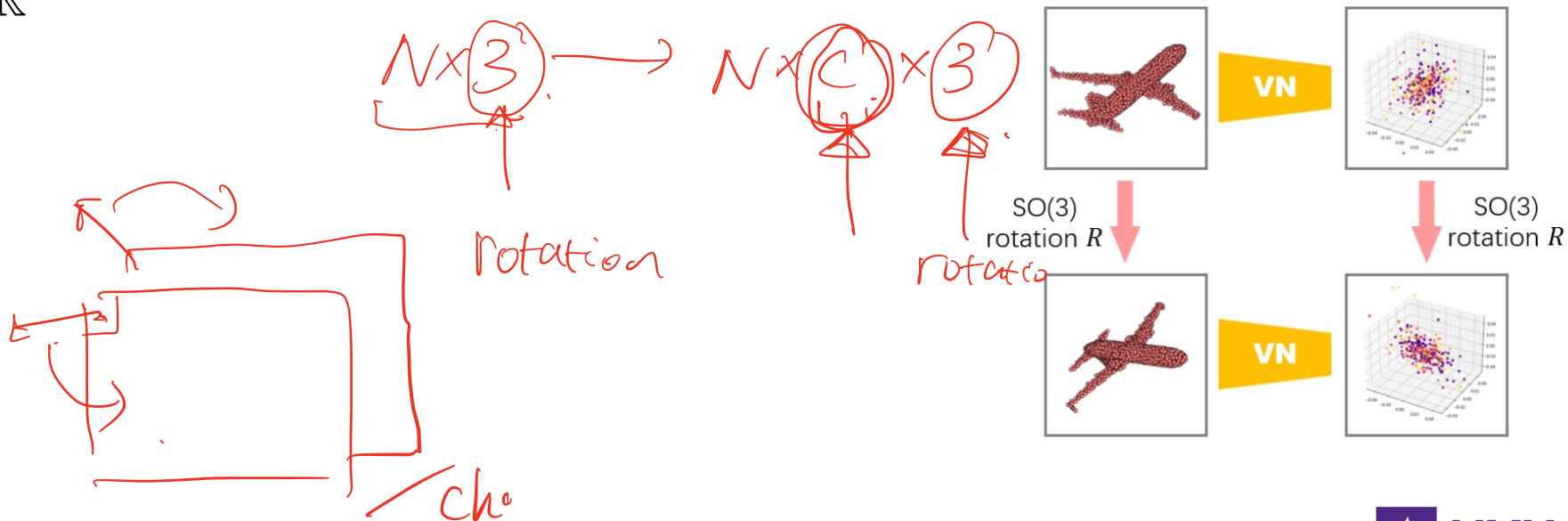
$$f(\mathcal{T}x) = \mathcal{T}' f(x)$$

- Useful property to constrain the hidden representation space.
- Convolution has the property of translational equivariance. Translated input = translated output. (up to discretization)
- 3D convolution handles translational equivariance.
- What about rotation?

3D Rotational Equivariance

- For object centered point clouds, rotation are more important.
- Imagine every feature has a vector / arrow.

$$V \in \mathbb{R}^{N \times C \times 3}$$

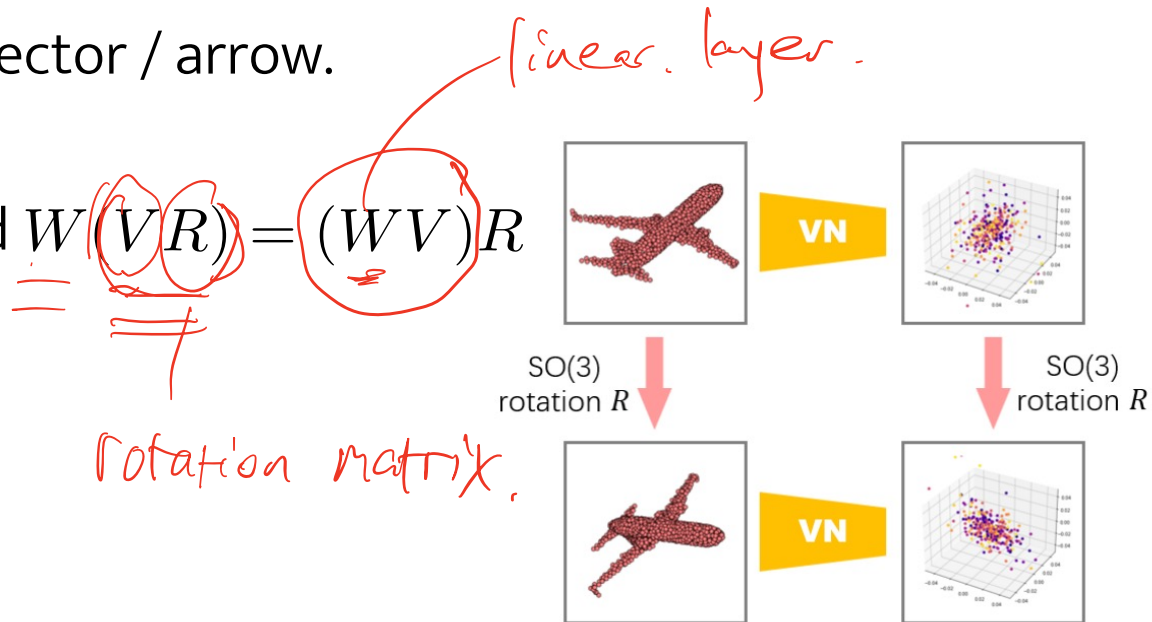


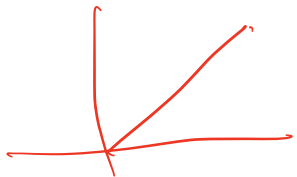
3D Rotational Equivariance

- For object centered point clouds, rotation are more important.
- Imagine every feature has a vector / arrow.

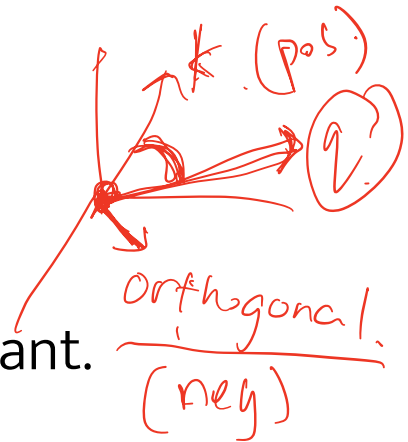
$$V \in \mathbb{R}^{N \times C \times 3}$$

- Linear layers are already good $W(VR) = (WV)R$





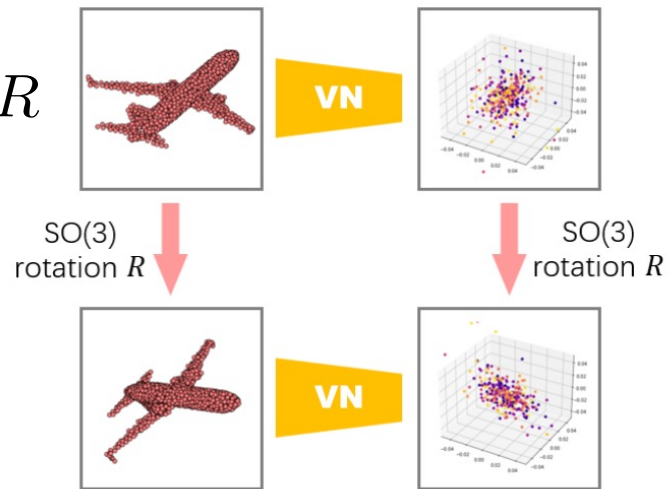
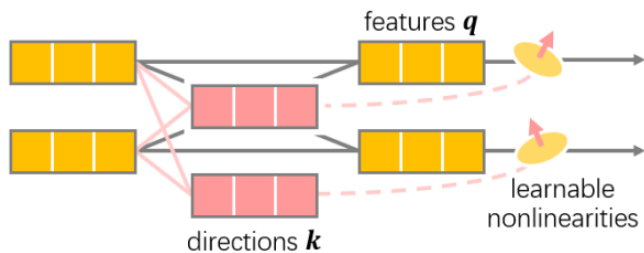
3D Rotational Equivariance



- For object centered point clouds, rotation are more important.
- Imagine every feature has a vector / arrow.
 $V \in \mathbb{R}^{N \times C \times 3}$
- Linear layers are already good $W(VR) = (WV)R$
- Need a different nonlinearity

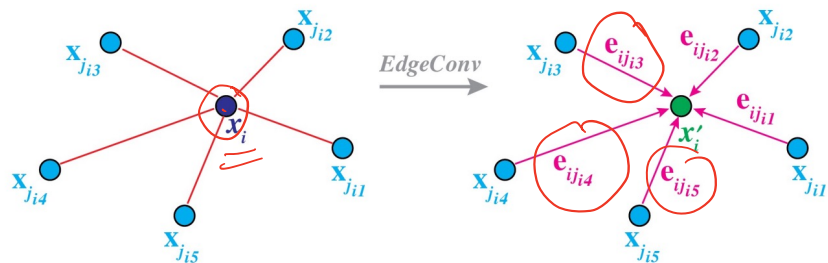
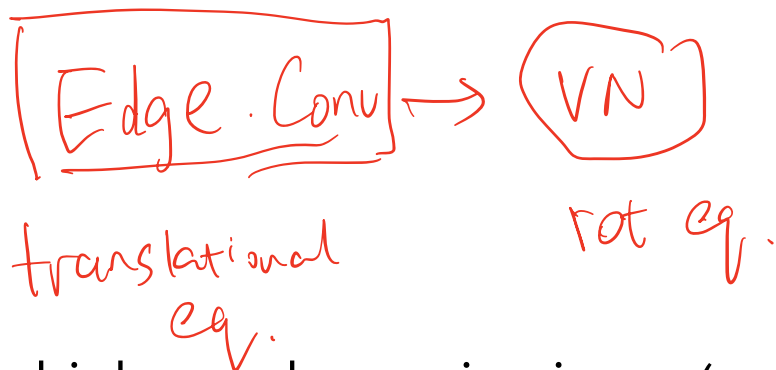
$$q = WV, \quad k = UV$$

$$v' = \begin{cases} q & \text{if } \langle q, k \rangle \geq 0 \\ q - \langle q, \frac{k}{\|k\|} \rangle \frac{k}{\|k\|} & \text{otherwise,} \end{cases}$$

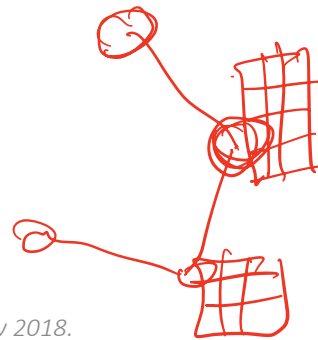


More General Equivariance

- Making VN also translational equivariant.



- For higher order equivariance (e.g. matrices, tensors), check out Tensor Field Networks.
- Convert features into spherical harmonics.

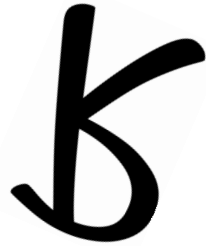
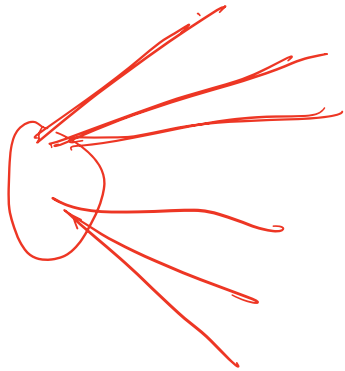


Challenges and Tradeoffs

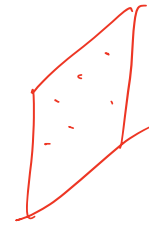
- Computational and memory efficiency vs. training data efficiency

Challenges and Tradeoffs

- Computational and memory efficiency vs. training data efficiency
- Mental rotation: Rotate in the latent space may not be as simple as a geometric transformation. The skill may still be experience driven.



Challenges and Tradeoffs



- Computational and memory efficiency vs. training data efficiency
- Mental rotation: Rotate in the latent space may not be as simple as a geometric transformation. The skill may still be experience driven.

b

- Neuroscience evidence: No spatial equivariance in connections. The LGN performs “pre-pre-training” from simulating retina waves [Blankenship & Feller, 2009].

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

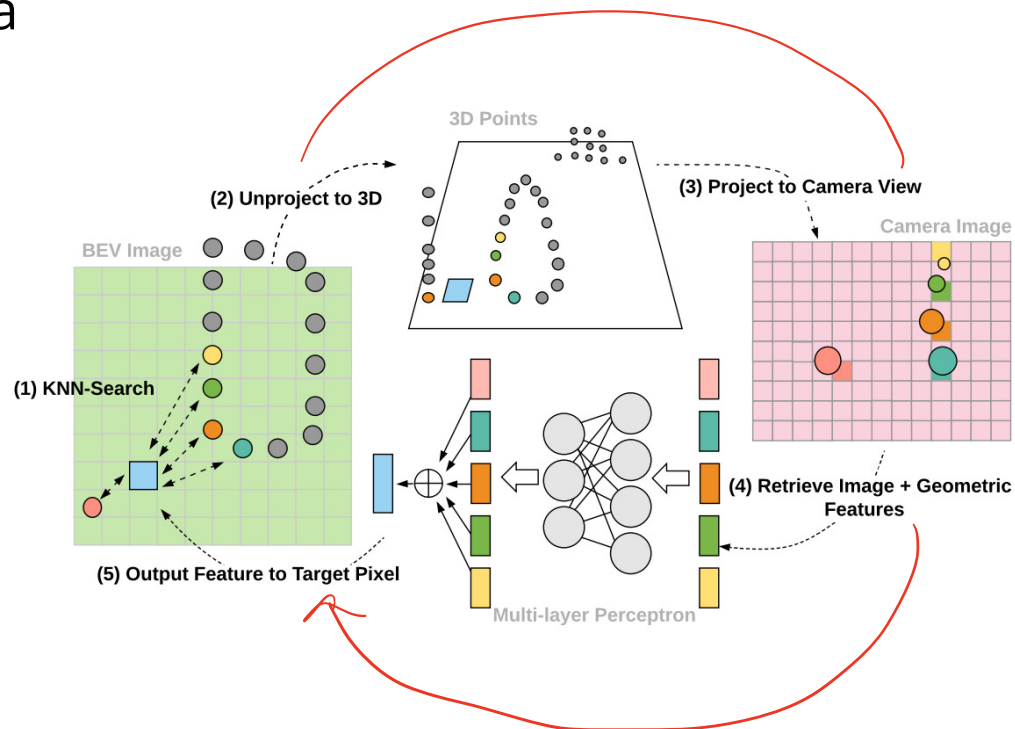
- LiDAR is precise in depth perception, but the point cloud format is sparse and non-uniform (dense around the ego-car and sparse in long distance.)
- Camera provides high resolution 2D view and good for long distance but lacks 3D. Can we achieve the best of both worlds?
- Late fusion: Generate proposals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

early fusion

mid fusion

Camera-LiDAR Projection

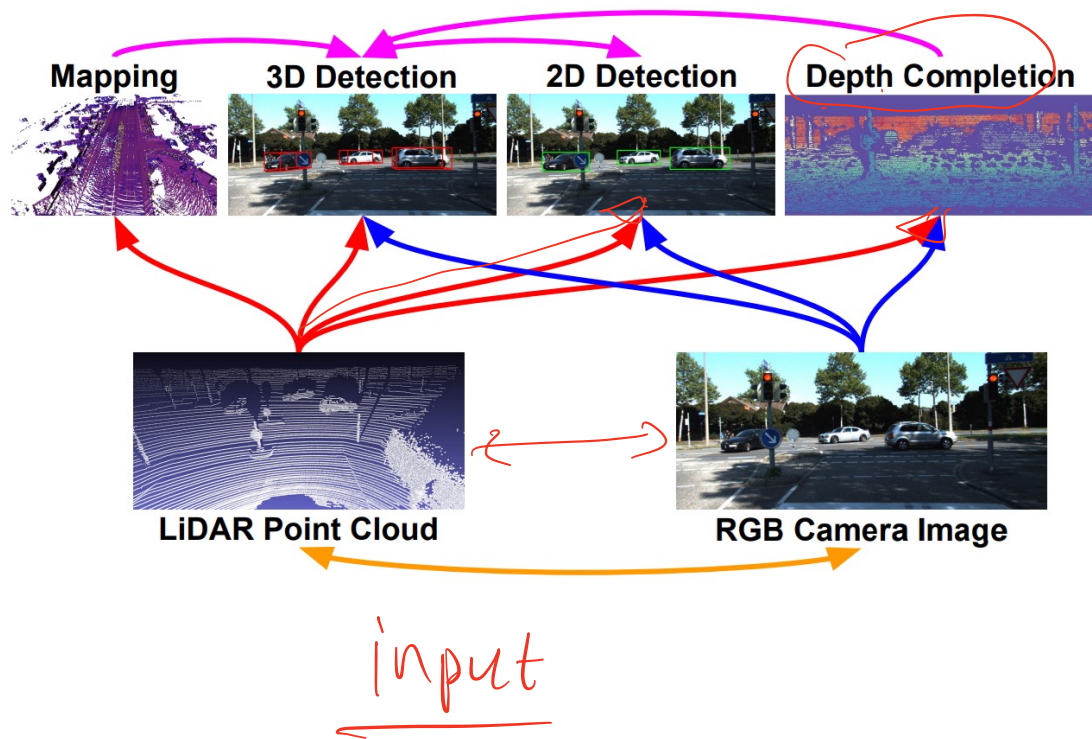
- Unproject LiDAR points to camera view (i.e. Range View)
- Query the closest camera RGB features for each LiDAR point.
- For empty space in BEV, we can interpolate from neighboring points using kNN.
- Continuous Fusion: $h_i = \sum_j MLP([f_j, x_j - x_i])$.



Supervised Dense Depth

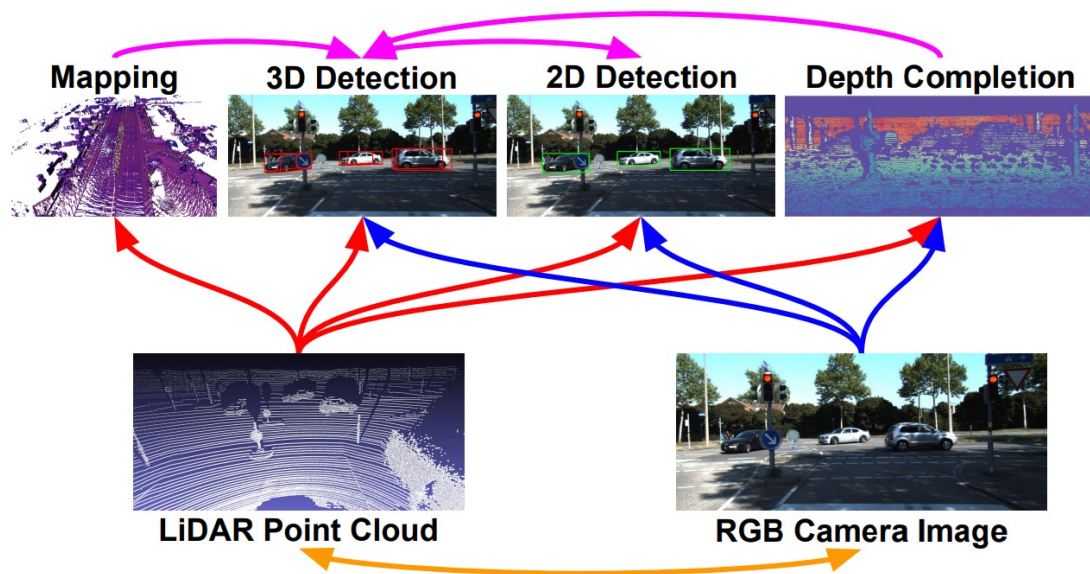
task.

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



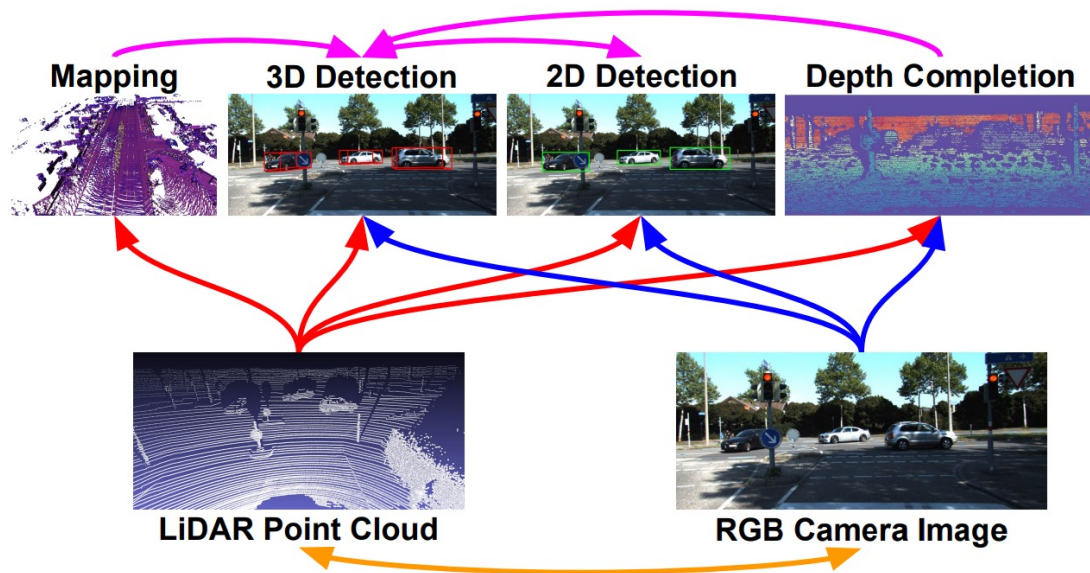
Supervised Dense Depth

- 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

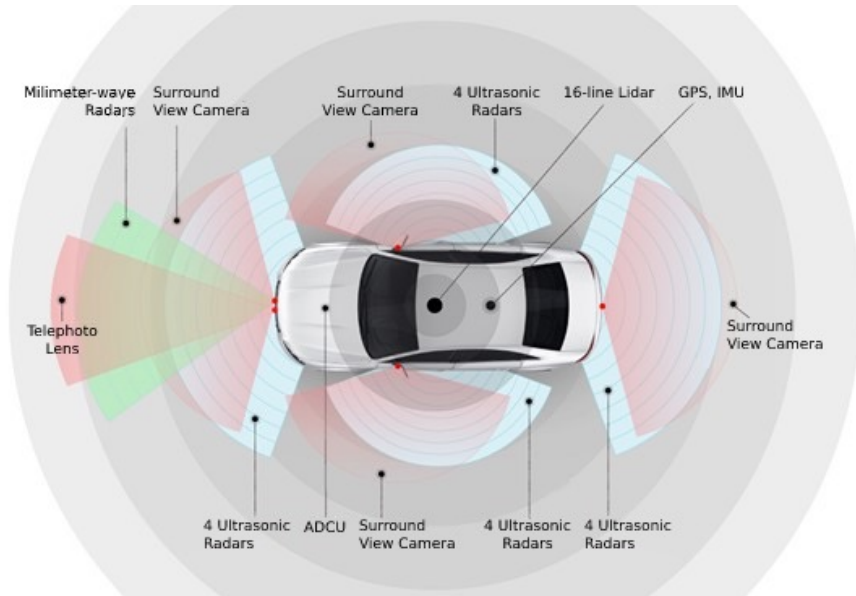
- With the ease of use of automatic differentiation libraries, we can compose a computation graph in millions of ways.
- We can design layers and operators to accommodate different types of inputs and outputs. 3D, point cloud, sparse data, etc.

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 accommodate 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

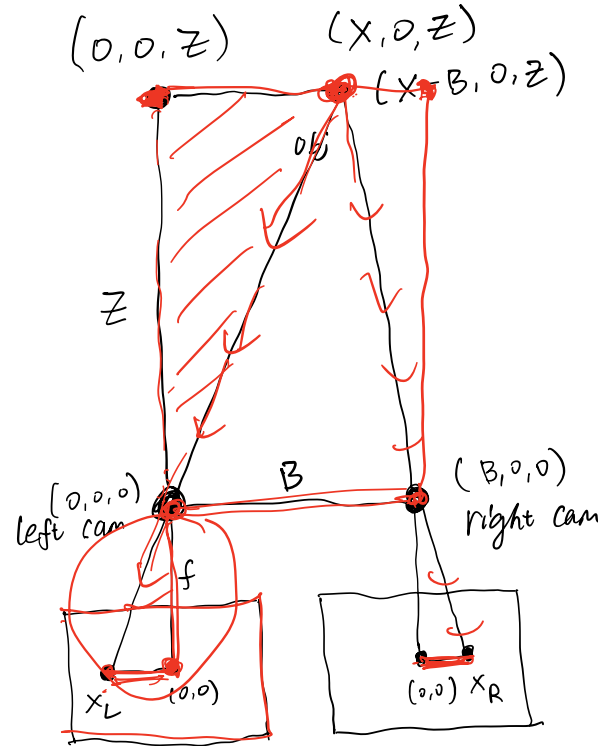
- Not all embodied agents have the luxury to have a full set of sensors.
- Can we infer the geometric structure with 2D perception?



$$(x, y) \rightarrow z$$

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.



$$x_L = \frac{x f}{z}$$

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

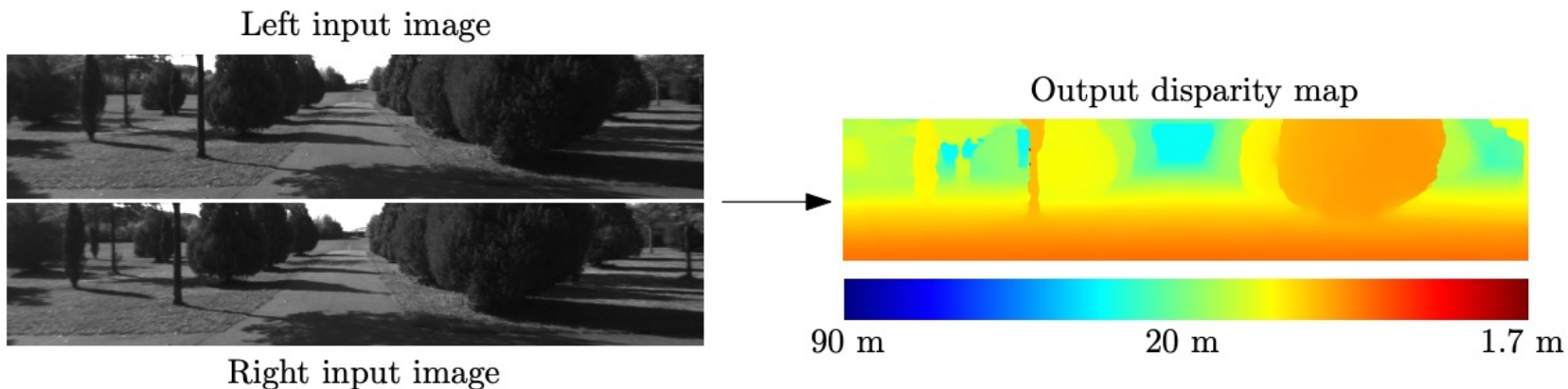
$$\begin{aligned} \delta &= x_L - x_R \\ &= \frac{[x - (x-B)] f}{z} \end{aligned}$$

$$= \frac{B f}{z}$$

$$z = \frac{B f}{\delta}$$

From 2D to 3D: Depth Network

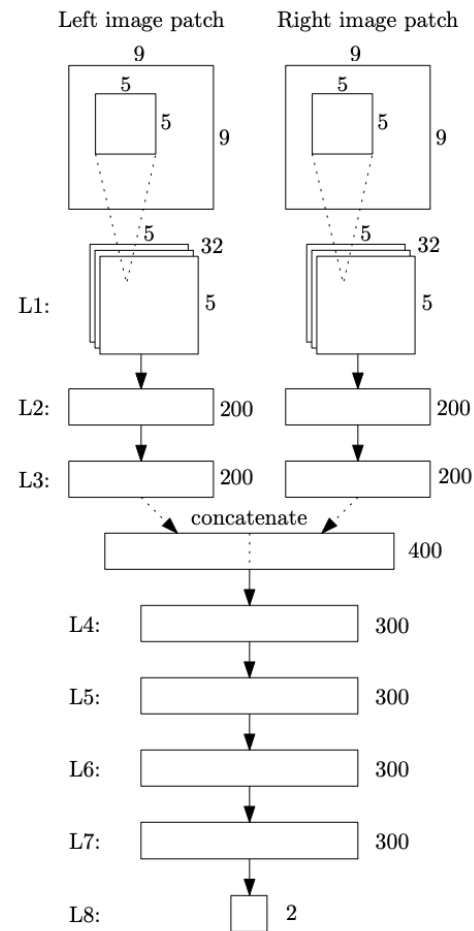
- A network that can output disparity.
- Using LiDAR or depth camera as groundtruth supervision.



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

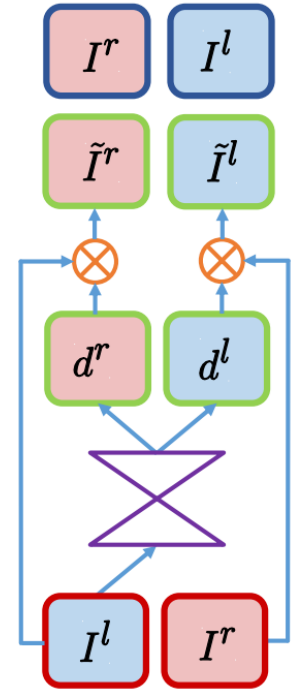
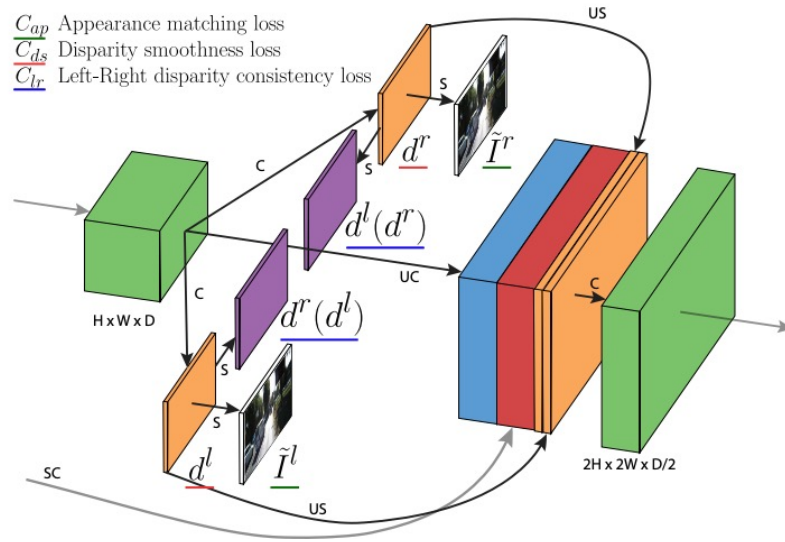
$$\begin{aligned}
 \text{Energy } E(D) = \sum_{\mathbf{p}} & \left(C_{\text{CBCA}}^4(\mathbf{p}, D(\mathbf{p})) \right. \\
 & + \sum_{\mathbf{q} \in \mathcal{N}_{\mathbf{p}}} P_1 \times 1\{|D(\mathbf{p}) - D(\mathbf{q})| = 1\} \\
 & \left. + \sum_{\mathbf{q} \in \mathcal{N}_{\mathbf{p}}} P_2 \times 1\{|D(\mathbf{p}) - D(\mathbf{q})| > 1\} \right), \\
 D(\mathbf{p}) = \operatorname{argmin}_d & C(\mathbf{p}, d).
 \end{aligned}$$



Self-Supervised Depth

- Appearance matching loss

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



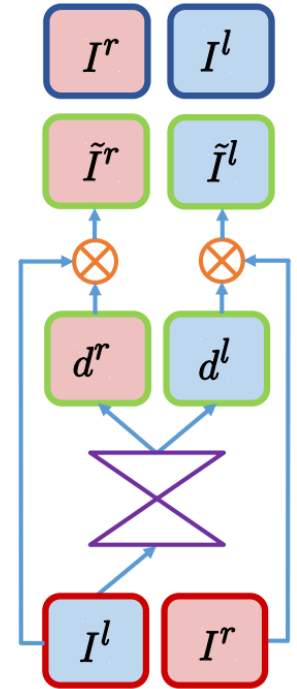
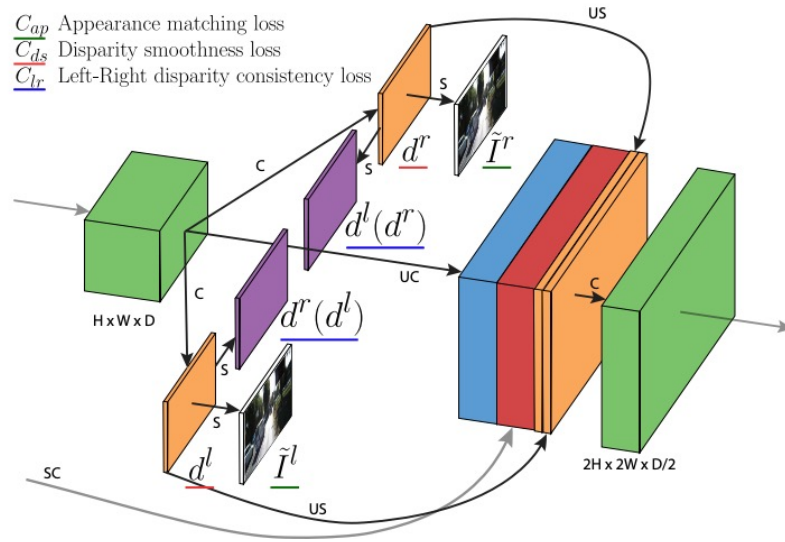
Self-Supervised Depth

- 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

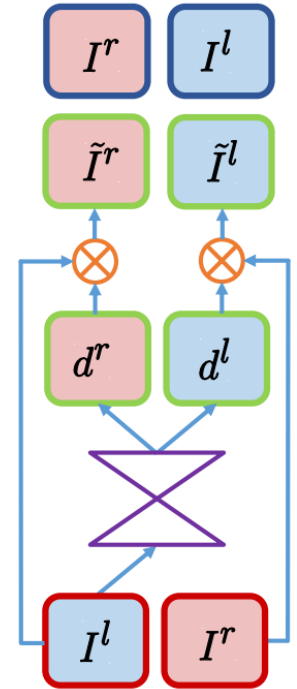
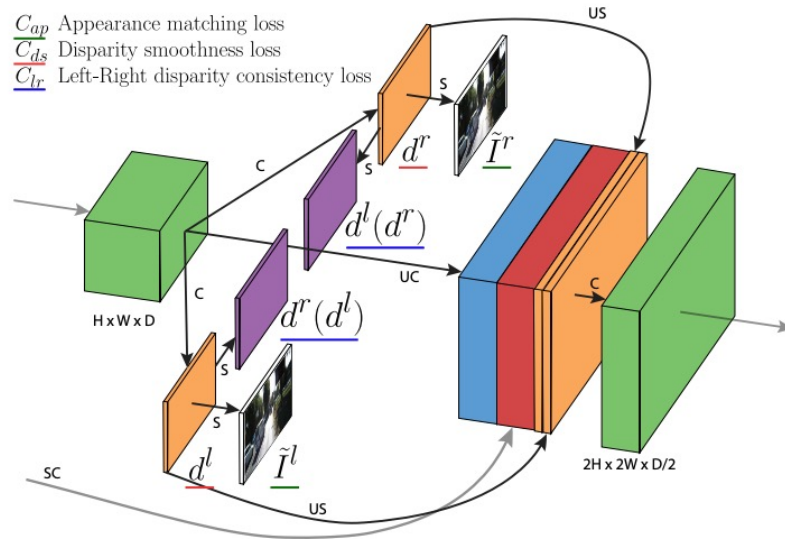
$$C_{ap}^l = \frac{1}{N} \sum_{i,j} \alpha \frac{1 - \text{SSIM}(I_{ij}^l, \tilde{I}_{ij}^l)}{2} + (1 - \alpha) \|I_{ij}^l - \tilde{I}_{ij}^l\|.$$

- 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

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



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 $I_x u + I_y v + I_t = 0.$
- Use a local patch and assume smooth motion

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

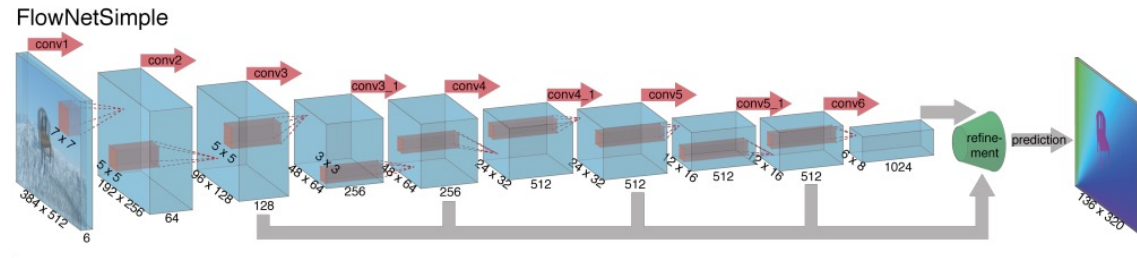
Classical Approach

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

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

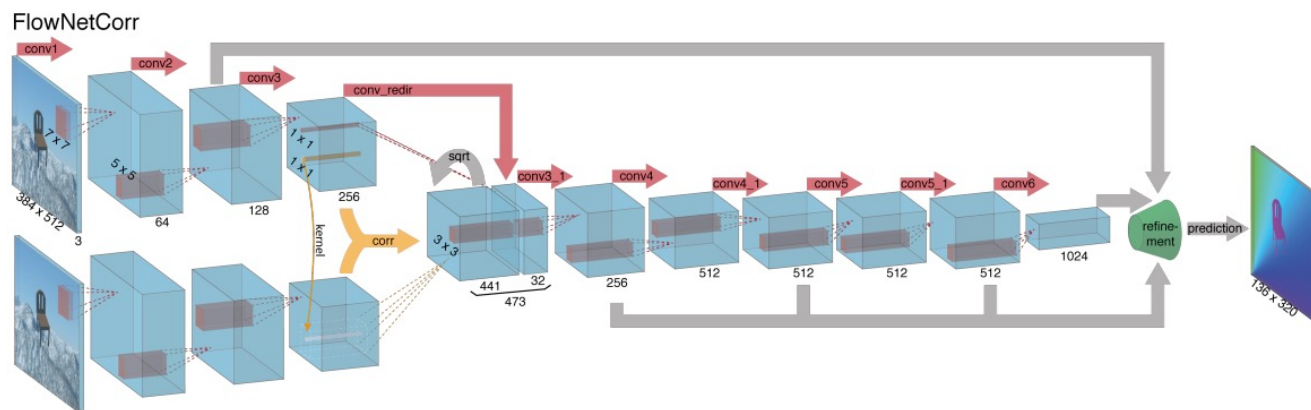
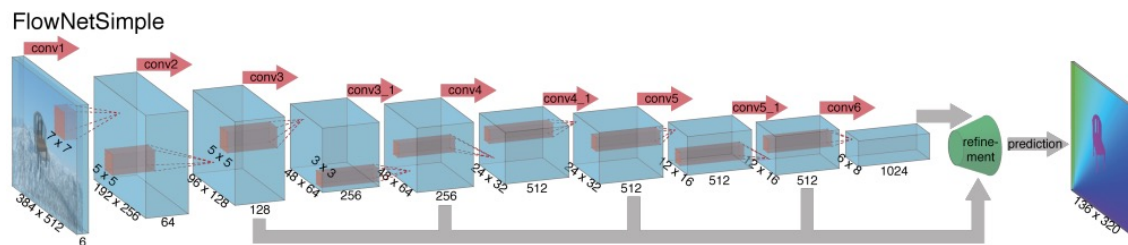
Correlation Volume Approach

- Simple Approach: Concatenate the two images together.

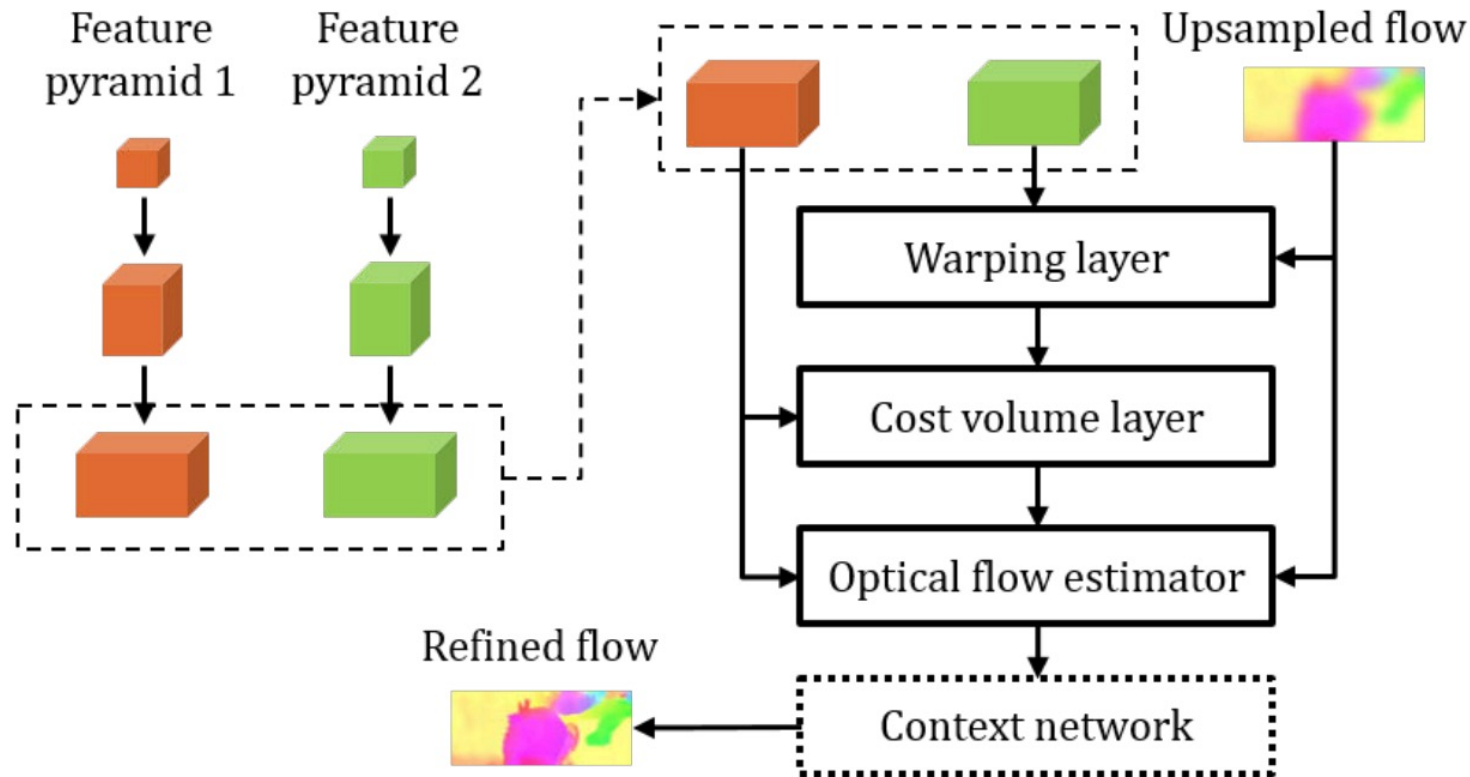


Correlation Volume Approach

- 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

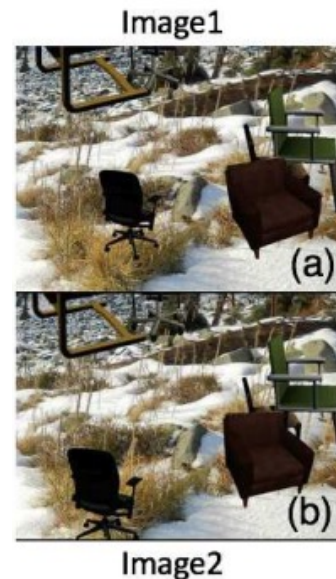


Unsupervised Flow

- Photometric Consistency (Appearance)

Unsupervised Flow

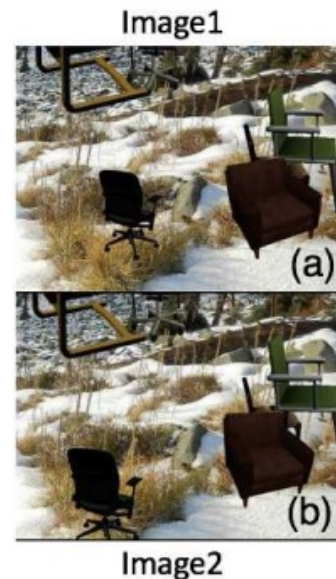
- Photometric Consistency (Appearance)
- Occlusion Estimation
 - Forward-backward consistency



Wang et al., 2018

Unsupervised Flow

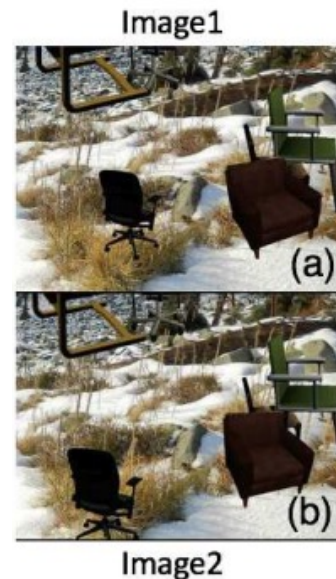
- Photometric Consistency (Appearance)
- Occlusion Estimation
 - Forward-backward consistency
- Smoothness



Wang et al., 2018

Unsupervised Flow

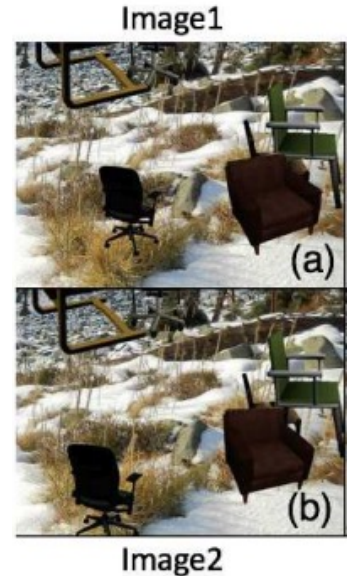
- Photometric Consistency (Appearance)
- Occlusion Estimation
 - Forward-backward consistency
- Smoothness
- Self-supervision: Ensure consistent flow at different augmentation (e.g. crops)



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?

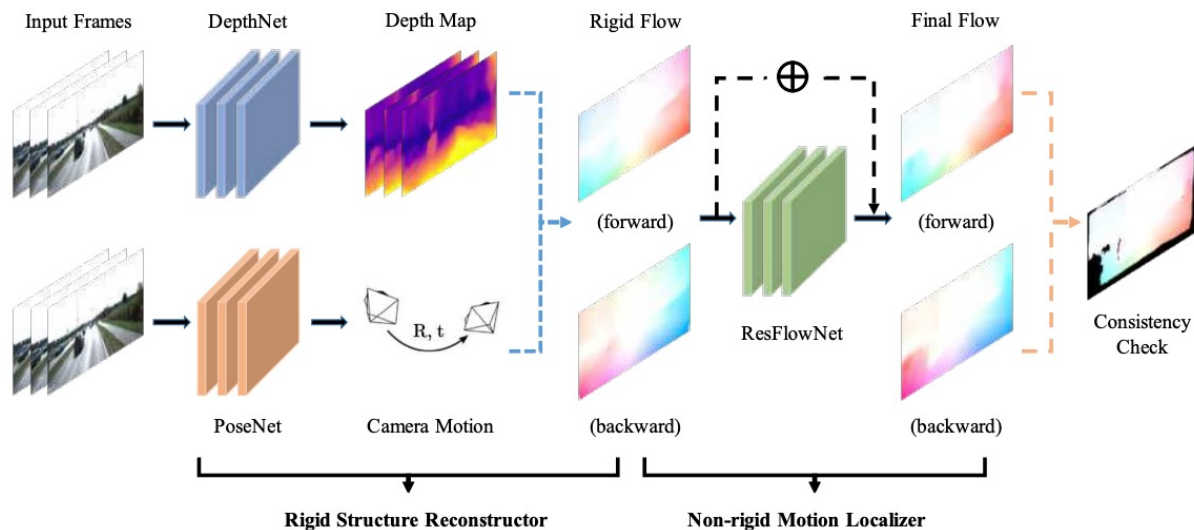


Wang et al., 2018

Depth, Flow, and Pose Movement

- The static objects follow rigid flow: determined by camera motion and depth.

$$f_{t \rightarrow s}^{rig}(p_t) = K T_{t \rightarrow s} D_t(p_t) K^{-1} p_t - p_t.$$



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:

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

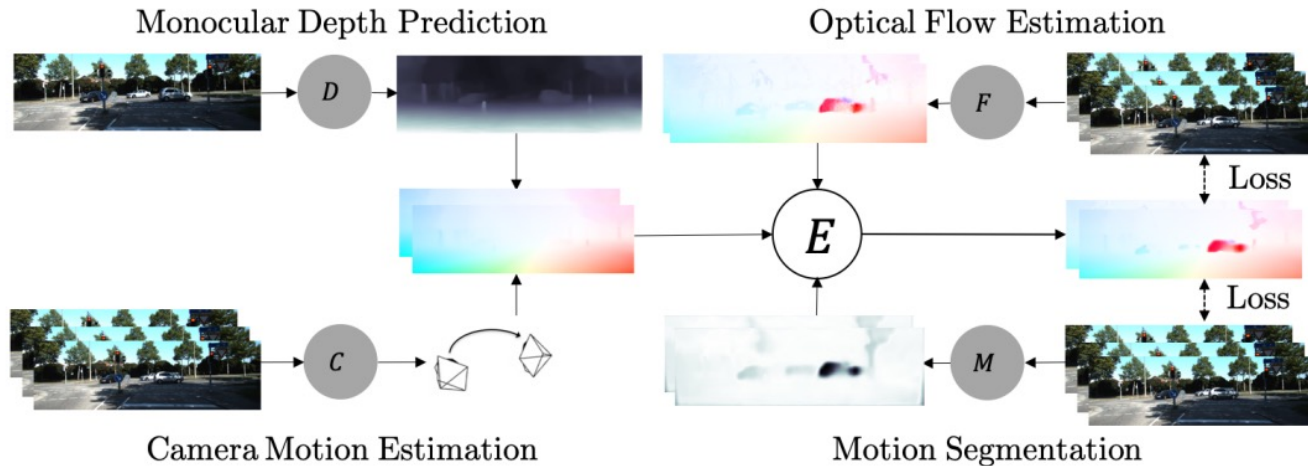
- Forward-Backward Consistency:

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

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

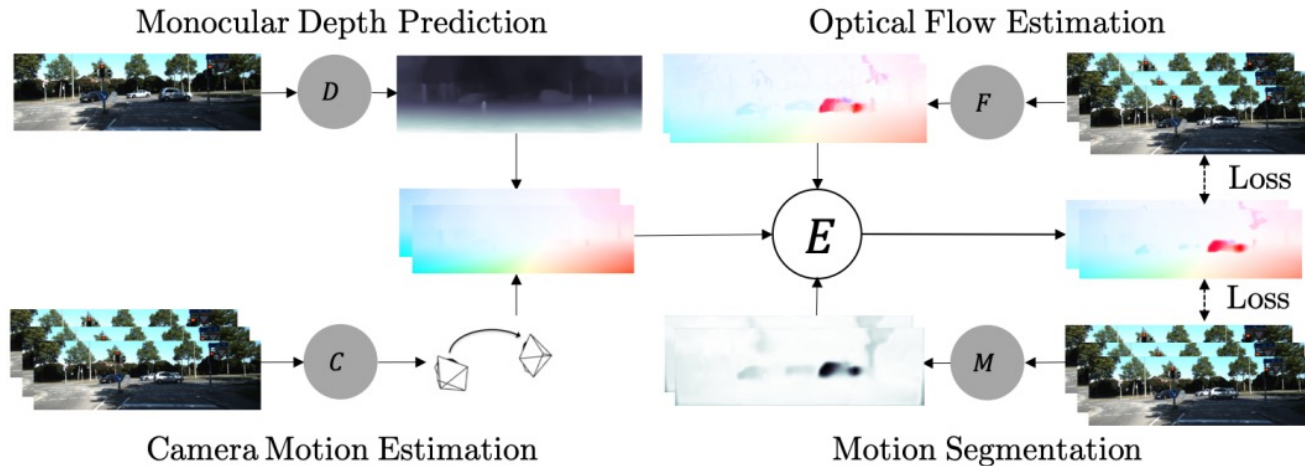
Summary

- Leverage cross correlation structure for spatial similarity matching.



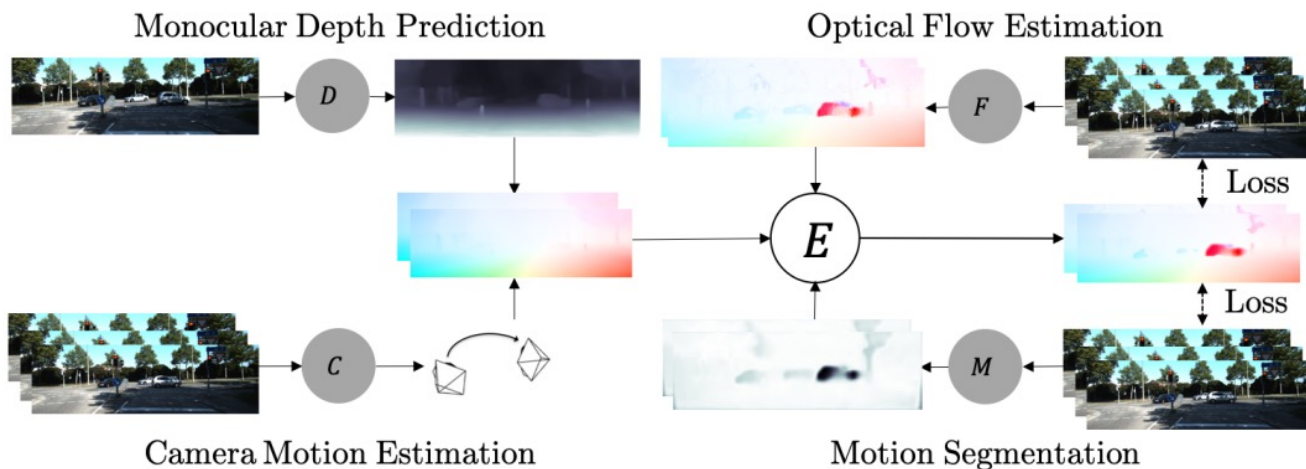
Summary

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



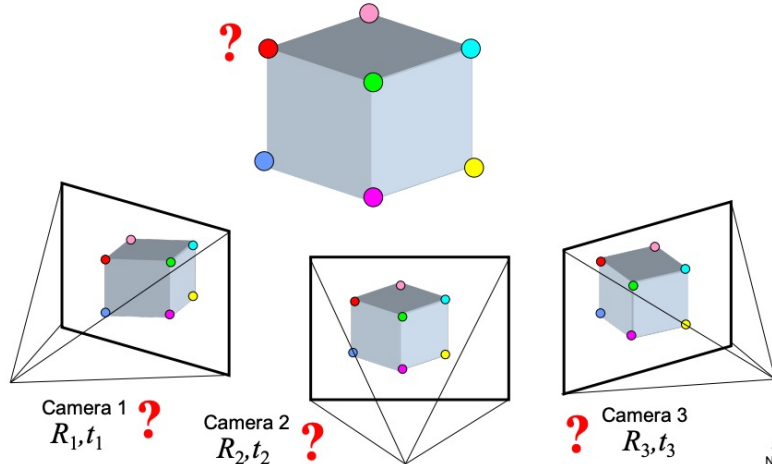
Summary

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

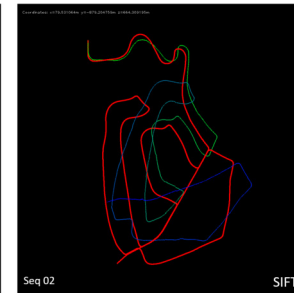
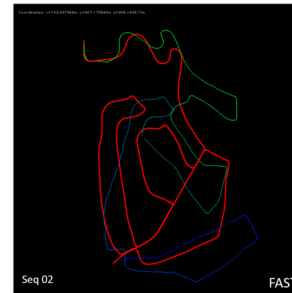
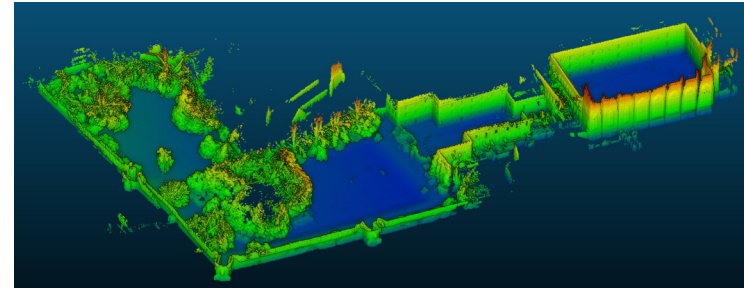


Classical Mapping

- 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



Garg & Jain

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.

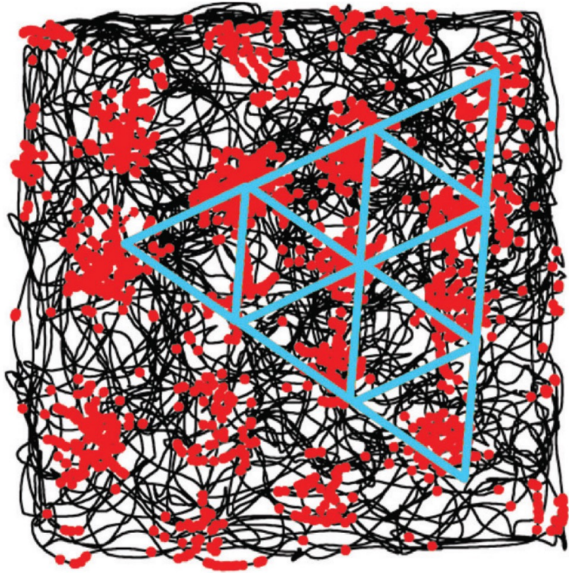
Common Drawbacks

- 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).

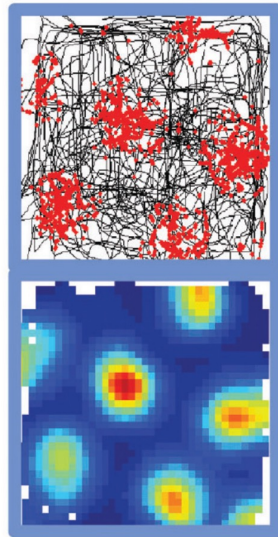
Common Drawbacks

- 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?

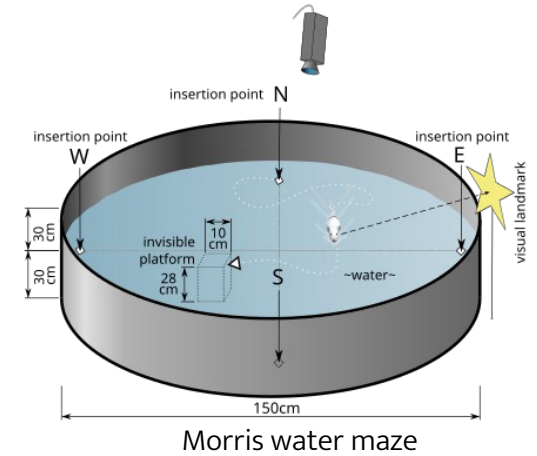
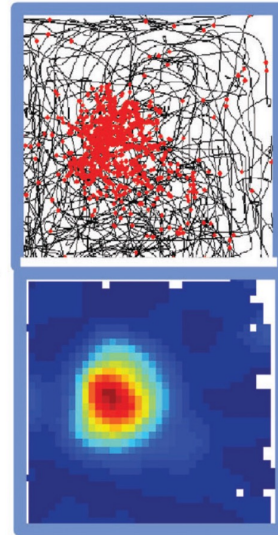
Mapping in the Brain: Grid and Place Cells



Grid



Place



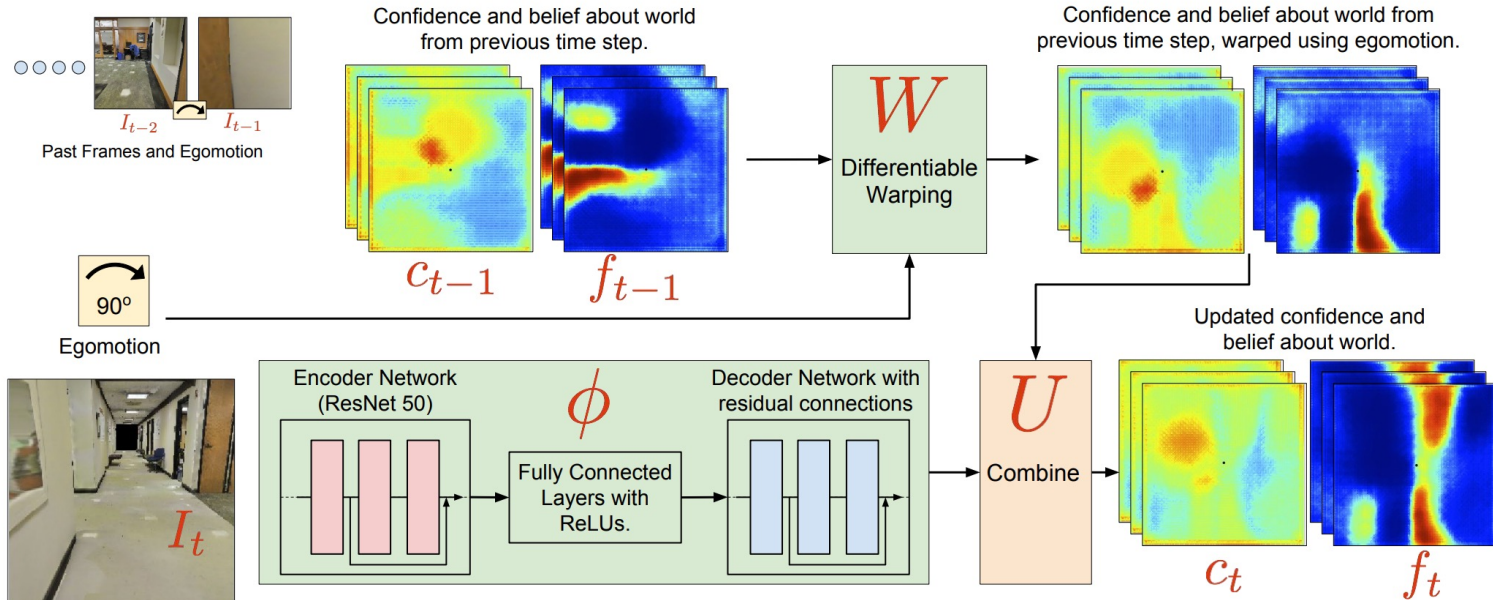
Morris water maze



Matthias Wandel, 2018

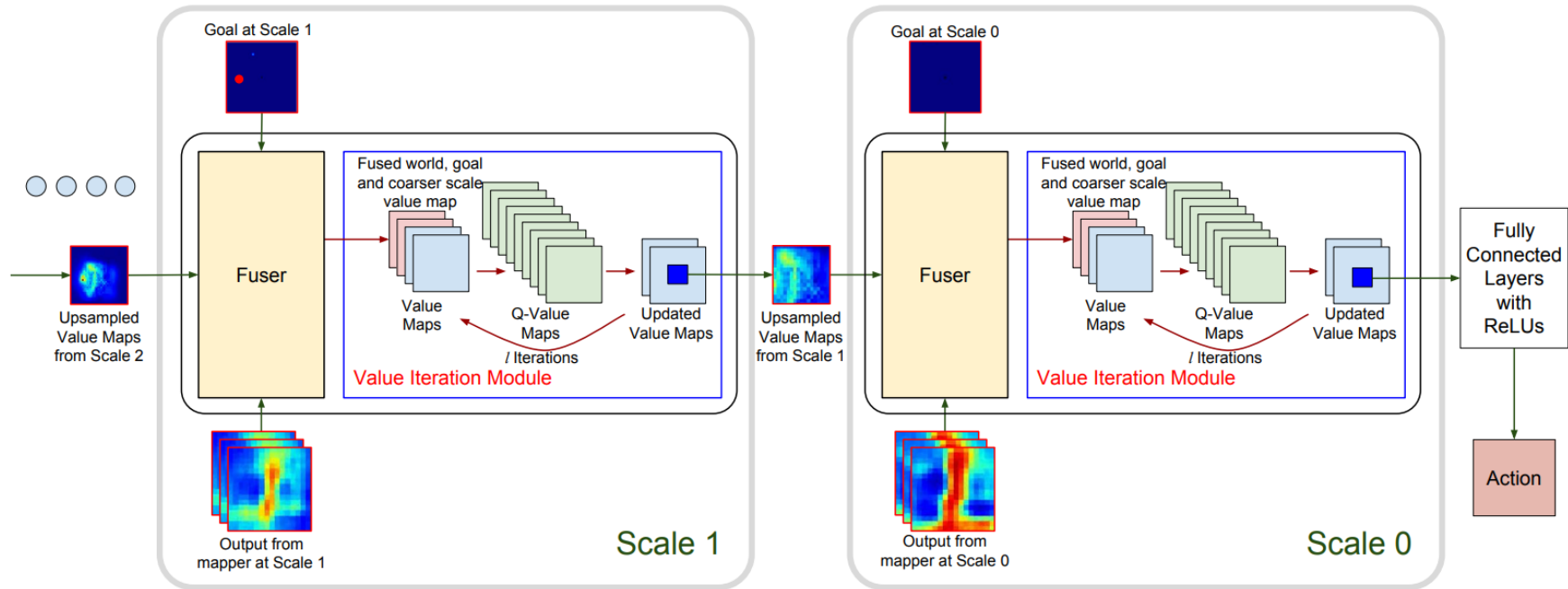
Neural Mapping

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

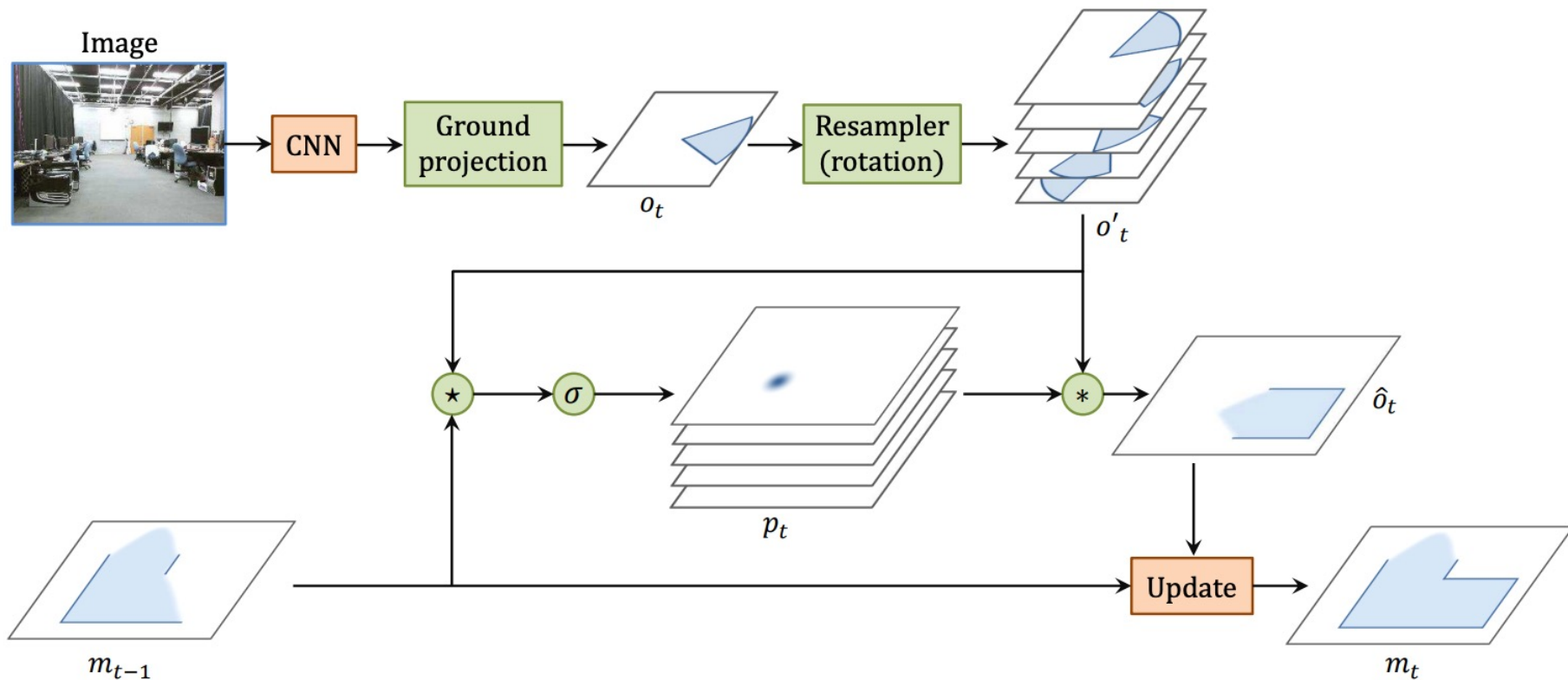


Hierarchical Planning

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



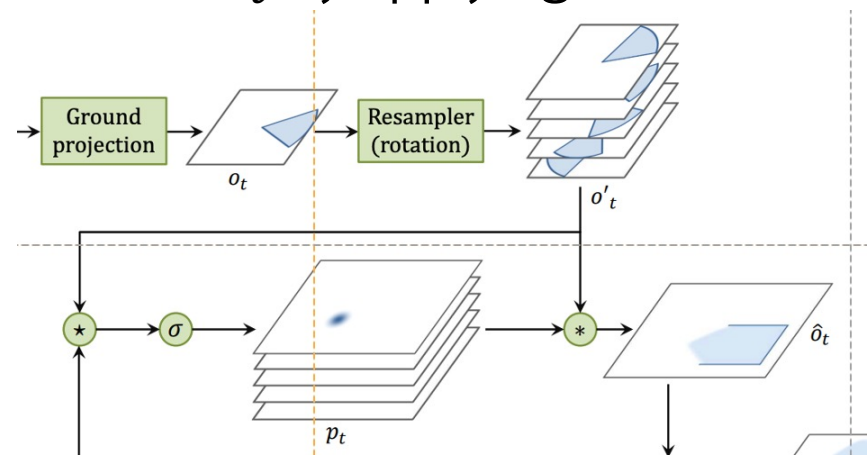
Simultaneous Localization and Registration



Simultaneous Localization and Registration

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

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



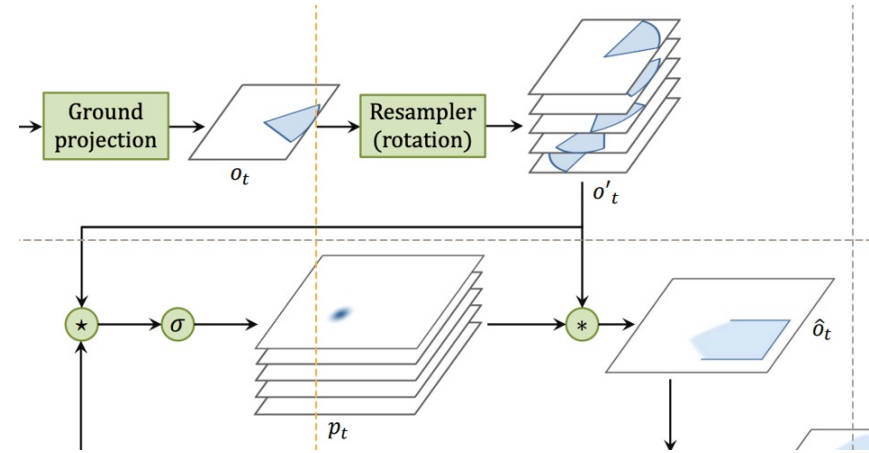
Simultaneous Localization and Registration

- 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.

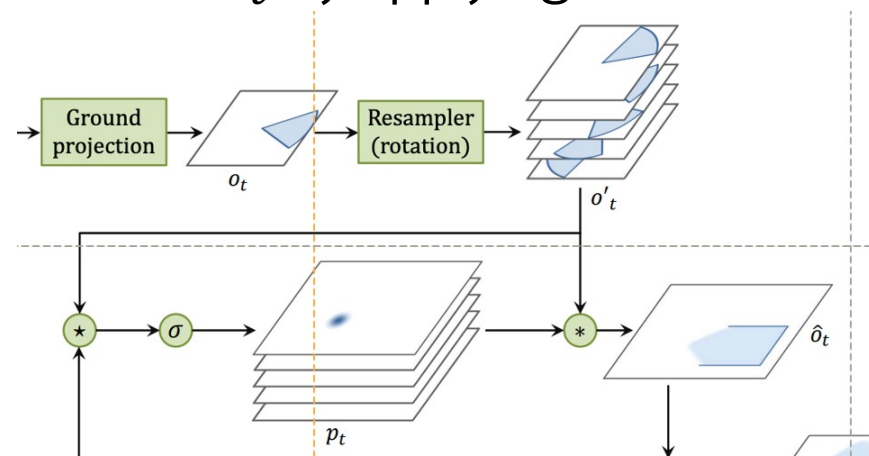
$$o'_{ijkl} = [R(o, 2\pi l/r)]_{ijk}.$$

- 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.

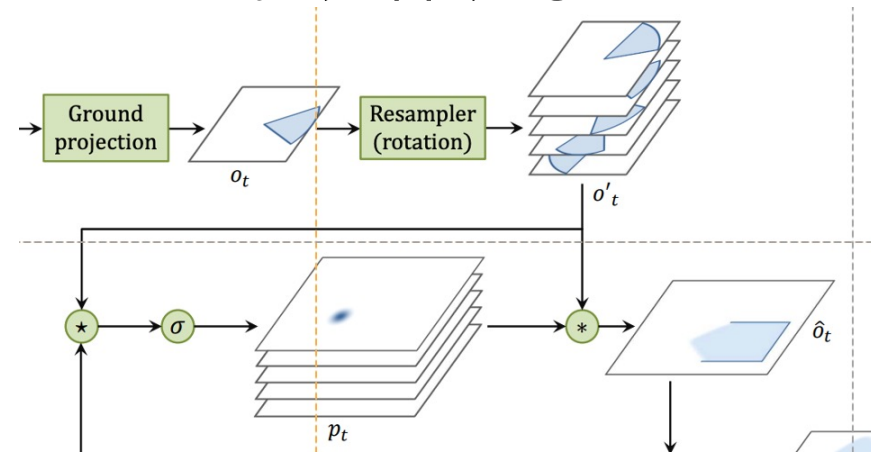
$$p_t = \text{Softmax}(m_{t-1} * o'_t).$$

- 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.

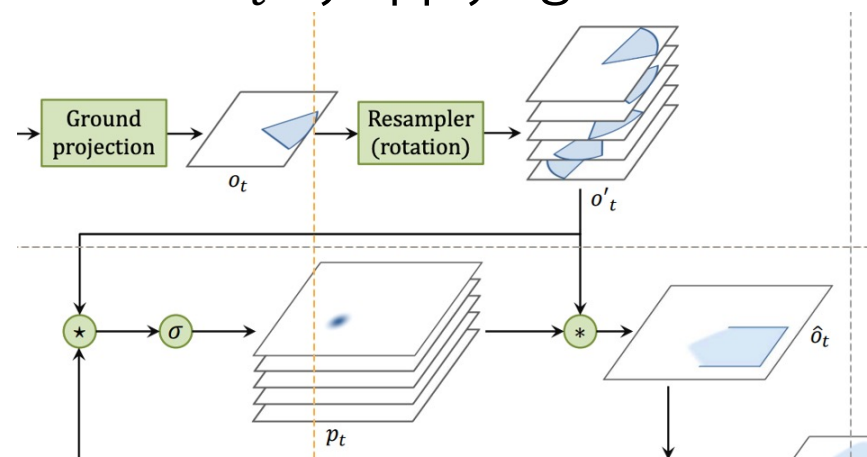
$$p_t = \text{Softmax}(m_{t-1} * o'_t).$$

- Transform observations into allocentric

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

- Update belief:

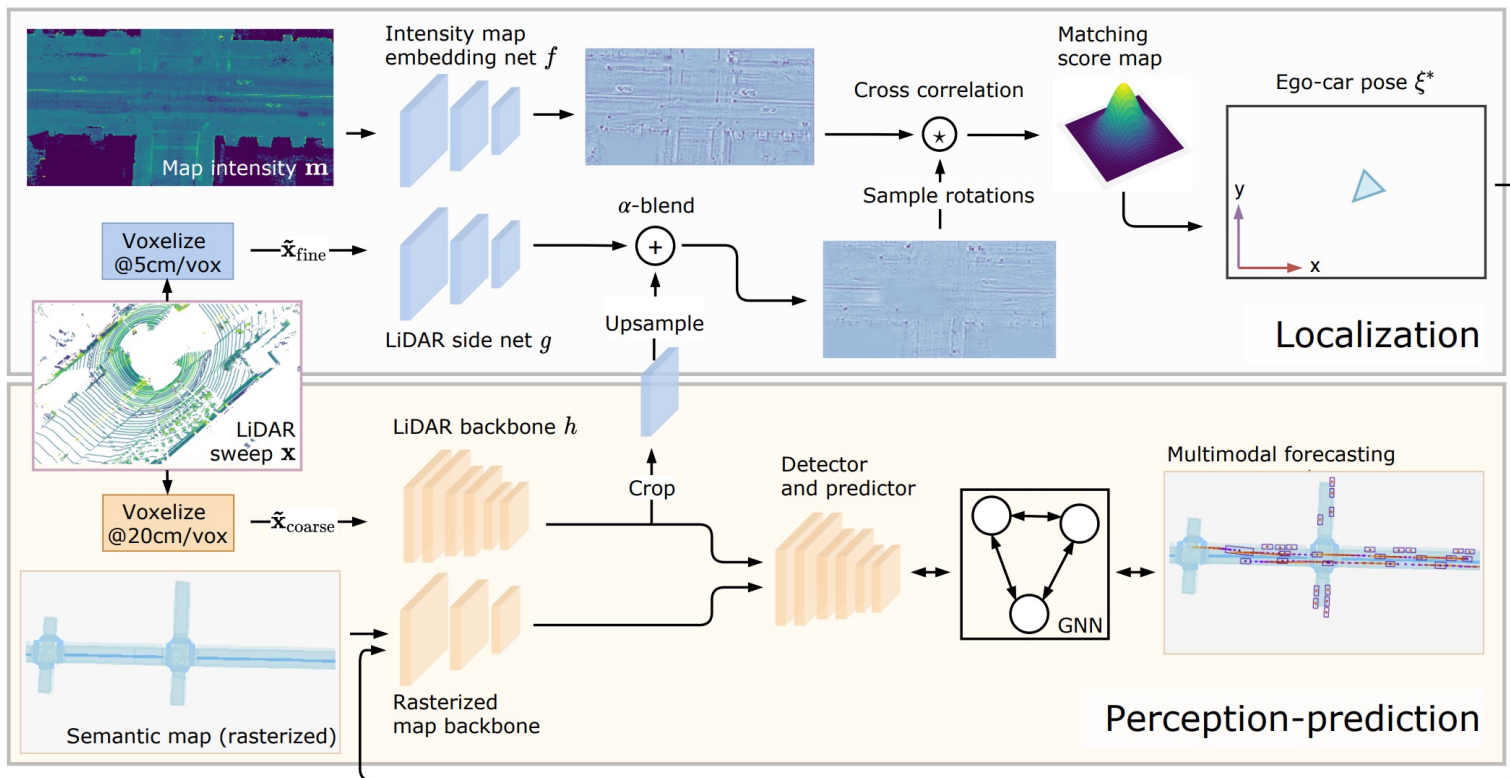
$$m_{i,j,t+1} = \text{LSTM}(m_{i,j,t}, \hat{o}_{i,j,t}).$$



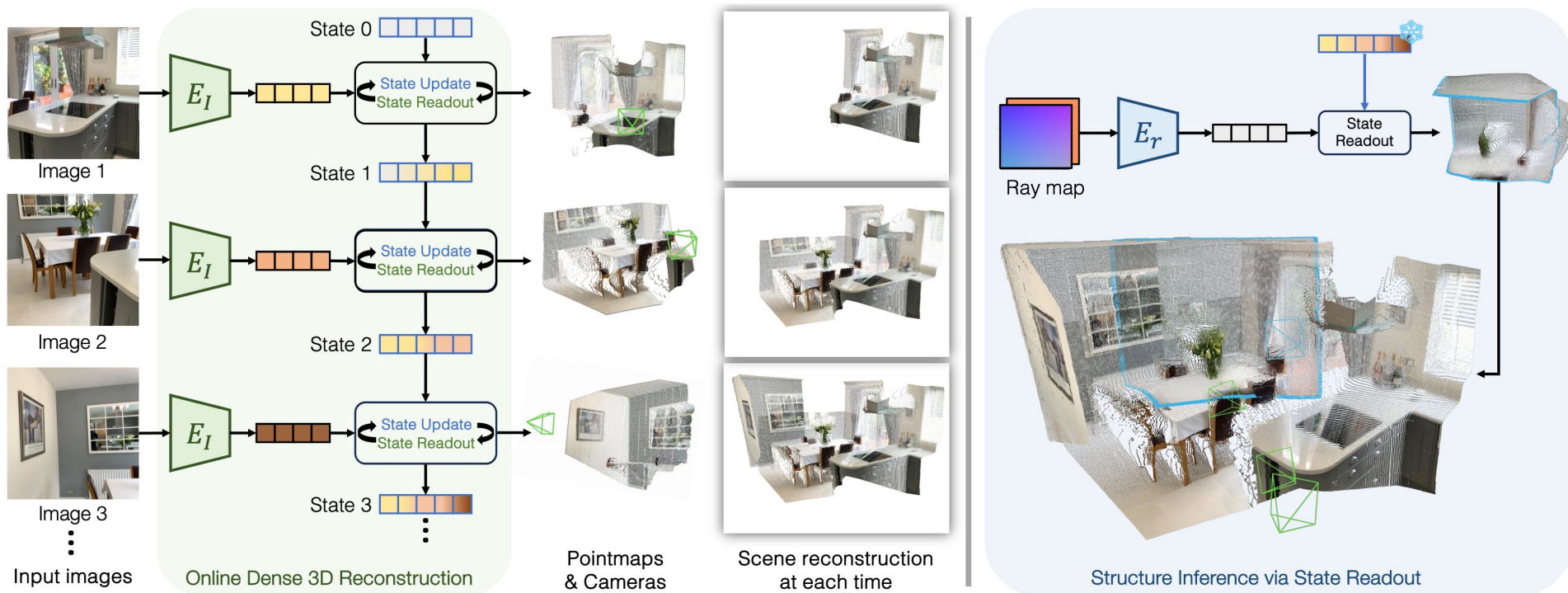
Loss:

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

Joint Localization, Perception and Prediction

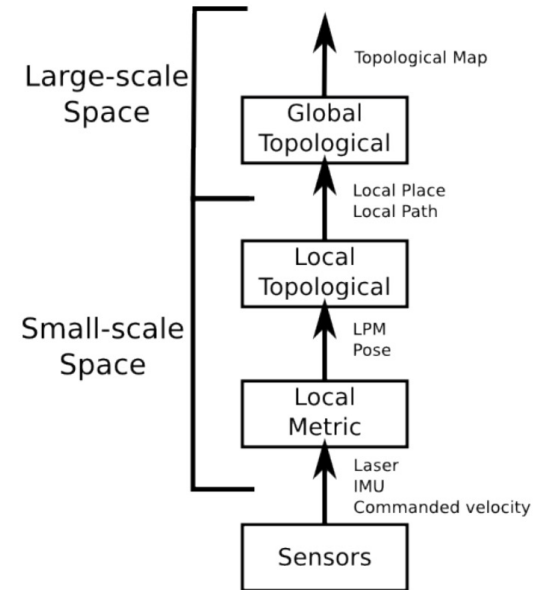
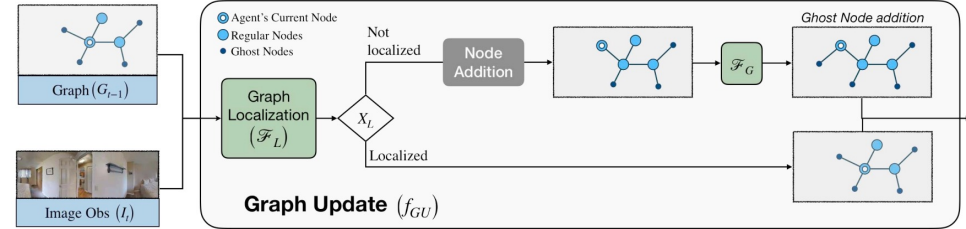
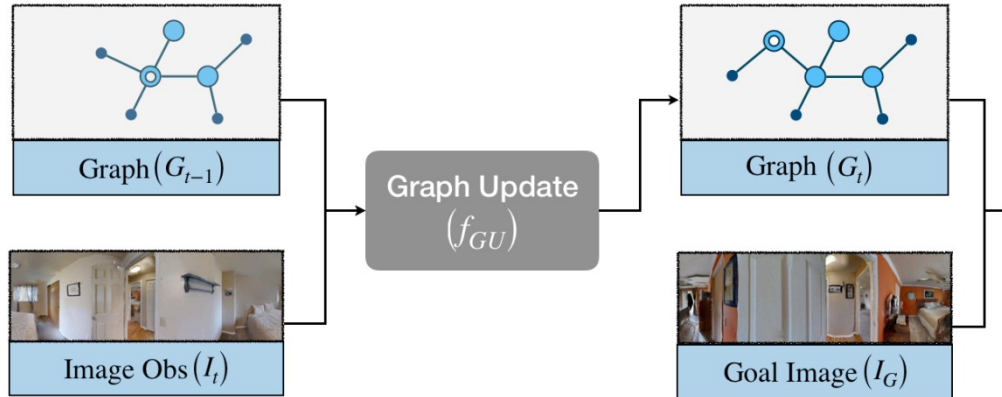


Continuous 3D Perception and Mapping



Topological Mapping

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



Summary - 3D Vision

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

Summary - 3D Vision

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

Summary - 3D Vision

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

Summary - 3D Vision

- Covers 3D, equivariance, 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 - 3D Vision

- Covers 3D, equivariance, 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 - 3D Vision

- Covers 3D, equivariance, 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 - 3D Vision

- Covers 3D, equivariance, 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).