# Grasping Pattern Estimation Based on Co-occurrence of Object and Hand Shape

**Takuya Kawakami**, Tadashi Matsuo, Yoko Ogawa, Nobutaka Shimada,
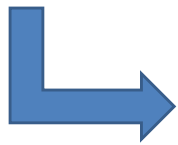
Ritsumeikan University

# Introduction

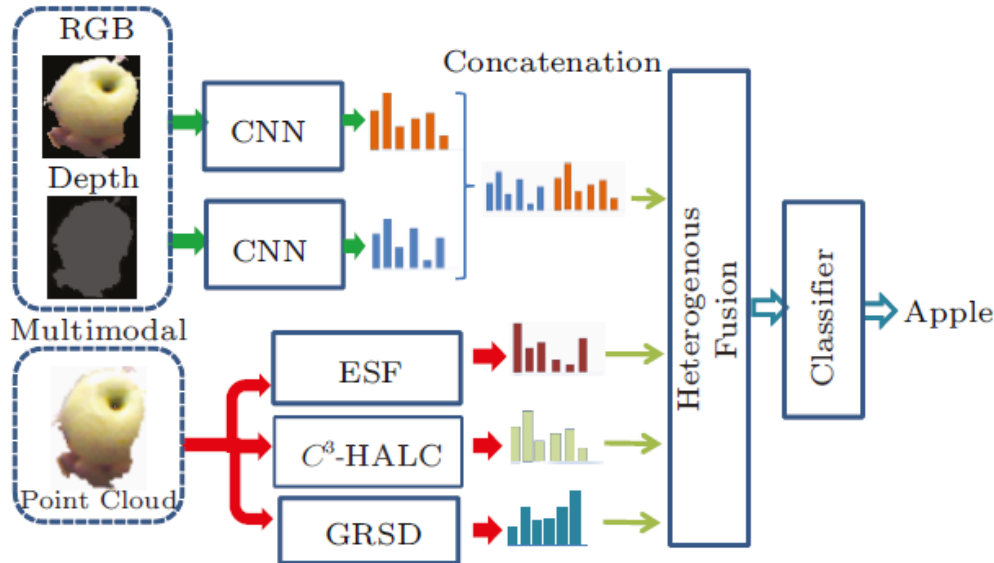- An object as a tool has its own function. The function is closely related to how a human grasp it [1].



**Can we estimate how to grasp an object from the object itself?**

It will be useful for object recognition and robot manipulation.

[1] N. Kamakura, "Shape of hand and Hand motion". Ishiyaku Publishers, 1989.

# Related work

- Xiong Lv et al., "RGB-D Hand-Held Object Recognition Based on Heterogeneous Feature Fusion",
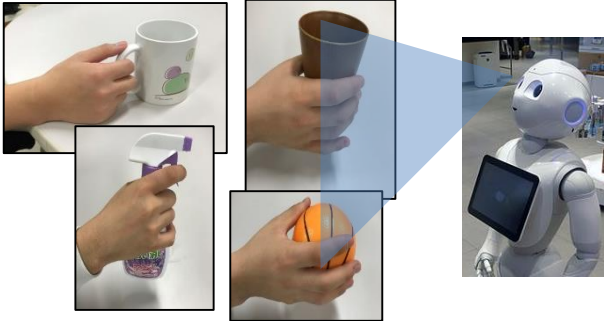  Journal of Computer Science and Technology(2015)



They achieved highly accurate classification by utilizing how to grasp an object, but...

- It estimates only an object label (not how to grasp it).
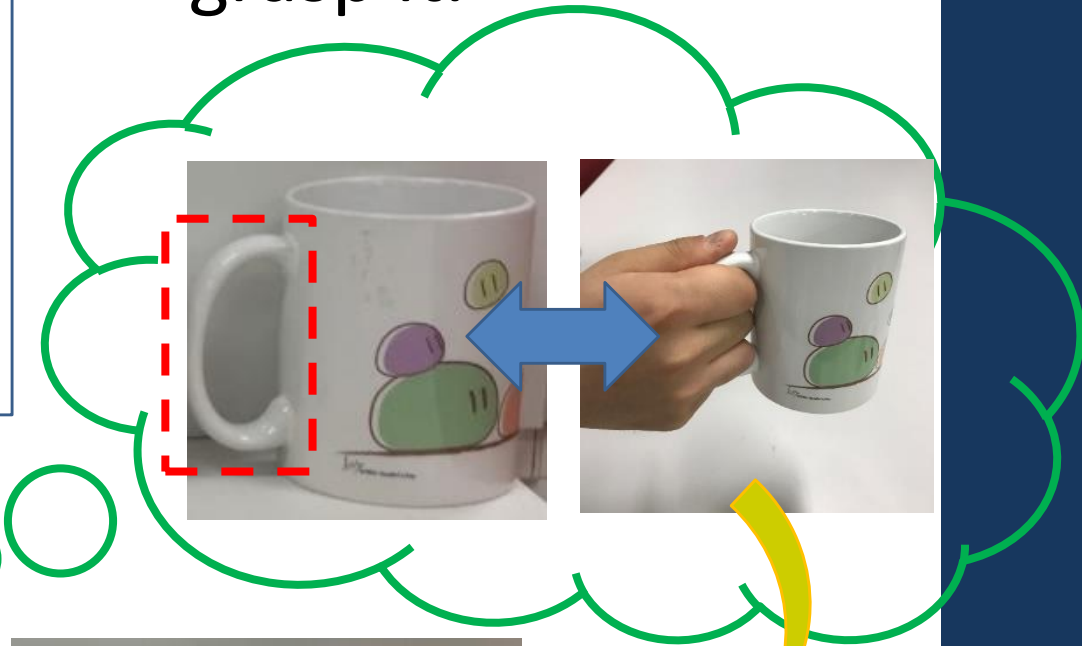- All teacher labels must be given manually.

# Our goal

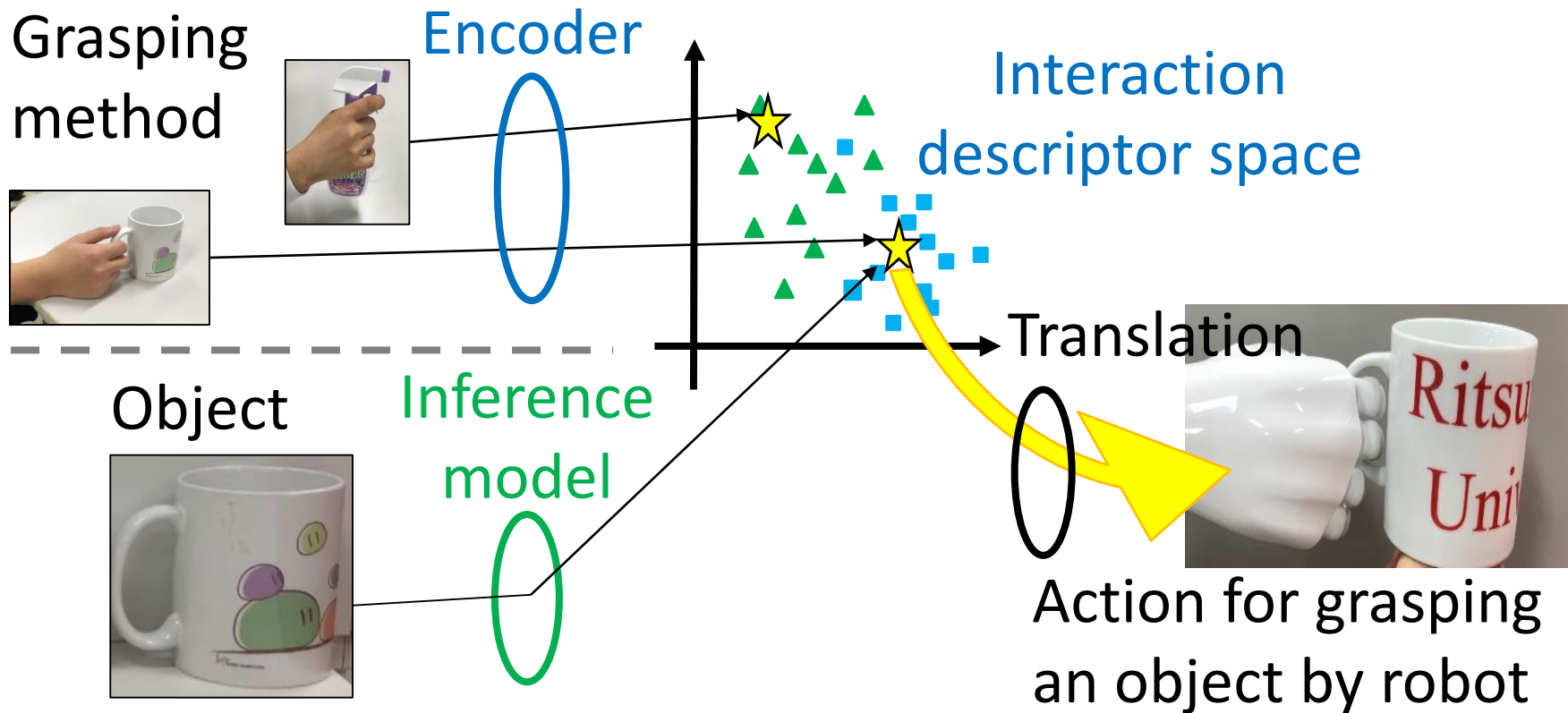Training

Learn human interactions without teacher labels.

Recall how a human grasp it.

Make action to grasp it.

# Proposed method

- We generate an interaction descriptor, a numeral representation of a human grasping method.

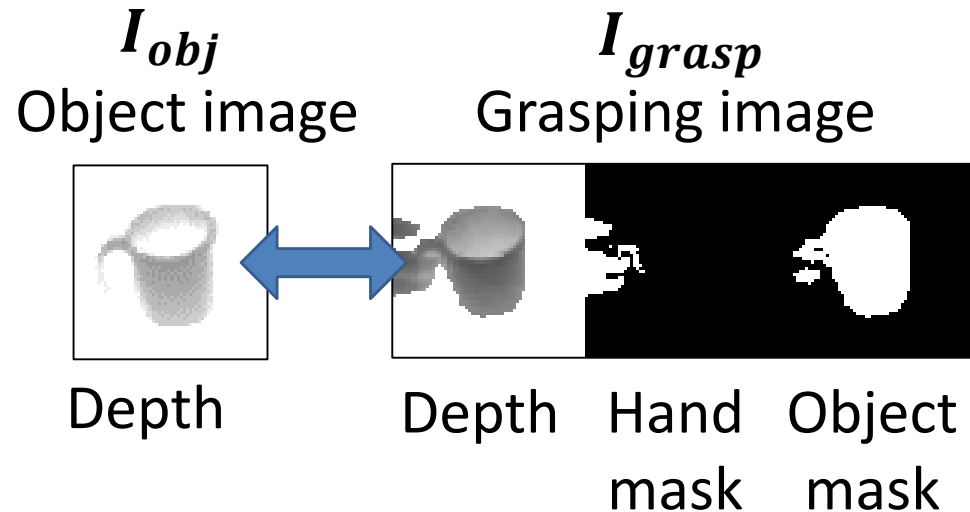- And then we make an inference model to learn the relation between object and grasping method.

Grasping method

Encoder

Interaction descriptor space

Translation

Object

Inference model

Action for grasping an object by robot

# Grasping image

Observing human grasping

Automatically collect Images for learning

$I_{obj}$
Object image

$I_{grasp}$
Grasping image

Depth

Depth Hand mask Object mask

Grasping method is represented as a grasping image. It consists of a depth image, hand mask and object mask.

It is paired with the corresponding object image.

# Capture of human's grasping scene
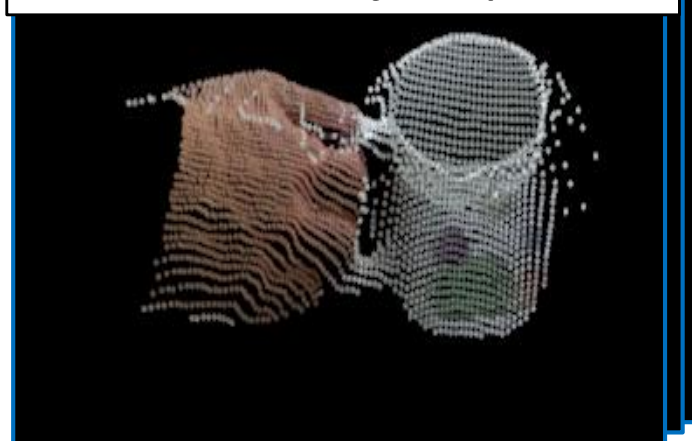
Observing human's grasping scene

RGB-D sensor

Captured point cloud
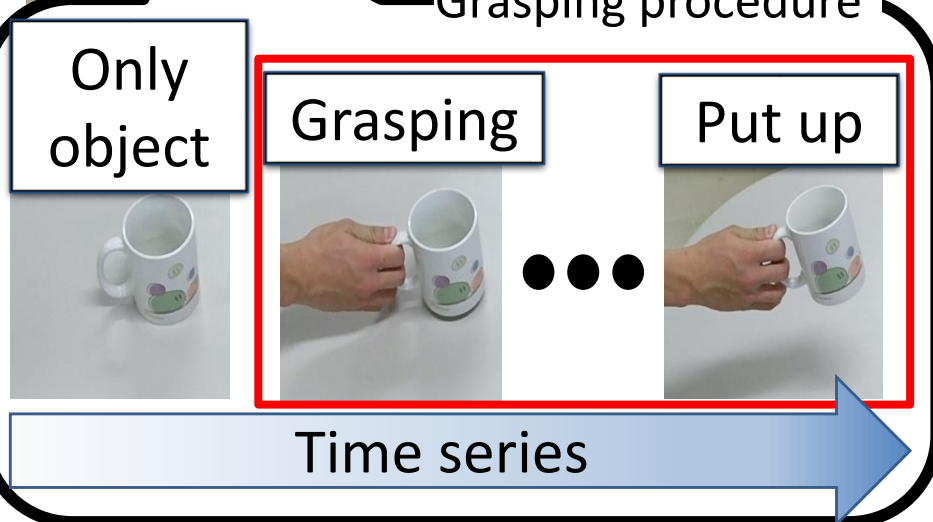
Remove unnecessary points

Hand and object points

Grasping procedure

Only object

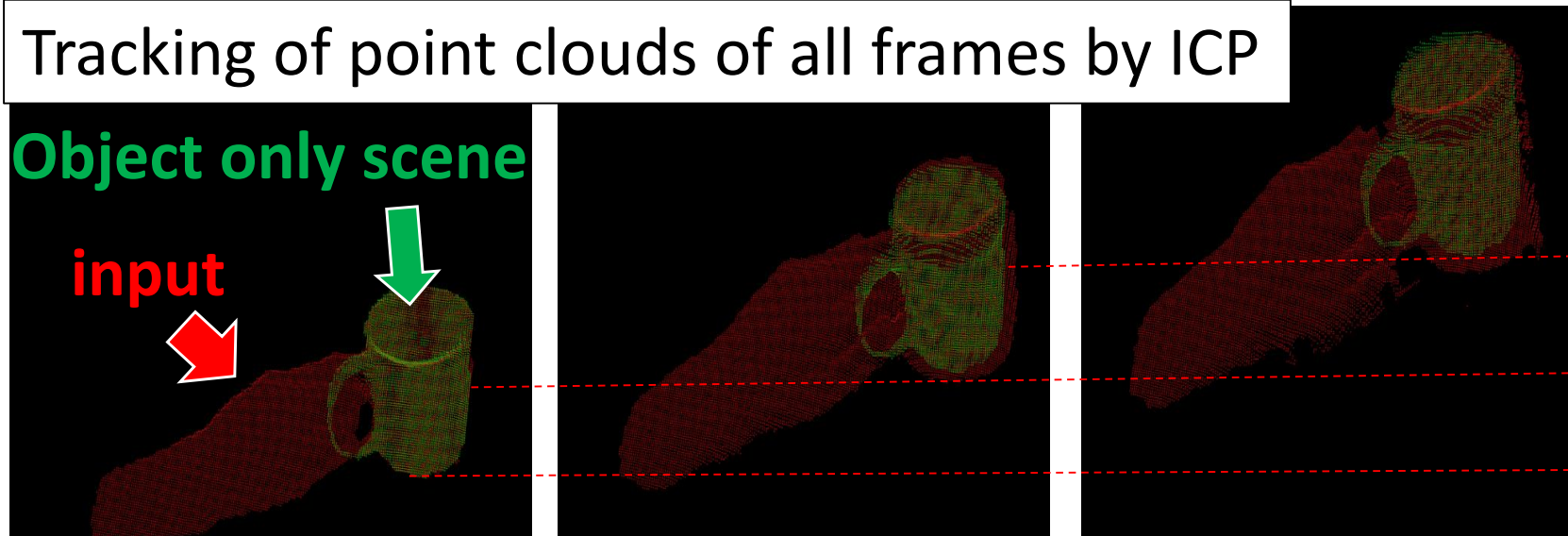Grasping

Put up

Time series

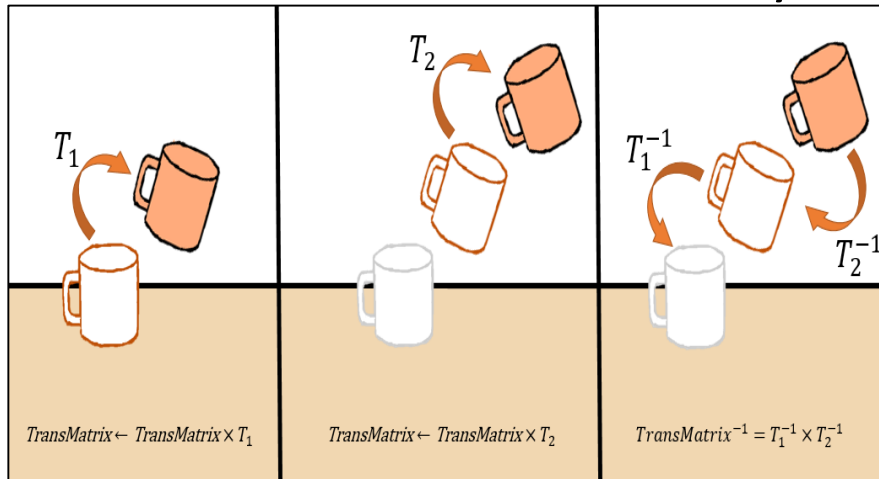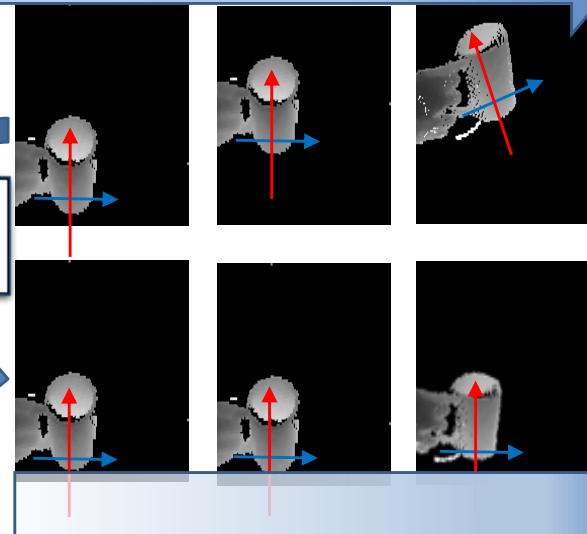# Alignment based on object

Tracking of point clouds of all frames by ICP



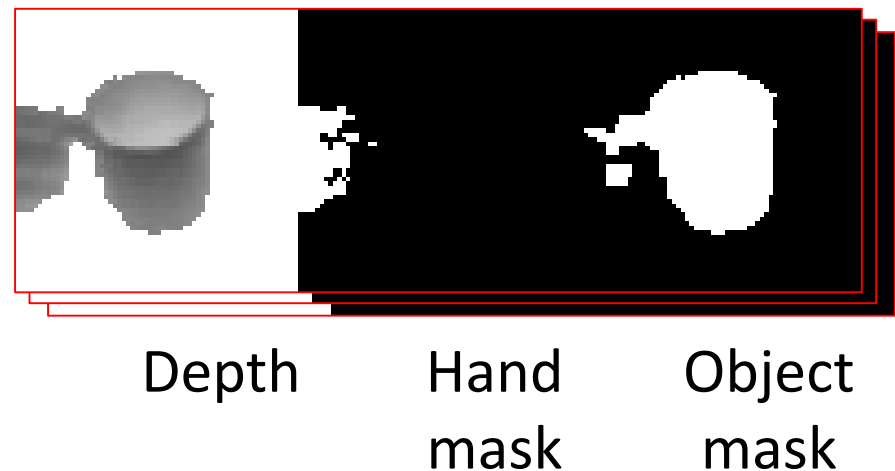**Object only scene**

**input**

Coordinate transformation by ICP



$T_1$

$T_2$

$T_1^{-1}$

$T_2^{-1}$

**Align**

$TransMatrix \leftarrow TransMatrix \times T_1$

$TransMatrix \leftarrow TransMatrix \times T_2$
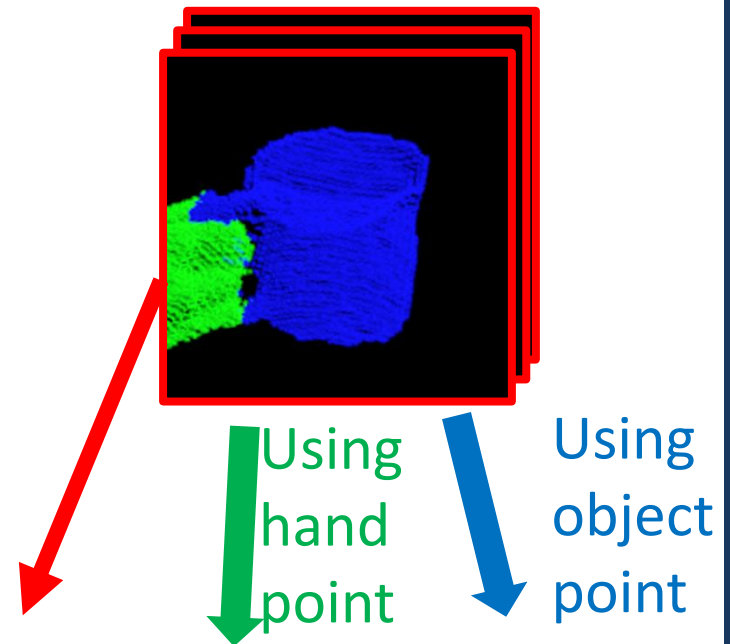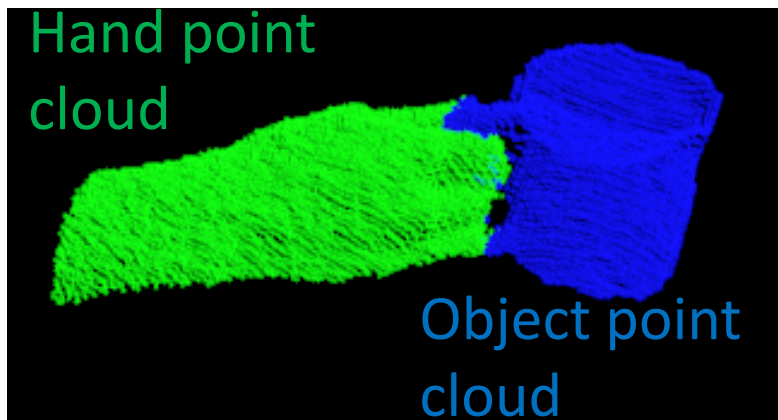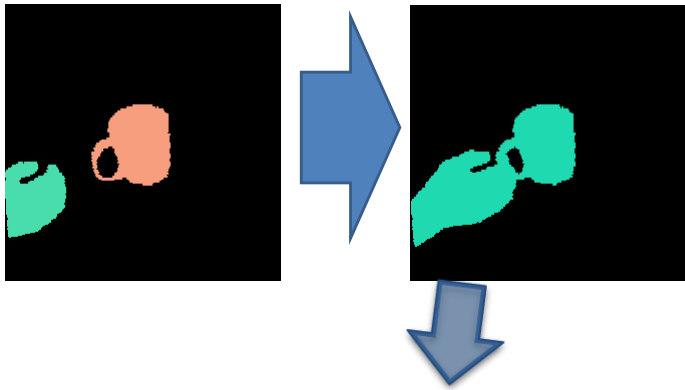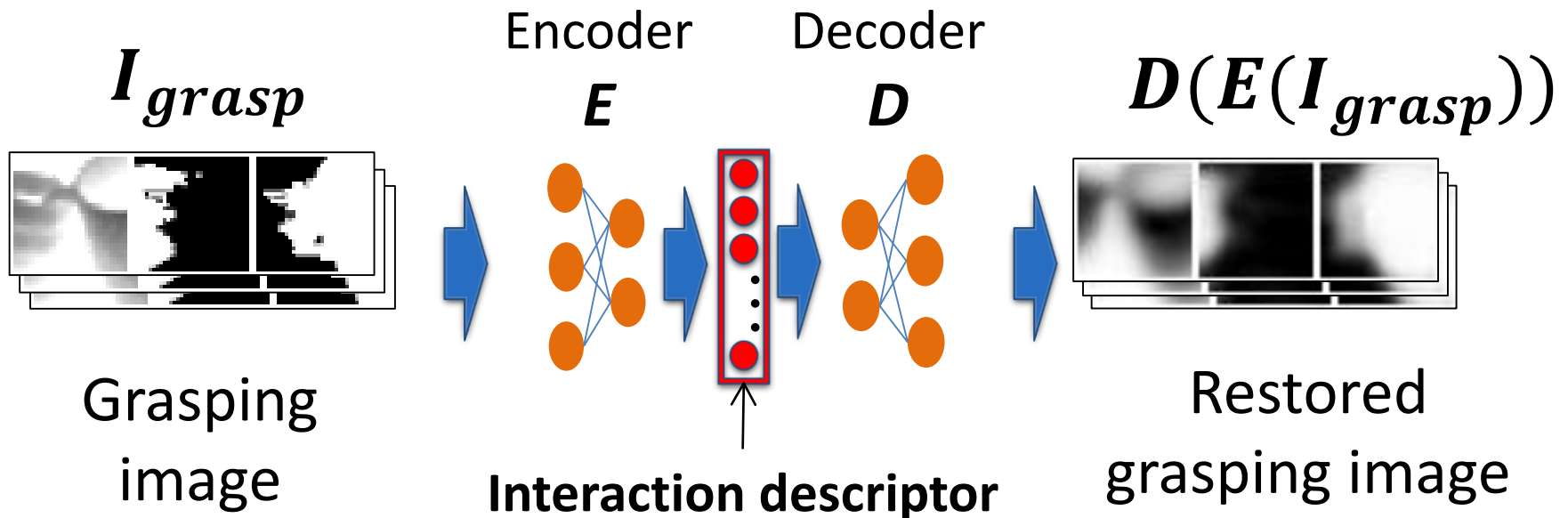
$TransMatrix^{-1} = T_1^{-1} \times T_2^{-1}$

# Segmentation of a hand and an object

Just before changing the number of regions, we segment hand points / object points.



Hand point cloud

Object point cloud

Using hand point

Using object point

Depth     Hand mask     Object mask

# Interaction descriptor

$I_{grasp}$

Encoder **E**

Decoder **D**

$D(E(I_{grasp}))$



Grasping image

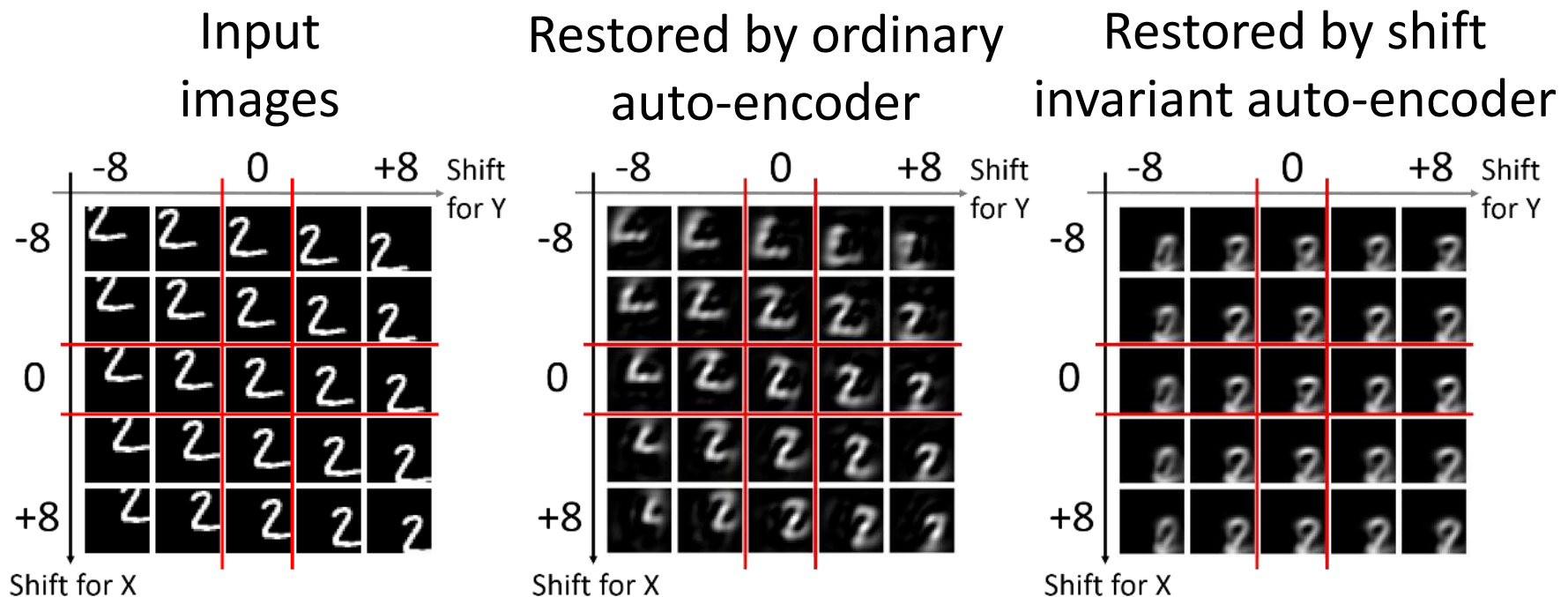**Interaction descriptor**

Restored grasping image

$E$ and $D$ are trained by minimizing restoration error
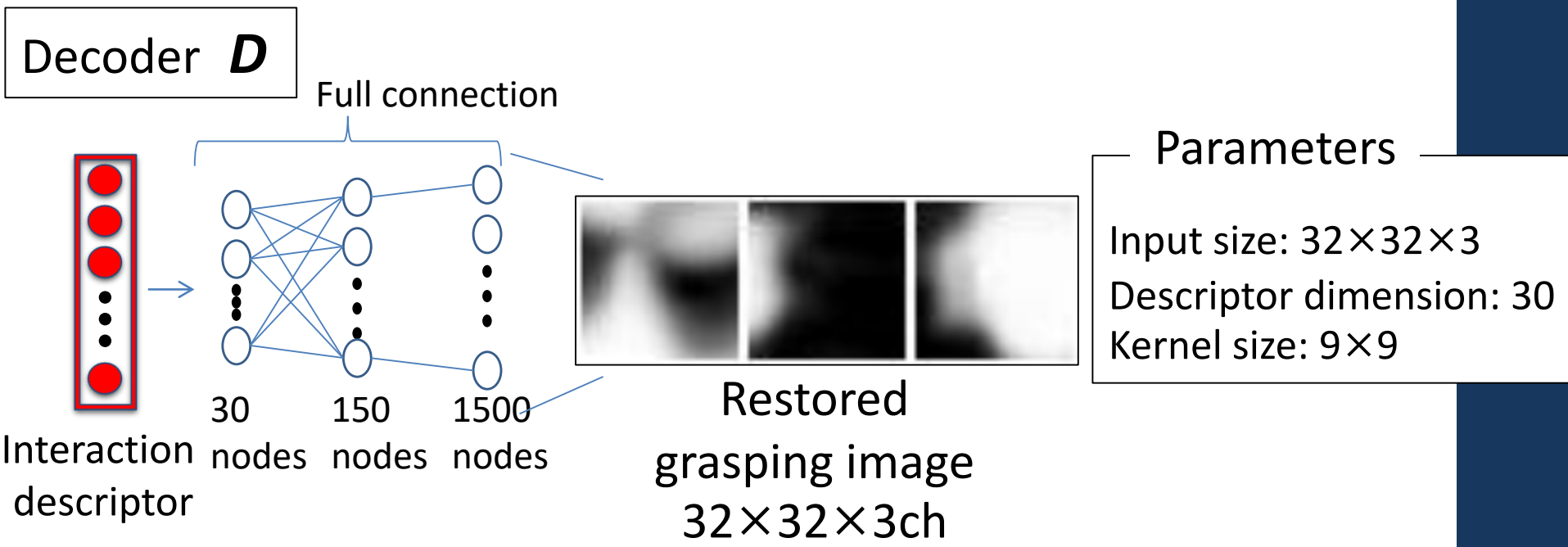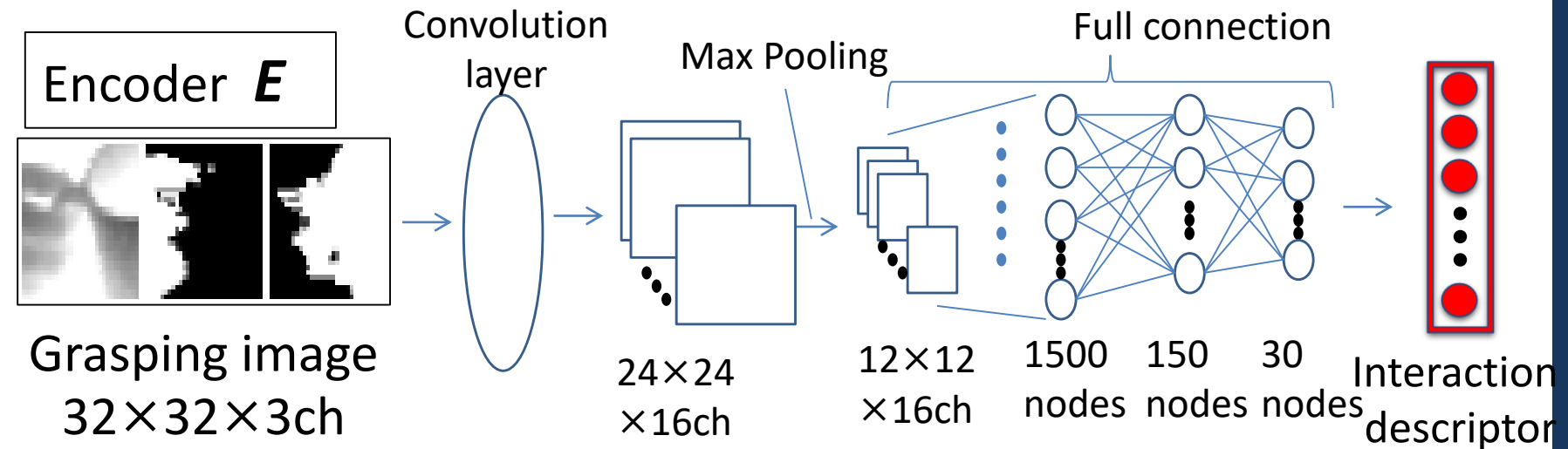without teacher labels.

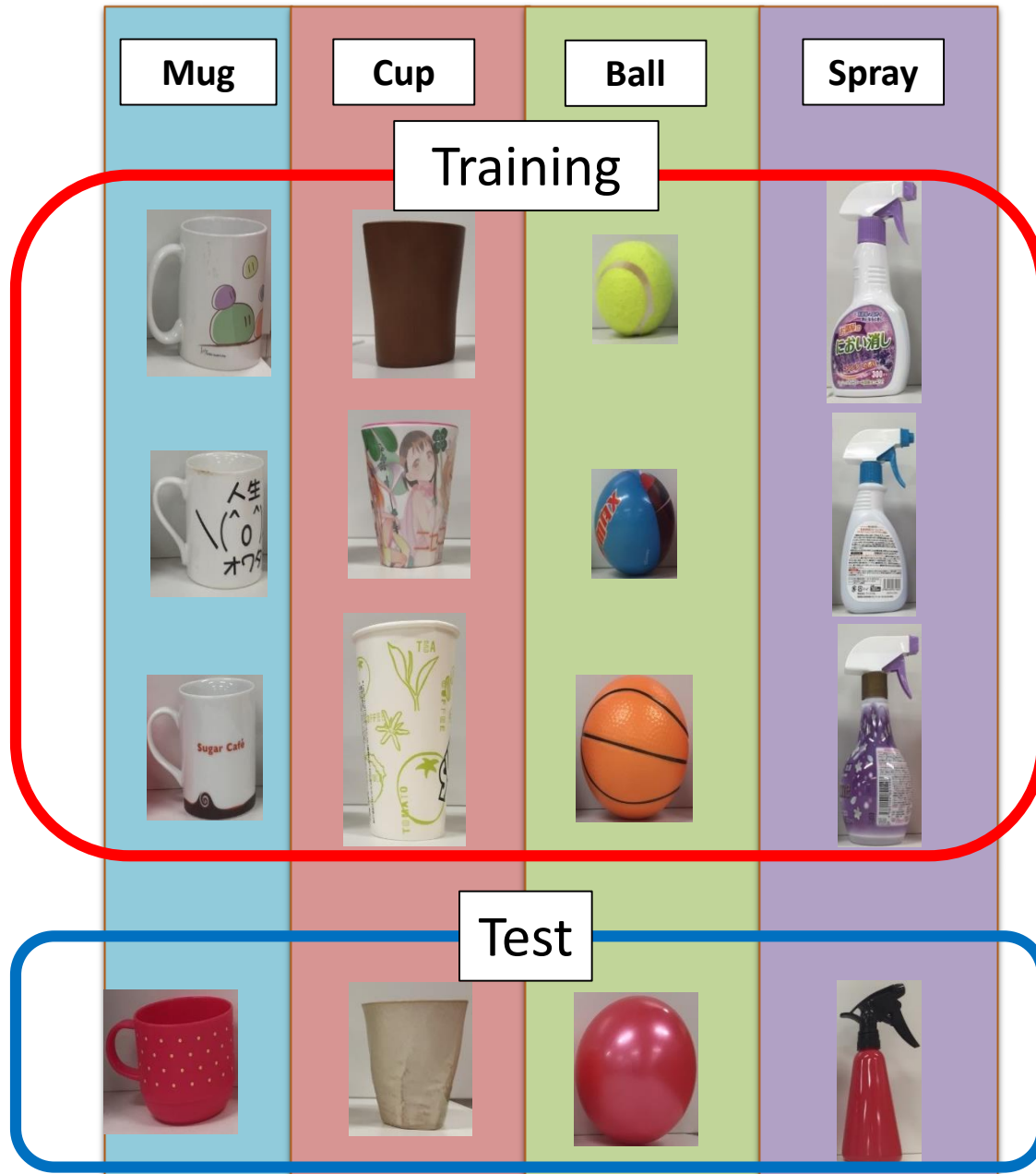A descriptor represents essence of an input.

# Shift invariant auto-encoder

An ordinary auto-encoder encodes shape and position.
But spatial shift in grasping images is not important.
We use shift invariant auto-encoder to encode a shape
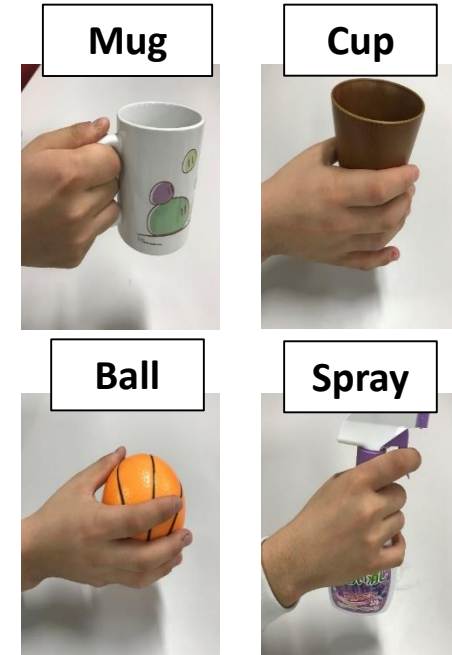itself. (descriptor includes shape information only)

Input
images

Restored by ordinary
auto-encoder

Restored by shift
invariant auto-encoder



T. Matsuo, et al., "Transform invariant auto-encoder," IROS 2017, https://doi.org/10.1109/IROS.2017.8206047

# Structure of auto-encoder



Encoder **E**

Convolution layer

Max Pooling

Full connection

Grasping image
32×32×3ch

24×24×16ch

12×12×16ch

1500 nodes

150 nodes

30 nodes

Interaction descriptor

Decoder **D**

Full connection

Interaction descriptor

30 nodes

150 nodes

1500 nodes

Restored grasping image
32×32×3ch

Parameters

Input size: 32×32×3
Descriptor dimension: 30
Kernel size: 9×9

# Objects and grasping types

# Restored grasping images

Input grasping images

Images restored from interaction descriptor



mug

cup

ball

spray

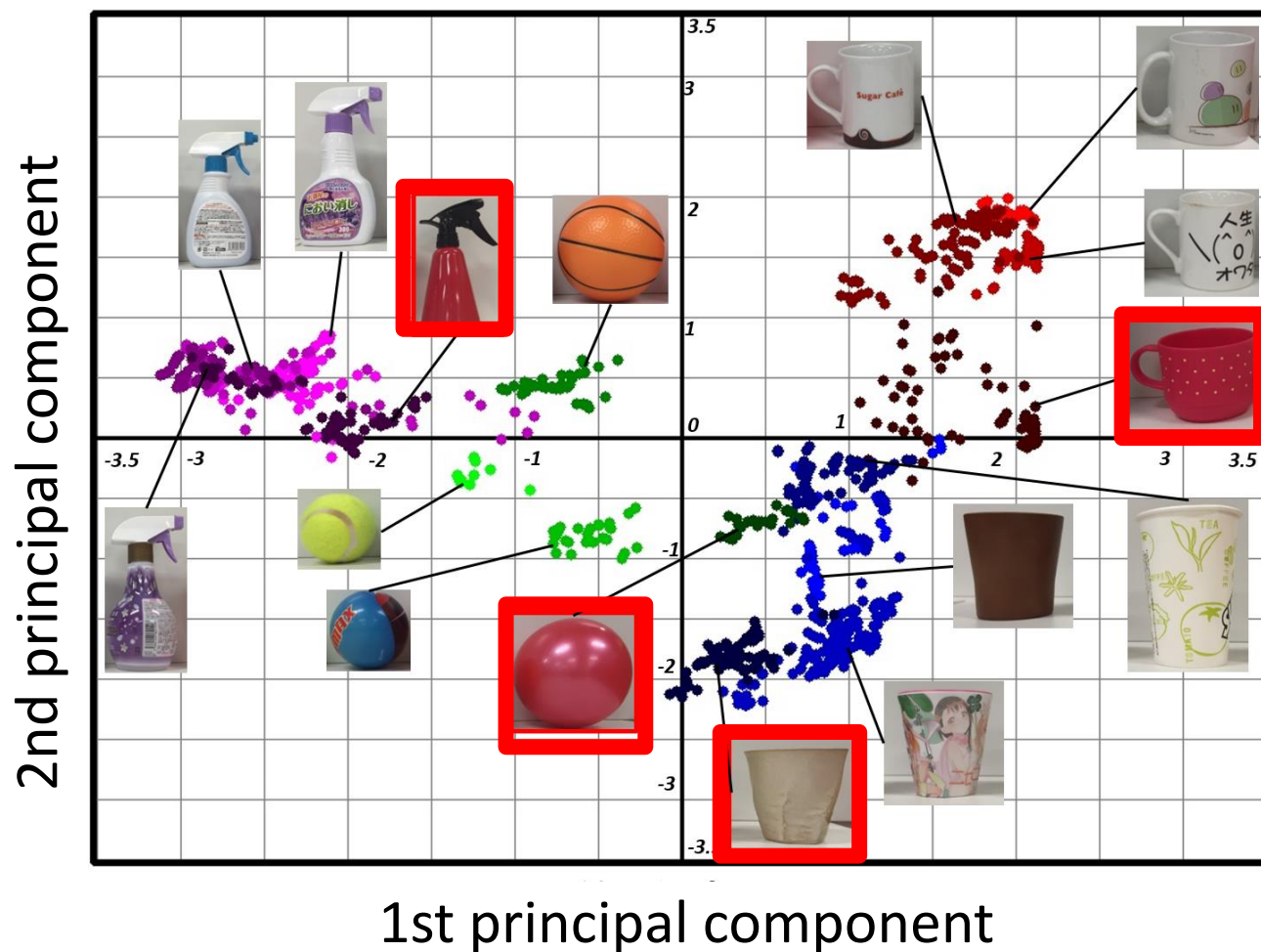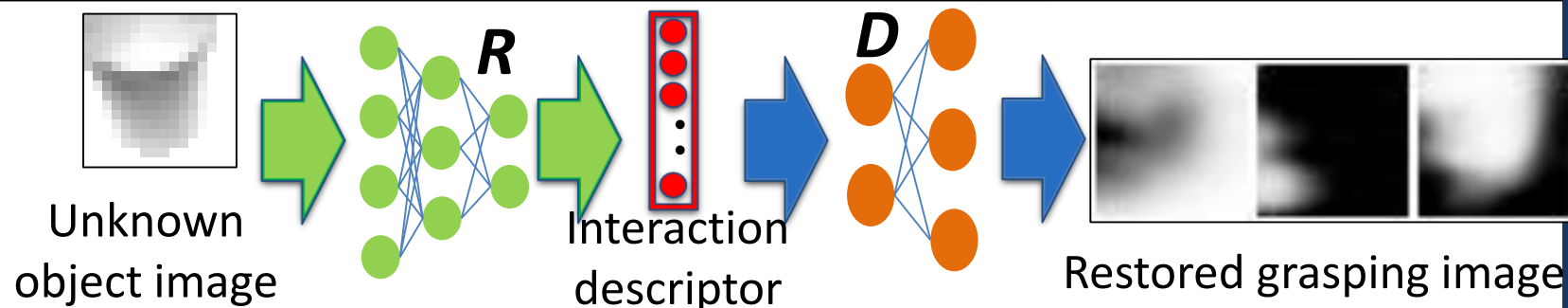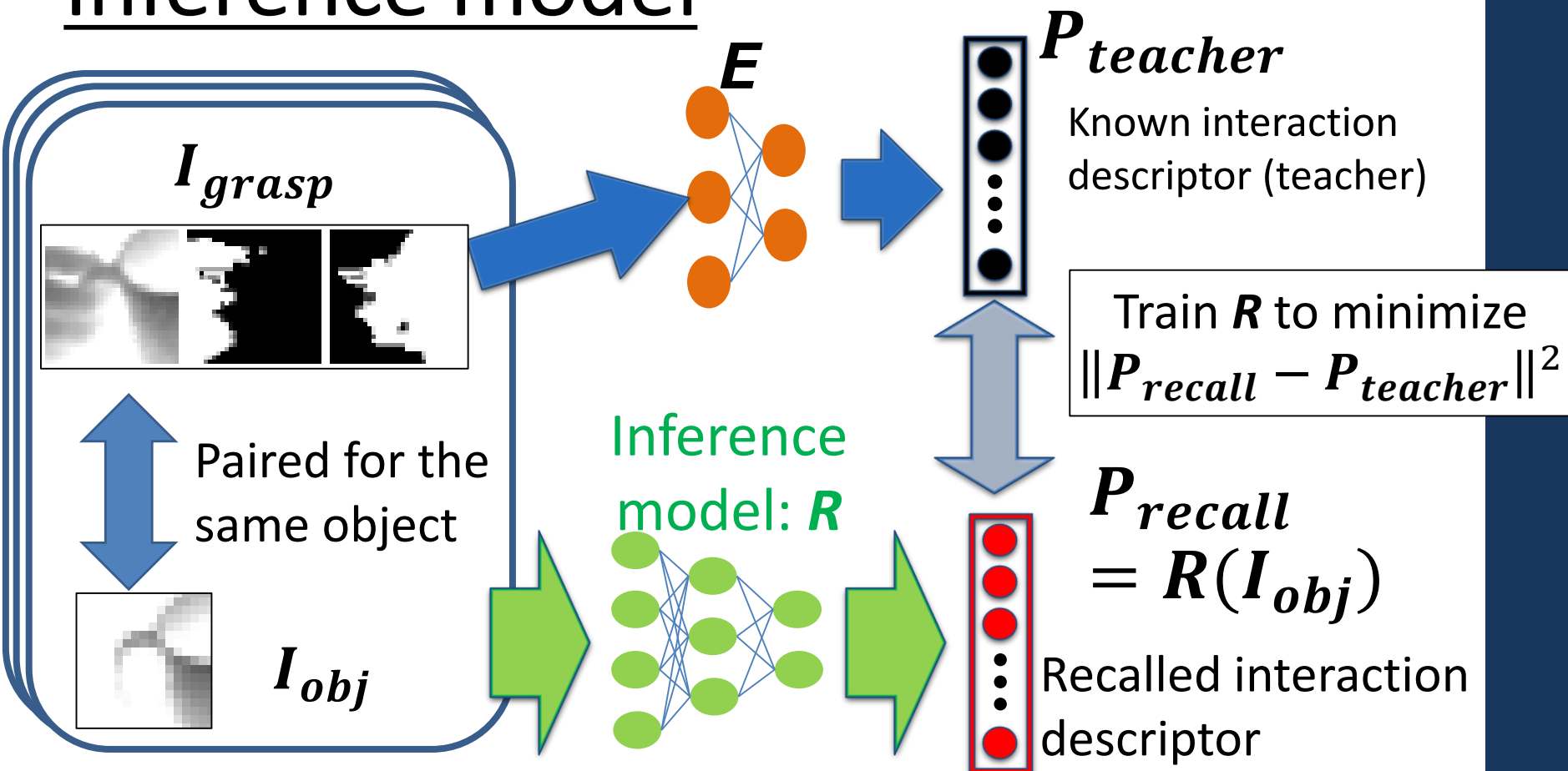Depth    Hand    Object

Depth    Hand    Object

Interaction descriptor has approximate shape information.

# Distribution of interaction descriptors



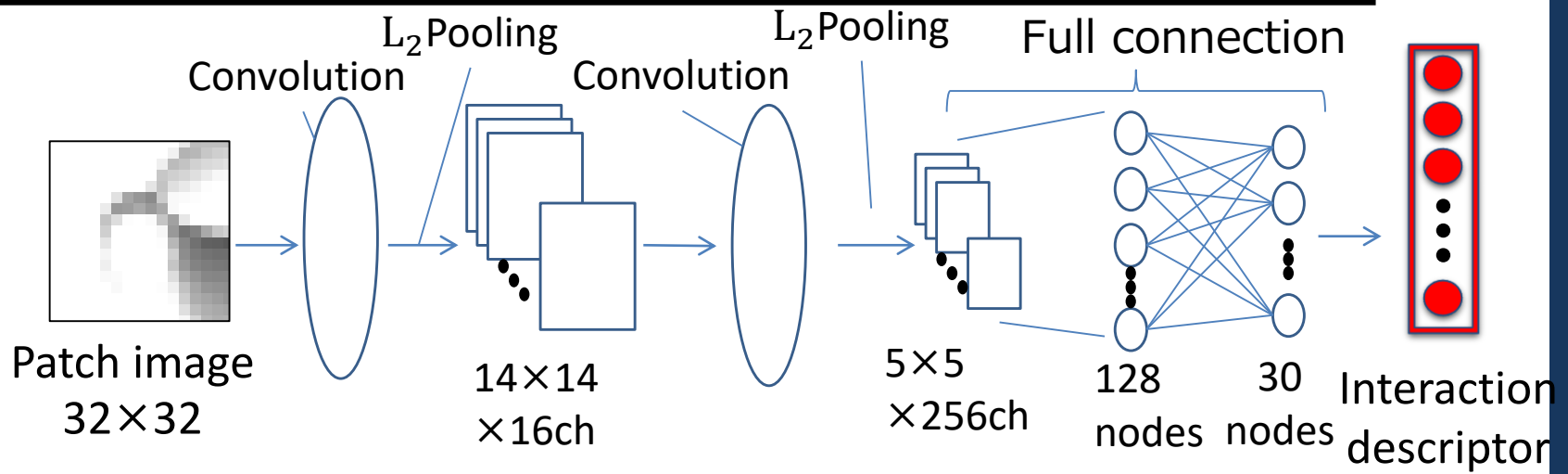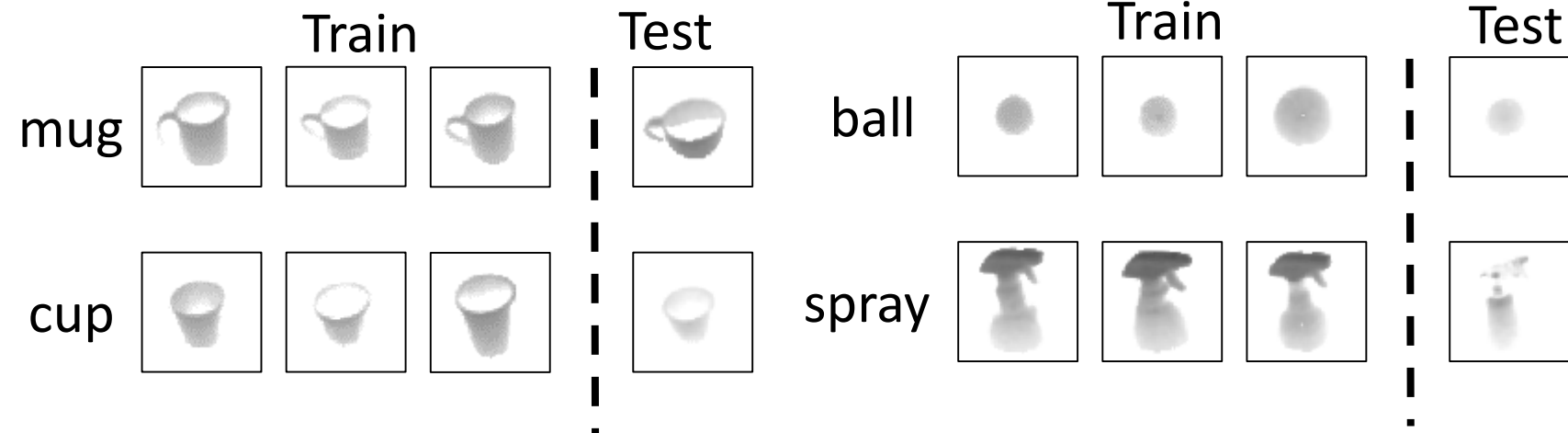Non-frame: Object for training

Red frame: Object for test

1st principal component

2nd principal component

# Inference model

$I_{grasp}$

$E$

$P_{teacher}$

Known interaction descriptor (teacher)

Paired for the same object

$I_{obj}$

Inference model: $R$

Train $R$ to minimize $\|P_{recall} - P_{teacher}\|^2$

$P_{recall} = R(I_{obj})$

Recalled interaction descriptor

$R$

$D$

Unknown object image

Interaction descriptor

Restored grasping image

# Structure of the inference model



Convolution  $L_2$Pooling  Convolution  $L_2$Pooling  Full connection

Patch image
32×32

14×14
×16ch

5×5
×256ch

128
nodes

30
nodes

Interaction
descriptor

Object images

Train    Test    Train    Test

mug

cup

ball

spray

# Recalled grasping images (train)

| Input | Recalled grasping image | Correct grasping image |
|-------|------------------------|------------------------|



mug

cup

ball

spray

Depth  Hand  Object  |  Depth  Hand  Object

The inference model successfully recalls grasping images.

# Recalled grasping images (test)



The model approximately recalls hand region masks.

# Recalled grasping images (test)



The model approximately recalls hand region masks.

# Recall from images with/without an important part
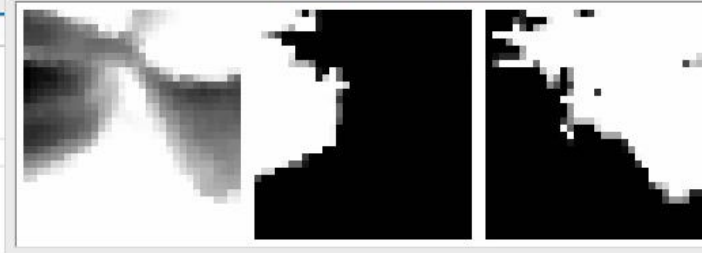
Recalled grasping image

with handle

Depth    Hand    Object

Mug

without handle

Cup

Cup

Similar!

# Recall from images with/without an important part



Object image 64×64

[1]物体画像

[3]把持画像のパッチ画像

Correct grasping image

[4]CNNで想起された把持画像

Recalled grasping image

Patch image 32×32

[2]

[5]AEで復元された把持画像

Restored grasping image

?

# Integration of recalled hand region masks



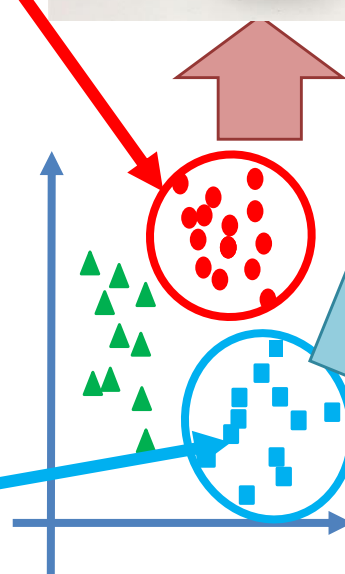Recalled image ($32 \times 32$)

Recalled image ($32 \times 32$)

Integrated hand region masks ($64 \times 64$)

Integrate descriptors in the same cluster

Interaction descriptor space

# Multiple grasping types for object

To see part-specific inference, we train auto-encoder and inference model with below grasping types.

# Integrated hand region mask

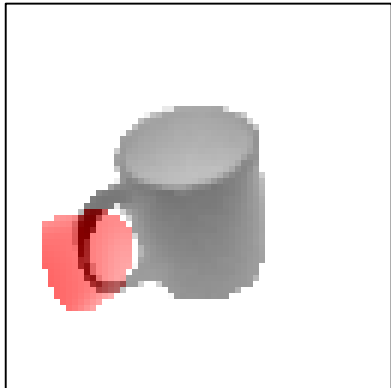**Integrated hand region mask** | **A real example of grasping**

Cluster 1



Cluster 2

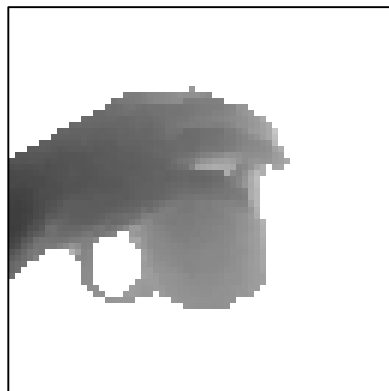The integrated hand mask for cluster $i$ is defined as:

$$P_i(x, y) = \frac{S_i(x, y)}{N_i(x, y)}$$

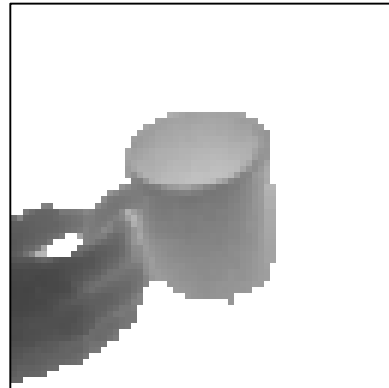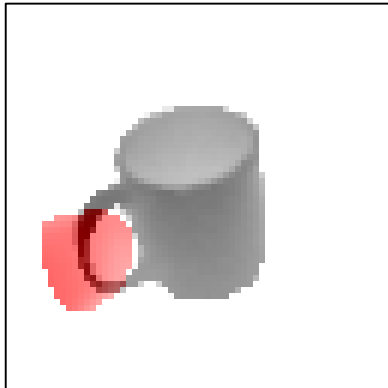$S_i(x, y)$: Sum of recalled hand mask in the $i$-th cluster

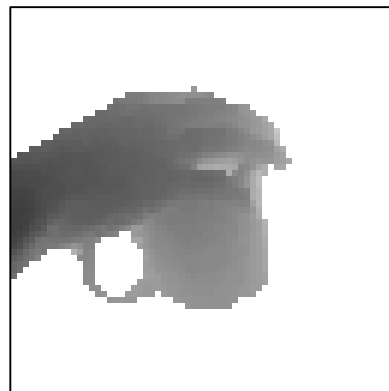$N_i(x, y)$: Number of non-zero at $(x, y)$ of recalled hand mask in the $i$-th cluster

Integrated hand region mask indicates hand region when human grasps the object.

# Integrated hand region mask

Integrated hand region mask | A real example of grasping

The integrated hand mask for cluster $i$ is defined as:

$$\text{(} \text{)} = \frac{S_i(x, y)}{N_i(x, y)}$$

Cluster 1

Cluster 2

$S_i(x, y)$: Sum of recalled ha... in the $i$-th cluster

$N$... Number of non-ze... $y$) of recalled hand mask in the $i$-th cluster

Integrated hand region mask indicate hand region when a human grasps the object.

# Conclusion

- We proposed a method to recall grasping method from an object. It is based on:

  - Interaction descriptor by shift invariant auto-encoder
    We can generate numeral representation of grasping method without teacher labels.

  - Inference model by CNN
    The relation between object shape and grasping method can be modeled by utilizing interaction descriptor.

- The proposed method can estimate hand region for grasping an object from the object itself.

- The proposed method will be useful for robot manipulator.

# Distribution of descriptors from shift invariant auto-encoder

We trained auto-encoders with shifted MNIST training images.
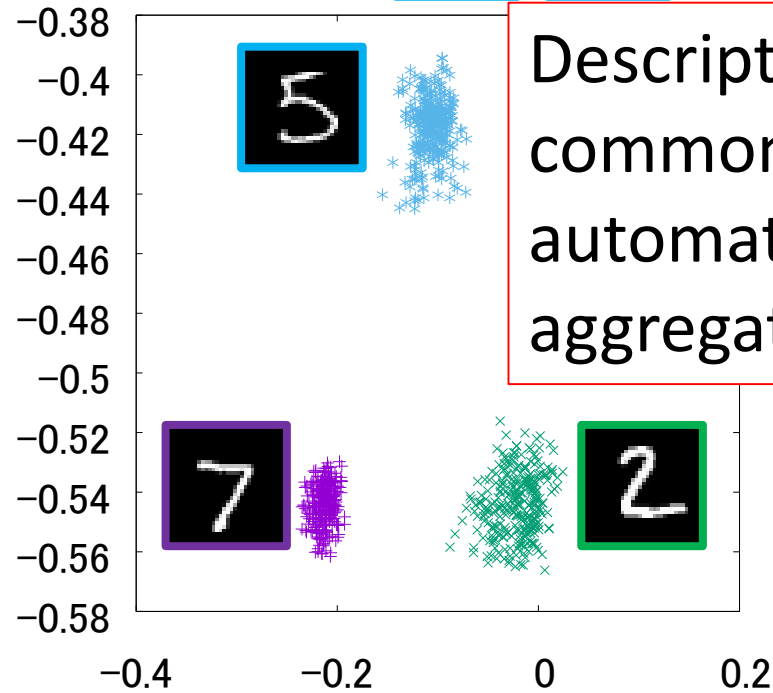
Input: $32 \times 32$
Descriptor dim: 30
Max shift width: 8

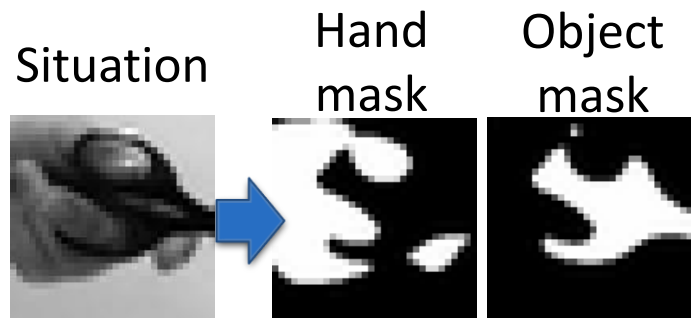Distributions of descriptors from shifted test images such as .

Descriptors from common shapes automatically aggregate.

Ordinary auto-encoder

Shift invariant auto-encoder

# Example for hand-object interaction

Situation
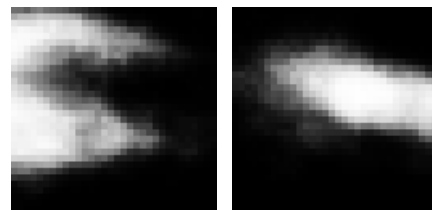
Hand mask

Object mask

Train AEs with 2-channel images consisting of hand/object masks.
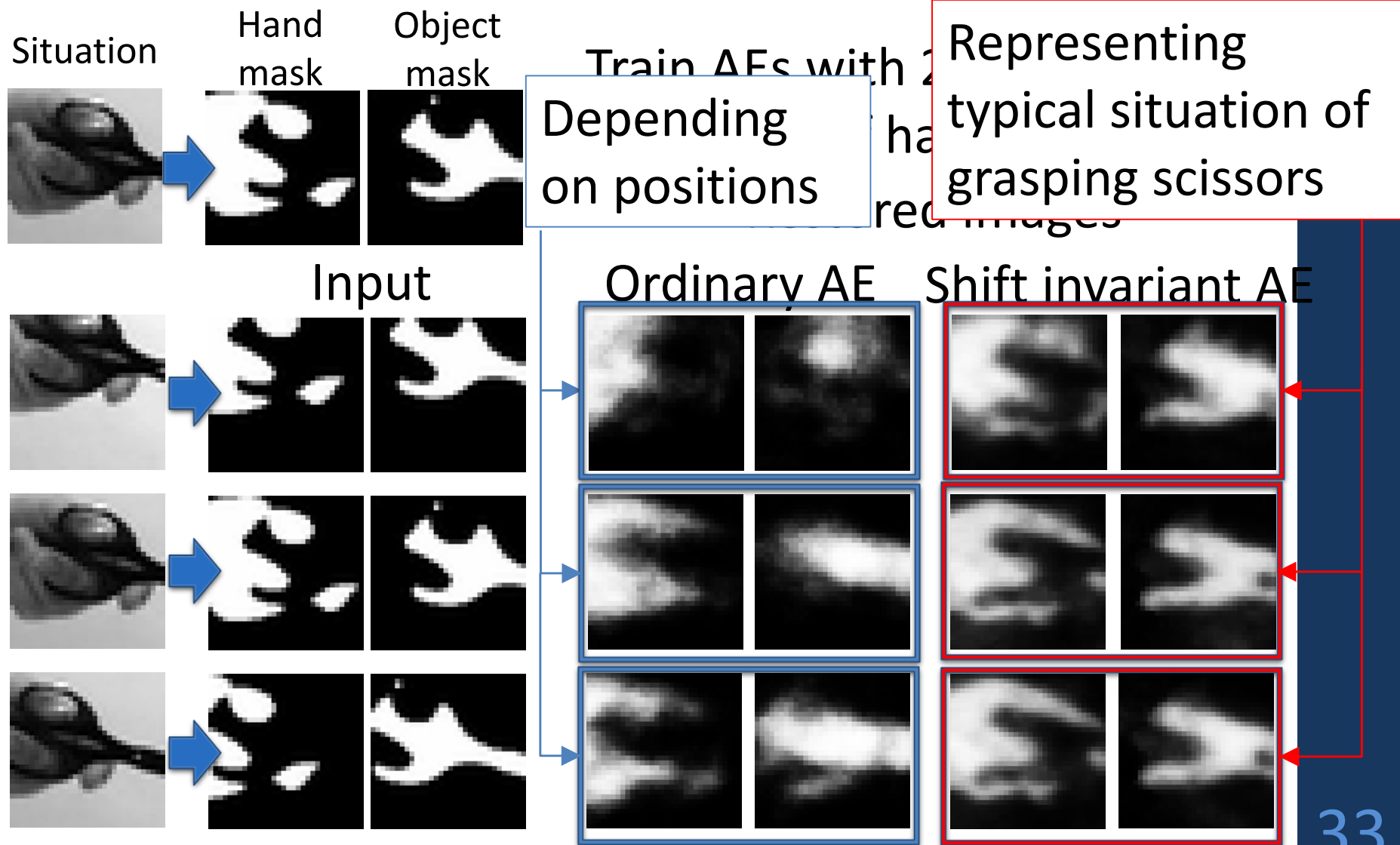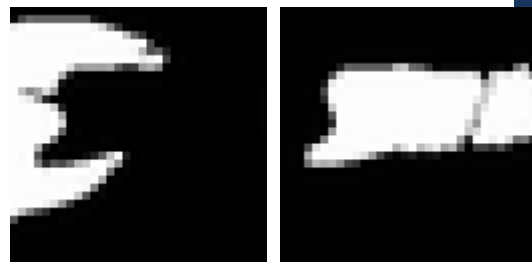
Restored images

Input

Ordinary AE
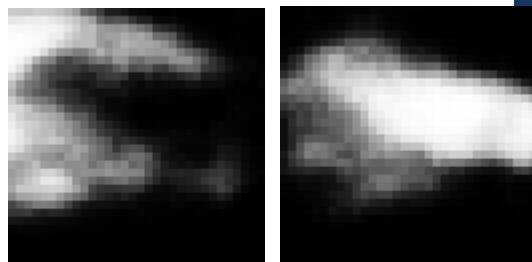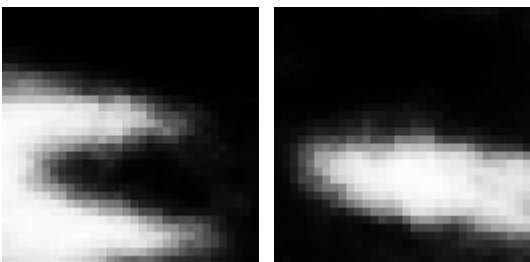
Shift invariant AE

# Example for hand-object interaction



Situation | Hand mask | Object mask

Train AEs with 2...

Depending on positions

Representing typical situation of grasping scissors

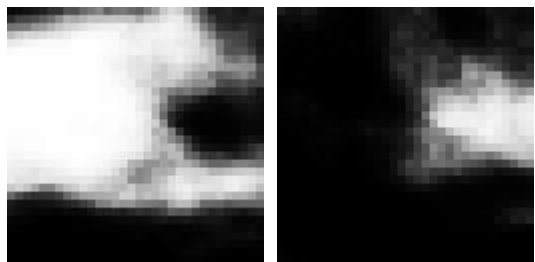Input | Ordinary AE | Shift invariant AE
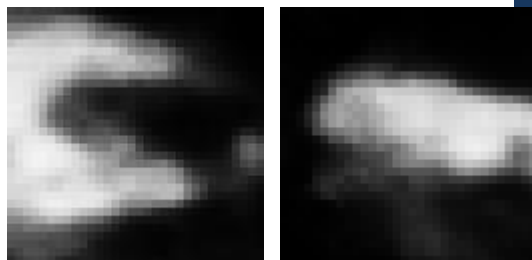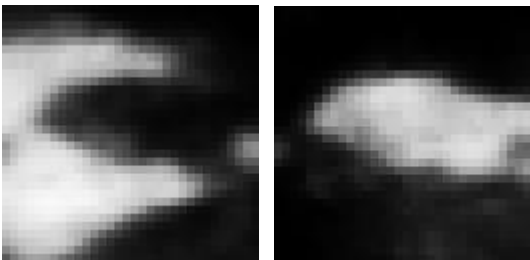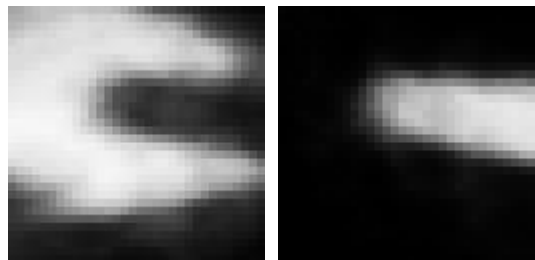
33

Input image $I$     Input image $I$     Input image $I$

Images restored by an ordinary auto-encoder

Images restored by a shift invariant auto-encoder