

인공지능 기반의 로봇 파지 계획 기술

2021. 06. 23(Wed.)

Sungkyunkwan University

Robotics and Intelligent System Engineering Lab.

Minseok Kang

연구의 최종목표

• Problem Statement •

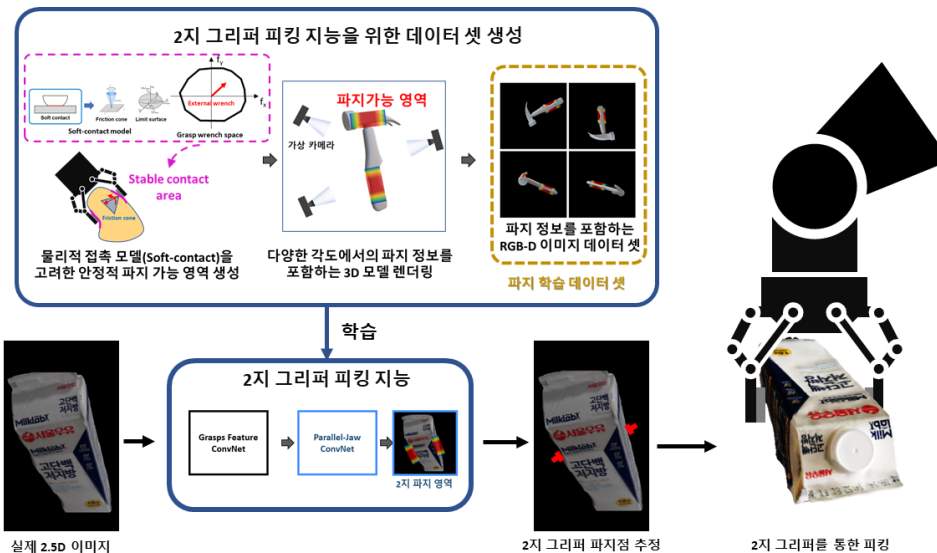
- 산업에서 성공적으로 사용되는 공압 및 2지 그리퍼로 **다양한 물체를 빠르고 안정적으로 피킹** 할 수 있는 **지능**은 무엇인가
- 시간과 노력을 최소화**하며 **양질의 학습 데이터**를 얻기 위한 방법은 무엇인가?
- 센싱, 제어에서 발생하는 **오차에 강인하도록 피킹 지능을 학습하기 위한 방법**은 무엇인가?

• Proposed Solution •

- 2지 및 공압 그리퍼 피킹 **지능의 융합**을 통해 **2지와 공압이 융합된 그리퍼로 물체를 피킹**할 수 있도록 지능 증강
- 실제 RGB-D 데이터와 파지 정보를 포함한 가상 RGB-D 데이터를 **Canonical-form 형태로 변환** 후 학습하여 **예측 성공률을 높임**
- 피킹 지능 학습 데이터의 안정적인 파지위치 라벨을 **파지점으로부터 파지영역으로 확장**

□ 그리퍼 피킹 지능의 학습을 위한 물체 파지 데이터 생성 기술 개발

□ 비전 센서, 3D 센서 등의 센서 데이터를 기반으로 물체 파지 가능 위치를 추정하는 그리퍼 피킹 지능 기술 개발



<2지 그리퍼 피킹을 위한 물체 파지 학습 데이터 생성 및 지능 개발>

What is good grasp?

- Find a gripper configuration that maximizes a success(or quality) metric

- Empirical methods

1) Human Label [1-2]

Cornell Grasp Dataset



Amazon Picking Challenging (MIT-Princeton)

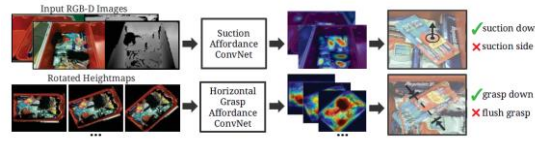
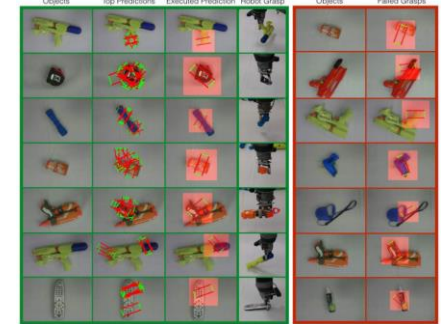
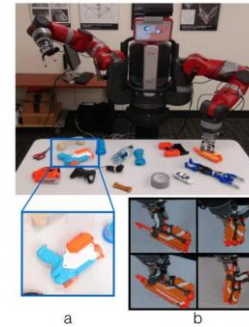


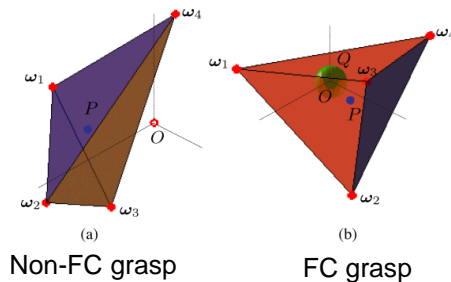
Fig. 5. Suction and grasp affordance prediction. Given multi-view RGB-D images, we estimate suction affordances for each image with a fully convolutional residual network. We then aggregate the predictions on a 3D point cloud, and generate suction down or suction side proposals based on surface normals. In parallel, we merge RGB-D images into an RGB-D heightmap, rotate it by 16 different angles, and estimate horizontal grasp for each heightmap. This effectively produces affordance maps for 16 different grasp angles, from which we generate the grasp down and flush grasp proposals.

2) Physical error [3]

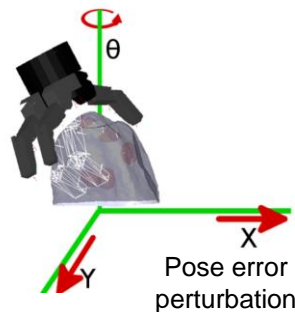


- Analytic method: Consider performance according to physical models

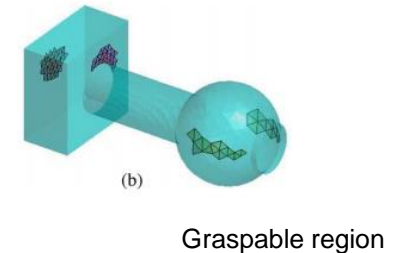
1) Grasp Wrench Space(GWS) [4]



2) Robust GWS [5]



3) Independent Contact Region [4]

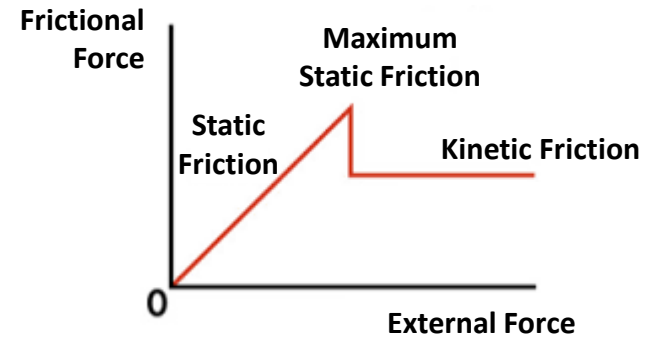
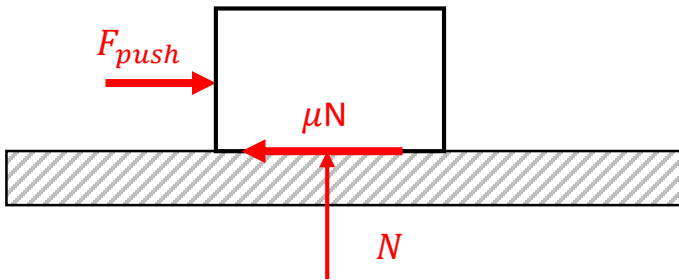


Graspable region

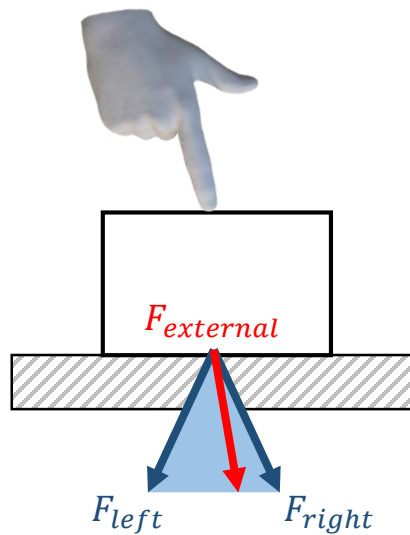
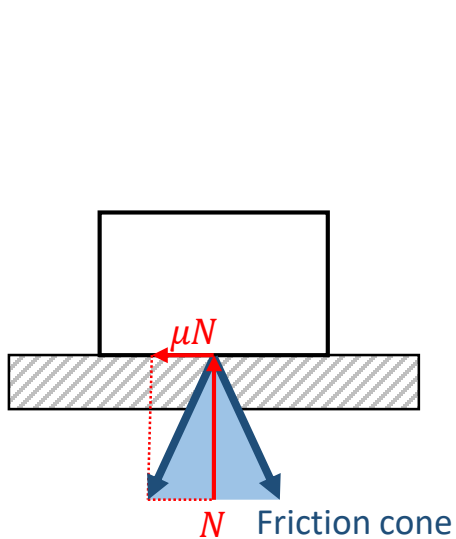
- [1] http://pr.cs.cornell.edu/grasping/rect_data/data.php
- [2] Zeng, Andy, et al. "Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching." *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018.
- [3] Pinto, Lerrel, and Abhinav Gupta. "Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours." *2016 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2016.
- [4] Roa, Máximo A., and Raúl Suárez. "Computation of independent contact regions for grasping 3-d objects." *IEEE Transactions on Robotics* 25.4 (2009): 839-850.
- [5] Weisz, Jonathan, and Peter K. Allen. "Pose error robust grasping from contact wrench space metrics." *2012 IEEE international conference on robotics and automation*. IEEE, 2012.

Grasp Quality

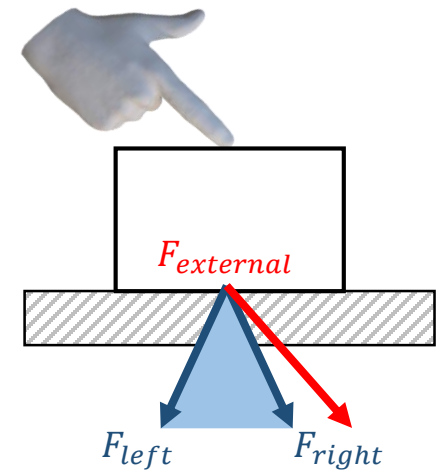
□ Coulomb friction



□ Friction cone



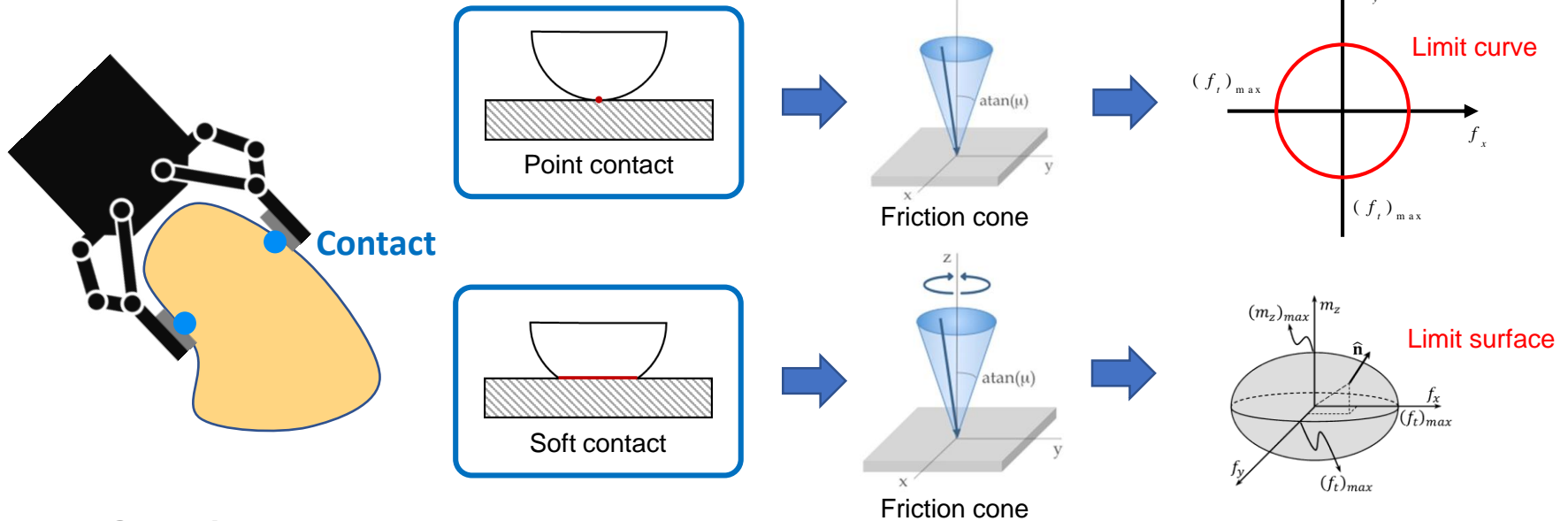
It can resist external force



It can not resist external force

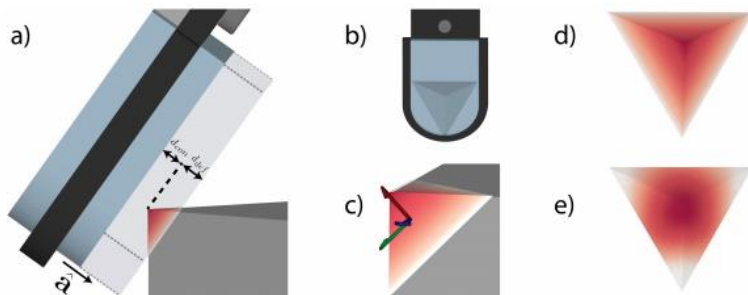
Grasp Quality

□ Contact model

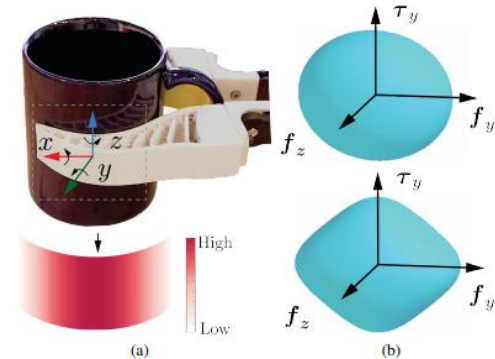


□ Complicate contact models

Assume gripper pad is deformable [1]



Nonplanar surface contact [2]

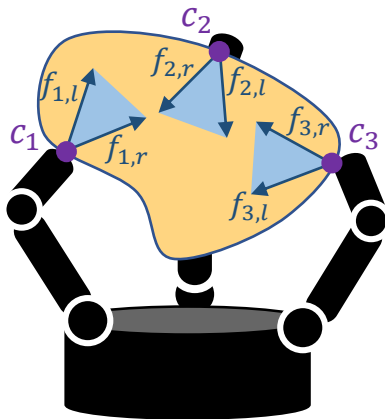


Grasp Quality

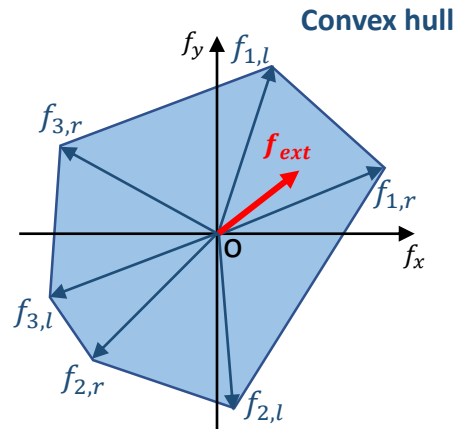
Grasp Wrench Space(GWS)

Assumption :

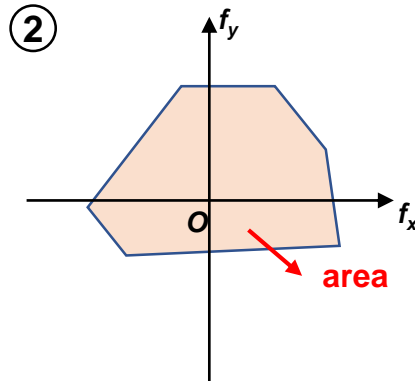
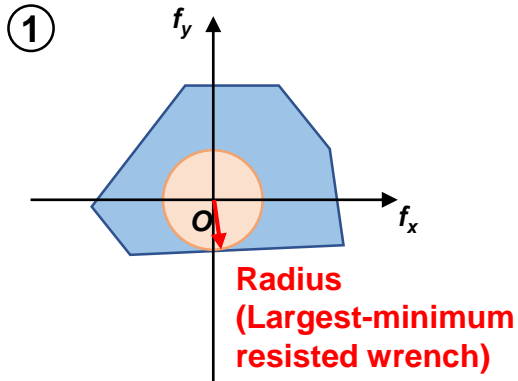
- 1) 2D plane object
- 2) No rotational motion. Thus, we don't consider torque



Construct GWS



Grasp Quality



...

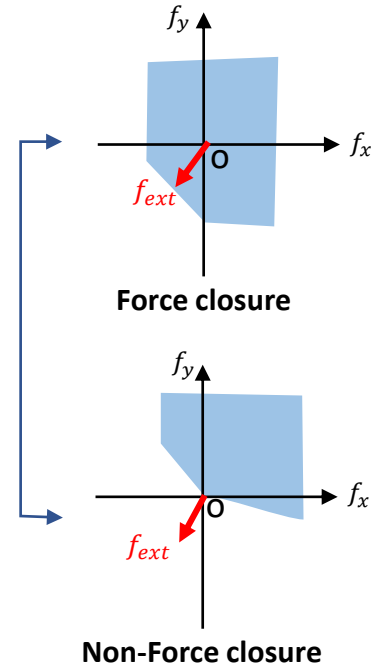
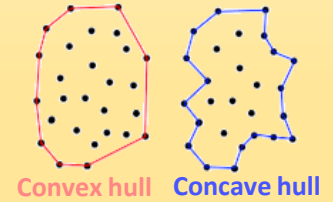
There are more grasp quality metrics.

Wrench?

Force and torque vectors

$$[f_x, f_y, f_z, \tau_x, \tau_y, \tau_z] \in \mathbb{R}^6$$

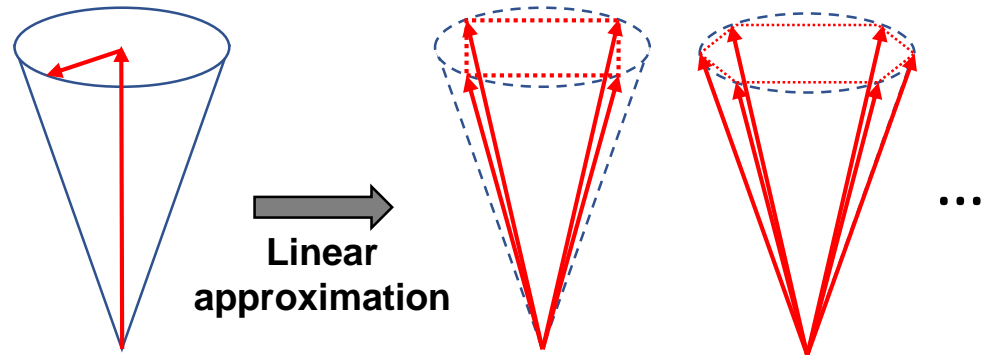
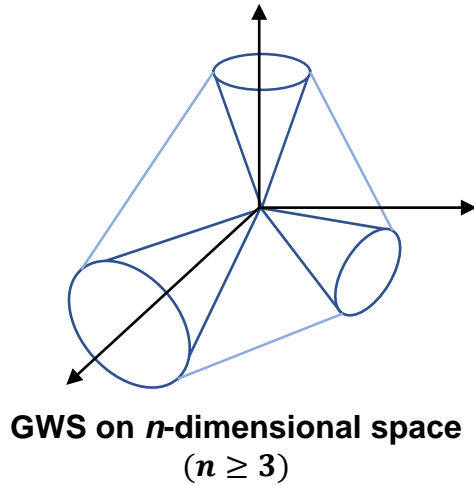
Convex Hull?



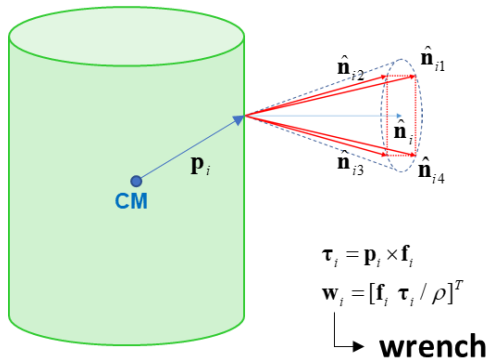
Grasp Quality

□ GWS on 3D object

✓ Approximation of friction cone



✓ GWS in 3D object



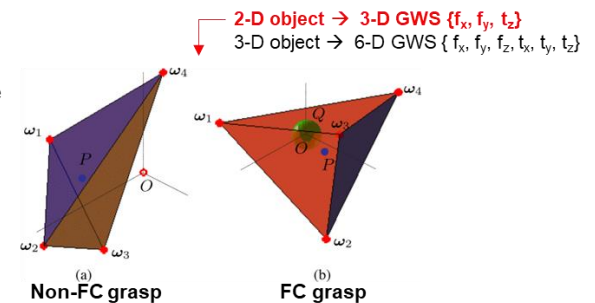
Wrenches applied at p_i

- 1) w_i : wrench generated by a unitary force f_i orthogonal to the object surface
 - 2) w_{ij} : wrench generated by a unitary force f_i along an edge of the linearized friction cone
- Primitive wrench

$$C = \{p_1, \dots, p_n\}$$

$$G = \{w_1, \dots, w_n\}$$

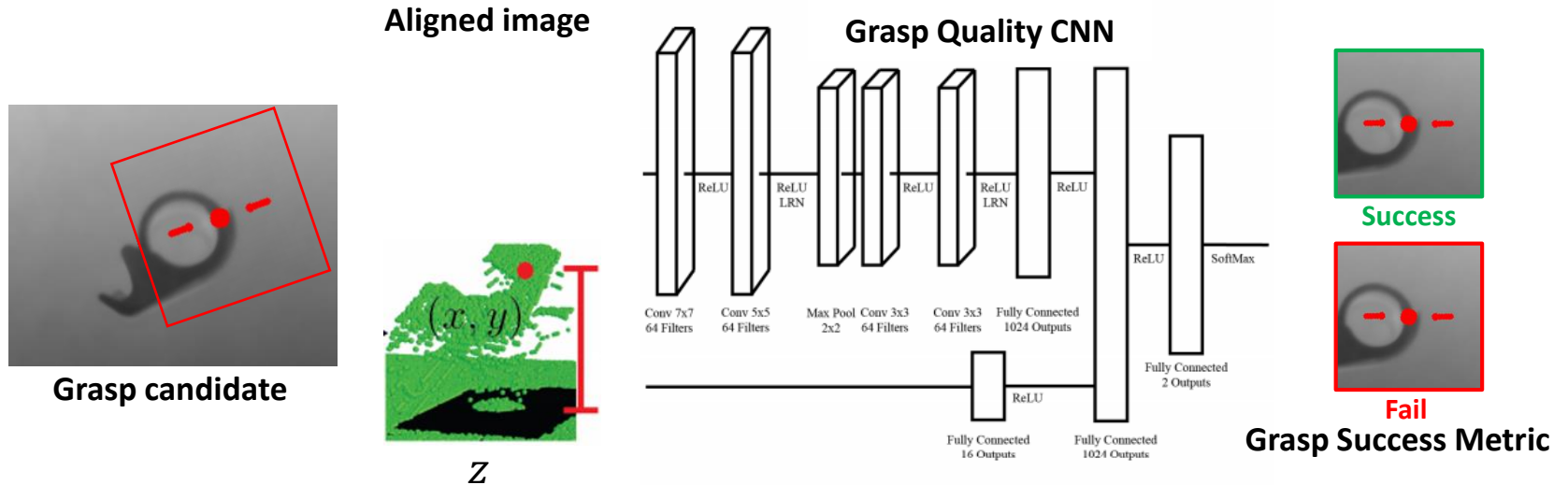
$$W = \{w_{11}, \dots, w_{1m}, \dots, w_{n1}, \dots, w_{nm}\}$$



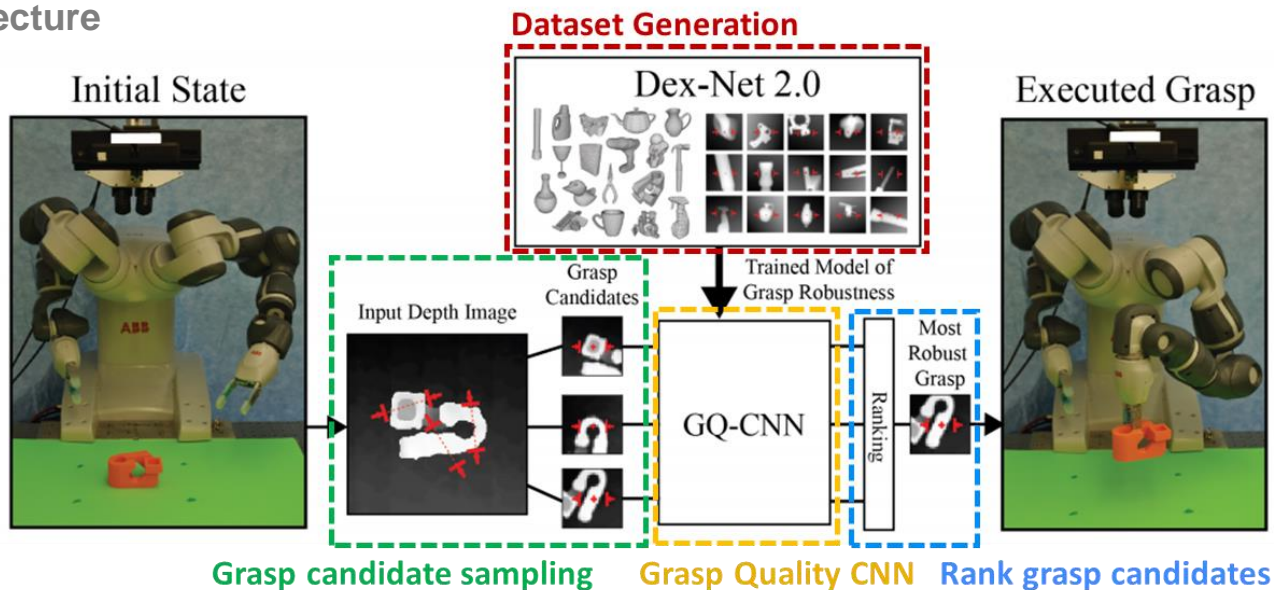
Dex-Net 2.0

Dex-Net 2.0

Grasp Quality Convolutional Neural Network(GQ-CNN)

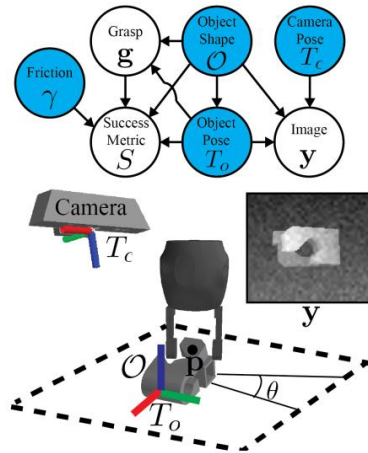


Architecture



Dex-Net 2.0

Problem Statements



Graphical model for robust parallel-jaw grasping of objects

States

$$\mathbf{x} = (O, T_o, T_c, \gamma)$$

obj info, obj pose, cam pose, friction coef

Point cloud

$$\mathbf{y} = \mathbb{R}_+^{H \times W}$$

Grasps

$$\mathbf{u} = (\mathbf{p}, \varphi) \in \mathbb{R}^3 \times S^1$$

Grasp metric
(robust analytic)

$$S(\mathbf{u}, \mathbf{x}) \in \{0, 1\}$$

state distribution

observation model

grasp candidate model

analytic model of grasp success

$$p(S, \mathbf{u}, \mathbf{x}, \mathbf{y}) = p(\mathbf{x}) p(\mathbf{y} | \mathbf{x}) p(\mathbf{u} | \mathbf{x}) p(S | \mathbf{u}, \mathbf{x})$$

Let $p(S, \mathbf{u}, \mathbf{x}, \mathbf{y})$ be a joint distribution on ... imprecision in sensing and control.

Let the *robustness* of a grasp given an observation: $Q(\mathbf{u}, \mathbf{y}) = E(S | \mathbf{u}, \mathbf{y})$

1) Goal is to find(learn) a **robustness function** : $Q_{\theta^*}(\mathbf{u}, \mathbf{y}) \in \{0, 1\}$

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} E_{p(S, \mathbf{u}, \mathbf{x}, \mathbf{y})} [\mathcal{L}(S, Q_{\theta}(\mathbf{u}, \mathbf{y}))]$$

NN parameters Cross entropy loss

2) Using in grasp planning

$$\pi_{\theta}(\mathbf{y}) = \underset{\mathbf{u} \in \mathcal{C}}{\operatorname{argmax}} Q_{\theta}(\mathbf{u}, \mathbf{y})$$

Grasp candidates

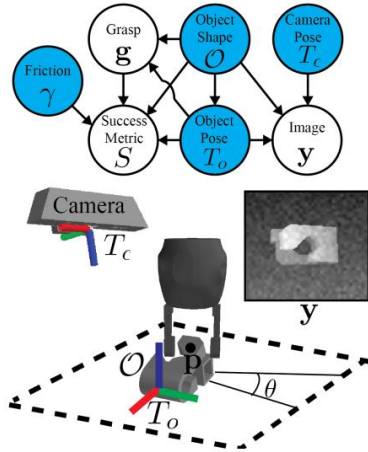
Cross Entropy

$$H(p, q) = - \sum_{x \in X} p(x) \log q(x)$$

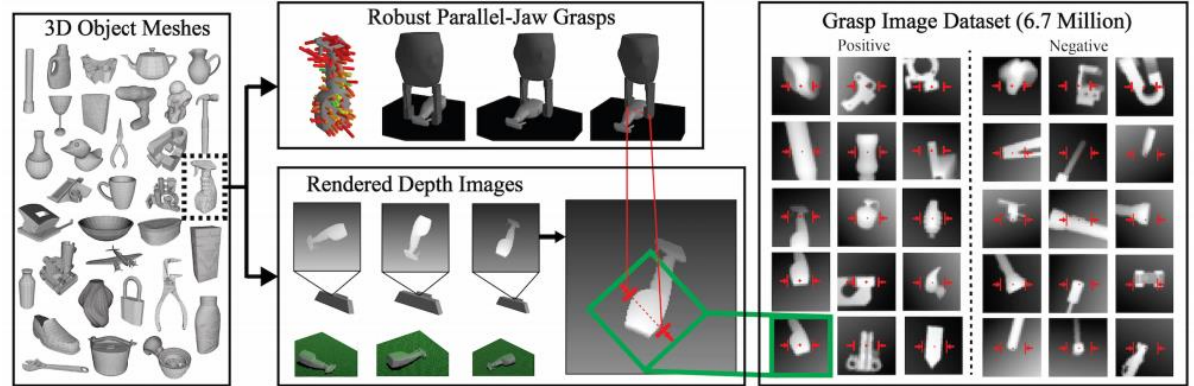
$p(x)$: true distribution

$q(x)$: estimate distribution

Dex-Net 2.0 – Dataset Generation



Graphical model for robust parallel-jaw grasping of objects



Dex-Net 2.0 pipeline for training dataset generation

Model $(S_1, \mathbf{u}_1, \mathbf{x}_1, \mathbf{y}_1), \dots, (S_N, \mathbf{u}_N, \mathbf{x}_N, \mathbf{y}_N) \sim p(S, \mathbf{u}, \mathbf{x}, \mathbf{y}) = p(\mathbf{x}) p(\mathbf{y} | \mathbf{x}) p(\mathbf{u} | \mathbf{x}) p(S | \mathbf{u}, \mathbf{x})$

i.i.d samples

grasp candidate model, analytic model of grasp success

1) state distribution $p(\mathbf{x})$

$$p(\mathbf{x}) = p(\gamma) p(\mathcal{O}) p(T_o | \mathcal{O}) p(T_c)$$

Distribution	Description
$p(\gamma)$	truncated Gaussian distribution over friction coefficients
$p(\mathcal{O})$	discrete uniform distribution over 3D object models
$p(T_o \mathcal{O})$	continuous uniform distribution over the discrete set of object stable poses and planar poses on the table surface
$p(T_c)$	continuous uniform distribution over spherical coordinates for radial bounds $[r_\ell, r_u]$ and polar angle in $[0, \delta]$

$$p(\gamma) \sim \mathcal{N}(0.5, 0.1) \text{ truncated to } [0, 1]$$

$$p(\mathcal{O}) \sim \mathcal{U}(\text{given 3D obj dataset})$$

$$p(T_o | \mathcal{O}) = p(T_o | T_s) p(T_s | \mathcal{O})$$

$$p(T_s | \mathcal{O}) \sim \mathcal{U}(\text{stable poses})$$

$$p(T_o | T_s) \sim \mathcal{U}([-0.1, 0.1] \times [-0.1, 0.1] \times [0, 2\pi])$$

$$p(T_c) \sim \mathcal{U}([0.65, 0.75] \times [0, 2\pi] \times [0.05\pi, 0.1\pi]) \quad (\square \ r, \theta, \varphi)$$

2) grasp candidate model $p(\mathbf{u} | \mathbf{x})$

Uniform distribution over 1) & 2) & 3)

- 1) Pairs of antipodal contact points on obj surface
- 2) Grasp axis \perp table plane
- 3) Reject no FC nor no parallel ($\mu=0.6$)

3) observation model $p(\mathbf{y} | \mathbf{x})$

$$\mathbf{y} = \alpha \hat{\mathbf{y}} + \varepsilon$$

$\hat{\mathbf{y}}$: depth image created using OSMesa offscreen rendering

α : Gamma random variable with shape=1000.0 and scale=0.001

ε : Gaussian Process noise over pixel coordinates

with measurement noise $\sigma=0.005$ and kernel bandwidth $\sqrt{2} p_x$

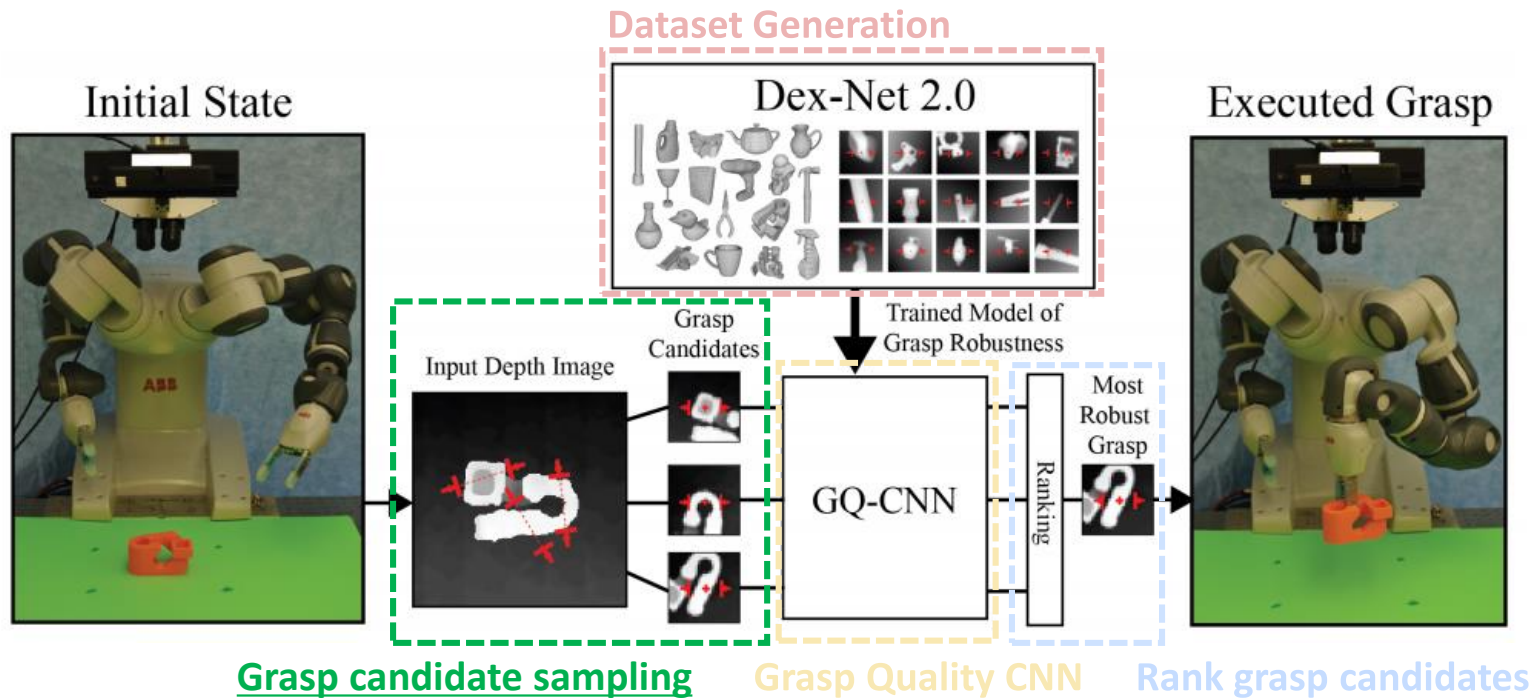
4) Analytic model of grasp success $p(S | \mathbf{u}, \mathbf{x})$

Robust epsilon quality

$$S(\mathbf{u}, \mathbf{x}) = \begin{cases} 1 & E_q > \delta \text{ and } \text{collfree}(\mathbf{u}, \mathbf{x}) \\ 0 & \text{otherwise} \end{cases}$$

E_q : robust epsilon quality

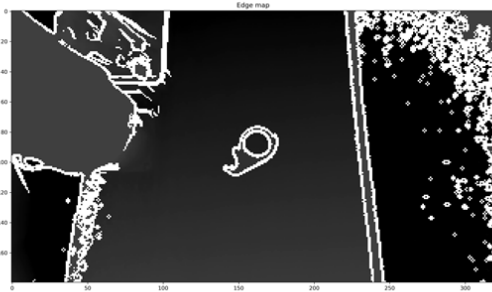
Dex-Net 2.0 – Grasp Candidate Sampling



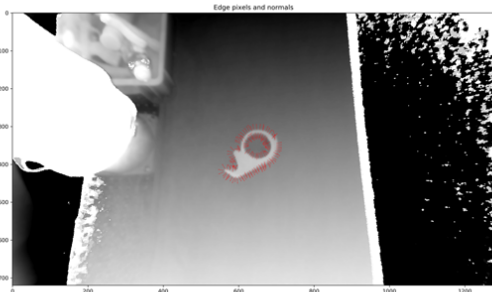
Dex-Net 2.0 – Grasp Candidate Sampling

Image based parallel-jaw grasp candidate sampling

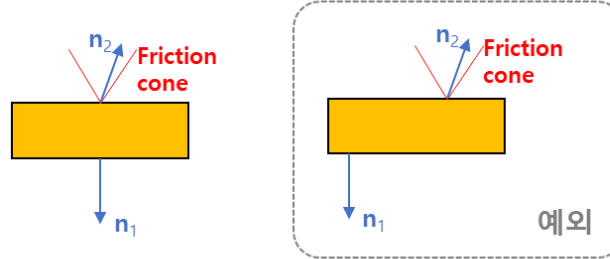
1. Get edge(gradient threshold)



2. Get normal vectors at edge pixel(gradient)



3-1. Get antipodal points(F.C)

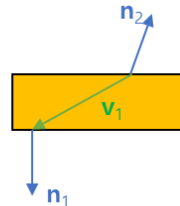


<conditions>

contact points dist(px) < grasp width(px)

$$\mathbf{n}_1 \cdot \mathbf{n}_2 < -\cos(\text{atan}(\mu))$$

3-2. Get antipodal points(F.C)



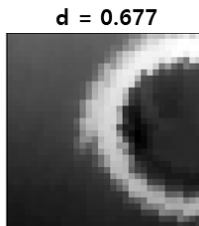
<conditions>

$$\angle(\mathbf{n}_1, \mathbf{v}) \leq \text{atan}(\mu)$$

$$\angle(\mathbf{n}_2, -\mathbf{v}) \leq \text{atan}(\mu)$$

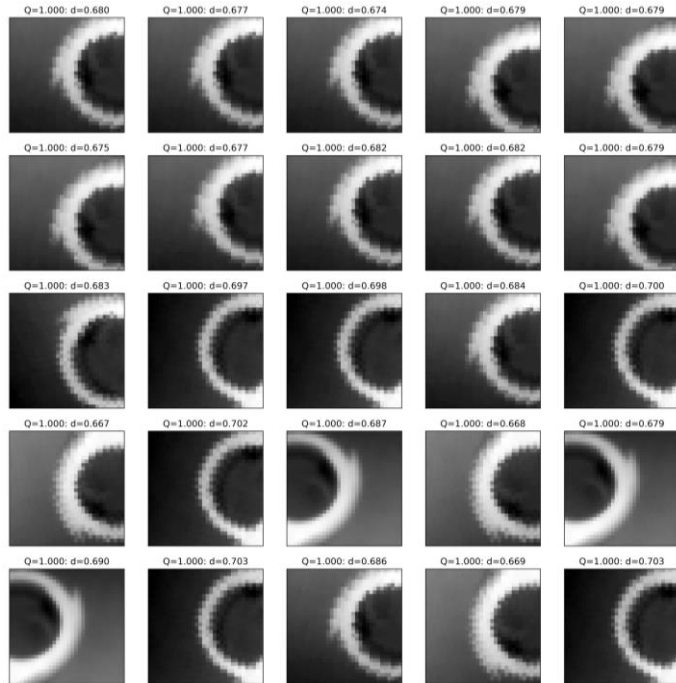
4. contact1, contact2 pixel → center pixel, theta → crop image

5. depth: depth at center pixel + offset(0.015~0.05)

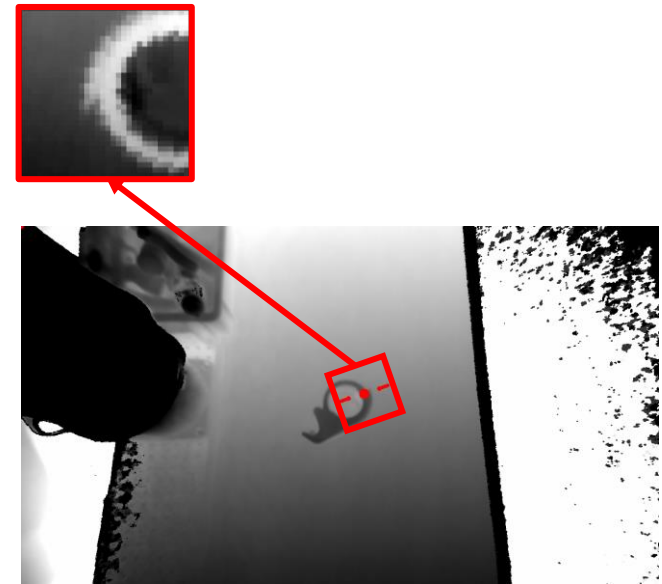


Crop image

Dex-Net 2.0 – Rank grasp candidates

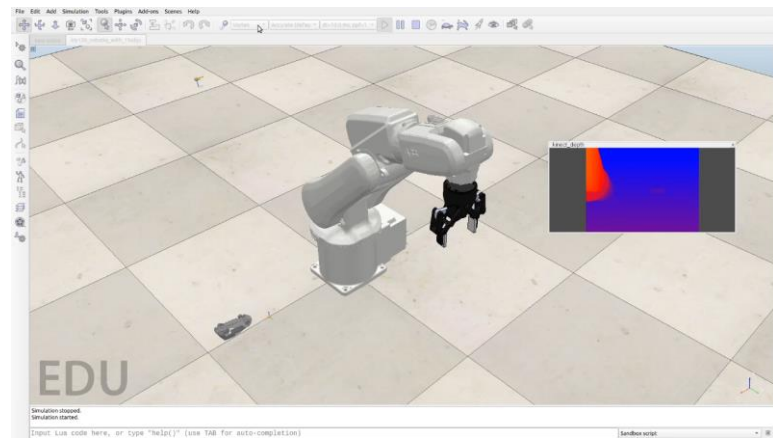
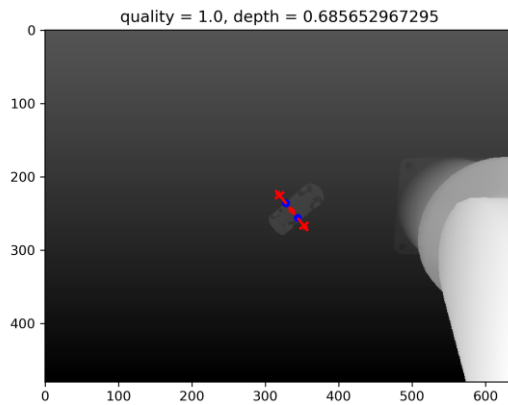
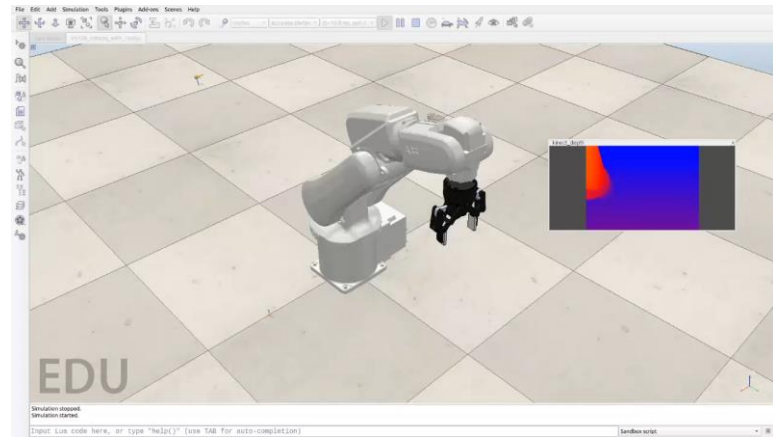
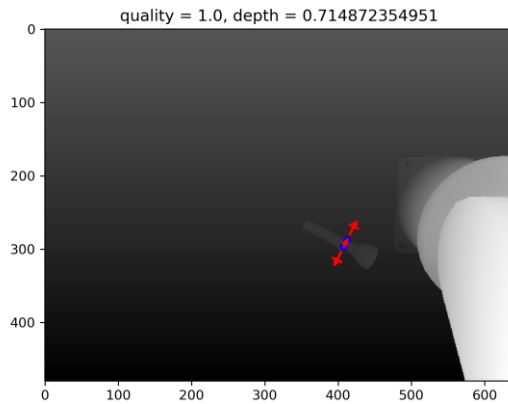


Grasp candidates



Top rank grasp

Dex-Net 2.0 – Simulation



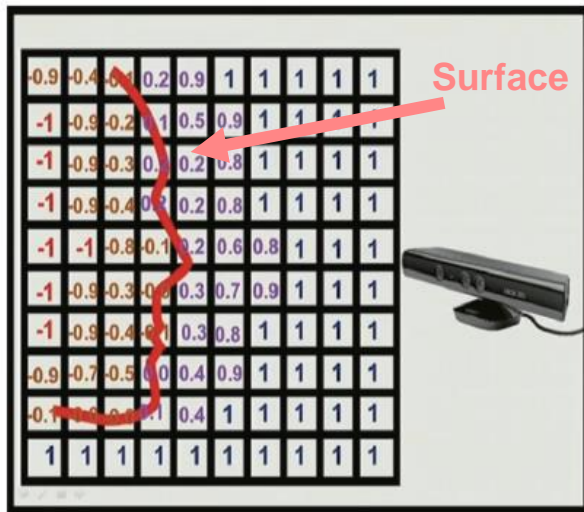
Simulation results

VGN

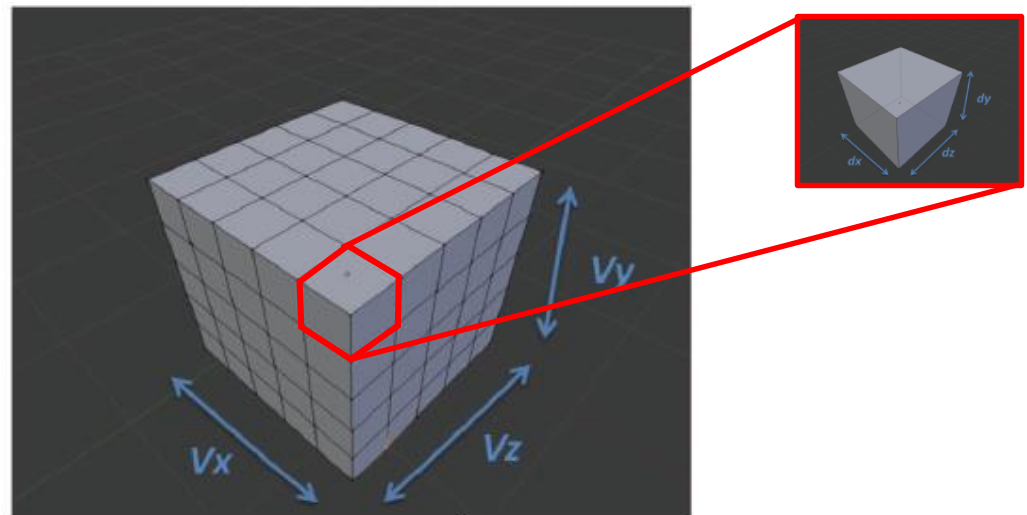
VGN – What is TSDF?

❑ Truncation Signed Distance Function (TSDF)

- ❑ The truncated distance value formed when the camera observes the surface of object



TSDF on 2D pixel grid

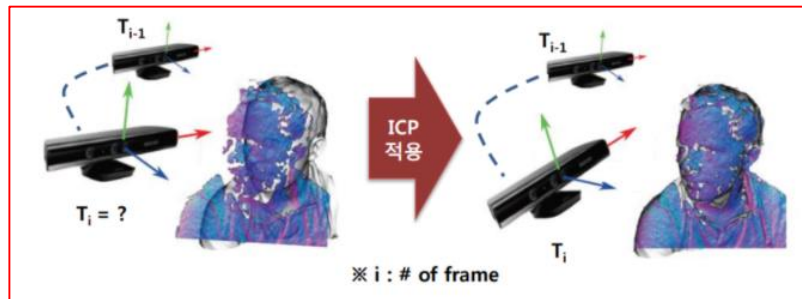
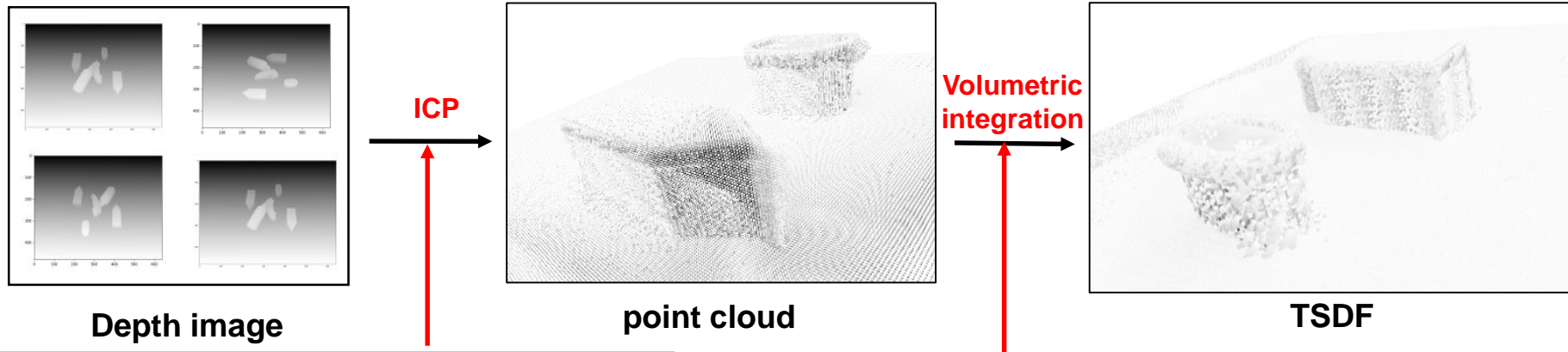


Voxel grid

$$TSDF = \begin{cases} (0, 1] & \text{(Outside of the surface)} \\ 0 & \text{(On the surface)} \\ [-1, 0) & \text{(Inside the surface)} \end{cases}$$

VGN – What is TSDF?

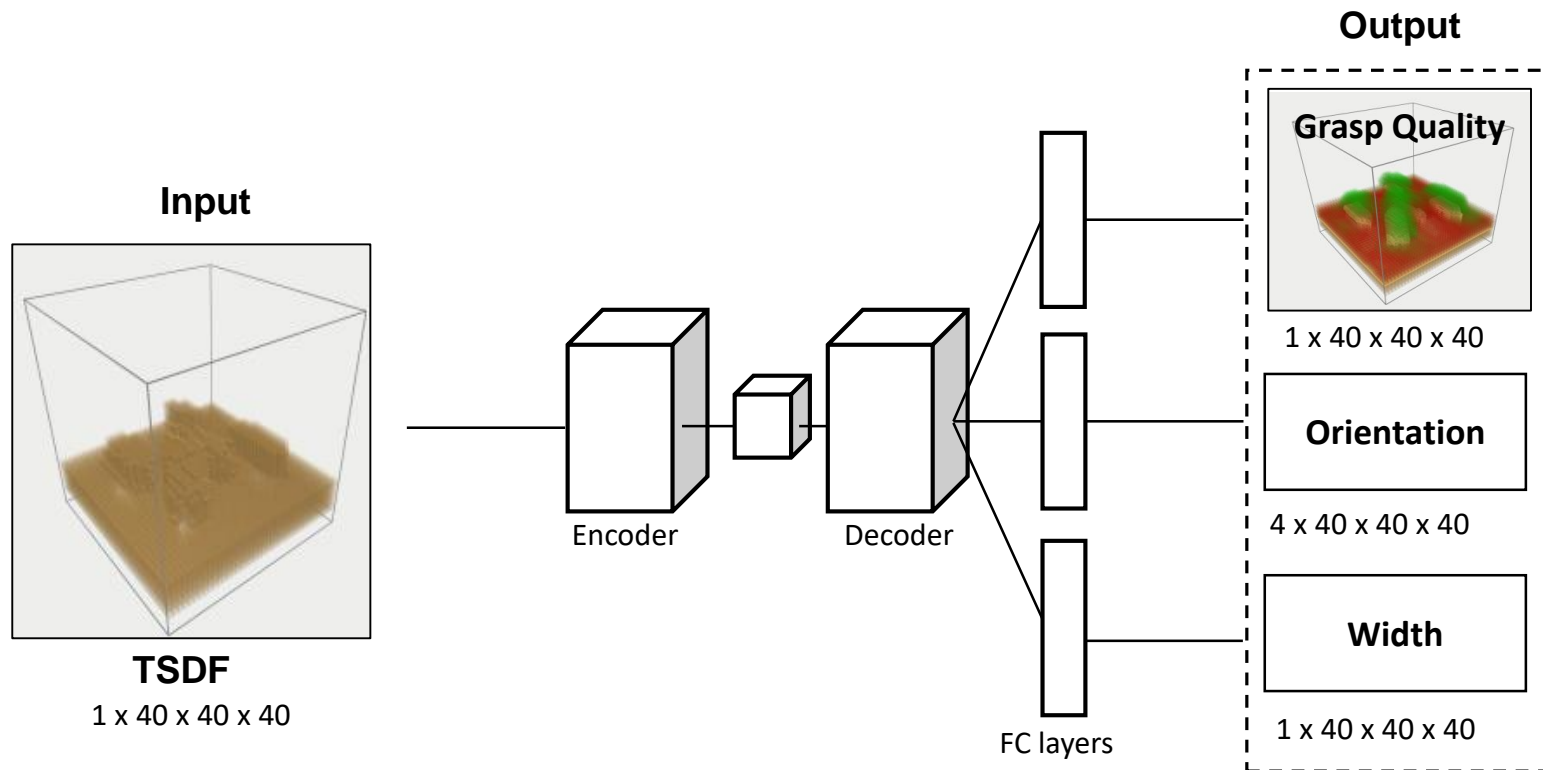
❑ 3D reconstruction with TSDF volume [1]



- Each Voxel of TSDF include distance function
- Reduces the noise of the sensor

VGN

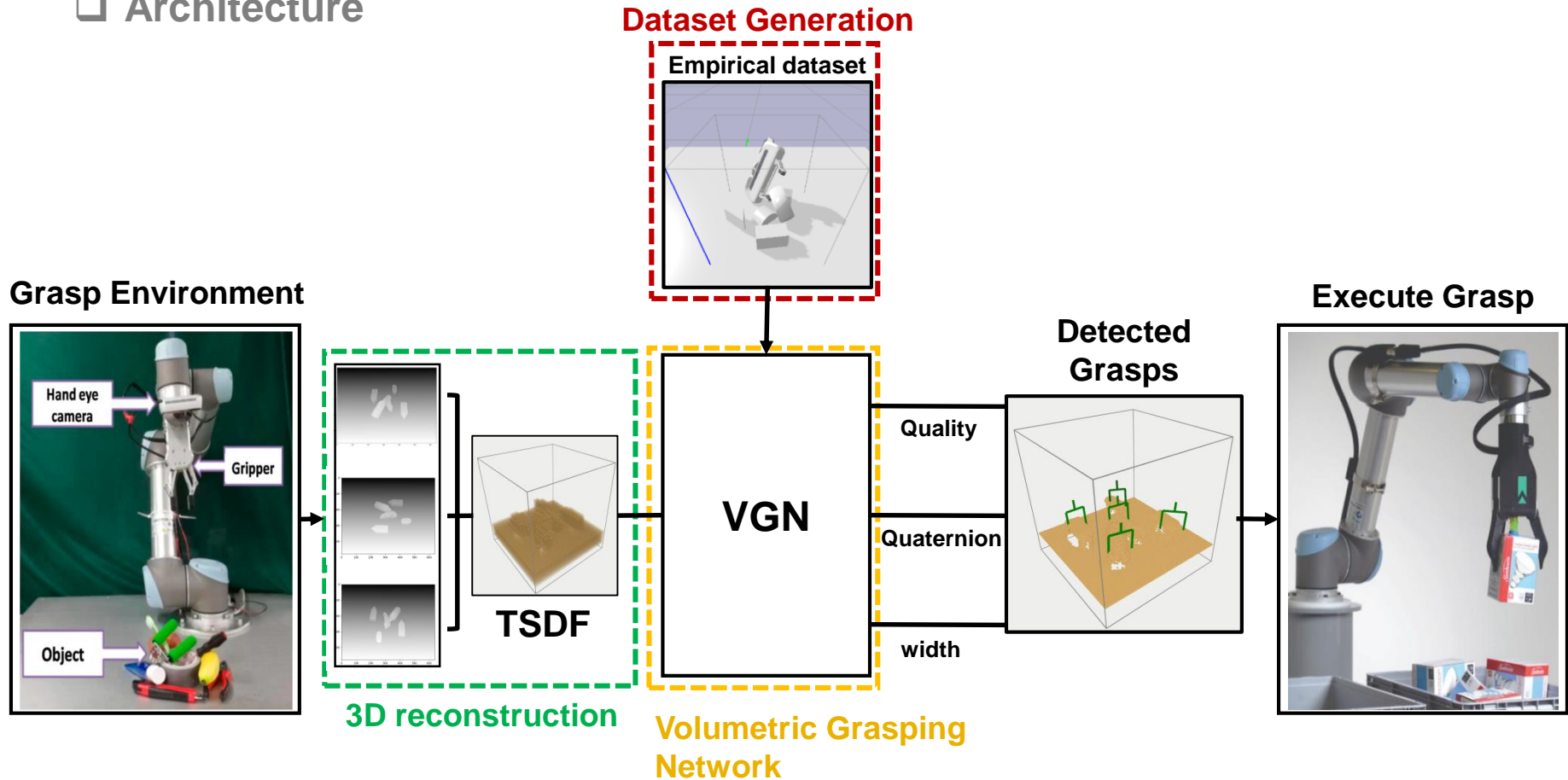
□ Volumetric Grasping Network (VGN)



Volumetric Grasping Network

VGN

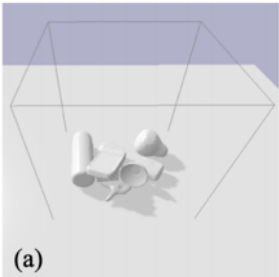
Architecture



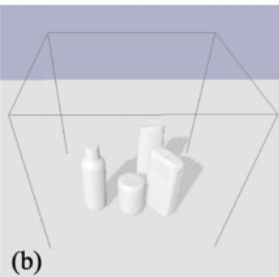
Dataset Generation

- ☐ Empirical methods : Physical Trials using Pybullet simulator
- ☐ Active Search : Spherical coordinate system (r, θ, φ)
: $N \sim \mathcal{U}(1,6)$, $r \sim \mathcal{U}(0.48, 0.72)$, $\theta \sim \mathcal{U}(0, \frac{\pi}{4})$, $\varphi \sim \mathcal{U}(0, 2\pi)$
- ☐ Random point Grasp

1. Select scenes

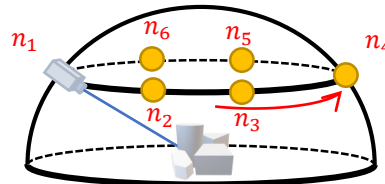


- Pile scenes : 4DOF

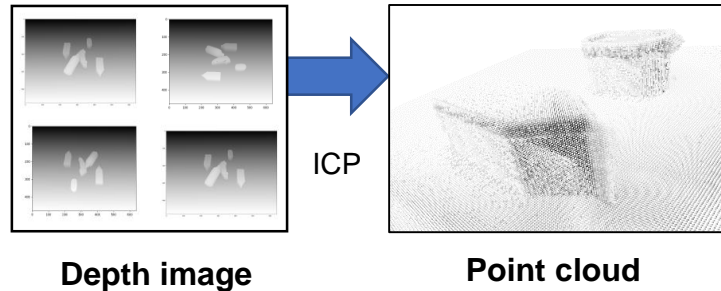


- Packed scenes : 6DOF

2. Object recognition

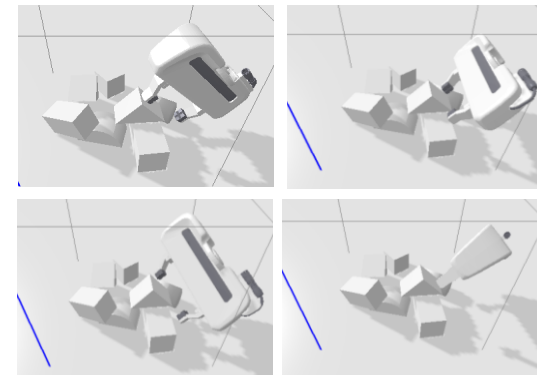


Active search



3. Random point grasp

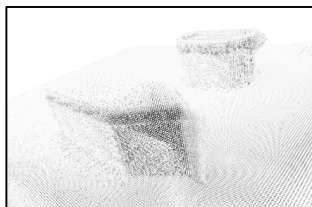
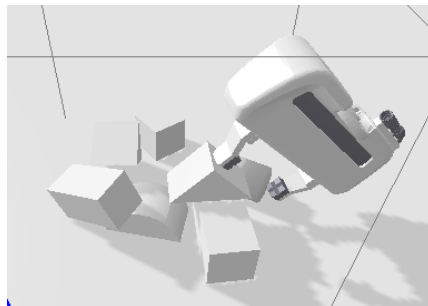
- Z axis rotation grasp
(Yaw $\sim \mathcal{U}(0, \pi)$)



Dataset Generation

- Transform from point cloud to TSDF
- Data labelling

4. Transform from point cloud to TSDF

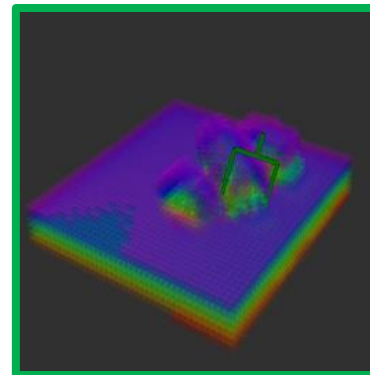


Point cloud

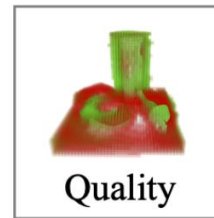
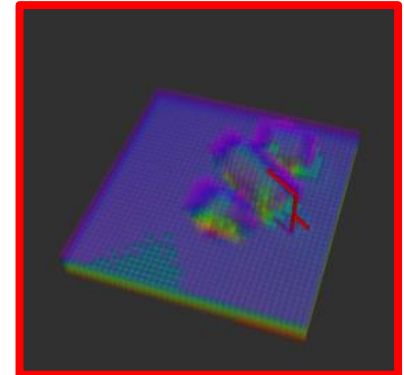


TSDF transform
& Result Labelling

Success



Fail



1×40×40×40

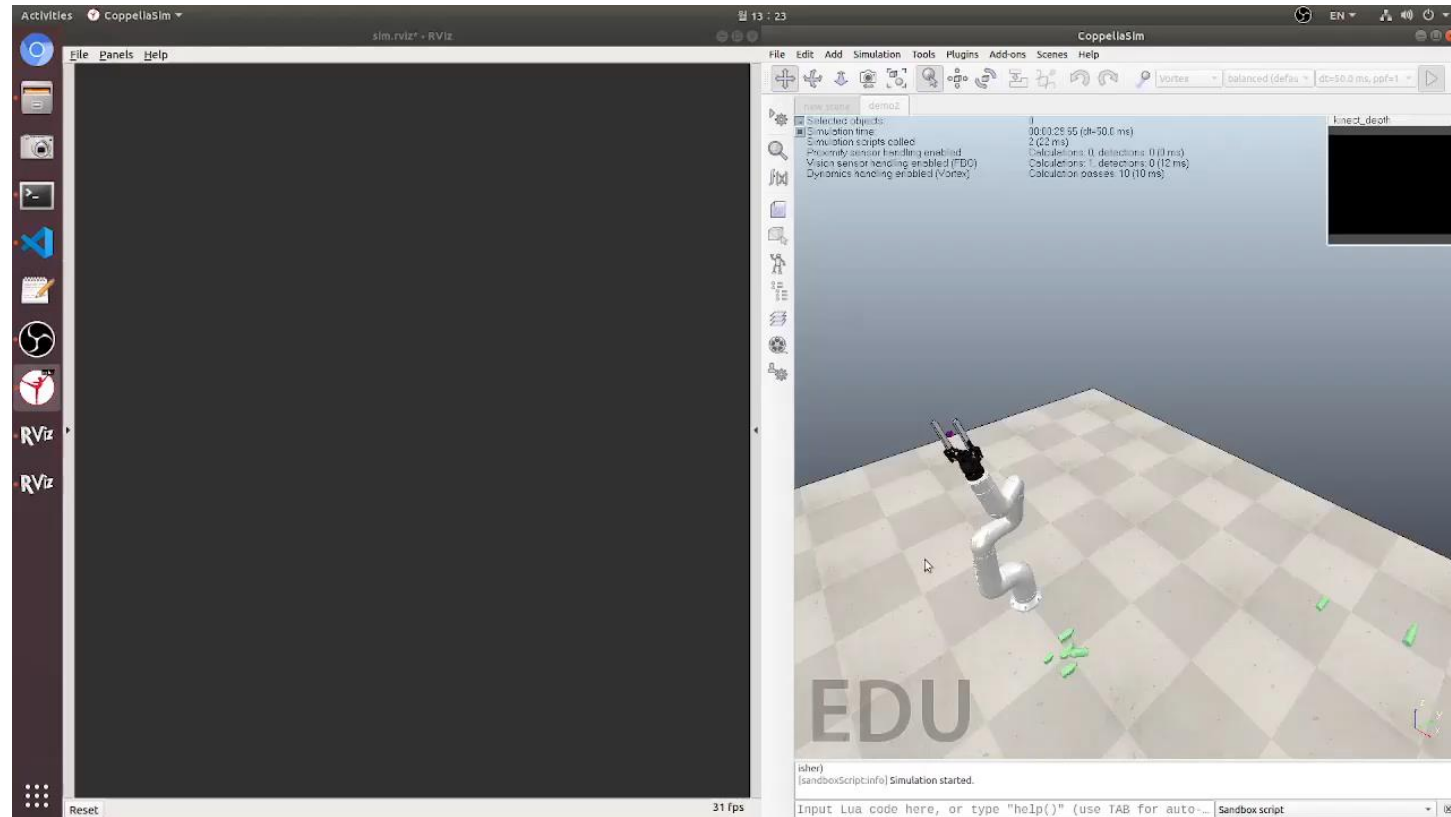
Orientation

4×40×40×40

Width

1×40×40×40

Simulation

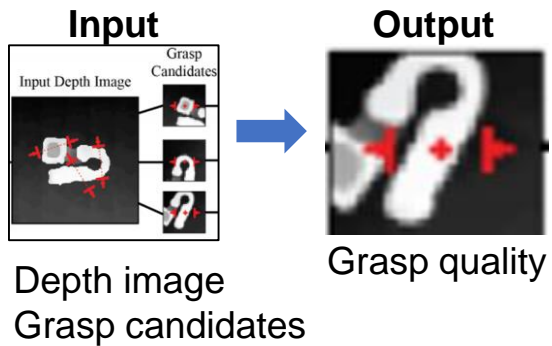


Simulation results

Conclusion

Dex-net 2.0 vs VGN Compare

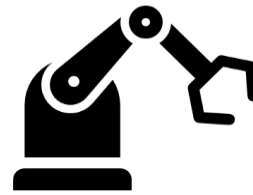
Dex-Net 2.0



Analytic method

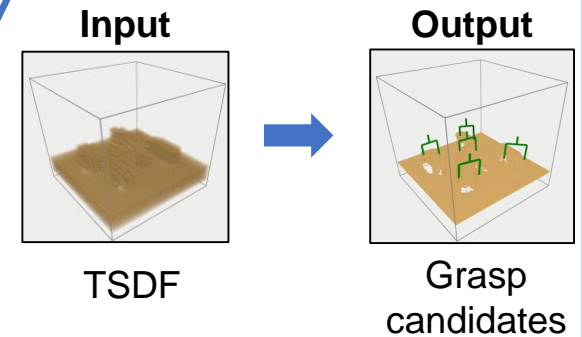
4DOF Grasp

Discriminative



Grasp

VGN



Empirical method

6DOF Grasp

Generative

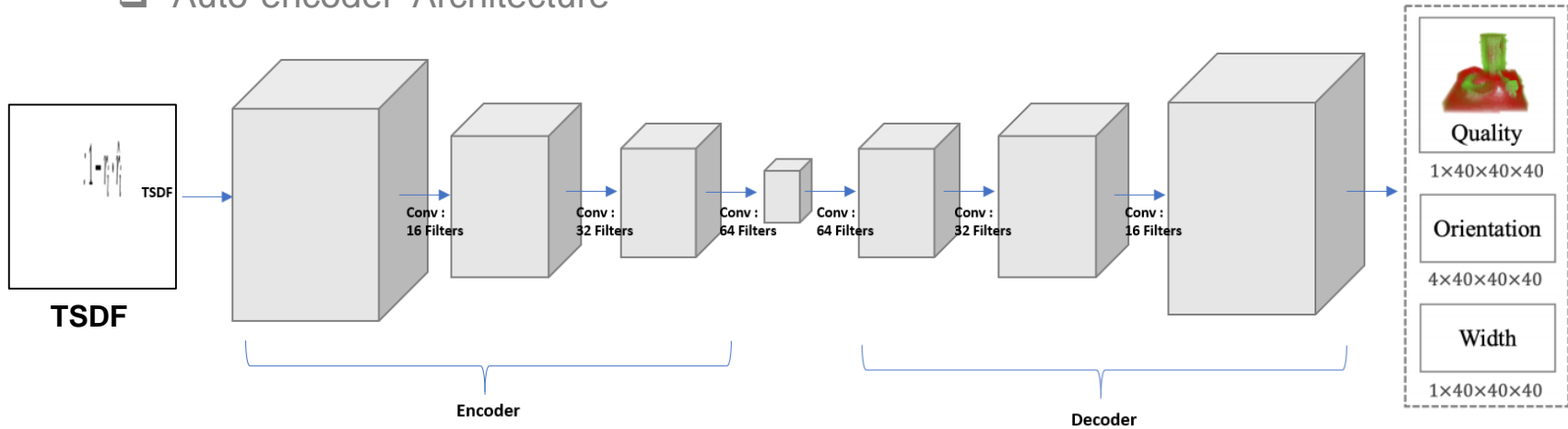
Thank you for your attention!

Appendix

VGN – 2) Volumetric Grasping Network

VGN- NET

Auto-encoder Architecture



Loss Function

$$\text{loss function} : \mathcal{L}(g_i, \hat{g}_i) = \mathcal{L}(q_i, \hat{q}_i) + q_i(\mathcal{L}(r_i, \hat{r}_i) + \mathcal{L}(w_i, \hat{w}_i))$$

voxel Grasp quality Grasp quaternions Grasp width

$$\begin{cases} \mathcal{L}(q_i, \hat{q}_i) & : \text{cross entropy loss function} \\ \mathcal{L}(r_i, \hat{r}_i) & : 1 - r_i \cdot \hat{r}_i \\ \mathcal{L}(w_i, \hat{w}_i) & : \text{MSE} \end{cases}$$

g_i : voxel

q_i : grasp label $\in \{0,1\}$

w_i : grasp quaternions

r_i : grasp width

\hat{q}_i : Ground truth grasp label

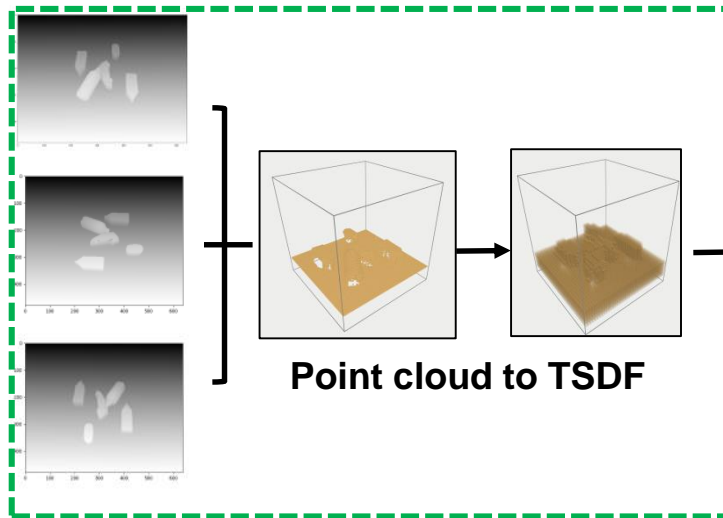
\hat{w}_i : Ground truth grasp quaternions

\hat{r}_i : Ground truth grasp width

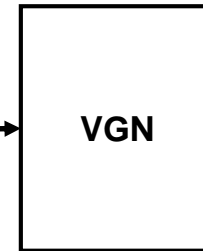
VGN – 3) Grasp Planning



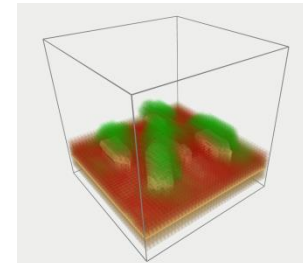
Simulation Environment



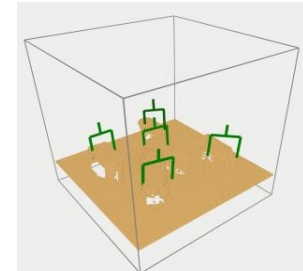
3) 3D reconstruction



4) Detected grasp



Grasp Quality

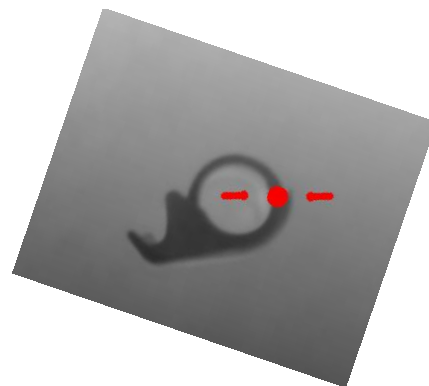
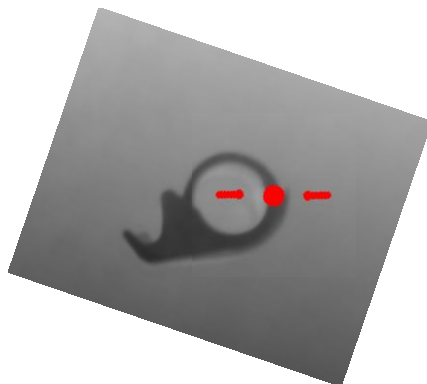
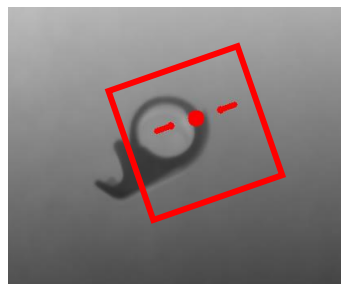


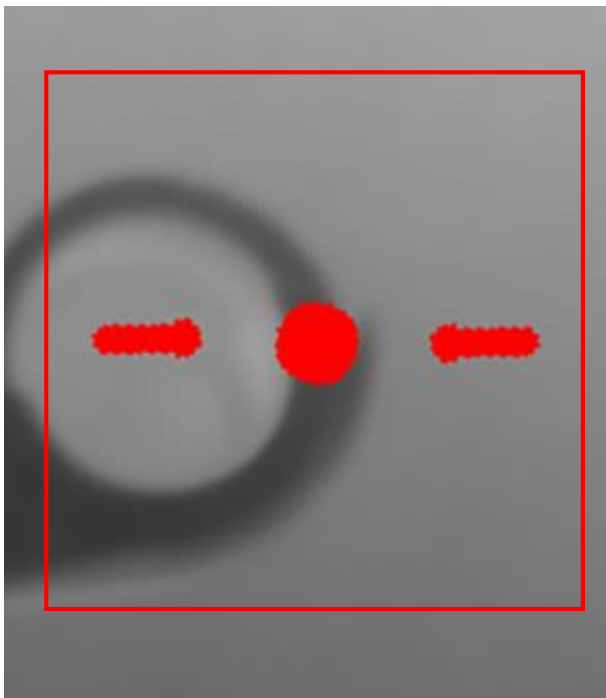
Select Grasp Candidate

$$f : V \rightarrow Q, R, W \xrightarrow{\text{red arrow}} t' = \frac{T_{RV}(t)}{v} \quad r' = T_{RV}(r) \quad w' = \frac{w}{v}$$

v : voxel size

T_{RV} : Transformation matrix between base, TSDF frame





Dex-net 2.0

Architecture

