

SE(3)-DiffusionFields:

Learning smooth cost functions for joint grasp
and motion optimization through diffusion

ICRA 2023

박지호

2024.03.12

SE(3)-DiffusionFields

Task

- Pick and Place: Grasp + Motion Planning

Contributions

- Diffusion for Grasp(SE(3))
- *Energy Prediction (not Score)* → Grasp & Motion Joint Optimization





Contents

1. Preliminaries
2. Method
3. Experiment
4. Discussion

Preliminaries

1. Motion Optimization
2. $SE(3)$ Lie Group
3. Score-Based Generative Model

1. Motion Planning

Methods

1) Data-Driven:

- Train Motion Generator → Sample top-k → eval & select

2) Optimization:

- Objective Function(ex. collision, target distance) → Trajectory Optimization

Common Strategy: *1) Data Driven*

- Advantages: Fast Inference, High performance,
- Disadvantages: Train Cost(Rely on Task Specific Generator!)

2. SE(3) LieGroup

"Lie 군 중 하나인 Special Euclidean 3 (SE(3))군은 3차원 공간 상에서 강체의 변환과 관련된 행렬과 이에 닫혀있는 연산들로 구성된 군을 의미한다."

Tangent Space (Lie Algebra)

$$\xi^{\wedge} = \begin{bmatrix} \omega^{\wedge} & \mathbf{v} \\ \mathbf{0} & 0 \end{bmatrix} \in se(3)$$

$$\begin{array}{c} \uparrow (\cdot)^{\wedge} \\ \downarrow (\cdot)^{\vee} \end{array}$$

$$\xi = \begin{bmatrix} \omega \\ \mathbf{v} \end{bmatrix} \in \mathbb{R}^6$$

Vector Space

$\exp(\cdot)$

$\log(\cdot)$

$\text{Exp}(\cdot)$

$\text{Log}(\cdot)$

Manifold (Lie Group)

$$\mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix} \in SE(3)$$

2. SE(3) LieGroup

- Rotation & Translation Representations

$$SE(3) = \left\{ \mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4} \mid \mathbf{R} \in SO(3), \mathbf{t} \in \mathbb{R}^3 \right\} \longleftrightarrow se(3) = \left\{ \xi^\wedge = \begin{bmatrix} \omega^\wedge & \mathbf{v} \\ \mathbf{0} & 0 \end{bmatrix} \in \mathbb{R}^{4 \times 4} \mid \xi = \begin{bmatrix} \omega \\ \mathbf{v} \end{bmatrix} \in \mathbb{R}^6 \right\}$$

Rotation(9-dim) + Translation(3-dim)
=> 12-dim

Rotation(3-dim) + Translation(3-dim)
=> 6-dim

2. SE(3) LieGroup

- Rotation & Translation Representations

$$SE(3) = \left\{ \mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4} \mid \mathbf{R} \in SO(3), \mathbf{t} \in \mathbb{R}^3 \right\}$$

$$\begin{bmatrix} x_2 \\ y_2 \\ z_2 \\ 1 \end{bmatrix} = \begin{bmatrix} R_{11} & R_{12} & R_{13} & T_x \\ R_{21} & R_{22} & R_{23} & T_y \\ R_{31} & R_{32} & R_{33} & T_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \\ z_1 \\ 1 \end{bmatrix}$$

Rotation(9-dim) + Translation(3-dim)
=> 16-dim

3. Score-Based Generative Model

- Probability Model: Gibbs(Boltzman) Distribution
- Score : Gradient of Log probability = Gradient of Energy
- Diffusion Model: Langevin Dynamics

3. Score-Based Generative Model

Gibbs(Boltzman) Distribution

- Probability Model: Energy(ϵ) of $i \rightarrow$ Probability of i

Physics

$$p_i = \frac{1}{Q} e^{-\epsilon_i/(kT)} = \frac{e^{-\epsilon_i/(kT)}}{\sum_{j=1}^M e^{-\epsilon_j/(kT)}}$$



Generative Model

$$p_{\theta}(\mathbf{x}) = \frac{\exp(-E_{\theta}(\mathbf{x}))}{Z(\theta)}$$

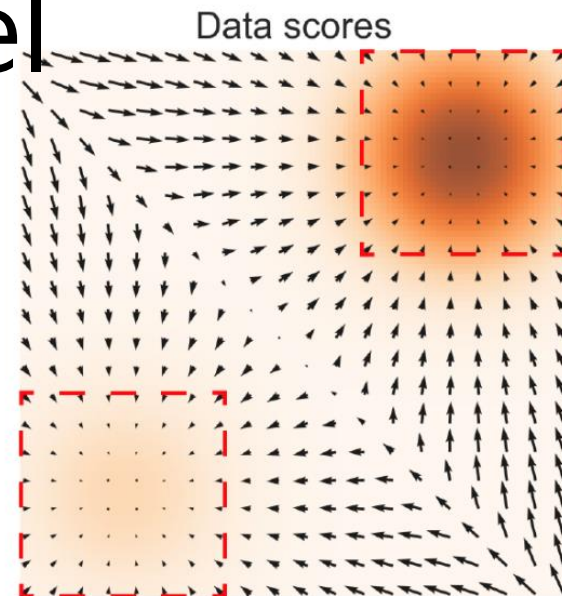
$$Z(\theta) = \int \exp(-E_{\theta}(\mathbf{x})) d\mathbf{x}$$

3. Score-Based Generative Model

Score: **Gradient of log Probability = Gradient of Energy**

$$\text{Score} = \nabla_x \log p(x)$$

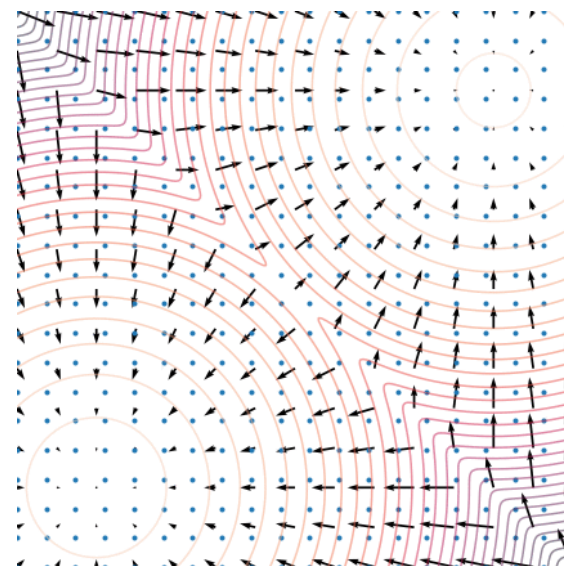
$$= \nabla_x \log \frac{1}{Z} e^{-E_\theta(x)} = -\nabla_x E_\theta(x) - \nabla_x \log Z = -\nabla_x E_\theta(x)$$



Diffusion Model: Score Predictor $s_\theta(x)$

$$\text{Loss : } \frac{1}{2} \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})p_{\text{data}}(\mathbf{x})} [\|\mathbf{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}} | \mathbf{x})\|_2^2].$$

(Denoising Score Matching)



3. Score-Based Generative Model

Langevin Dynamics → Diffusion Sampling

- mathematical modeling of the dynamics of molecular systems
- Nature prefers to go to Low Energy (Low Energy = High Probability)

Physics

$$M \ddot{\mathbf{X}} = -\nabla U(\mathbf{X}) - \gamma M \dot{\mathbf{X}} + \sqrt{2 M \gamma k_B T} \mathbf{R}(t),$$



Generative Model

$$\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} + \frac{\epsilon}{2} \nabla_{\mathbf{x}} \log p(\tilde{\mathbf{x}}_{t-1}) + \sqrt{\epsilon} \mathbf{z}_t,$$

Method

1. **Grasp**: SE(3)-DiffusionField

- Object Dependent Grasp(SE(3)) Prediction

2. **Grasp + Motion Joint Optimization** (via Diffusion Inverse Sampling)

- No training, Only Optimization (w/ pretrained SE(3)-Diff)

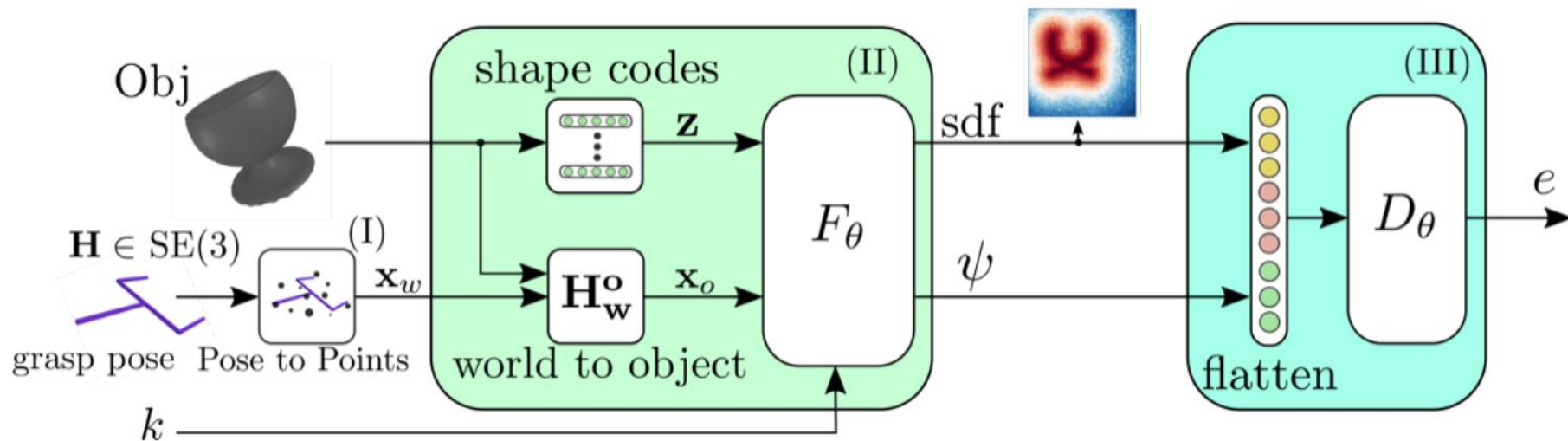
1. Grasp: SE(3)-DiffusionField

Train Dataset

- Acronym: Grasp Dataset(Simulation)

SE(3)-Diff

- Input: Object Pose & id, Grasp(6DoF SE(3)), Diffusion Timestep(k)
- Output: Energy of Grasp(H)



1. Grasp: SE(3)-DiffusionField

If Model Predicts **Energy** instead of **Score**...

How can we Train Diffusion Denoising?

$$\rightarrow s_{\theta}(\mathbf{H}, k) = -DE_{\theta}(\mathbf{H}, k)/D\mathbf{H}$$

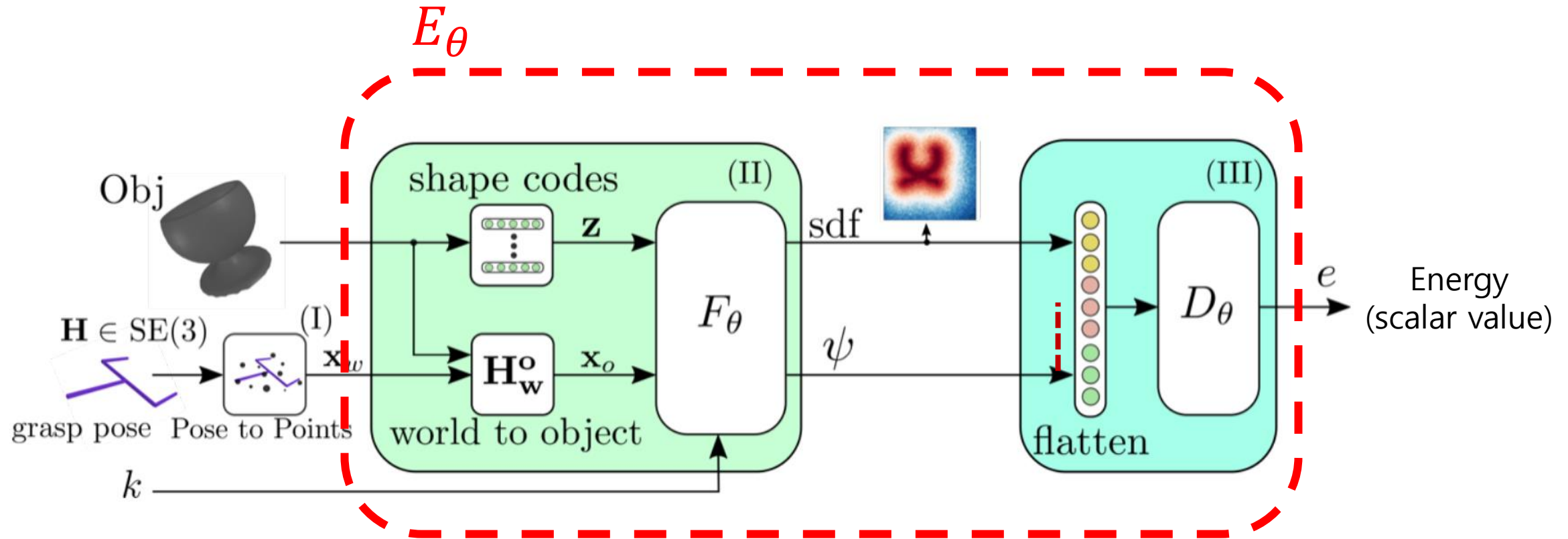
Score Matching Loss

$$\mathcal{L}_{\text{dsm}} = \frac{1}{L} \sum_{k=0}^L \mathbb{E}_{\mathbf{x}, \hat{\mathbf{x}}} [\|s_{\theta}(\hat{\mathbf{x}}, k) - \nabla_{\hat{\mathbf{x}}} \log \mathcal{N}(\hat{\mathbf{x}}|\mathbf{x}, \sigma_k^2 \mathbf{I})\|],$$

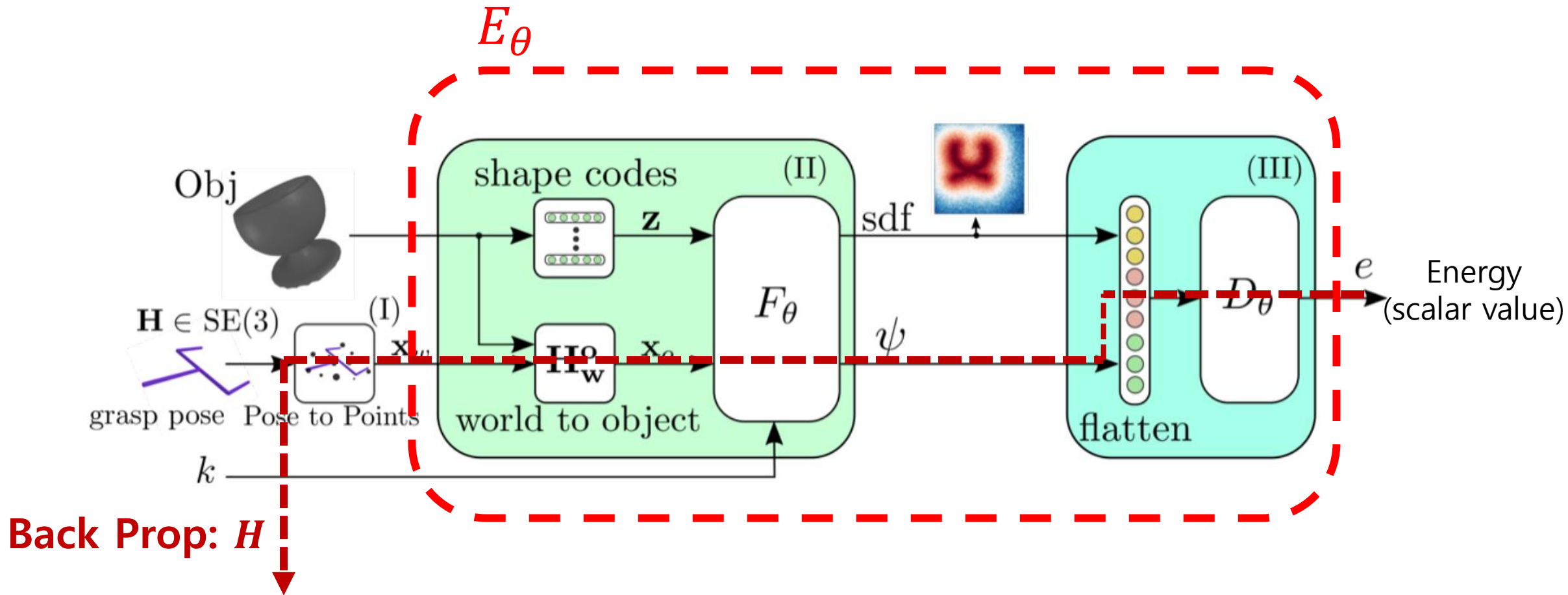
Diffusion Sampling(Denoising)

$$\mathbf{x}_{k-1} = \mathbf{x}_k + \frac{\alpha_k^2}{2} s_{\theta}(\mathbf{x}_k, k) + \alpha_k \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

1. Grasp: SE(3)-DiffusionField

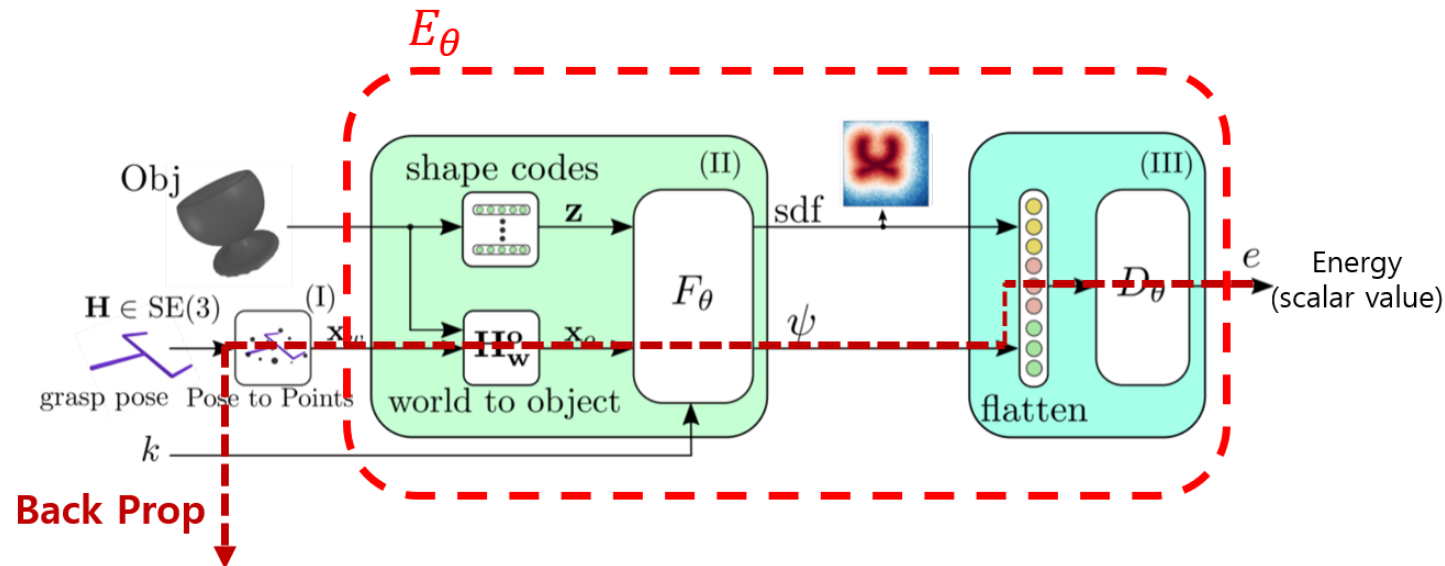


1. Grasp: SE(3)-DiffusionField



$$s_\theta(\mathbf{H}, k) = -DE_\theta(\mathbf{H}, k)/D\mathbf{H}$$

1. Grasp: SE(3)-DiffusionField

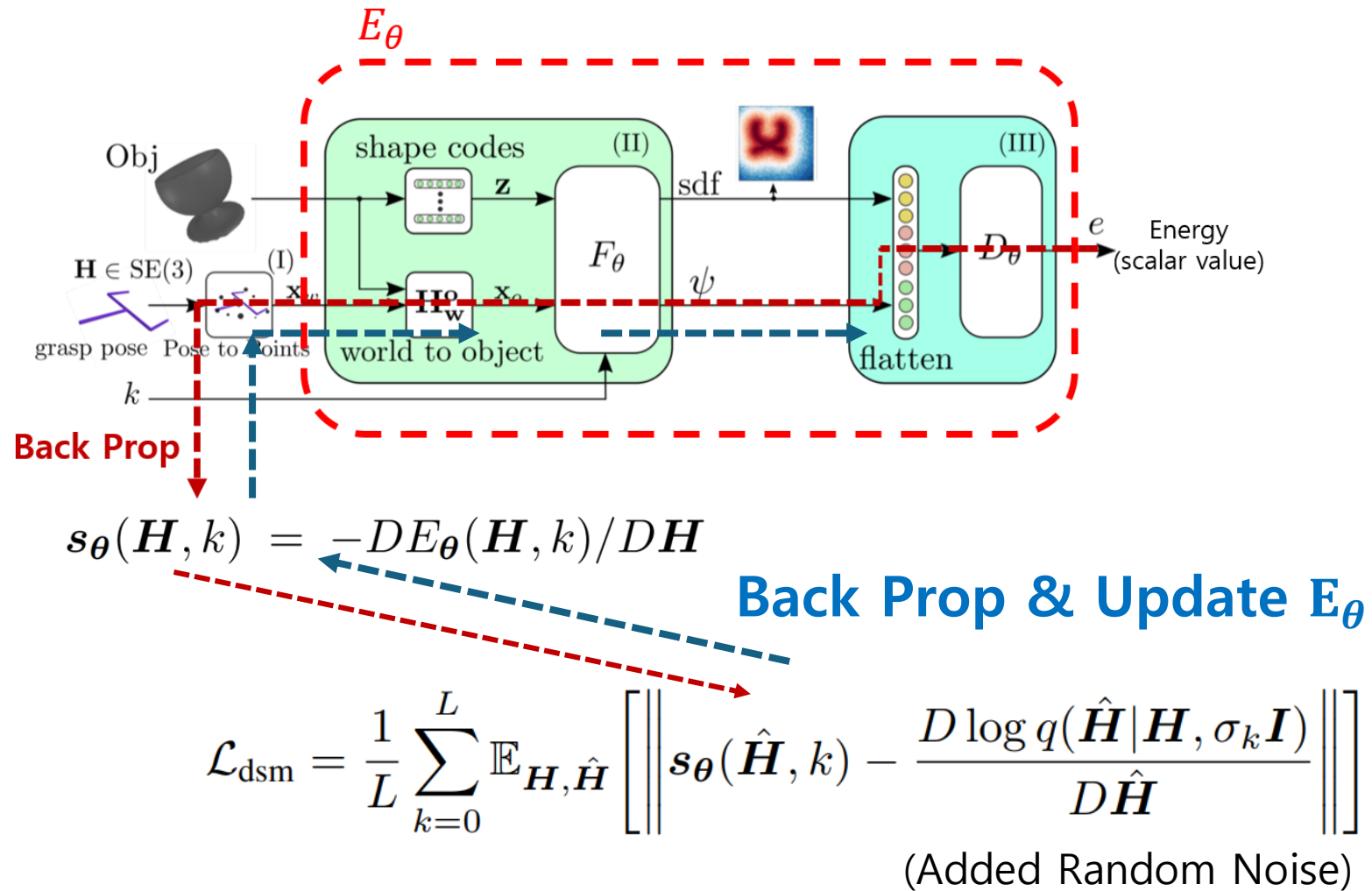


$$\mathbf{s}_\theta(\mathbf{H}, k) = -DE_\theta(\mathbf{H}, k)/D\mathbf{H}$$

$$\mathcal{L}_{\text{dsm}} = \frac{1}{L} \sum_{k=0}^L \mathbb{E}_{\mathbf{H}, \hat{\mathbf{H}}} \left[\left\| \mathbf{s}_\theta(\hat{\mathbf{H}}, k) - \frac{D \log q(\hat{\mathbf{H}} | \mathbf{H}, \sigma_k \mathbf{I})}{D\hat{\mathbf{H}}} \right\| \right],$$

(Added Random Noise)

1. Grasp: SE(3)-DiffusionField



1. Grasp: SE(3)-DiffusionField

Train Loss

Diffusion Loss + SDF Loss
(3D Edge of Object)

Algorithm 1: Grasp SE(3)-DiF Training

Given: θ_0 : initial params for z, F_θ, D_θ ;
Datasets: $\mathcal{D}_o : \{m, H_w^o\}$, object ids and poses,
 $\mathcal{D}_{sdf}^m : \{x, sdf\}$, 3D positions x and sdf for object m ,
 $\mathcal{D}_g^m : \{H\}$ succesful grasp poses for object m ;

```
1 for  $s \leftarrow 0$  to  $S - 1$  do
2    $k, \sigma_k \leftarrow [0, \dots, L]$ ; // sample noise scale
3    $m, H_w^o \in \mathcal{D}_o$ ; // sample objects ids and poses
4    $z = \text{shape codes}(m)$ ; // get shape codes
5   SDF train
6    $x, sdf \in \mathcal{D}_{sdf}^m$ ; // get 3D points and sdf for obj.  $m$ 
7    $\hat{sdf}, _ = F_\theta(H_w^o x, z, k)$ ; // get predicted sdf
8    $L_{sdf} = \mathcal{L}_{mse}(\hat{sdf}, sdf)$ ; // compute sdf error
9   Grasp diffusion train
10   $H \sim \mathcal{D}_g^m$ ; // Sample success grasp poses for obj.  $m$ 
11   $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_k I)$ ; // sample white noise on  $k$  scale
12   $\hat{H} = H \text{Expmap}(\epsilon)$ ; // perturb grasp pose Eq. (4)
13   $x_n^o = \hat{H} x_n$ ; // Transform to N 3d points (see Figure 3)
14   $sdf_n, \psi_n = F_\theta(x_n^o, z_b, k)$ ; // get features
15   $\Psi = \text{Flatten}(\hat{sdf}_n, \psi_n)$ ; // Flatten the features
16   $e = D_\theta(\Psi)$ ; // compute energy
17   $L_{dsm} = \mathcal{L}_{dsm}(e, \hat{H}, H, \sigma_k)$ ; // Compute dsm loss Eq. (5)
18  Parameter update
19   $L = L_{dsm} + L_{sdf}$ ; // Sum losses
20   $\theta_{s+1} = \theta_s - \alpha \nabla_\theta L$ ; // Update parameters
21 return  $\theta^*$ ;
```

2. Grasp + Motion Optimization

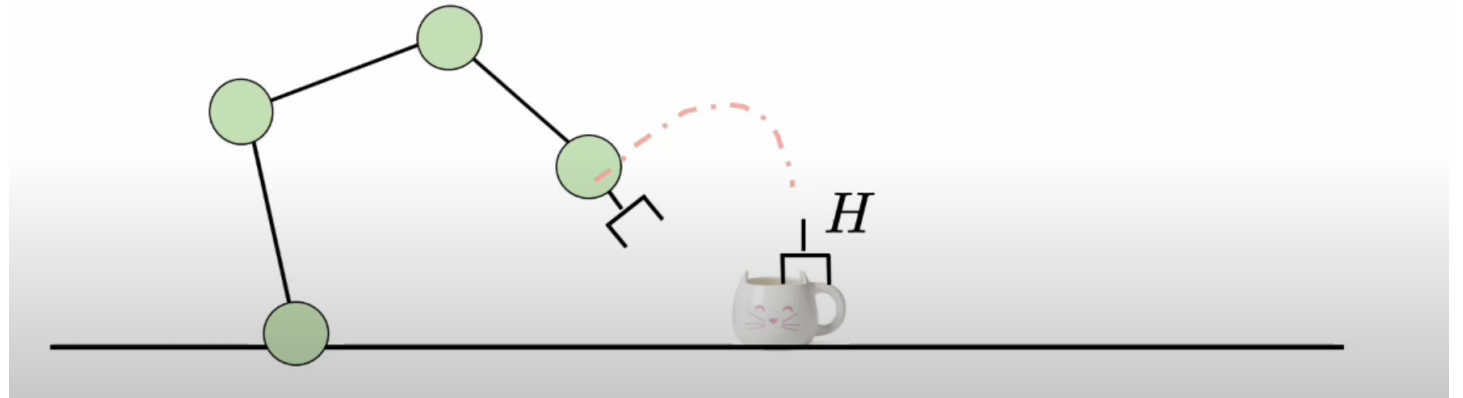
q_t : Robot Joints

$\tau: \{q_t\}_{t=1}^T$: Trajectory

\mathcal{J} : Objective Function

Grasp Pose selection --> Motion Planning

$$\tau^* = \arg \max_{\tau} \sum_{t=0}^T c_{coll}(\mathbf{q}_t) + \cdots + c_{reach}(\mathbf{q}_T, H)$$



Motion Optimization

$$\tau^* = \arg \min_{\tau} \mathcal{J}(\tau) = \arg \min_{\tau} \sum_j \omega_j c_j(\tau)$$

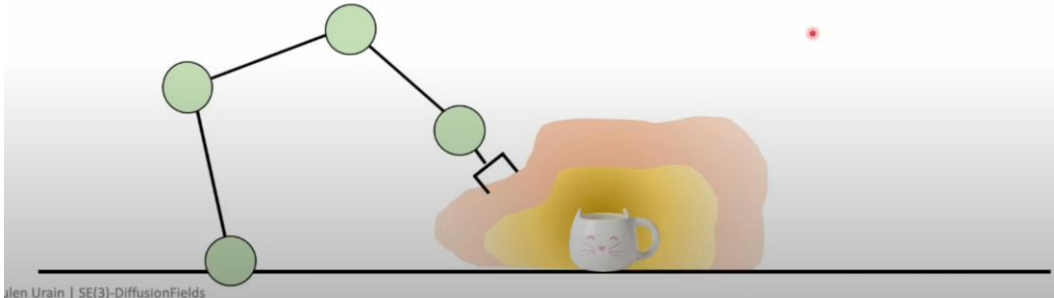
2. Grasp + Motion Optimization

1) Add Energy in \mathcal{J}

- $c_{grasp} = E_{\theta}(q_t)$

- Avoid two steps (Grasp Selection -> Motion planning) and only ONE STEP

$$\tau^* = \arg \max_{\tau} \sum_{t=0}^T c_{coll}(\mathbf{q}_t) + \dots + c_{grasp}(\mathbf{q}_T)$$



2) Diffusion Sampling

- Diffusion Inverse Process
(Instead of Gradient Descent)

$$\tau_{k-1} = \tau_k + 0.5 \alpha_k^2 \nabla_{\tau_k} \log q(\tau|k) + \alpha_k \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

$$q(\tau|k) \propto \exp(-\mathcal{J}(\tau, k))$$

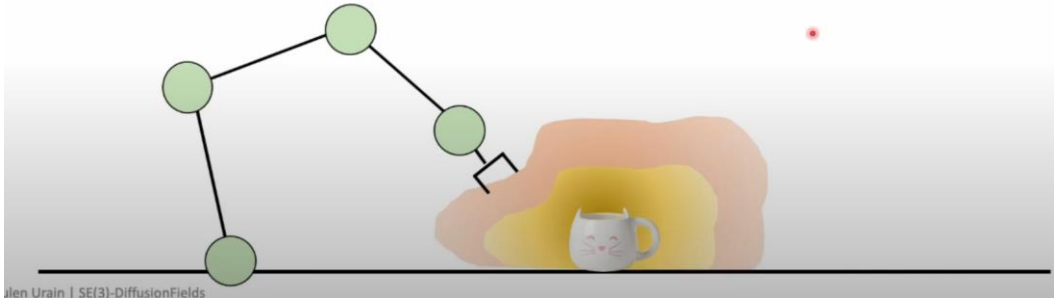
2. Grasp + Motion Optimization

1) Add Energy in \mathcal{J}

- $c_{grasp} = E_{\theta}(q_t)$

- Avoid two steps (Grasp Selection -> Motion planning) and only ONE STEP

$$\tau^* = \arg \max_{\tau} \sum_{t=0}^T c_{coll}(\mathbf{q}_t) + \dots + c_{grasp}(\mathbf{q}_T)$$



2) Diffusion Sampling

- Diffusion Inverse Process
(Instead of Gradient Descent)

$$\tau_{k-1} = \tau_k + 0.5 \alpha_k^2 \nabla_{\tau_k} \log q(\tau|k) + \alpha_k \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

$$q(\tau|k) \propto \exp(-\mathcal{J}(\tau, k))$$

**In Fact,
Adding Gradient of motion objectives, in Grasp Diffusion!**

Experiments

1) Grasp Evaluation (Simulation)

2) Grasp + Motion Evaluation (Simulation)

- Comparing with other models

3) Real World

- 20 trials for each experiment(4)
- Results: 100%, 90%, 95%, 100%

Experiments

1) Grasp Evaluation

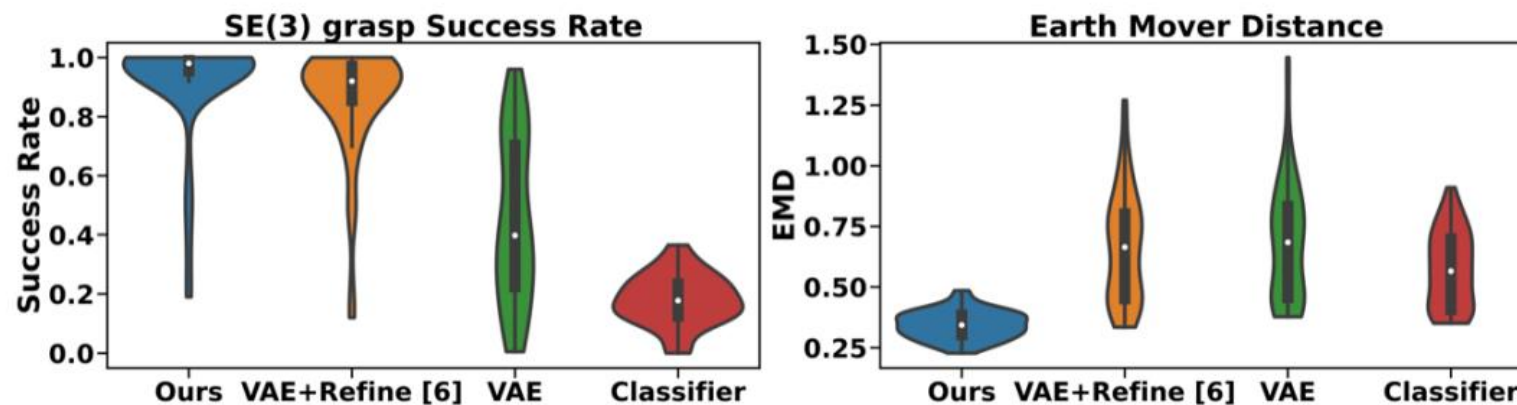


Fig. 4: 6D grasp pose generation experiment. Left: Success rate evaluation. Right: Earth Mover Distance (EMD) evaluation metrics (lower is better).

Experiments

2) Grasp + Motion Evaluation

- joint(class): Grasp Classifier + Motion Objective → Optimization

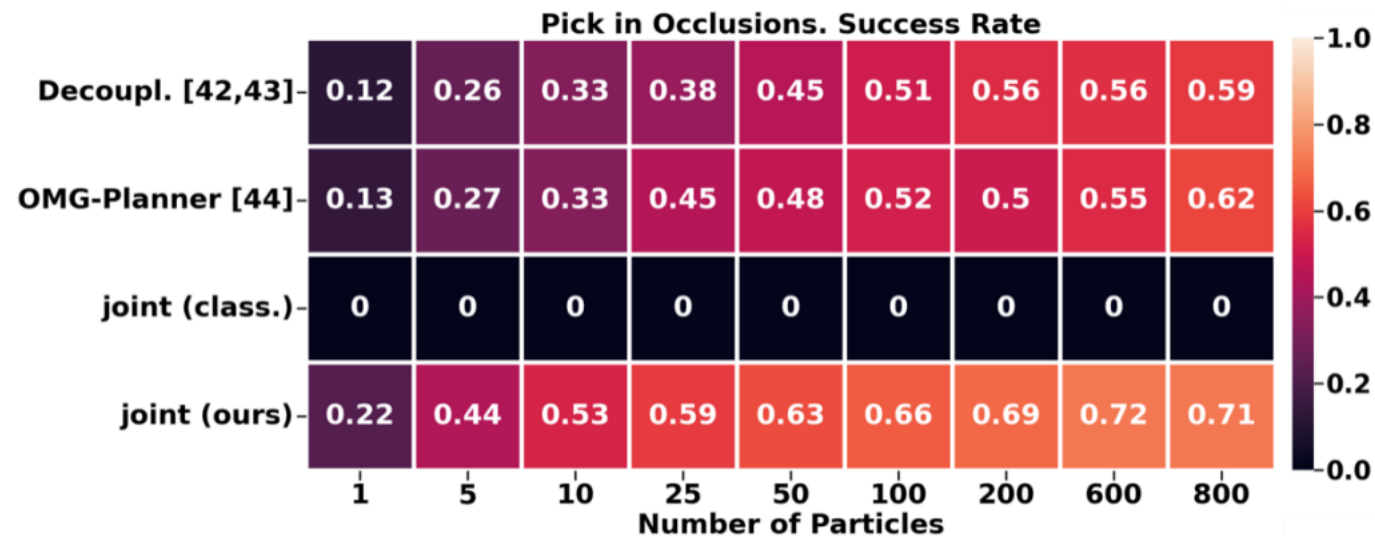
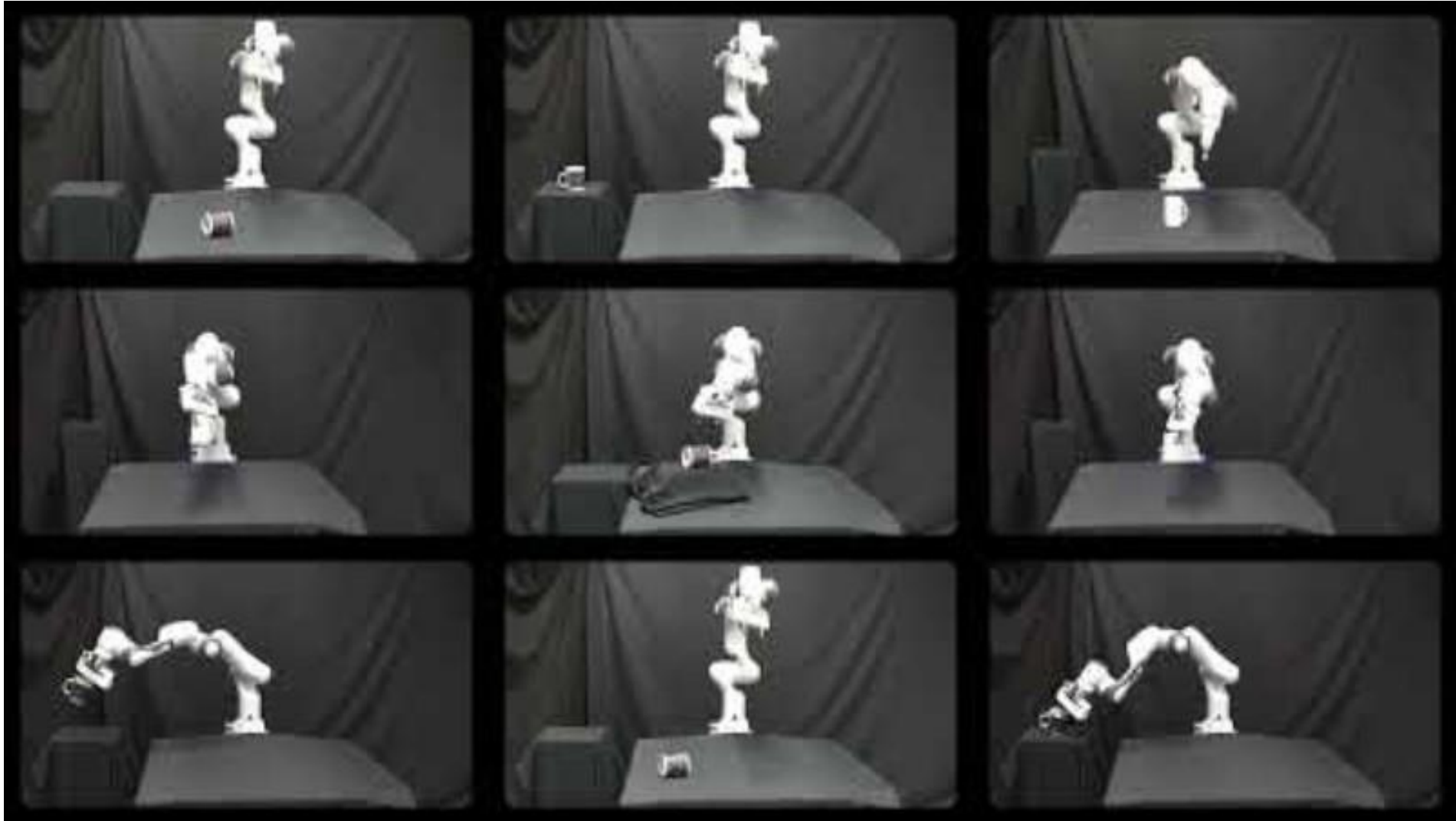
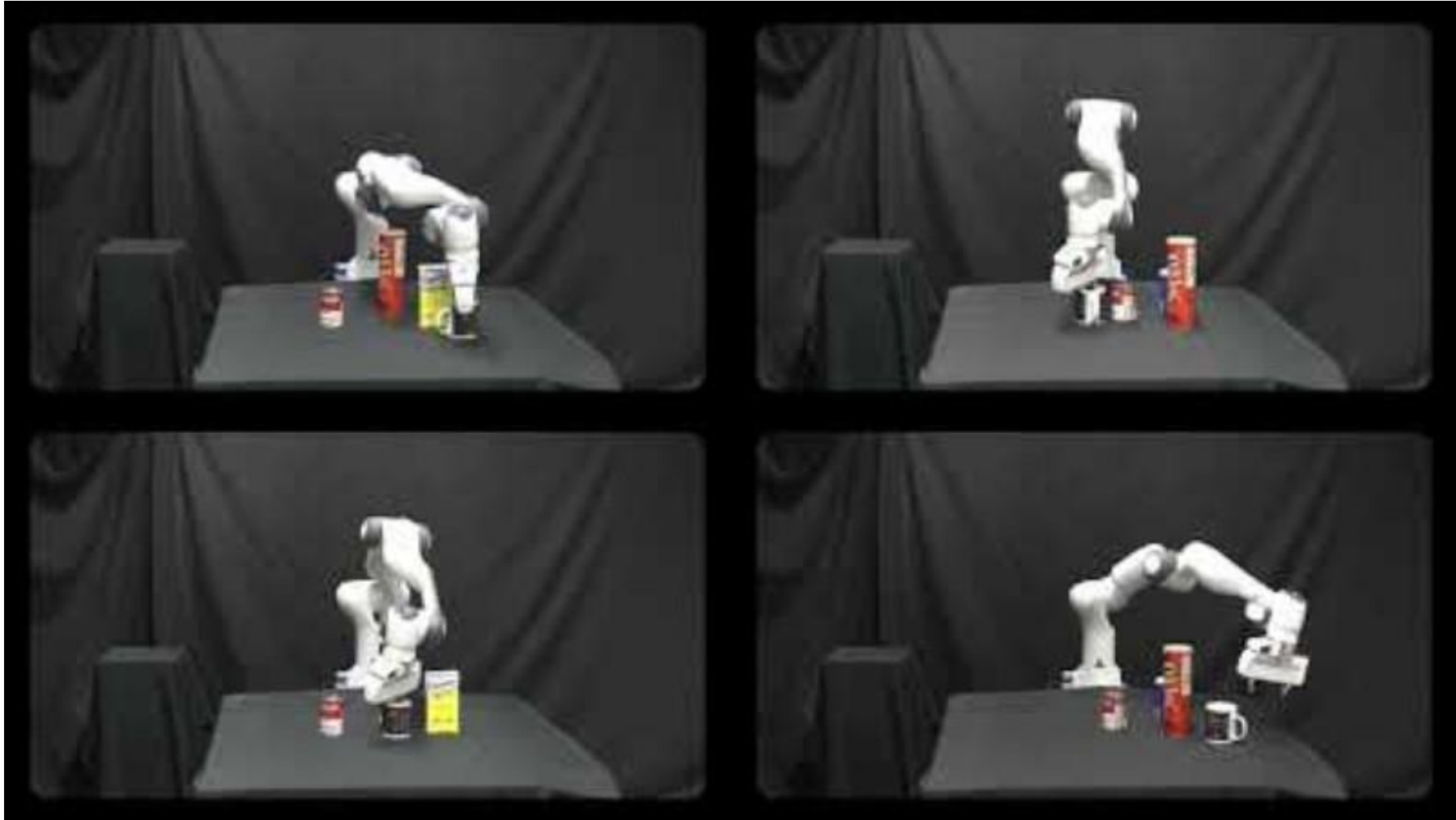


Fig. 5: Evaluation Pick in occlusion. We measure the success rate of 4 different methods based on different number of initializations.

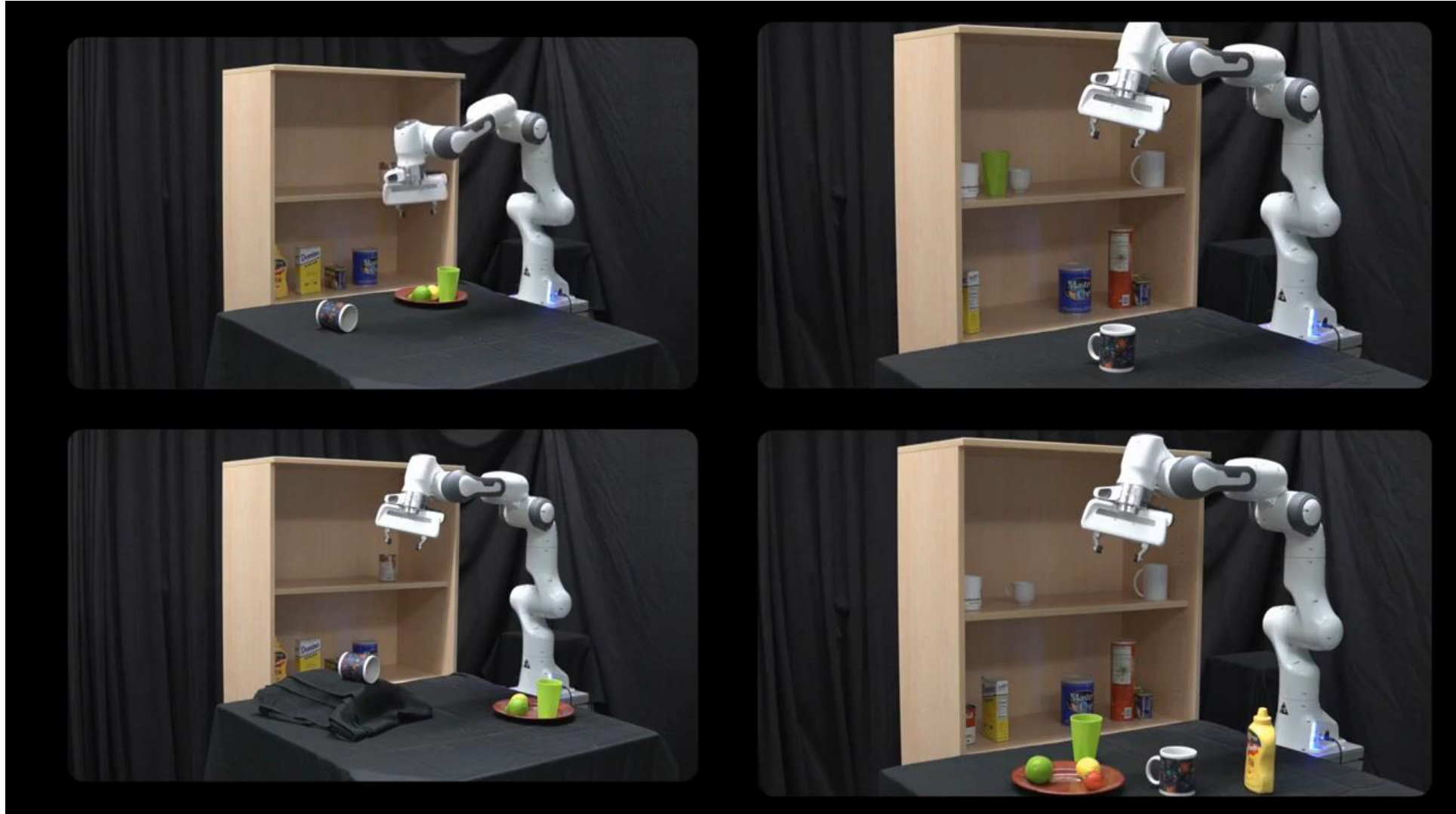
Experiments



Experiments



Experiments



Discussion

Is Predicting Energy Better?

- SE(3) dimension이 낮아서 energy-based modeling이 가능한 것일까?
- Score prediction 해도 Joint Optimization은 똑같이 적용이 가능한데, Unified Objective Function 스토리를 위한 것일까?



Thank You!