# DenseFusion

KIST
**송명하**

# Dataset



**RGB**



**Depth**



**Label(seg)**



**Point cloud**



**Pose(Matrix)**

**Bounding Box**

# Dataset (Symmetric Object)



**Bowl**

**Wood_block**

**Large_clamp**

**Extra_large_clamp**

**Foam_brick**

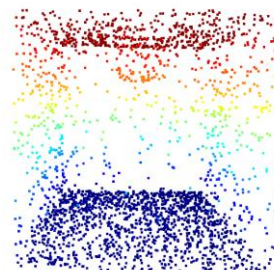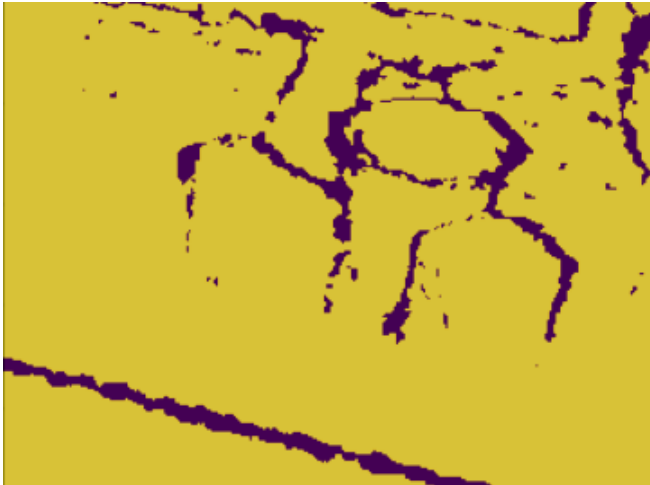Use the Segmentation label(PoseCNN result) to remove the background

-> Data augmentation(adding noise)

**In the Depth image, only nonzero parts are made True and then multiplied by the semantic label of the object. The result is shown in the image below**

1. **Use the bound box's Corrdinates(PoseCNN result) to cut out the image and the depth mask**
2. **Only a nonzero value is selected after flattening the depth mask value**
3. **Randomly sampling 1000 out of the values created in step2, The location of the selected value is stored in a variable named choose**
4. **In the depth masks, indexing is done using the choose variable**
5. **Create a variable named xmap and ymap indexing it in the same way, and create a pointcloud.**

**The number of points in the object – 500 randomly sampled. Remove from model point and save 500 points as model point.**

**Multiply the rotation matrix at the created model point and add translation.**

# Network Architecture

**RGB** — **Resnet 18** — **PSP Module** — **PSP Up Sample Module** — **Flatten sampling**

# Network Architecture

# Network Architecture

RGB ── Resnet 18 ── PSP Module ── PSP Up Sample Module ── Flatten sampling

# Network Architecture

RGB — **Resnet 18** — **PSP Module** — **PSP Up Sample Module** — **Flatten sampling**

# Network Architecture

PSP Module

| Adap Avg Pool (1,1) | Conv 2d 512 | Up sample bilinear |
| Adap Avg Pool (2,2) | Conv 2d 512 | Up sample bilinear |
| Adap Avg Pool (3,3) | Conv 2d 512 | Up sample bilinear |
| Adap Avg Pool (6,6) | Conv 2d 512 | Up sample bilinear |

Concat 2560

Conv2d 1024

Center for Robotics Research

# Network Architecture

RGB — **Resnet 18** — **PSP Module** — **PSP Up Sample Module** — **Flatten sampling**

# Network Architecture

**Concat 2560**

**Conv2d 1024**

**PSP Upsample Module**

**Drop Out 0.3** | **Up sample bilinear** | **Conv 2d 512** | **PRelu** | **Drop Out 0.15** | **Up sample bilinear** | **Conv 2d 512** | **PRelu** | **Drop Out 0.15** | **Up sample bilinear** | **Conv 2d 512** | **PRelu** | **Conv 2d 32** | **Log softmax**

# Network Architecture

RGB — **Resnet 18** — **PSP Module** — **PSP Up Sample Module** — **Flatten sampling**
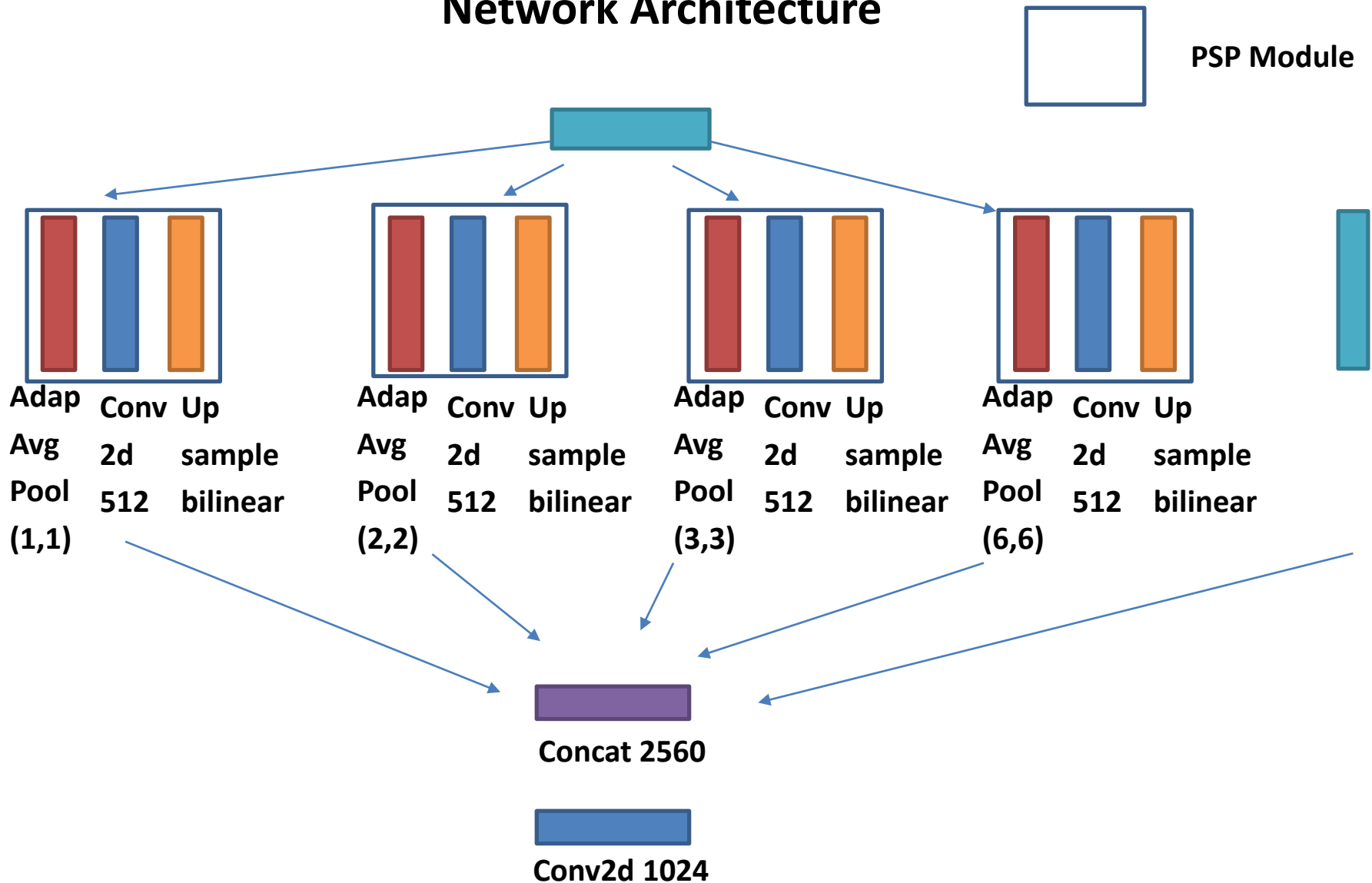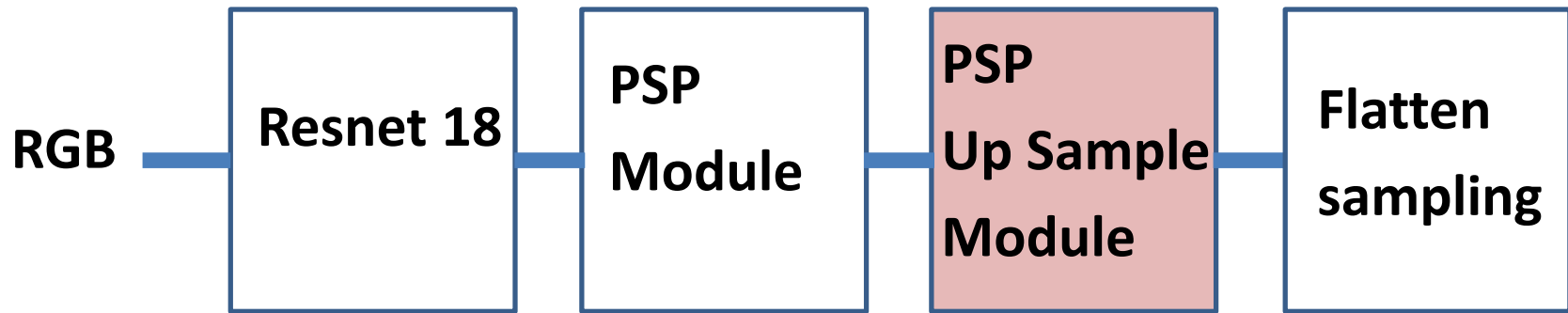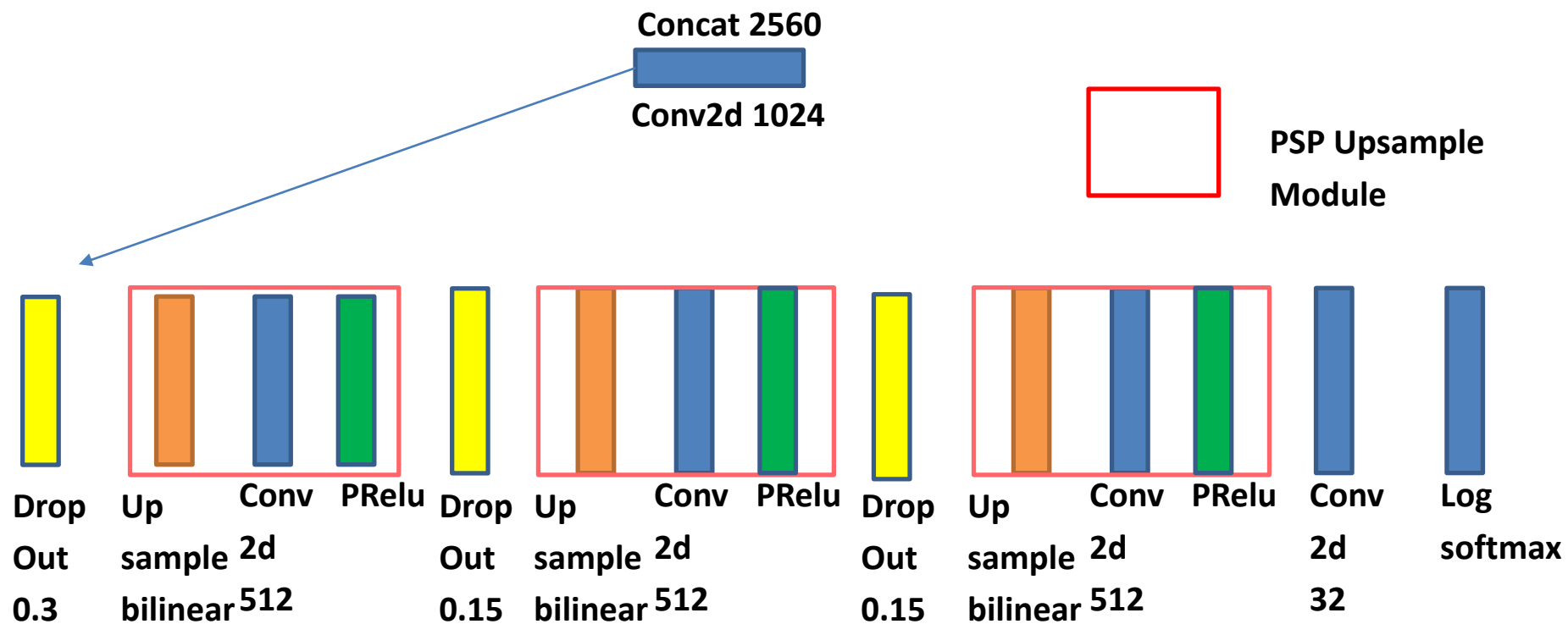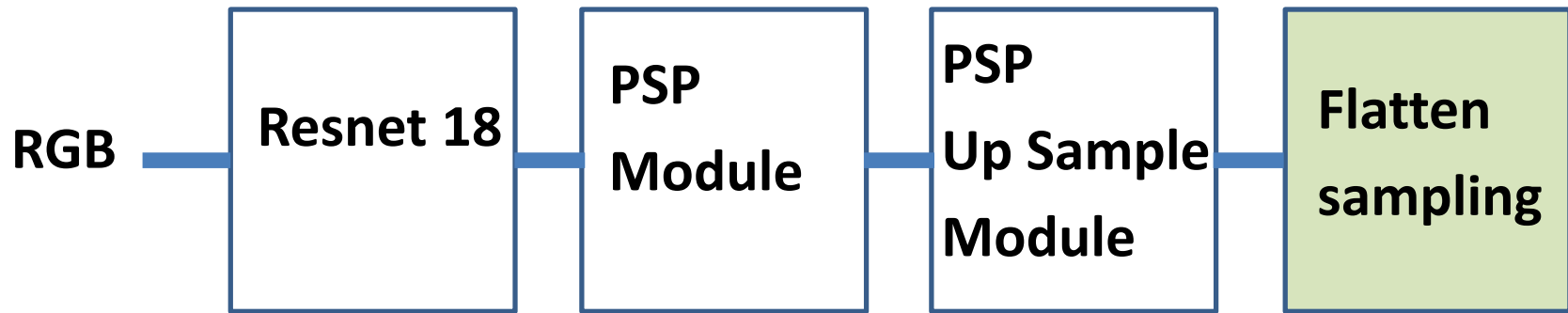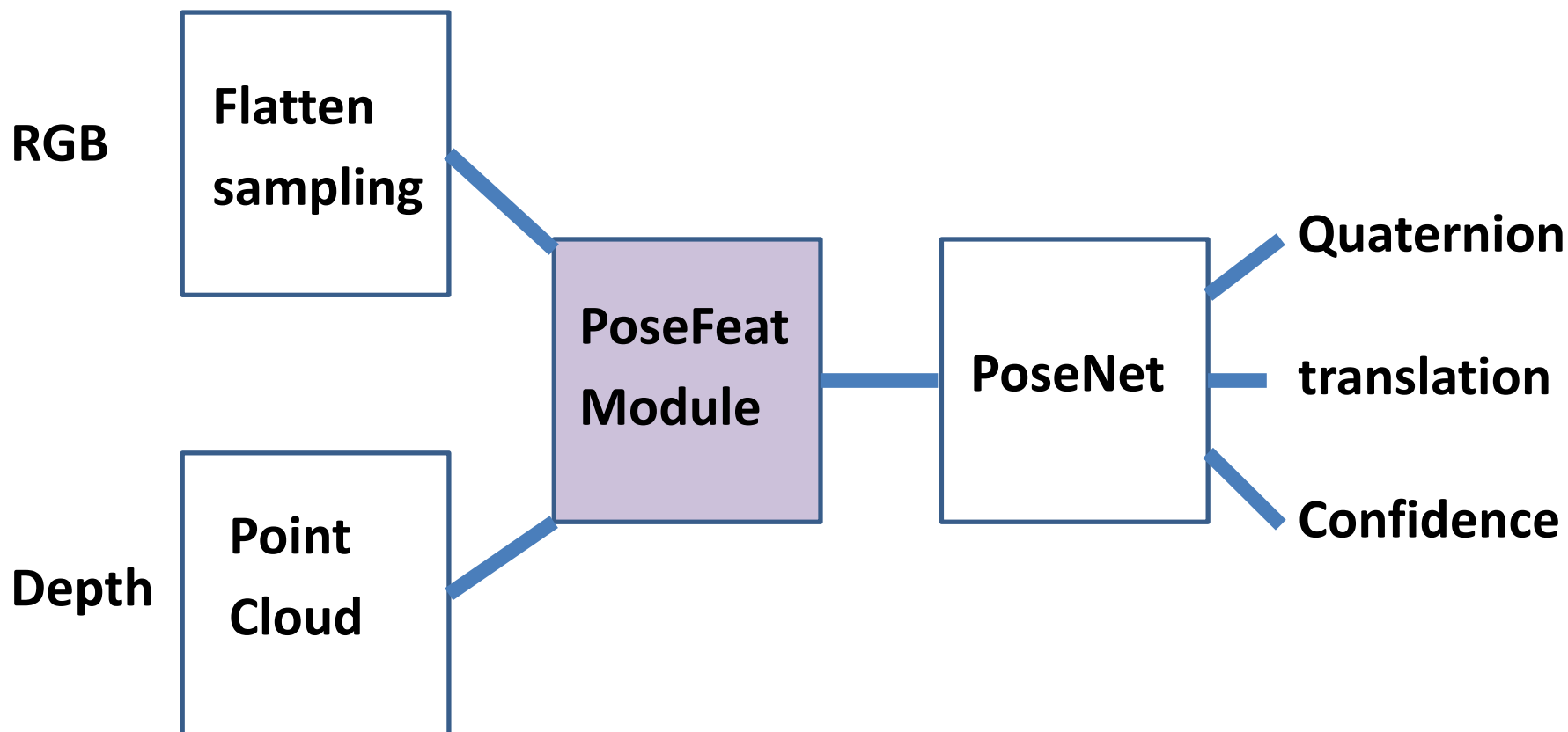
# Network Architecture

# Network Architecture
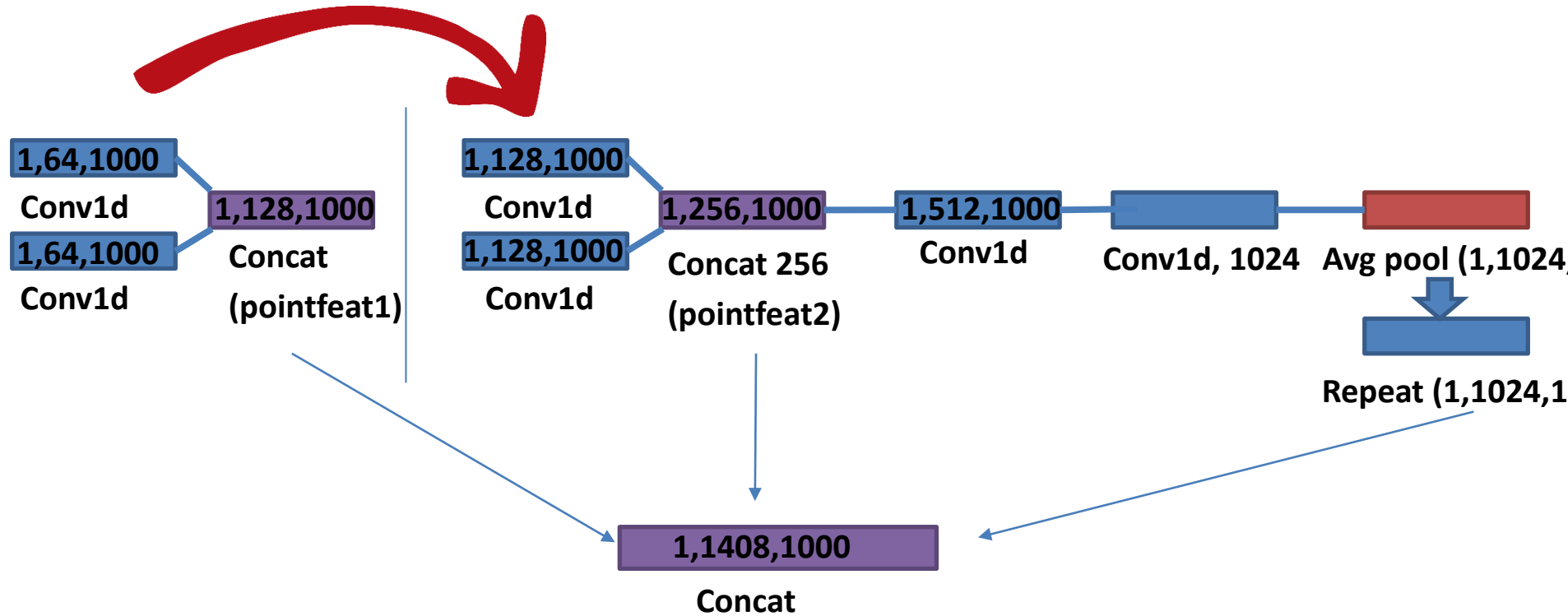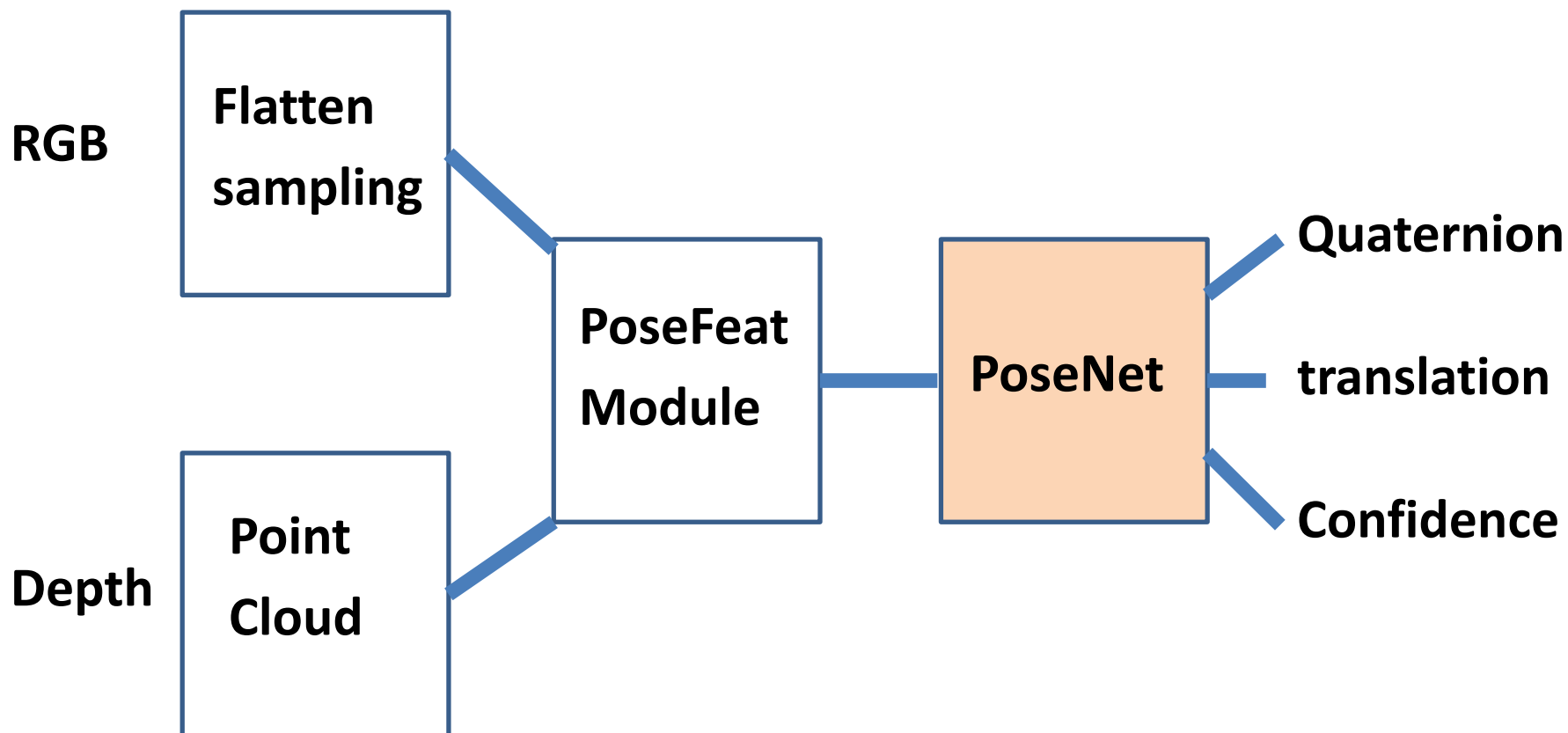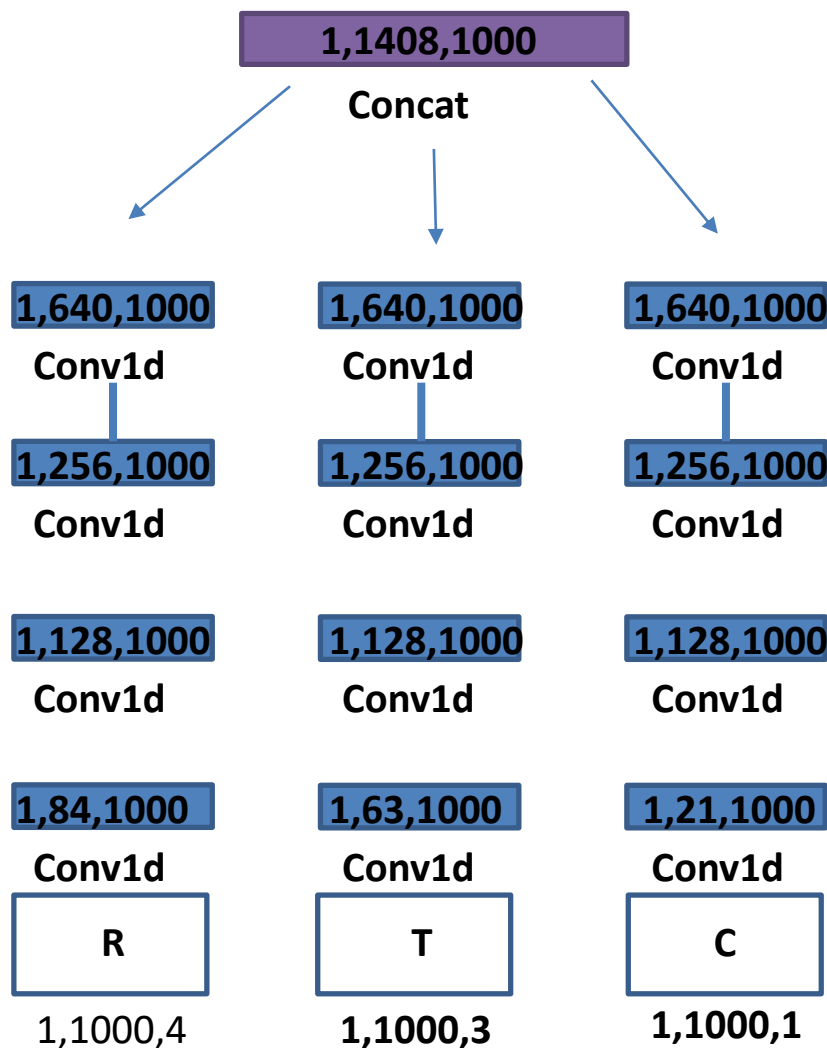
# Network Architecture

# Network Architecture



1,1408,1000

**Concat**

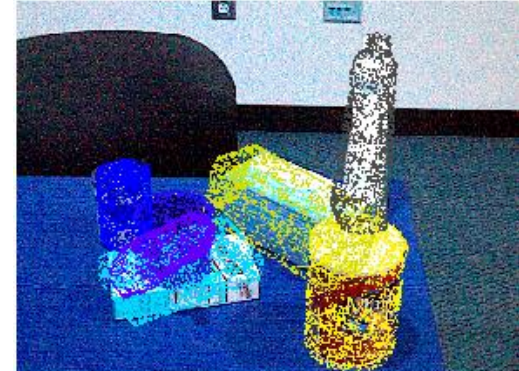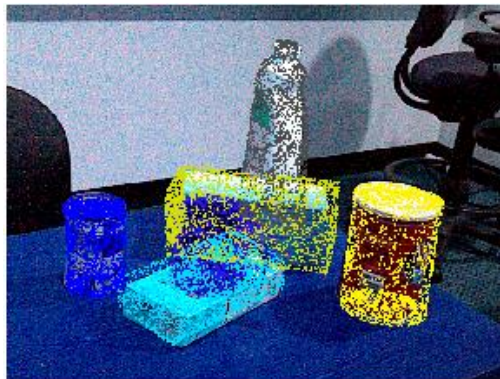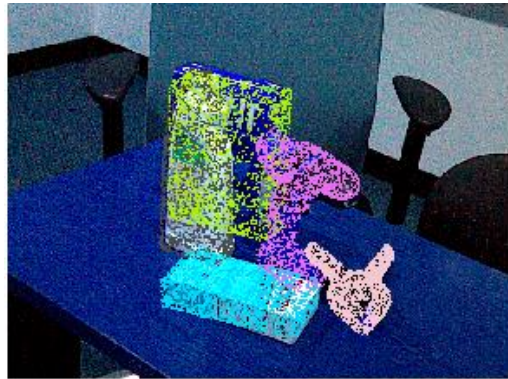| 1,640,1000 | 1,640,1000 | 1,640,1000 |
|---|---|---|
| Conv1d | Conv1d | Conv1d |
| 1,256,1000 | 1,256,1000 | 1,256,1000 |
| Conv1d | Conv1d | Conv1d |
| 1,128,1000 | 1,128,1000 | 1,128,1000 |
| Conv1d | Conv1d | Conv1d |
| 1,84,1000 | 1,63,1000 | 1,21,1000 |
| Conv1d | Conv1d | Conv1d |
| R | T | C |
| 1,1000,4 | 1,1000,3 | 1,1000,1 |

## Loss

$$L_i^p = \frac{1}{M} \sum_j \|(Rx_j + t) - (\hat{R}_i x_j + \hat{t}_i)\| \qquad (1)$$

$$L = \frac{1}{N} \sum_i (L_i^p c_i - w \log(c_i)),$$

# Result

# 속도

Table 3. Runtime breakdown (second per frame on YCB-Video Dataset). Our method is approximately 200x faster than PoseCNN+ICP. Seg means Segmentation, and PE means Pose Estimation.

| PoseCNN+ICP [40] | | | | Ours | | | |
|------|------|------|------|------|------|--------|------|
| Seg | PE | ICP | ALL | Seg | PE | Refine | ALL |
| 0.03 | 0.17 | 10.4 | 10.6 | 0.03 | 0.02 | 0.01 | 0.06 |

**16FPS, about 5 objects in each frame (16.6fps)**

# End