



# Take a ~~Deeper~~ Look at Object 6D Pose Estimation

**Siyu ZHANG**

Research Engineer

ZJU-SenseTime Joint Lab of 3D Vision

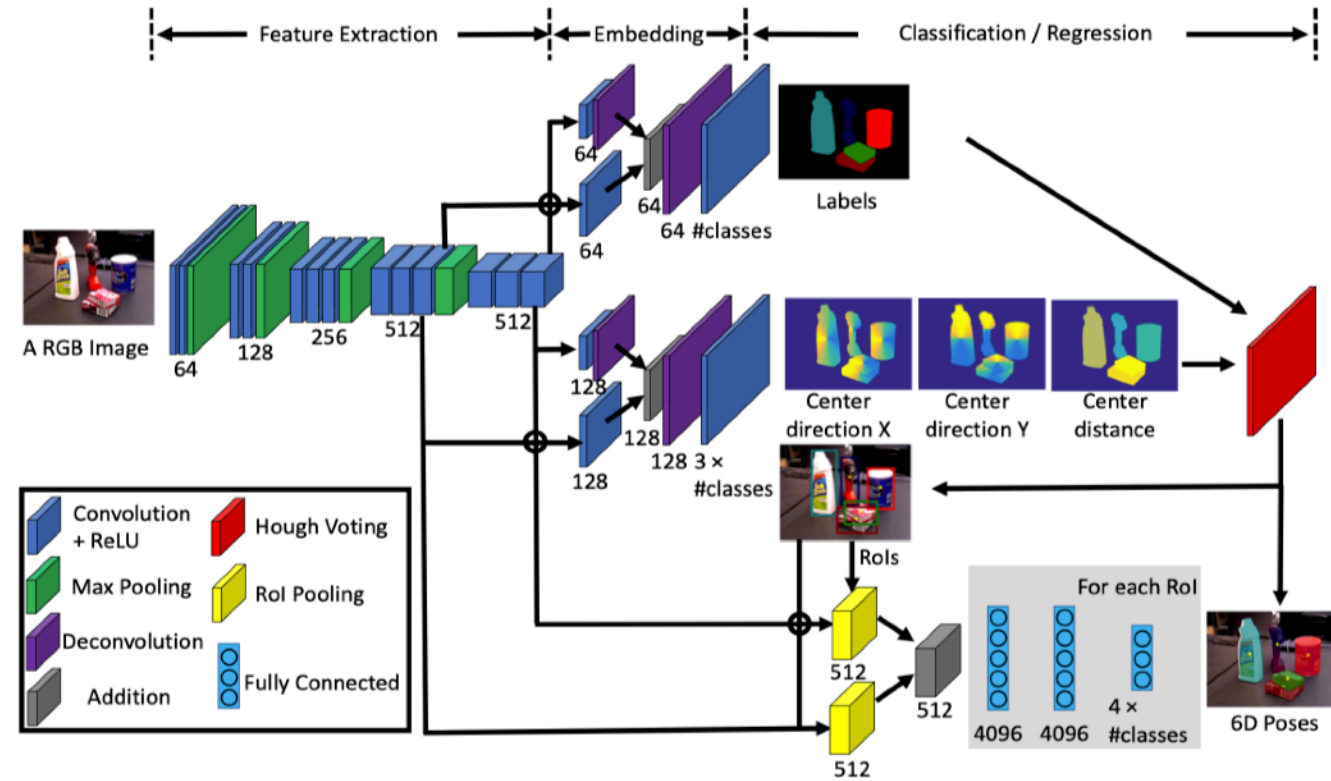
# Overview

---

- Some terminologies ...
  - Pose Estimation
    - Given: Image<sub>t</sub>
    - Target: Pose<sub>t</sub>
  - Pose Refinement
    - Given: Image<sub>t</sub>, Pose<sub>t\_init</sub>
    - Target: Pose<sub>t\_refined</sub>
  - Pose Tracking
    - Given: Image<sub>t-1</sub>, Pose<sub>t-1</sub>, Image<sub>t</sub>
    - Target: Pose<sub>t</sub>

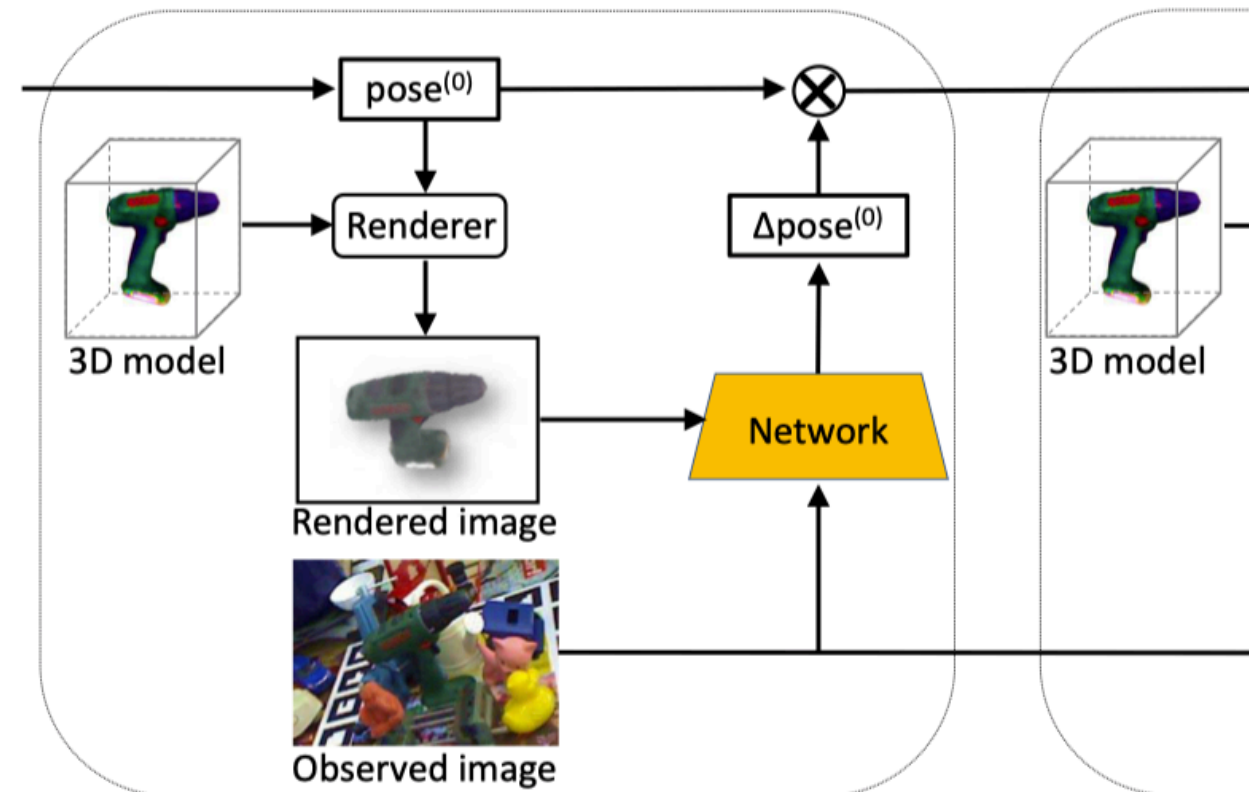
# Overview

- As a **regression** problem
  - Pose Estimation: direct regression
    - Extract feature from image
    - Directly regress the translation and orientation of target object



# Overview

- As a **regression** problem
- Pose Tracking: render and regression
  - Render the image of target object with initial pose and object 3D model
  - Feed the rendered image and input image into Neural Network to extract feature
  - Regress the additive pose for target object



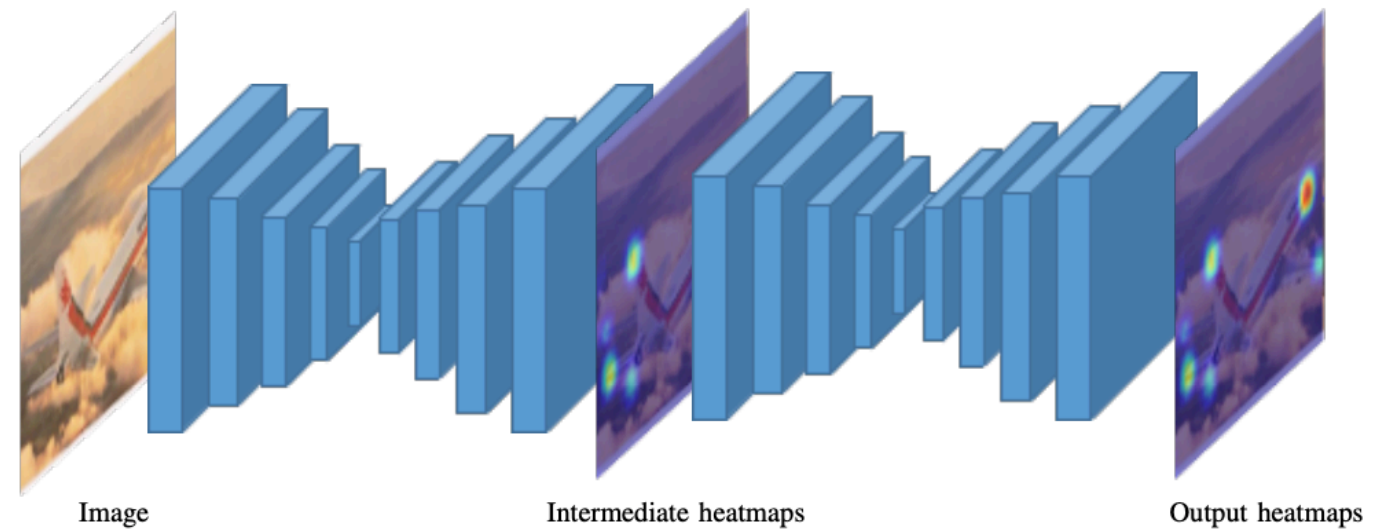
# Overview

---

- As a **regression** problem
  - Pose Estimation: direct regression from image feature
  - Pose Tracking: render, compare with observed image, then regress relative (additive) pose
- As a **matching** problem

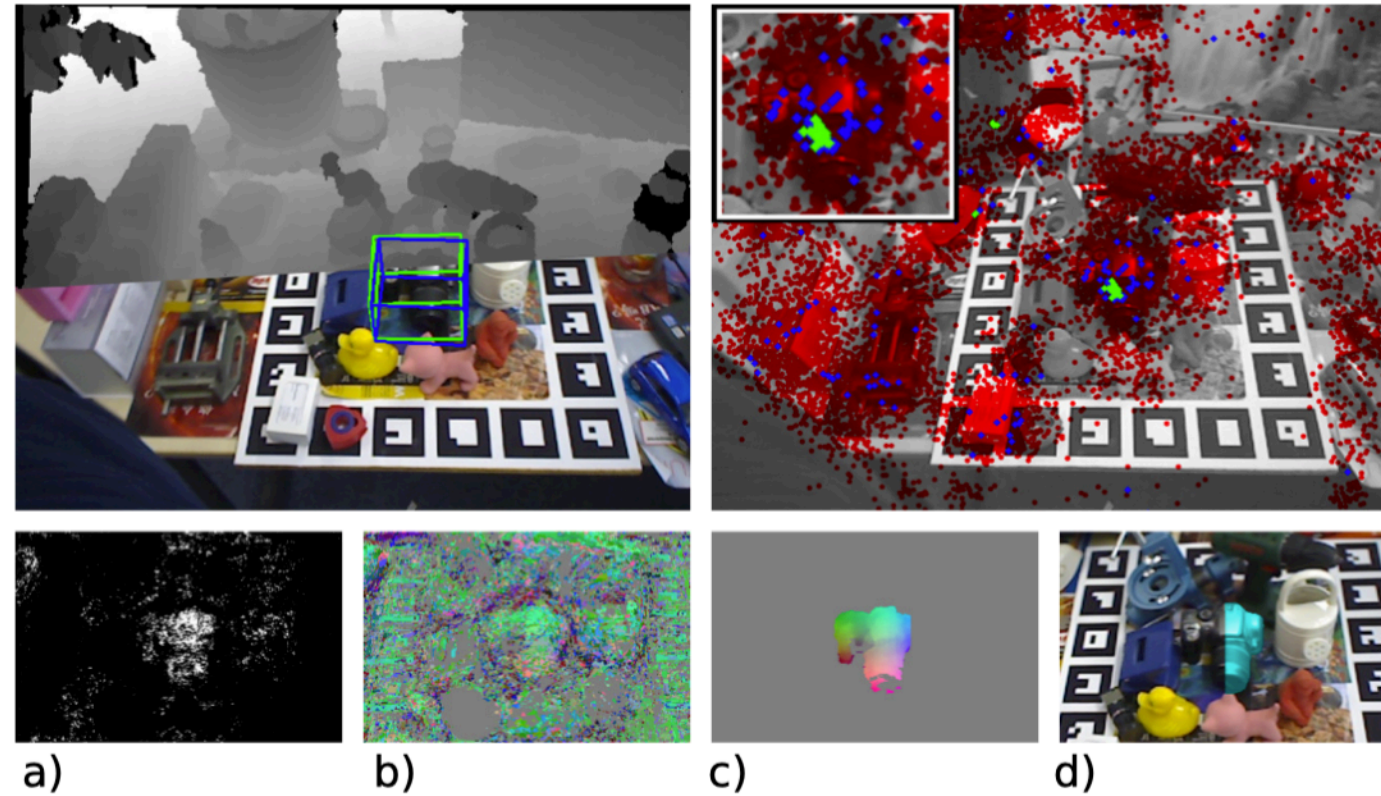
# Overview

- As a **matching** problem
- Pose Estimation by matching sparse correspondences
  - Extract semantic 2D keypoint from the image
  - Solve for PnP with corresponding 3D keypoints in object frame



# Overview

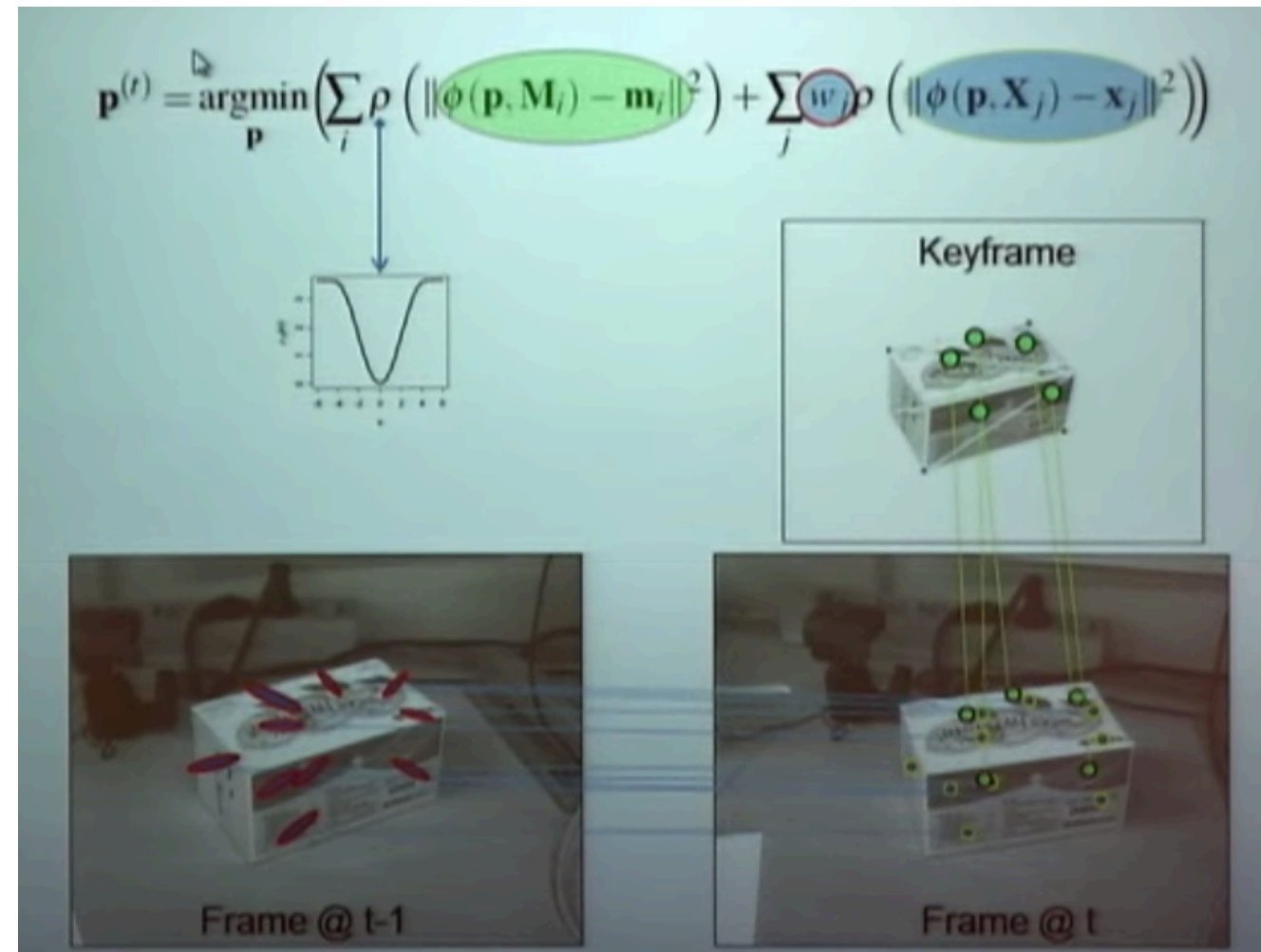
- As a **matching** problem
  - Pose Estimation by matching dense correspondences
    - Estimate pixel-wise object coordinates for all foreground pixels
    - Apply RANSAC + PnP to solve for object pose





# Overview

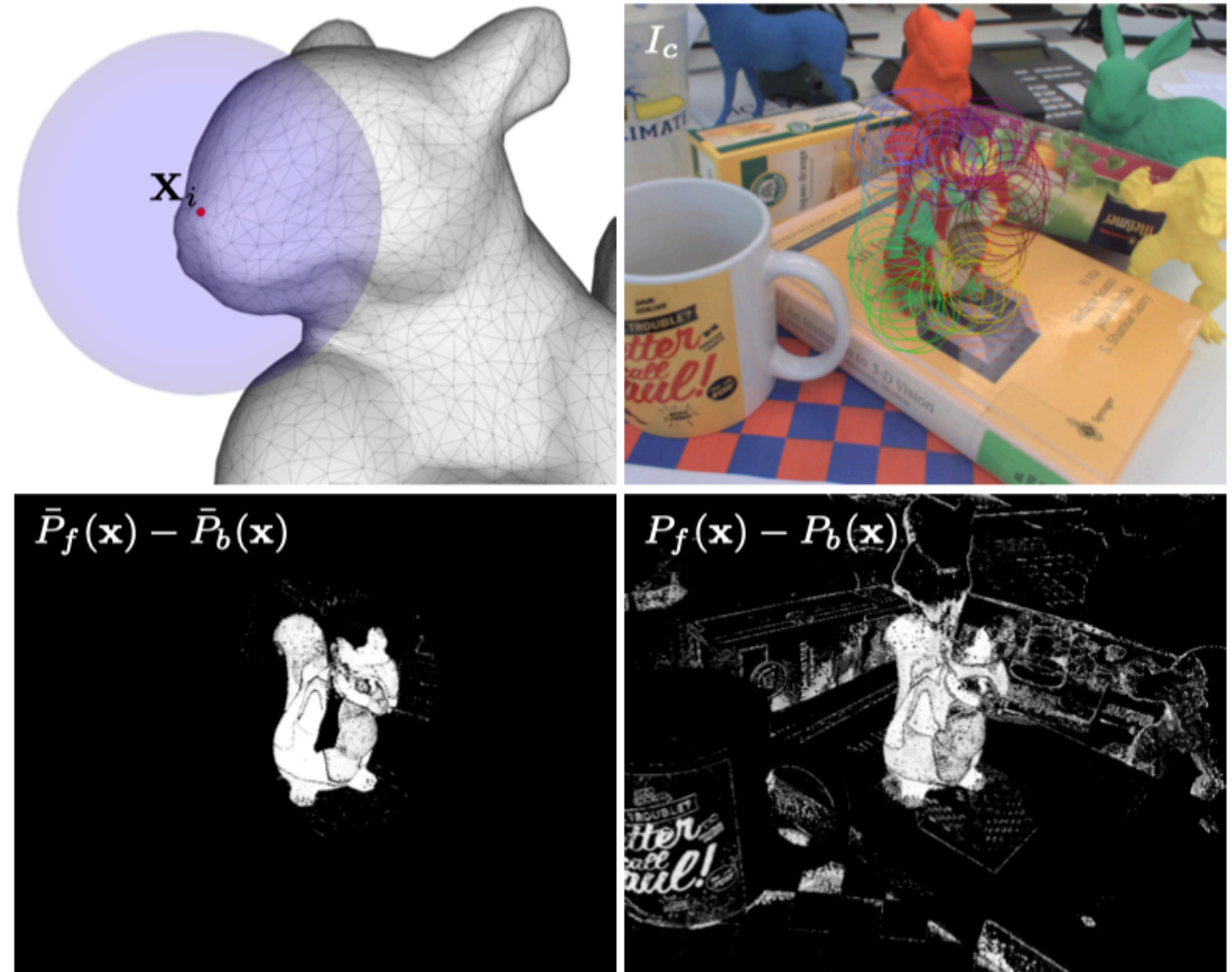
- As a **matching** problem
- Pose Tracking by keypoint tracking
  - Extract keypoint and descriptor from input images of different frames
  - Using matched keypoint to estimate relative pose between different frames





# Overview

- As a **matching** problem
- Pose Tracking by silhouette tracking:
  - Project object mesh with initial object pose and extract silhouette
  - Compute residual of current silhouette according to local foreground-background similarity
  - Optimize object pose to minimize the residual by gradient-based approach



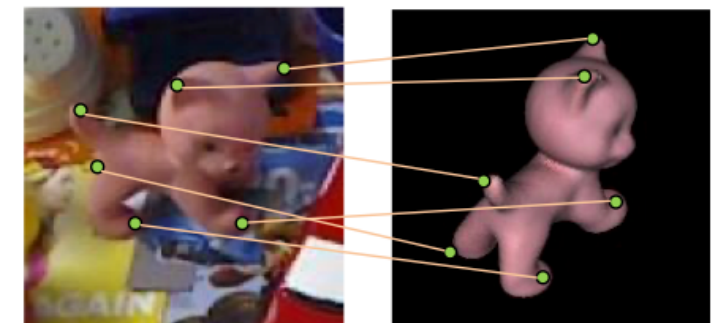
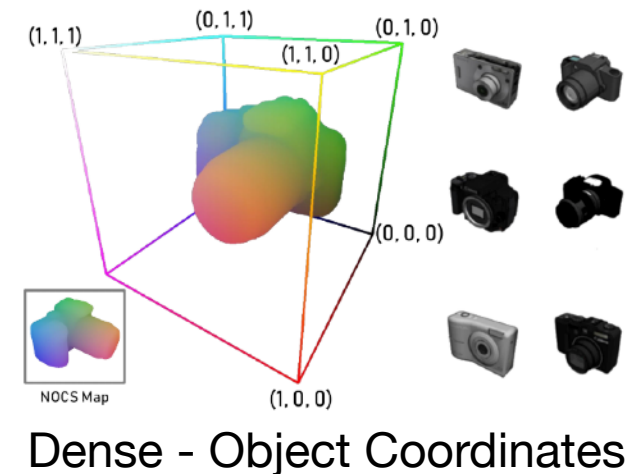
# Overview

---

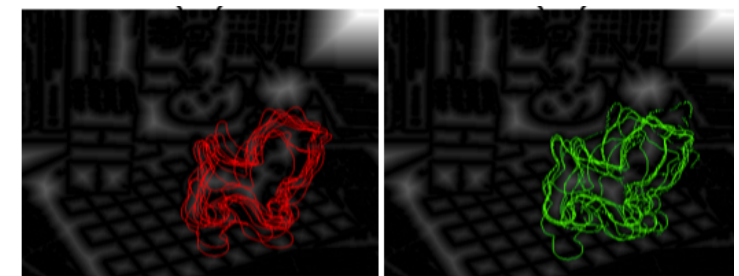
- As a **regression** problem
  - Pose Estimation: direct regression
  - Pose Tracking: render and regression
- As a **matching** problem
  - Pose Estimation: matching from image pixels to points in object frame
  - Pose Tracking: matching between frames
- Msic.
  - Tracking by Detection: Single frame estimator + filtering (smoothing)
  - Coarse-to-fine estimation: Coarse pose initialization (or sampling) + Iterative refinement

# Overview

- Pose estimation is a solved problem when:
  - Using sparse keypoint as representation
  - Given adequate well-annotated data and precise object model
- Some recent trends
  - sparse representation to denser representation
  - instance-level to category-level
  - model-based to model-free



Sparse - Key Points



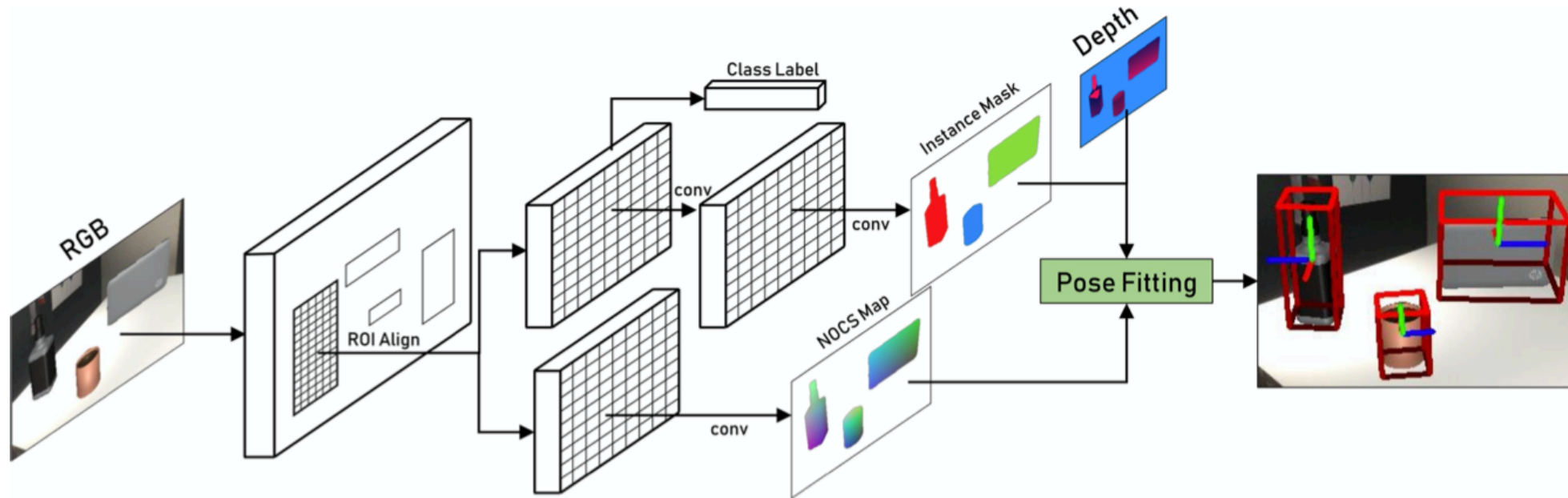
Semi-dense - Silhouette/Edges

# Overview

- Paper to cover

| Paper Name   | Conference | Model-Free                | Type           | Representation                 |
|--|------------|---------------------------|----------------|--------------------------------|
| NOCS   | CVPR19     | at inference              | category-level | NOCS                           |
| DPOD   | ICCV19     | No                        | instance-level | UV Map                         |
| Hybrid Pose  | CVPR20     | No                        | instance-level | hybrid of geometric primitives |
| 6-Pack   | ICRA20     | at training and inference | category-level | keypoints                      |
| Category Level Object Pose Estimation via Neural Analysis-by-Synthesis     | ECCV20     | at training and inference | category-level | latent vector                  |
| Shape Prior Deformation for Categorical 6D Object Pose and Size Estimation | ECCV20     | at inference              | category-level | NOCS                           |
| LatentFusion   | CVPR20     | at training and inference | unrestricted   | latent volume                  |
| Reconstruct Locally, Locally Globally                                      | CVPR20     | at training and inference | instance-level | object coordinates             |

# Pose Estimation with Various Representations

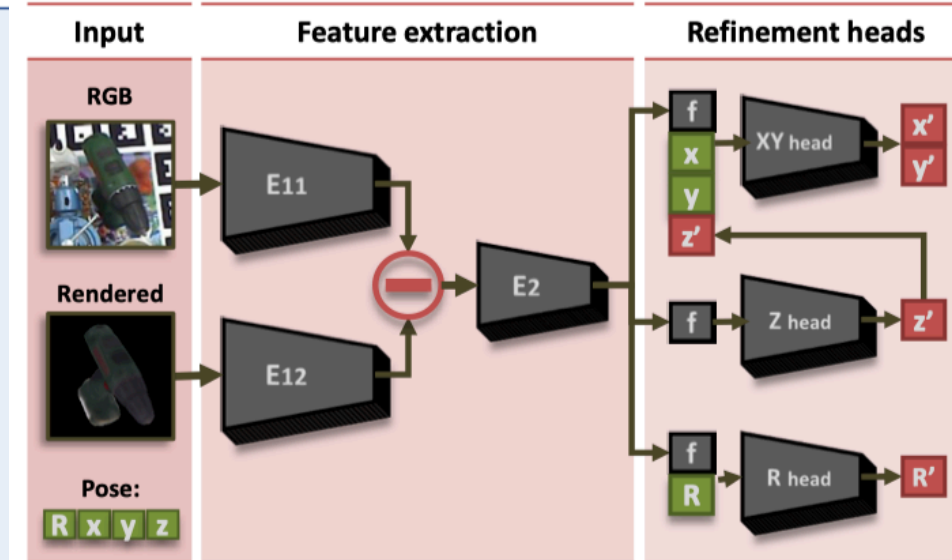
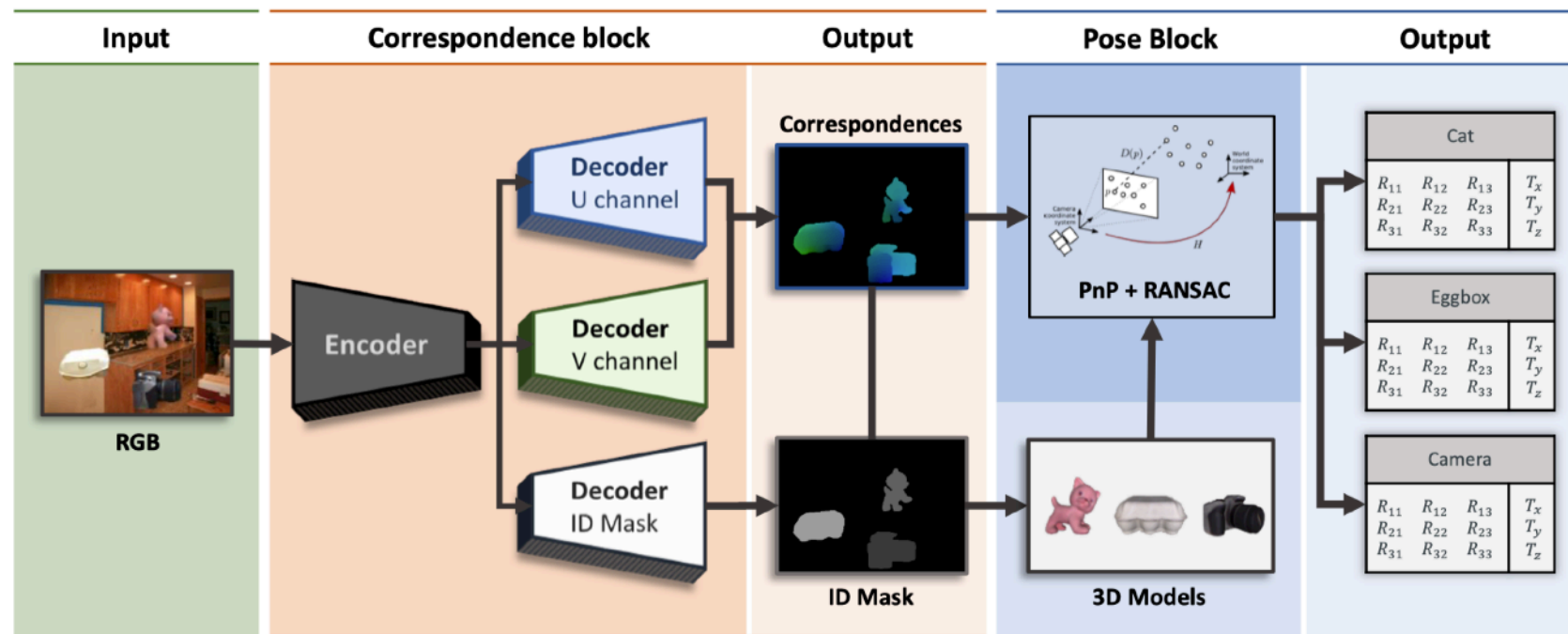


- NOCS

- Representation: Normalize Object Coordinate Space
- Approach:
  - Estimate NOCS Map, lift to original scale with depth
  - Obtain object point cloud with depth and instance mask
  - Estimate pose with transformation from object coordinate to object point cloud



# Pose Estimation with Various Representations



- DPOD
  - Representation: UV Map
  - Approach:
    - Estimate UV Map and instance mask
    - Estimate pose using RANSAC + PnP
    - Refine pose with render and compare

# Pose Estimation with Various Representations

- DPOD

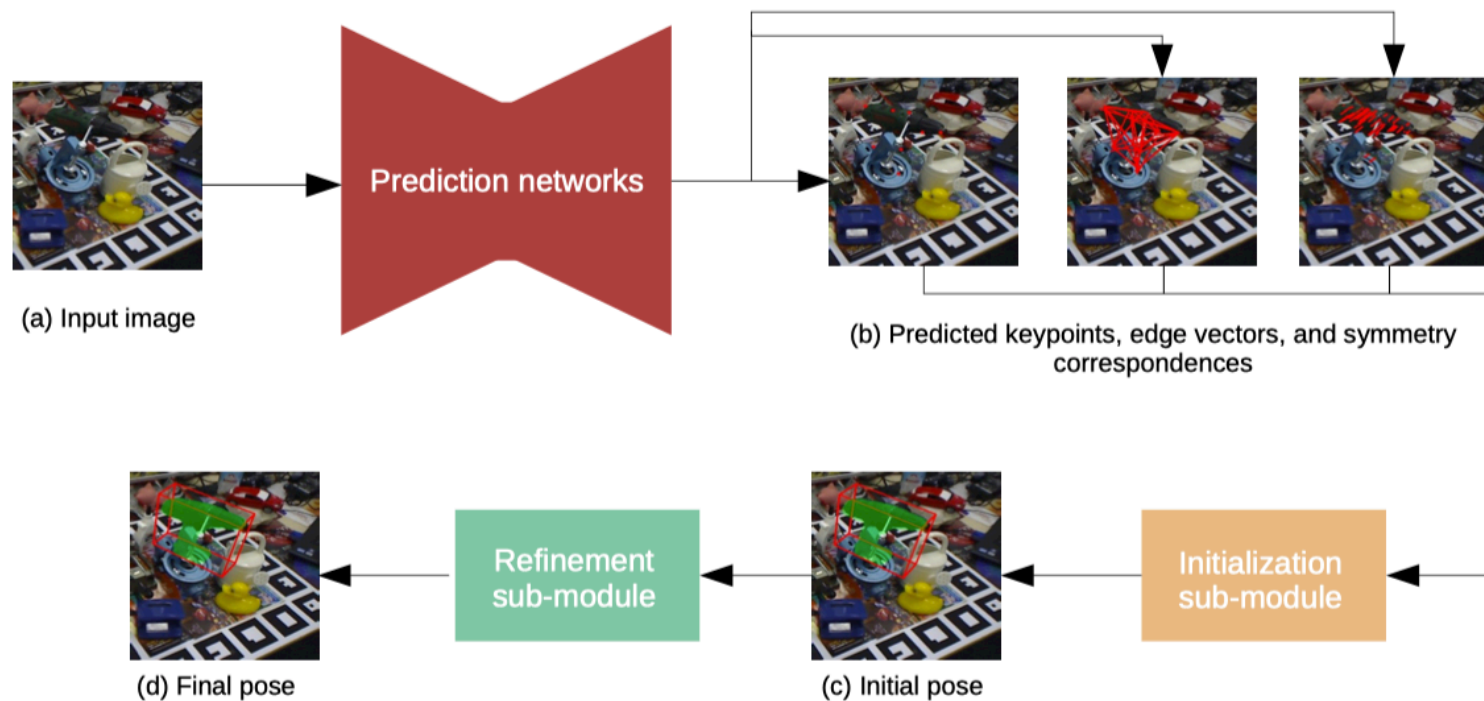
- Experiment result on LineMOD
- Discussion: UV Map VS. NOCS

| w/o Refinement |              | + Refinement |              |
|----------------|--------------|--------------|--------------|
| PVNet [25]     | Ours         | DeepIM [18]  | Ours         |
| 43.62          | <b>53.28</b> | 77.0         | <b>87.73</b> |
| <b>99.90</b>   | 95.34        | 97.5         | <b>98.45</b> |
| 86.86          | <b>90.36</b> | 93.5         | <b>96.07</b> |
| <b>95.47</b>   | 94.10        | 96.5         | <b>99.71</b> |
| <b>79.34</b>   | 60.38        | 82.1         | <b>94.71</b> |
| 96.43          | <b>97.72</b> | 95.0         | <b>98.80</b> |
| 52.58          | <b>66.01</b> | 77.7         | <b>86.29</b> |
| 99.15          | <b>99.72</b> | 97.1         | <b>99.91</b> |
| <b>95.66</b>   | 93.83        | <b>99.4</b>  | 96.82        |
| <b>81.92</b>   | 65.83        | 52.8         | <b>86.87</b> |
| 98.88          | <b>99.80</b> | 98.3         | <b>100</b>   |
| <b>99.33</b>   | 88.11        | <b>97.5</b>  | 96.84        |
| <b>92.41</b>   | 74.24        | 87.7         | <b>94.69</b> |
| <b>86.27</b>   | 82.98        | 88.6         | <b>95.15</b> |



# Pose Estimation with Various Representations

- Hybrid-Pose



- In a nut shell: estimate multiple geometric primitives with network, then optimize over them

- Geometric primitives:

- Pixel coordinate Keypoint
- 2D edges between pairwise key points
- Pairs of symmetric points

$$\bar{\mathbf{r}}_{R,t}^{\mathcal{K}}(\mathbf{p}_k) := \hat{\mathbf{p}}_k \times (R\bar{\mathbf{p}}_k + \mathbf{t}), \quad (1)$$

$$\bar{\mathbf{r}}_{R,t}^{\mathcal{E}}(\mathbf{v}_e, \mathbf{p}_{e_s}) := \hat{\mathbf{v}}_e \times (R\bar{\mathbf{p}}_{e_t} + \mathbf{t}) + \hat{\mathbf{p}}_{e_s} \times (R\bar{\mathbf{v}}_e) \quad (2)$$

$$\mathbf{r}_{R,t}^{\mathcal{S}}(\mathbf{q}_{s,1}, \mathbf{q}_{s,2}) := (\hat{\mathbf{q}}_{s,1} \times \hat{\mathbf{q}}_{s,2})^T R\bar{\mathbf{n}}_r. \quad (3)$$

Residual for initial pose estimation

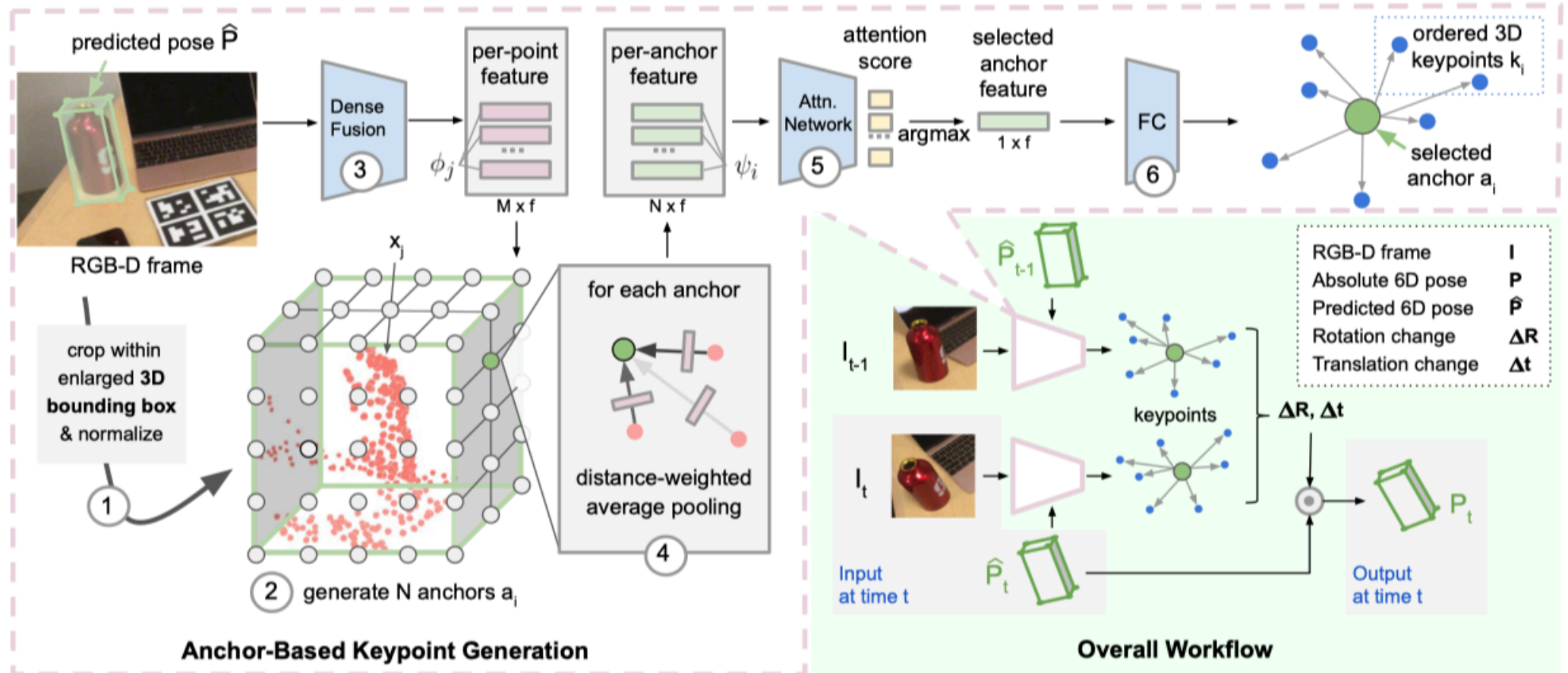
$$\begin{aligned} \min_{R,t} & \sum_{k=1}^{|\mathcal{K}|} \rho(\|\mathbf{r}_{R,t}^{\mathcal{K}}(\mathbf{p}_k)\|, \beta_{\mathcal{K}}) \|\mathbf{r}_{R,t}^{\mathcal{K}}(\mathbf{p}_k)\|_{\Sigma_k}^2 \\ & + \frac{|\mathcal{K}|}{|\mathcal{E}|} \sum_{e=1}^{|\mathcal{E}|} \rho(\|\mathbf{r}_{R,t}^{\mathcal{E}}(\mathbf{v}_e)\|, \beta_{\mathcal{E}}) \|\mathbf{r}_{R,t}^{\mathcal{E}}(\mathbf{v}_e)\|_{\Sigma_e}^2 \\ & + \frac{|\mathcal{K}|}{|\mathcal{S}|} \sum_{s=1}^{|\mathcal{S}|} \rho(r_{R,t}^{\mathcal{S}}(\mathbf{q}_{s,1}, \mathbf{q}_{s,2}), \beta_{\mathcal{S}}) \end{aligned} \quad (9)$$

Residual for pose refinement

- Different formulations are used for pose estimation and refinement
  - For pose estimation, the target is to form an  $\mathbf{Ax}=\mathbf{b}$  linear system
  - For pose refinement, the gradient-based optimization approach is used

# Category-Level Pose Estimation

## • 6-Pack



## • Approach:

- 3D anchor generation (3D grid)
- Point-wise feature aggregated to anchor feature
- Anchor scoring and 3D keypoint regression
- Intuition: coarse anchor selection + fine-grained keypoint selection to enlarge search space

$$L_{anc} = \frac{1}{N} \sum_i c_i (\|a_i - o_{gt}\|_2 - \beta)$$

anchor  
scoring loss

$$L_{mvc} = \frac{1}{K} \sum_i \|k_i^t - [\Delta R_t^{gt} | \Delta t_t^{gt}] \cdot k_i^{t-1}\|$$

multi-view  
consistency  
loss

$$L_{tra} = \|(\bar{k}^t - \bar{k}^{t-1}) - \Delta t_t^{gt}\|$$

pose  
estimation loss

$$L_{rot} = 2 \arcsin\left(\frac{1}{2\sqrt{2}} \|\Delta \hat{R}_t - \Delta R_t^{gt}\|\right)$$

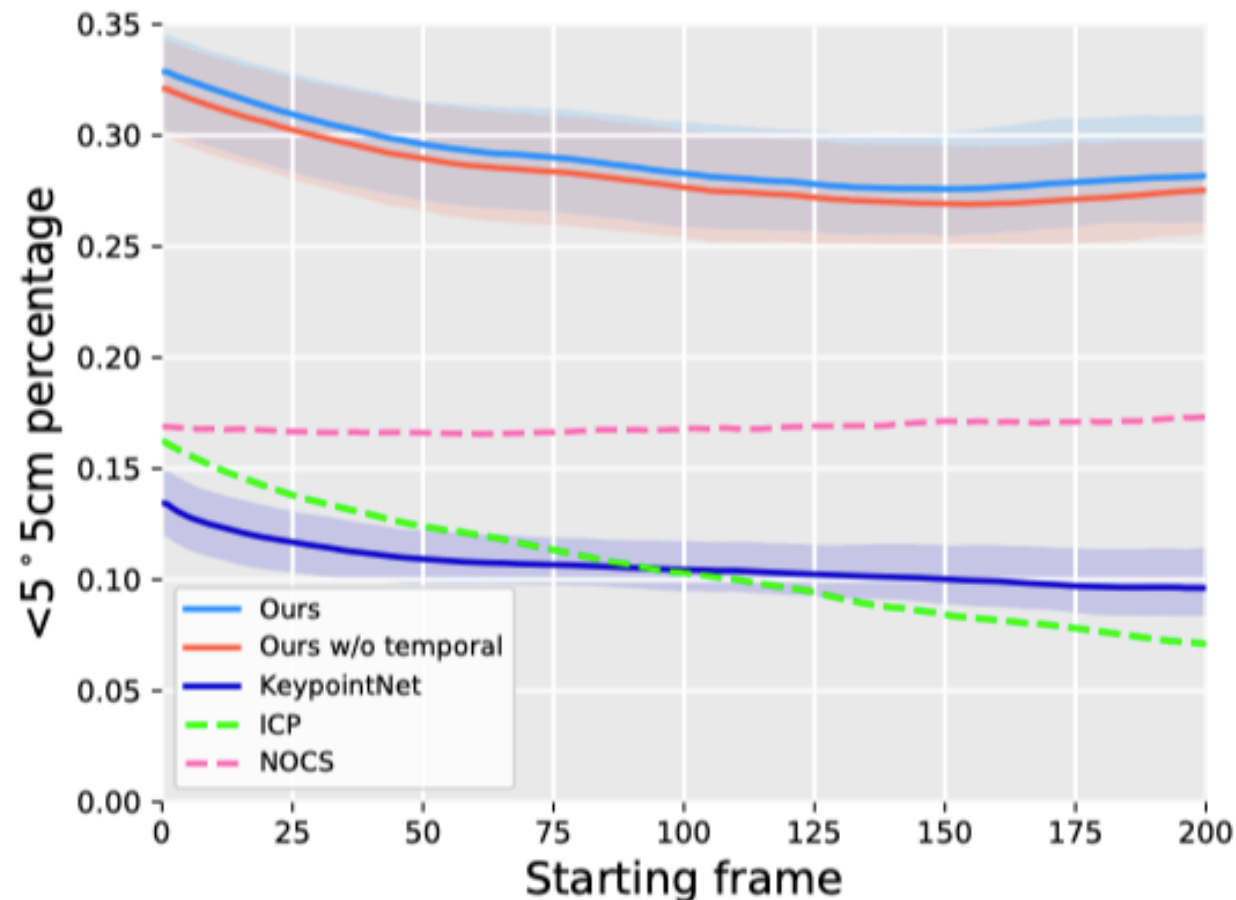
# Category-Level Pose Estimation

- 6-Pack

- Results on NOCS-REAL275

- Accurate and stable

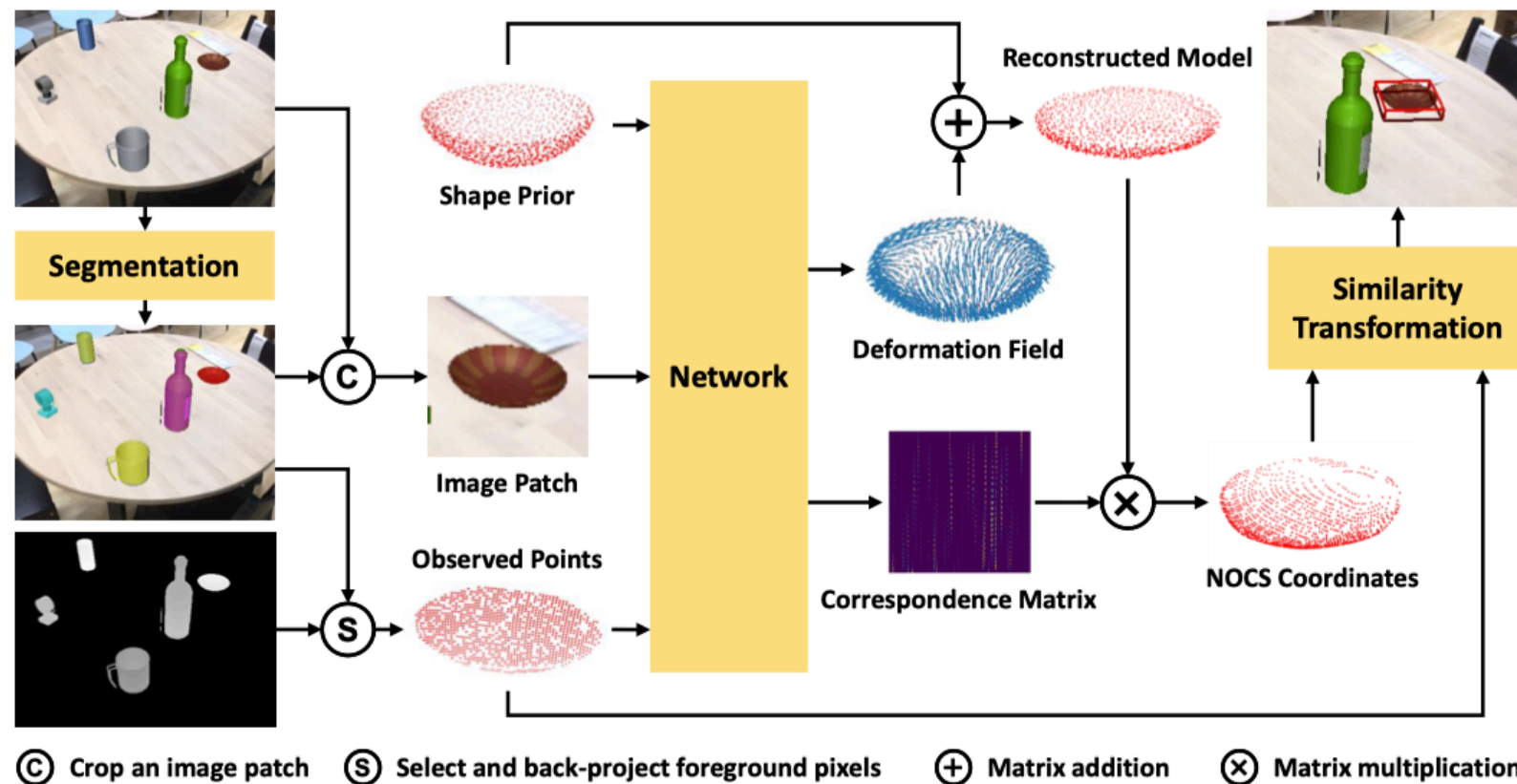
- Run at 10 HZ with GTX1070



|         |                  | NOCS<br>[46] | ICP<br>[50] | Keypoint<br>Net [41] | Ours w/o<br>temporal | Ours         |
|---------|------------------|--------------|-------------|----------------------|----------------------|--------------|
| bottle  | 5°5cm            | 5.5          | 10.1        | 5.9                  | 23.7                 | <b>24.5</b>  |
|         | IoU25            | 48.7         | 29.9        | 23.1                 | <b>92.0</b>          | 91.1         |
|         | R <sub>err</sub> | 25.6         | 48.0        | 28.5                 | 15.7                 | <b>15.6</b>  |
|         | T <sub>err</sub> | 14.4         | 15.7        | 9.5                  | 4.2                  | <b>4.0</b>   |
| bowl    | 5°5cm            | <b>62.2</b>  | 40.3        | 16.8                 | 53.0                 | 55.0         |
|         | IoU25            | 99.6         | 79.7        | 74.7                 | <b>100.0</b>         | <b>100.0</b> |
|         | R <sub>err</sub> | <b>4.7</b>   | 19.0        | 9.8                  | 5.3                  | 5.2          |
|         | T <sub>err</sub> | <b>1.2</b>   | 4.7         | 8.2                  | 1.6                  | 1.7          |
| camera  | 5°5cm            | 0.6          | <b>12.6</b> | 1.8                  | 8.4                  | 10.1         |
|         | IoU25            | 90.6         | 53.1        | 30.9                 | <b>91.0</b>          | 87.6         |
|         | R <sub>err</sub> | <b>33.8</b>  | 80.5        | 45.2                 | 43.9                 | 35.7         |
|         | T <sub>err</sub> | <b>3.1</b>   | 12.2        | 8.5                  | 5.5                  | 5.6          |
| can     | 5°5cm            | 7.1          | 17.2        | 4.3                  | <b>25.0</b>          | 22.6         |
|         | IoU25            | 77.0         | 40.5        | 42.6                 | 89.9                 | <b>92.6</b>  |
|         | R <sub>err</sub> | 16.9         | 47.1        | 28.8                 | <b>12.5</b>          | 13.9         |
|         | T <sub>err</sub> | <b>4.0</b>   | 9.4         | 13.1                 | 5.0                  | 4.8          |
| laptop  | 5°5cm            | 25.5         | 14.8        | 49.2                 | 62.4                 | <b>63.5</b>  |
|         | IoU25            | 94.7         | 50.9        | 94.6                 | 97.8                 | <b>98.1</b>  |
|         | R <sub>err</sub> | 8.6          | 37.7        | 6.5                  | 4.9                  | <b>4.7</b>   |
|         | T <sub>err</sub> | <b>2.4</b>   | 9.2         | 4.4                  | 2.5                  | 2.5          |
| mug     | 5°5cm            | 0.9          | 6.2         | 3.1                  | 22.4                 | <b>24.1</b>  |
|         | IoU25            | 82.8         | 27.7        | 52.0                 | <b>100.0</b>         | 95.2         |
|         | R <sub>err</sub> | 31.5         | 56.3        | 61.2                 | <b>20.3</b>          | 21.3         |
|         | T <sub>err</sub> | 4.0          | 9.2         | 6.7                  | <b>1.8</b>           | 2.3          |
| Overall | 5°5cm            | 17.0         | 16.9        | 13.5                 | 32.5                 | <b>33.3</b>  |
|         | IoU25            | 82.2         | 47.0        | 53.0                 | <b>95.1</b>          | 94.2         |
|         | R <sub>err</sub> | 20.2         | 48.1        | 30.0                 | 17.1                 | <b>16.0</b>  |
|         | T <sub>err</sub> | 4.9          | 10.5        | 8.4                  | <b>3.4</b>           | 3.5          |

# Category-Level Pose Estimation

- Shape Prior Deformation for Categorical 6D Object Pose and Size Estimation

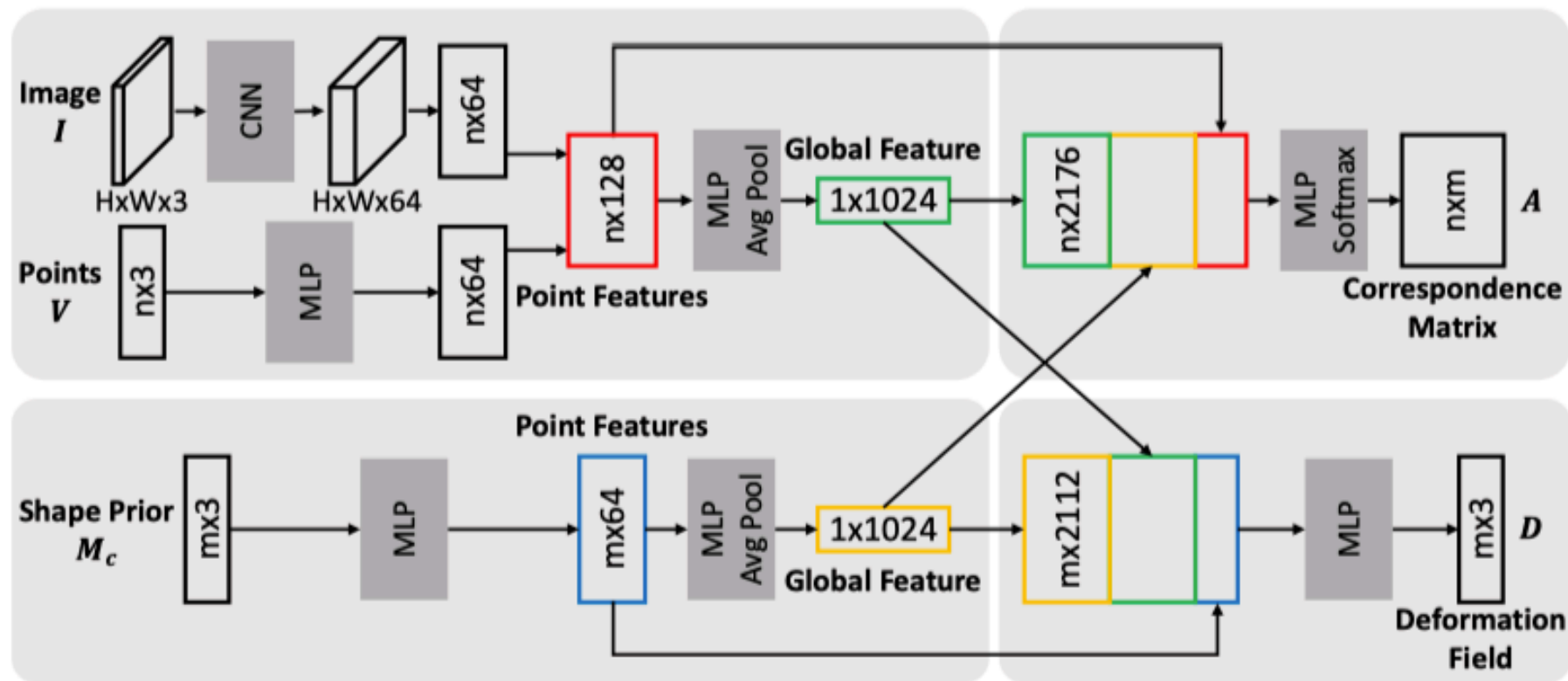


- In a nut shell: estimate NOCS with explicit category-level shape prior
  - Shape prior: mean (normalized) model point cloud by decoding mean latent vector for each category
  - Fuse feature from model point cloud, generated deformation field and correspondence.



# Category-Level Pose Estimation

- Shape Prior Deformation for Categorical 6D Object Pose and Size Estimation



- Loss

- For autoencoder & deformation field: reconstruction loss (Chamfer distance)

$$d_{CD}(M_c^i, \hat{M}_c^i) = \sum_{x \in M_c^i} \min_{y \in \hat{M}_c^i} \|x - y\|_2^2 + \sum_{y \in \hat{M}_c^i} \min_{x \in M_c^i} \|x - y\|_2^2.$$

- For  $A$ : correspondence loss (smooth L1 with gt NOCS)

$$L_{\text{corr}}(P, P_{gt}) = \frac{1}{N_v} \sum_{\mathbf{x} \in P} \sum_{i=1,2,3} \begin{cases} 5(x_i - y_i)^2, & \text{if } |x_i - y_i| \leq 0.1 \\ |x_i - y_i| - 0.05, & \text{otherwise} \end{cases},$$

# Category-Level Pose Estimation

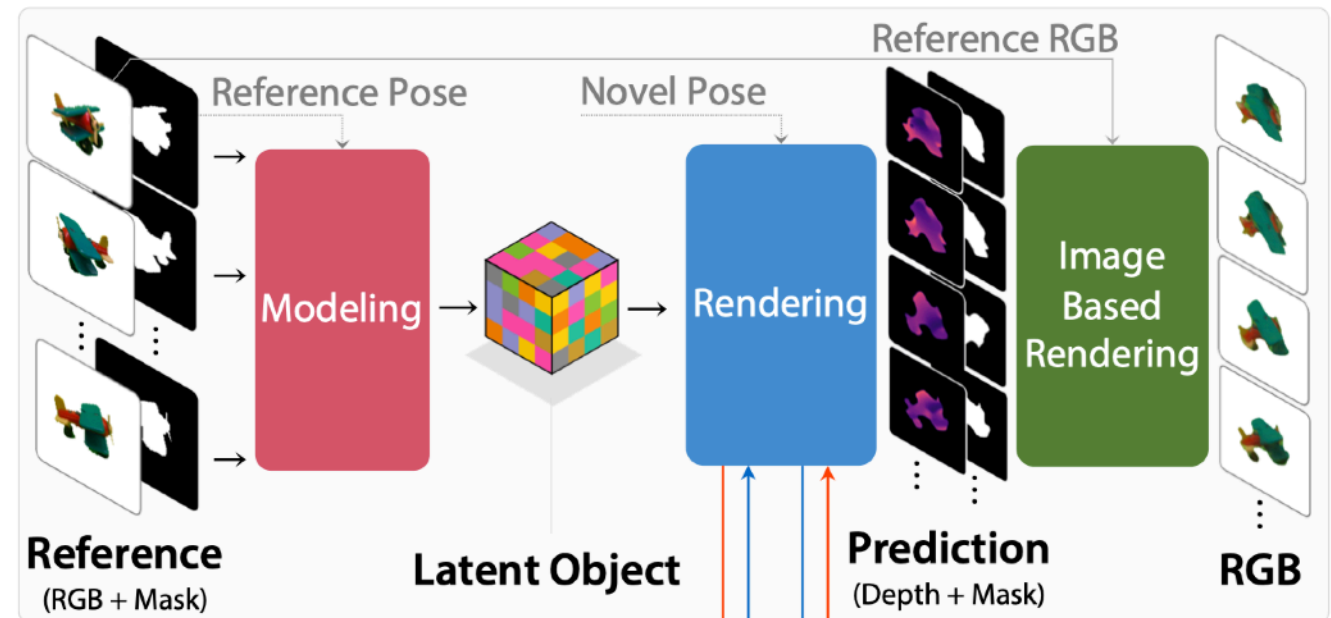
- Shape Prior Deformation for Categorical 6D Object Pose and Size Estimation
- Experiment on NOCS dataset

| Data     | Method        | mAP              |                  |             |             |             |             |
|----------|---------------|------------------|------------------|-------------|-------------|-------------|-------------|
|          |               | 3D <sub>50</sub> | 3D <sub>75</sub> | 5° 2cm      | 5° 5cm      | 10° 2cm     | 10° 5cm     |
| CAMERA25 | Baseline [33] | 83.9             | 69.5             | 32.3        | 40.9        | 48.2        | 64.6        |
|          | Ours (RGB)    | 93.1             | <b>84.6</b>      | 50.2        | 54.5        | 70.4        | 78.6        |
|          | Ours (RGB-D)  | <b>93.2</b>      | 83.1             | <b>54.3</b> | <b>59.0</b> | <b>73.3</b> | <b>81.5</b> |
| REAL275  | Baseline [33] | <b>78.0</b>      | 30.1             | 7.2         | 10.0        | 13.8        | 25.2        |
|          | Ours (RGB)    | 75.2             | 46.5             | 15.7        | 18.8        | 33.7        | 47.4        |
|          | Ours (RGB-D)  | 77.3             | <b>53.2</b>      | <b>19.3</b> | <b>21.4</b> | <b>43.2</b> | <b>54.1</b> |

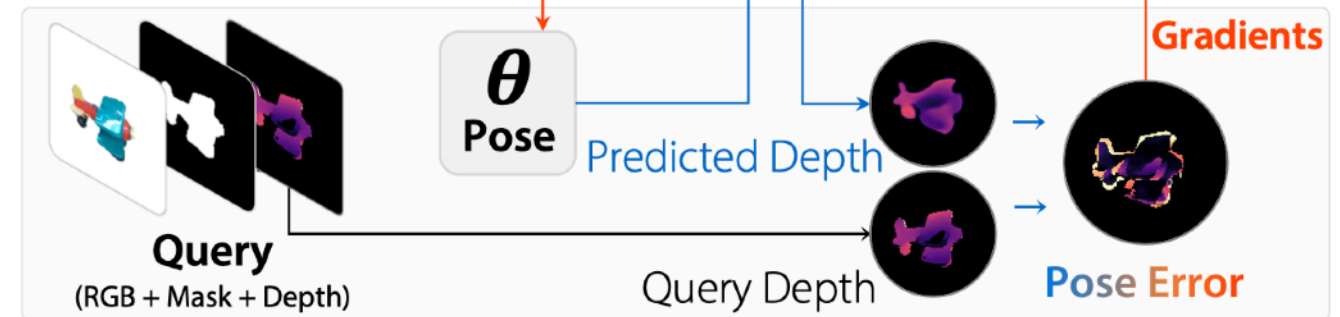
# Model-Free Pose Estimation

- Latent Fusion
  - Reconstruct with reference frames
    - Estimate 2D feature, lift to 3D feature voxels
    - Feature aggregation across views
  - Inference with optimization
    - Estimate initial translation, sample rotation
    - Refine pose by optimization over depth error and latent error

## 1. Modeling and Rendering



## 2. Pose Estimation





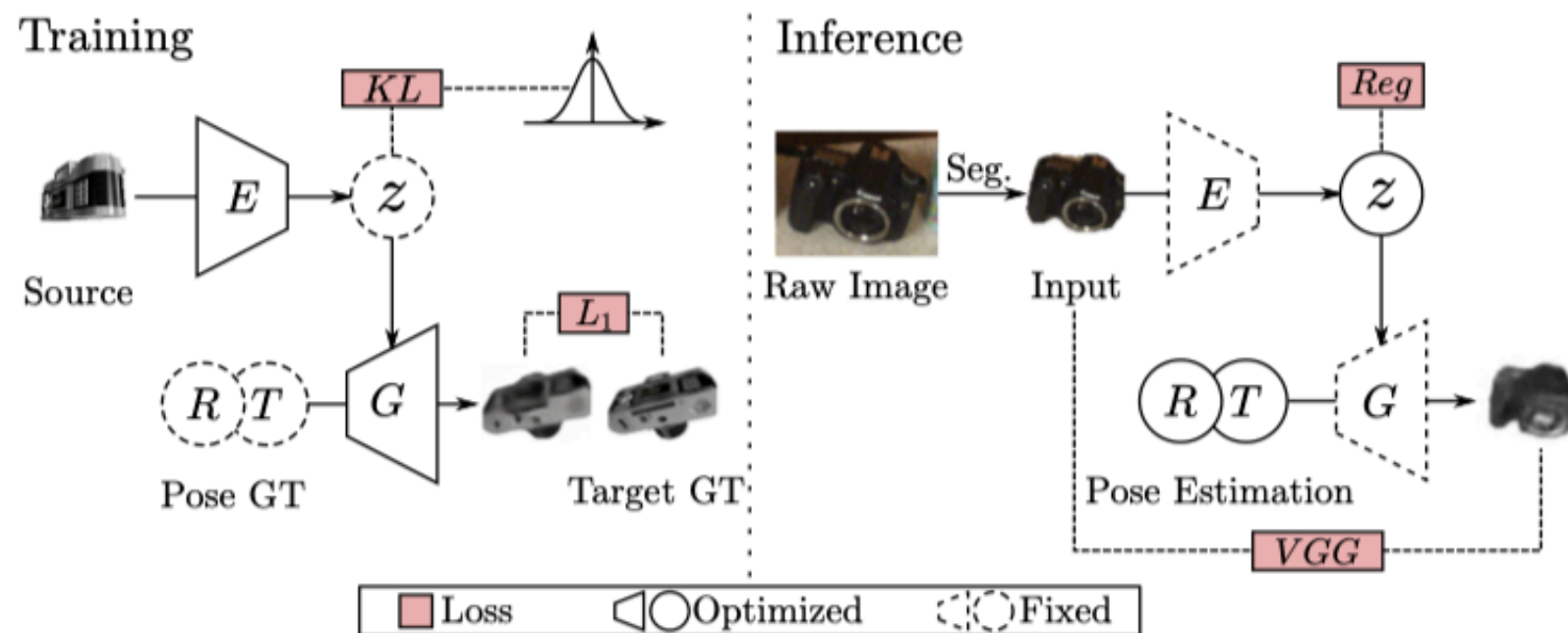
# Model-Free Pose Estimation

- Latent Fusion
- Experiment on MOPED dataset
- Comparable results with model-based approach

|               | PoseRBPF [6]     |              |              | IBR-LD   |              |              |
|---------------|------------------|--------------|--------------|--|--------------|--------------|
| Input         | Textured 3D Mesh |              |              |  |              |              |
| Training      | Yes              |              |              |  |              |              |
| # Networks    | Per-Object       |              |              |  |              |              |
| Pose Loss     | -                |              |              | $\mathcal{L}_{\text{latent}} + \mathcal{L}_{\text{depth}}$ |              |              |
|               | ADD              | ADD-S        | Proj.2D      | ADD  | ADD-S        | Proj.2D      |
| black_drill   | 59.78            | <b>82.94</b> | 49.80        | 56.67  | 79.06        | 53.77        |
| cheezit       | 57.78            | 82.45        | 48.47        | <b>61.31</b>   | <b>91.63</b> | <b>55.24</b> |
| duplo_dude    | 56.91            | 82.14        | 47.11        | 74.02  | 89.55        | 52.49        |
| duster        | <b>58.91</b>     | 82.78        | <b>46.66</b> | 49.13  | 91.56        | 19.33        |
| graphics_card | 59.13            | 83.20        | 49.85        | <b>80.71</b>   | <b>91.25</b> | <b>67.71</b> |
| orange_drill  | <b>58.23</b>     | 82.68        | <b>49.08</b> | 51.84  | 70.95        | 46.12        |
| pouch         | 57.74            | 82.16        | 49.01        | <b>60.43</b>   | <b>89.60</b> | <b>49.80</b> |
| remote        | 56.87            | 82.04        | <b>48.06</b> | 55.38  | 94.80        | 37.73        |
| rinse_aid     | 57.74            | 82.53        | <b>48.13</b> | 65.63  | 92.58        | 28.61        |
| toy_plane     | <b>62.41</b>     | 85.10        | 49.81        | 60.18  | <b>90.24</b> | <b>51.70</b> |
| vim_mug       | <b>58.09</b>     | <b>82.38</b> | <b>48.08</b> | 30.11  | 80.76        | 14.38        |
| mean          | 58.51            | 82.76        | <b>48.55</b> | 58.67  | 87.45        | 43.35        |

# Model-Free Pose Estimation

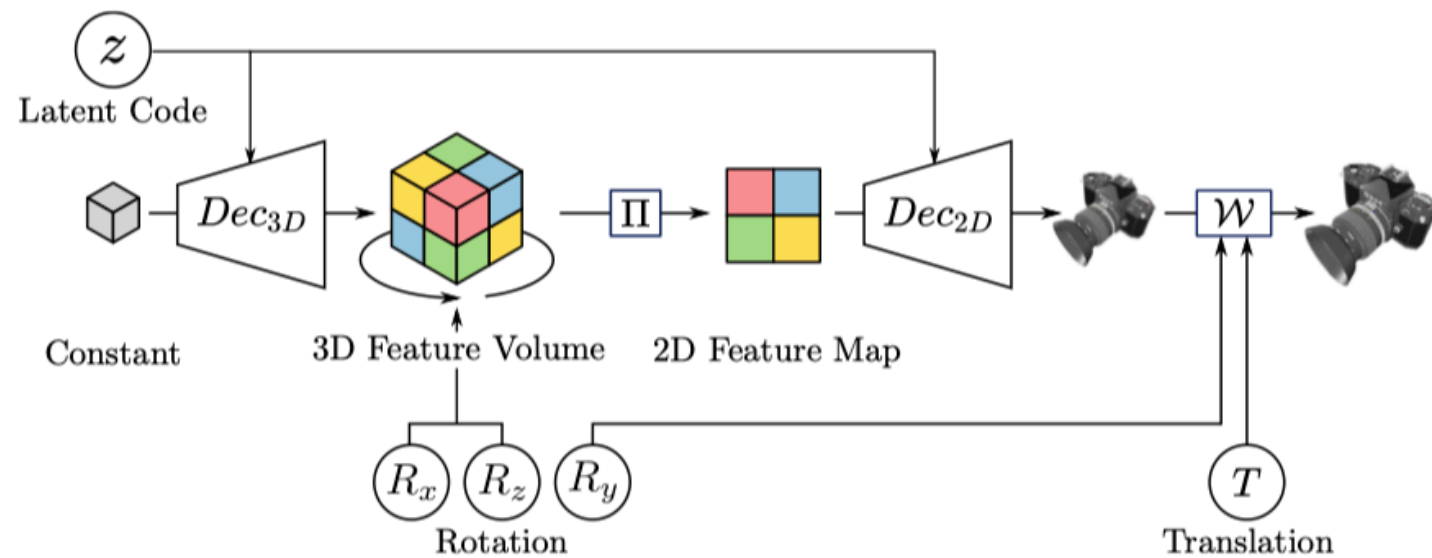
- Category Level Object Pose Estimation via Neural Analysis-by-Synthesis



- Approach:
  - Generate synthetic image given object pose and latent code
  - Compare feature metric error between synthesized and observed images to optimize object pose and object shape (latent code)

# Model-Free Pose Estimation

- Category Level Object Pose Estimation via Neural Analysis-by-Synthesis



$$\mathcal{W}(T, R_z) : \begin{bmatrix} u \\ v \end{bmatrix} \mapsto \frac{f}{t_z} \cdot \left( R_z \begin{bmatrix} u \\ v \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \end{bmatrix} \right)$$

$$\hat{I} = G(R, T, z) = \mathcal{W}(T, R_z) \circ G_{3D}(R_x, R_y, z)$$

- Details:

- Train generator using **synthetic data only**, use some tricks to reduce the need of network capacity

- Initial pose are randomly sampled...  $E(I, R, T, z) = \|F_{vgg}(I) - F_{vgg}(\hat{I})\|_2 + \|z\|_2,$

- Optimize with deep feature (pre-trained VGG) + regularization on latent code

# Model-Free Pose Estimation

- Category Level Object Pose Estimation via Neural Analysis-by-Synthesis

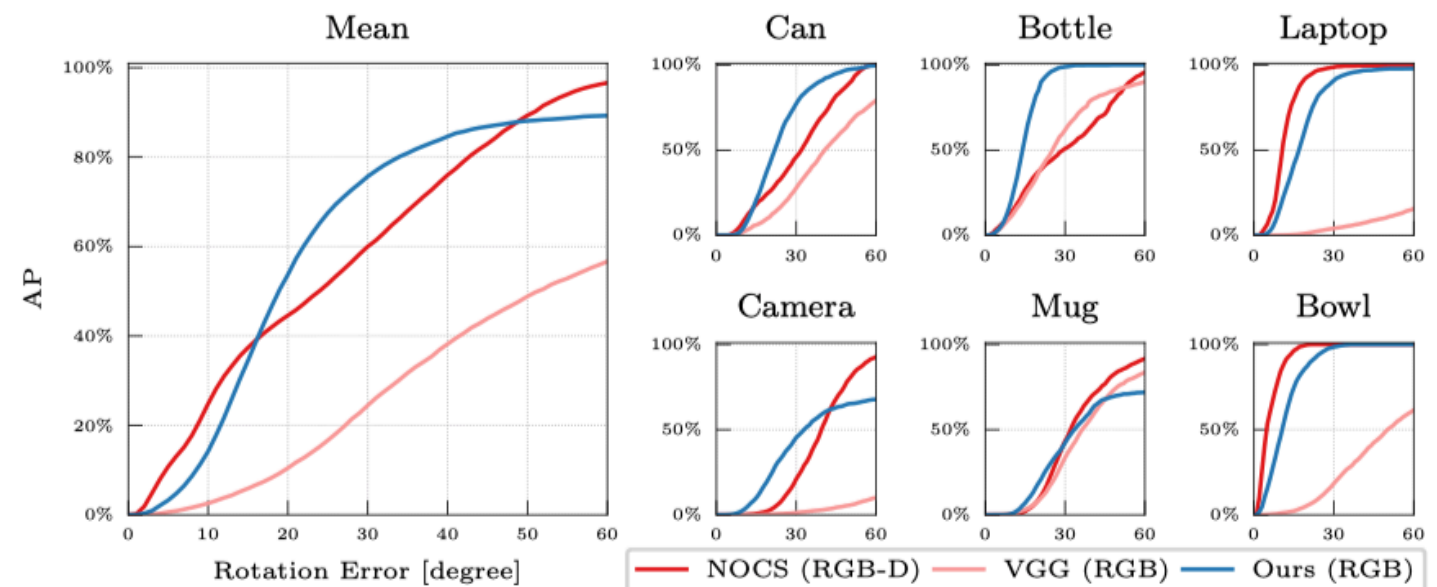
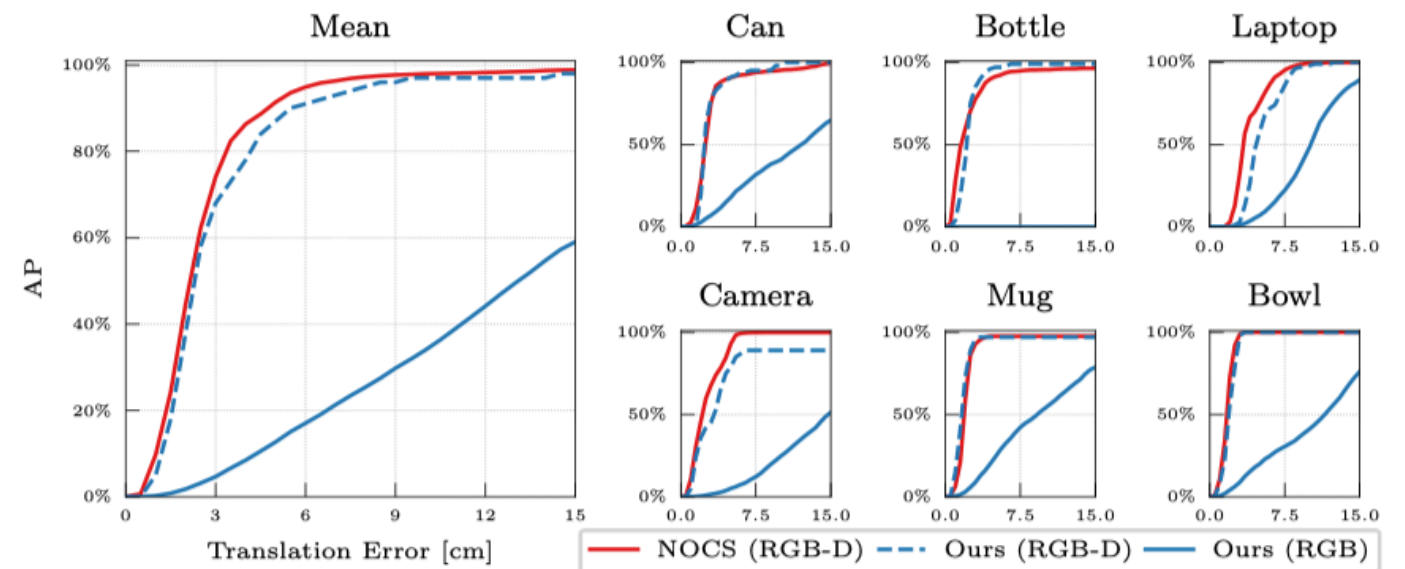
- Experiment on NOCS dataset

- Comparable to NOCS when using RGB-D input

- Does not require pose annotation

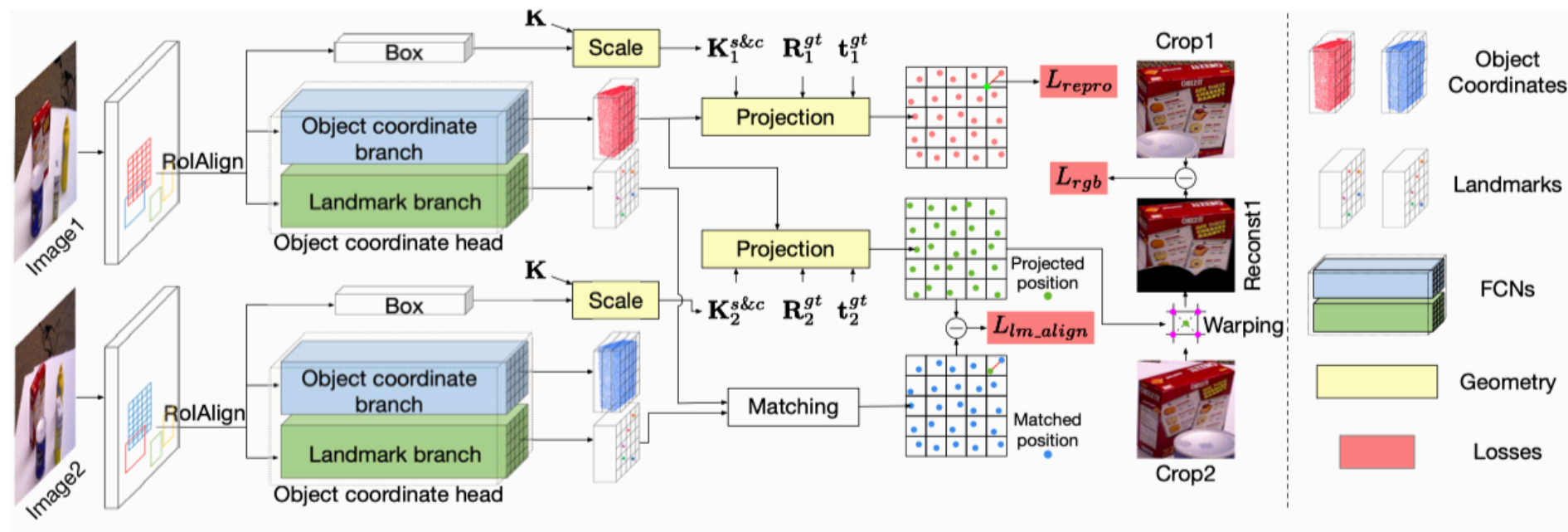
- Compare with *LatentFusion*

- Pros: No reference frames needed
- Cons: can it extend to arbitrary unseen object?



# Model-Free Pose Estimation

- Reconstruct Locally, Localize Globally

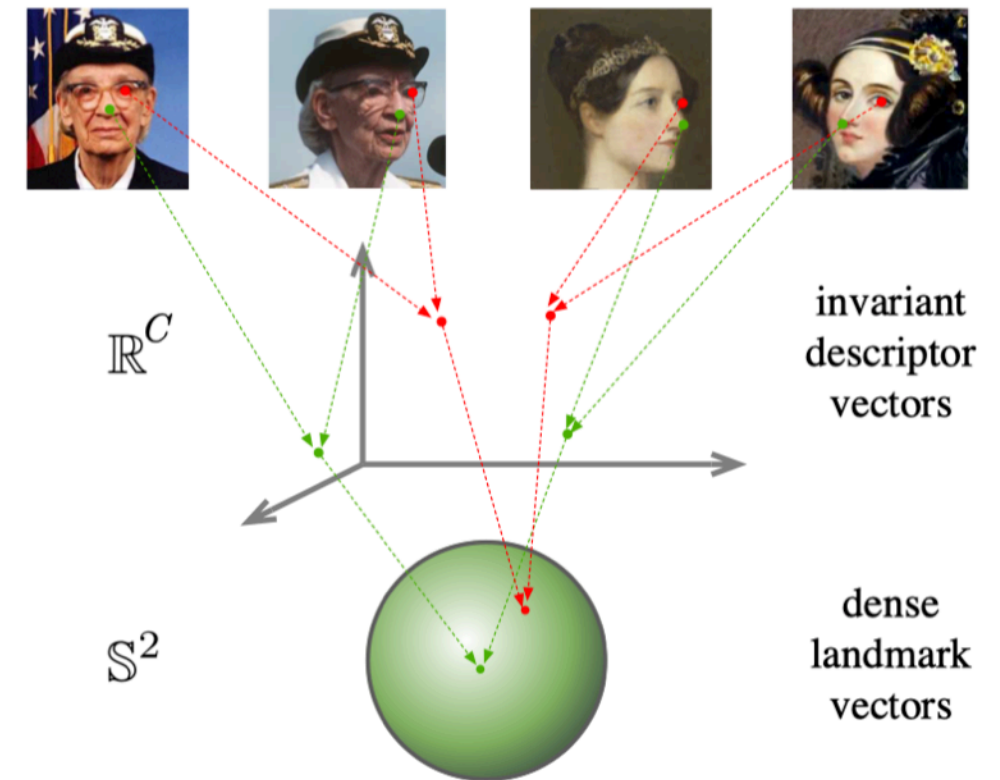


- Overall target: learn object coordinate estimation without CAD model
- Approach:
  - predict mask, object coordinates and landmark
  - Single-frame + cross frame consistency supervision



# Model-Free Pose Estimation

- Reconstruct Locally, Locally Globally
  - What is a *landmark*?
    - Descriptor with fewer channels.
    - Descriptors that gain ‘robustness’ to intra-class variations
  - In this work: make no difference with using descriptor...
  - May be used to extend the work to category-level (?)



marks [45, 9, 42]. A dense descriptor associates to each image pixel a  $C$ -dimensional vector, whereas a dense landmark detector associates to each pixel a 2D vector, which is the index of the landmark in a  $(u, v)$  parameterisation of the object surface. Thus we can interpret a landmark as a tiny 2D descriptor. Due to its small dimensionality, a landmark loses the ability to encode instance-specific details of the appearance, but gains robustness to intra-class variations.

# Model-Free Pose Estimation

- Reconstruct Locally, Locally Globally
- Results: Comparable with SoTA approaches

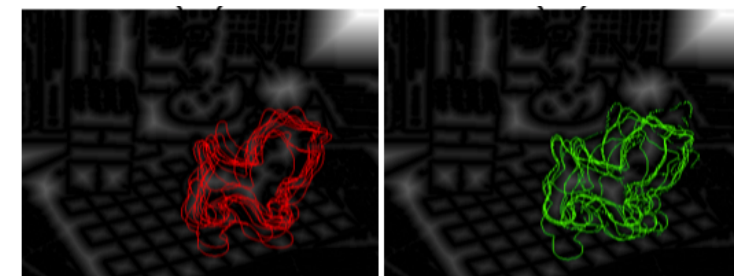
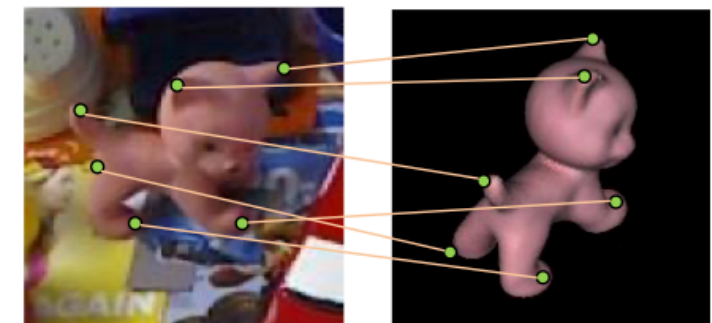
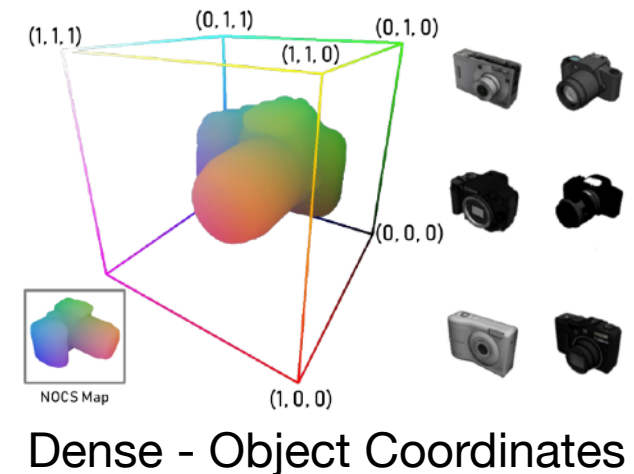
| method      | w/ CAD model |             |                        |               |                        |                          |                       |                       | w/o CAD model  |                |       |
|-------------|--------------|-------------|------------------------|---------------|------------------------|--------------------------|-----------------------|-----------------------|----------------|----------------|-------|
|             | BB8<br>[41]  | BB8<br>w/ r | SSD-6D<br>w/ r<br>[22] | Tekin<br>[48] | DeepIM<br>w/ r<br>[27] | Dense-<br>Fusion<br>[52] | Pix2-<br>Pose<br>[38] | PVNet<br>w/ r<br>[40] | SSD-6D<br>[22] | LieNet<br>[11] | Ours  |
| ape         | 27.9         | 40.4        | 65                     | 21.62         | 77.0                   | 92                       | 58.1                  | 43.62                 | 0.00           | 38.8           | 52.91 |
| benchwise   | 62.0         | 91.8        | 80                     | 81.80         | 97.5                   | 93                       | 91.0                  | 99.90                 | 0.18           | 71.2           | 96.51 |
| cam         | 40.1         | 55.7        | 78                     | 36.57         | 93.5                   | 94                       | 60.0                  | 86.86                 | 0.41           | 52.5           | 87.84 |
| can         | 48.1         | 64.1        | 86                     | 68.80         | 96.5                   | 93                       | 84.4                  | 95.47                 | 1.35           | 86.1           | 86.81 |
| cat         | 45.2         | 62.6        | 70                     | 41.82         | 82.1                   | 97                       | 65.0                  | 79.34                 | 0.51           | 66.2           | 67.30 |
| driller     | 58.6         | 74.4        | 73                     | 63.51         | 95.0                   | 87                       | 76.3                  | 96.43                 | 2.58           | 82.3           | 88.70 |
| duck        | 32.8         | 44.3        | 66                     | 27.23         | 77.7                   | 92                       | 43.8                  | 52.58                 | 0.00           | 32.5           | 54.74 |
| eggbox*     | 40.0         | 57.8        | 100                    | 69.58         | 97.1                   | 100                      | 96.8                  | 99.15                 | 8.90           | 79.4           | 94.74 |
| glue*       | 27.0         | 41.2        | 100                    | 80.02         | 99.4                   | 100                      | 79.4                  | 95.66                 | 0.00           | 63.7           | 91.98 |
| holepuncher | 42.4         | 67.2        | 49                     | 42.63         | 52.8                   | 92                       | 74.8                  | 81.92                 | 0.30           | 56.4           | 75.41 |
| iron        | 67.0         | 84.7        | 78                     | 74.97         | 98.3                   | 97                       | 83.4                  | 98.88                 | 8.86           | 65.1           | 94.59 |
| lamp        | 39.9         | 76.5        | 73                     | 71.11         | 97.5                   | 95                       | 82.0                  | 99.33                 | 8.20           | 89.4           | 96.64 |
| phone       | 35.2         | 54.0        | 79                     | 47.74         | 87.7                   | 93                       | 45.0                  | 92.41                 | 0.18           | 65.0           | 89.24 |
| average     | 43.6         | 62.7        | 79                     | 55.95         | 88.6                   | 94                       | 72.4                  | 86.27                 | 2.42           | 65.2           | 82.88 |

Table 2. **LineMOD: Percentages of correct pose estimates in ADD-10.** \* denotes that the object is symmetric and is evaluated in ADD-S. w/r means the pose is refined with 3D model.



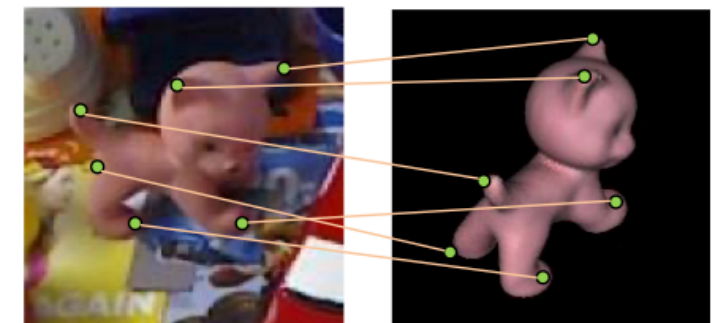
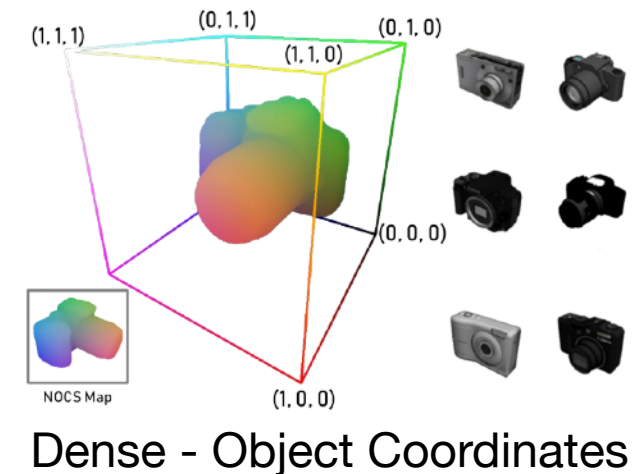
# Discussion

- Solved problem
  - Pose estimation with sparse keypoint set
  - Instance-level pose estimation
  - Given accurate CAD model and pose annotation
- Some recent trends
  - sparse representation to denser representation
  - instance-level to category level
  - model-based to model-free

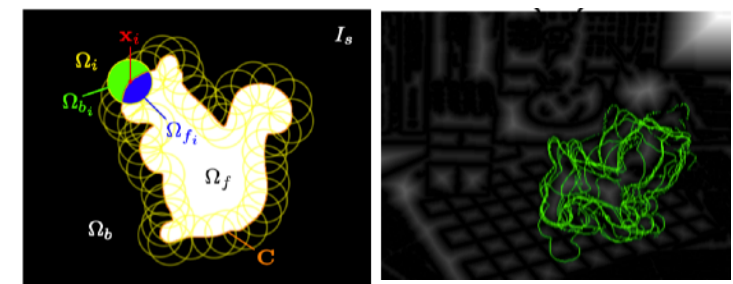


# Discussion

- Optimal representation of an object?
  - Preferred properties:
    - Available on weakly-textured object
    - Generalizable beyond instance-level
    - Available without accurate geometric models
    - Trackable across time
  - Potential solution:
    - Sparse keypoint
    - Dense coordinate map
    - Silhouette + appearance cue
    - Latent representation
    - Hybrid of geometric/appearance primitives



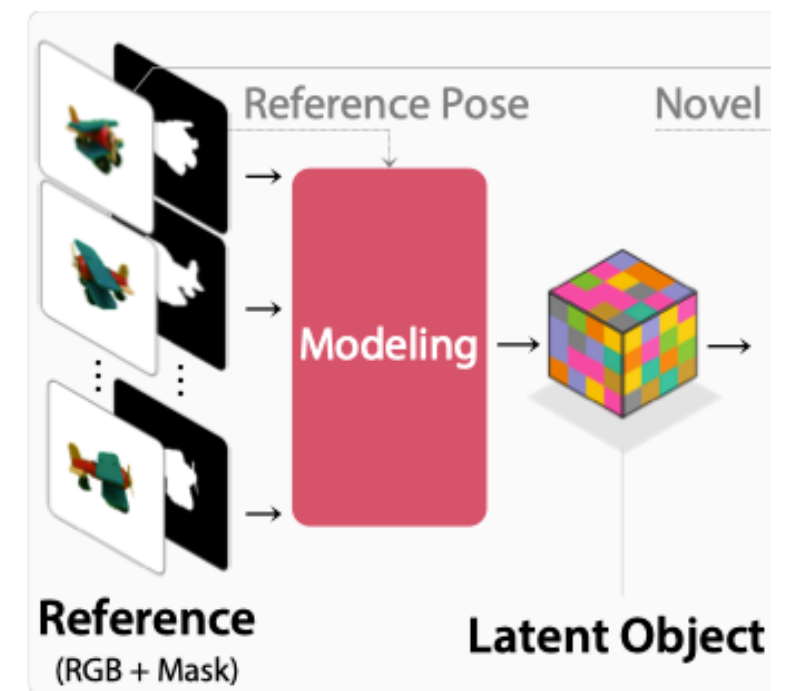
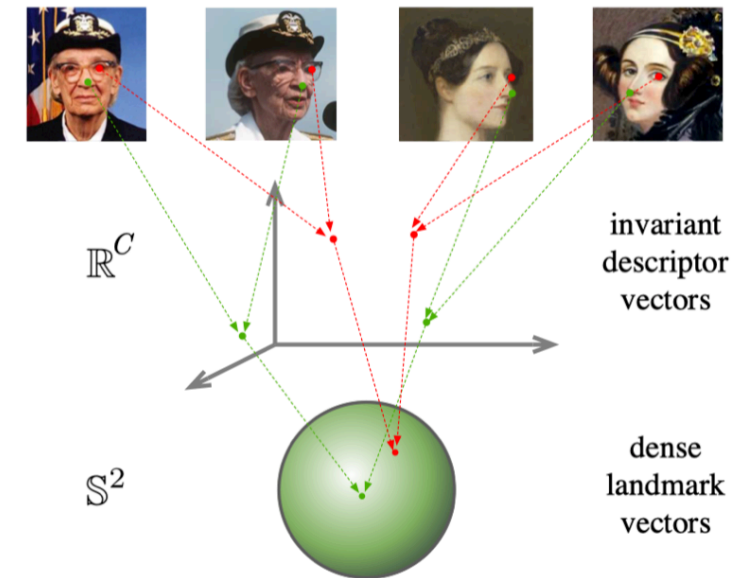
Sparse - Keypoints



Hybrid - Silhouette + TCLC Hist.

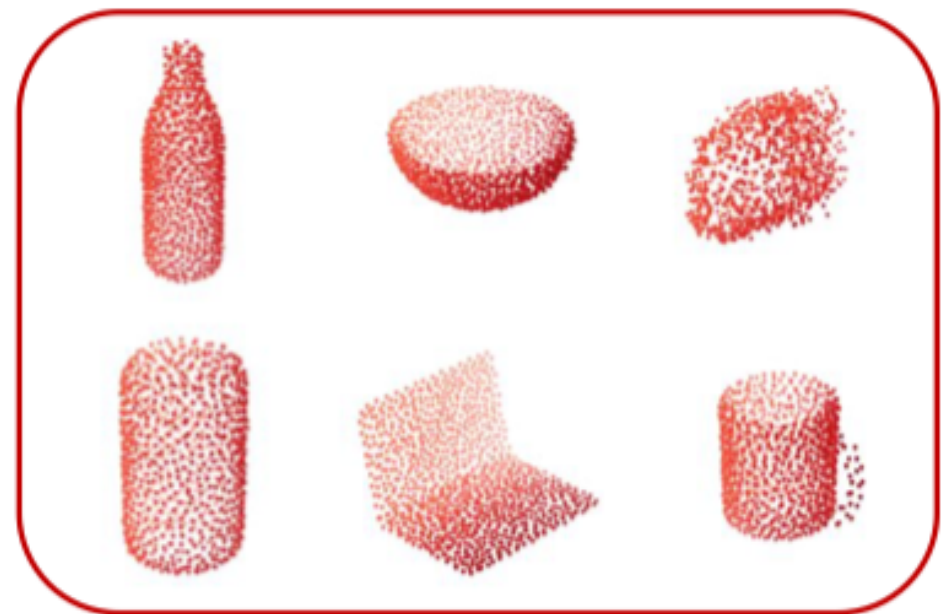
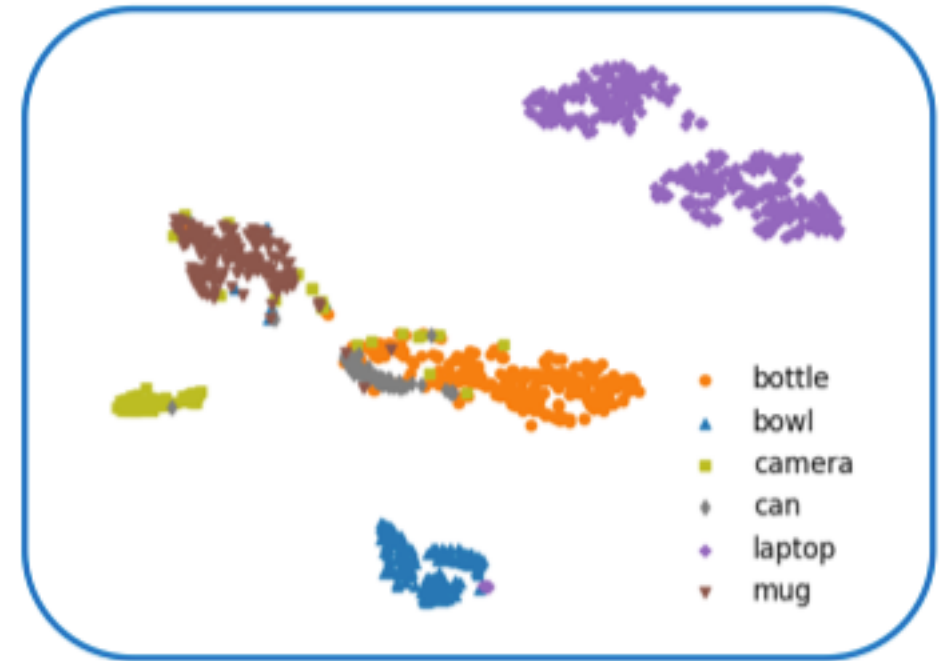
# Discussion

- Optimal representation of an object?
- How to achieve model-free pose estimation?
  - Learn to reconstruct geometry without accurate model
  - Neural synthesis to generate RGB image (and depth) for later optimization



# Discussion

- Optimal representation of an object?
- How to achieve model-free pose estimation?
- How to achieve category-level pose estimation?
  - With intra-category shape prior, either implicit (encoded in network) or explicit (mean shape)
  - Generalizable neural reconstruction



# From academia to industry

---

- Needs for 6DoF Pose Estimation in real application
  - Accuracy VS. Stability
  - Data VS. Algorithm
  - Scalability intra (or even inter) categories

# From academia to industry

---

- Hierarchy of problem to solve
  - Instance-level object 6DoF pose estimation in varied scenes and on varied sensor without fine-tune and adaptation
    - more of an engineering problem of scalable data collection
  - Category-level (with unseen instance) object 6DoF pose estimation on some common categories
    - worth working as academic problem
  - Estimate the pose of arbitrary unseen object
    - zero-shot learning, final target.. but hard if not impossible



# The Solution by Google Media Pipe

- MobilePose

- Approach

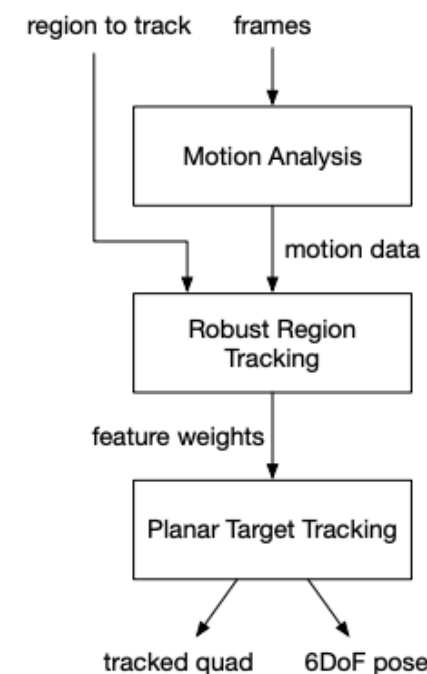
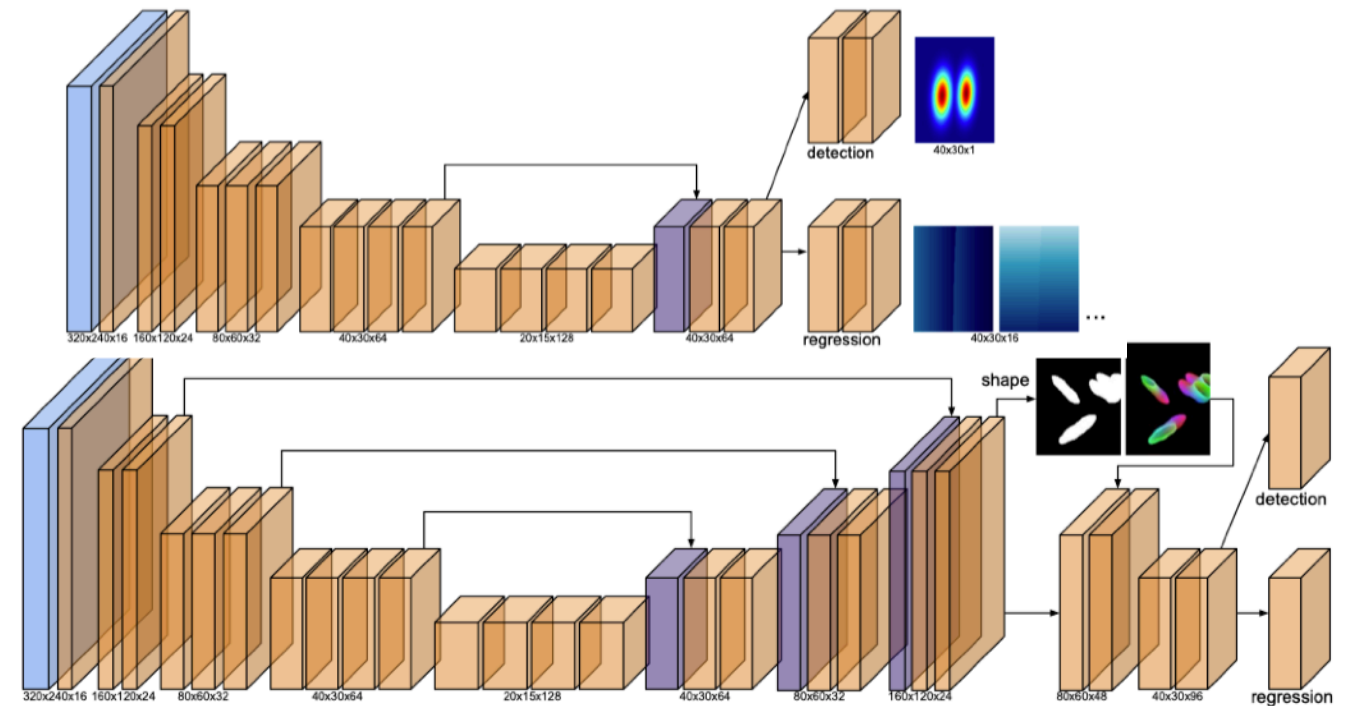
- Joint 2D detection and regression (for 3D box corners and center)
    - Estimate mask and object coordinate if available to augment feature

- Data pipeline

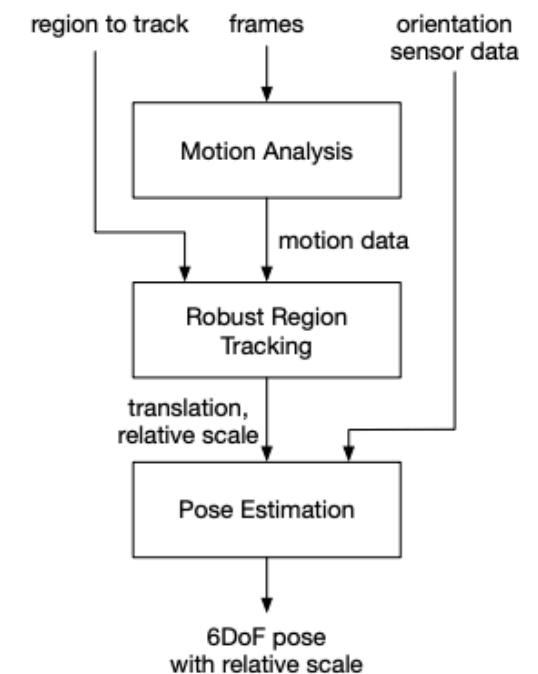
- Hand-annotated for the first frame
    - Propagate along the sequence using camera pose from ARCore

- Instant motion tracking

- Track the 9 keypoints by motion analysis



(a) Planar target tracking



(b) 6DoF tracking





# Thanks for your Attention

**Siyu ZHANG**

Research Engineer

ZJU-SenseTime Joint Lab of 3D Vision

[zhangsiyu1@sensetime.com](mailto:zhangsiyu1@sensetime.com)