# SA6D: Self-Adaptive Few-Shot 6D Pose Estimator for Novel and Occluded Objects

Ning Gao[1,2]   Ngo Anh Vien[1]   Hanna Ziesche[1]   Gerhard Neumann[2]

[1]Bosch Center for Artificial Intelligence   [2]Autonomous Learning Robots, KIT
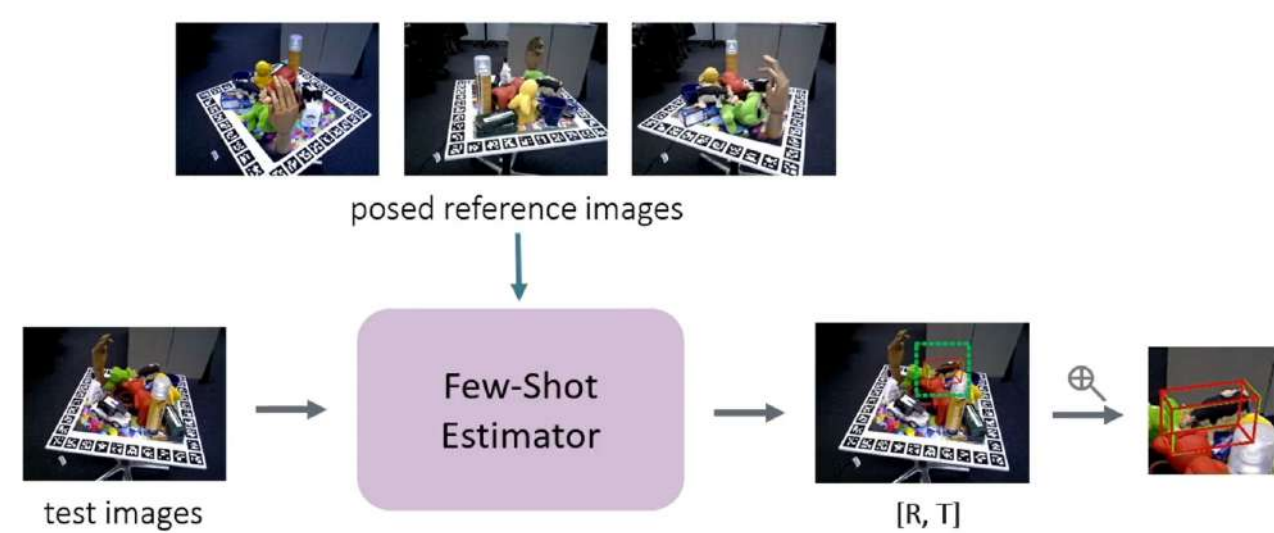
## ABSTRACT

**Goal:**
Exploit generic few-shot 6D pose estimation of novel objects from unseen categories, especially under severe clutter.

**Contributions:**
- SA6D increases the performance and robustness against heavy occlusions while *not* requiring any object information (3D model, object diameter or ground-truth mask) or object-centric images, while only requiring a small set of posed RGB-D reference images with known poses of the novel object.
- SA6D employs an online-adaptive segmentation module to identify the target object during inference.
- SA6D utilizes pretrained models from prior work *without* any retraining processes.
- SA6D significantly outperforms current state-of-the-art methods against occlusion in real-world scenarios while trained entirely on synthetic data.
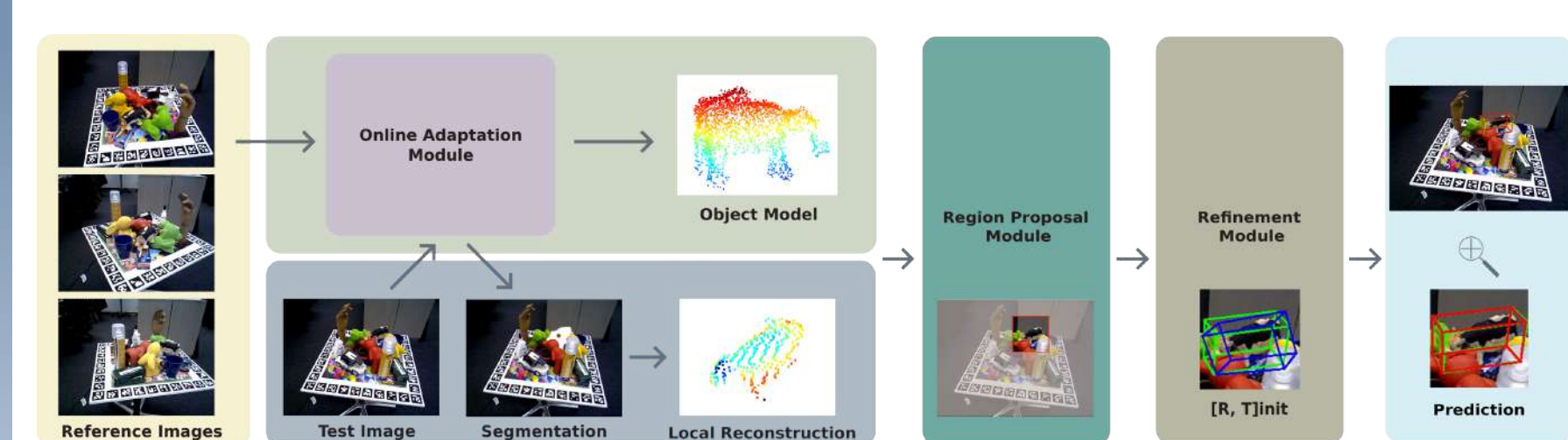
## TASK DESIGN

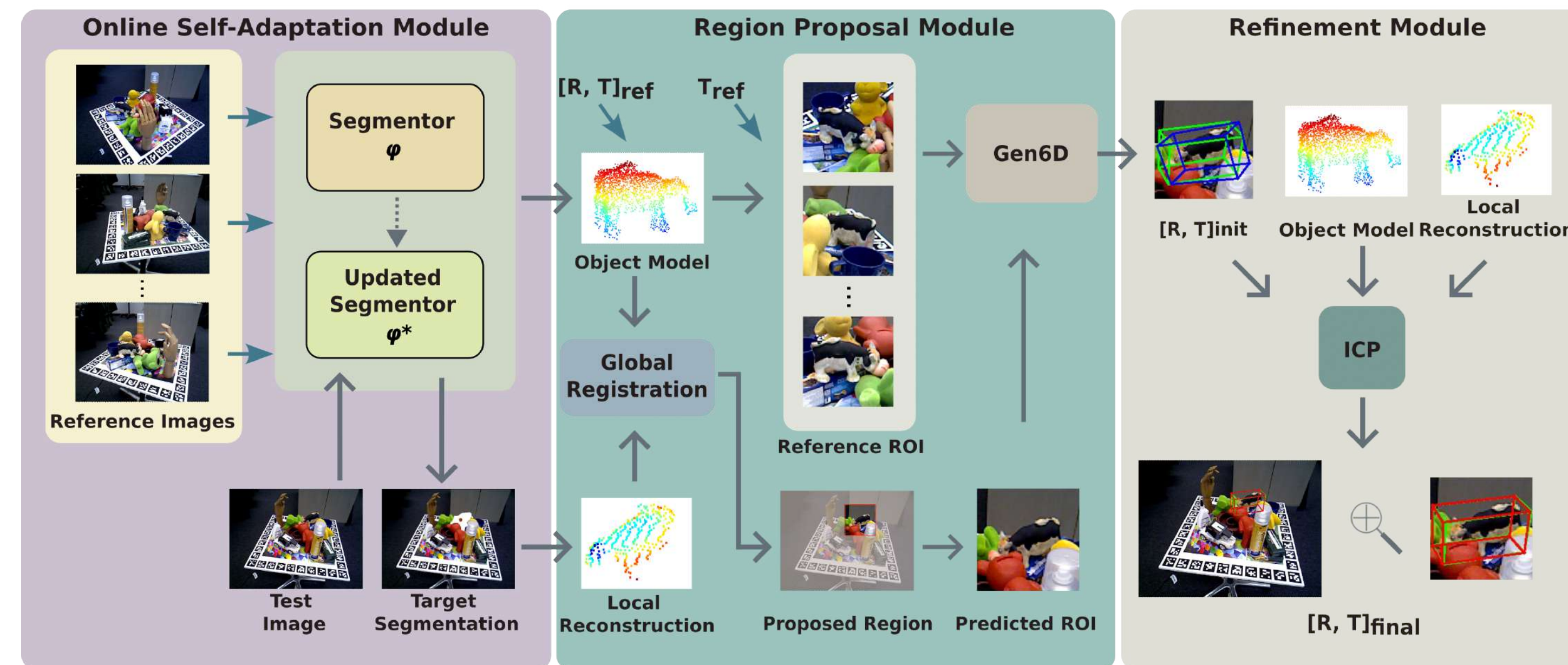Few-Shot 6D Novel Object Pose Estimation:



Predict 6D pose of a novel object in a new image (test image) from a few reference images with the known pose of the object (reference images).
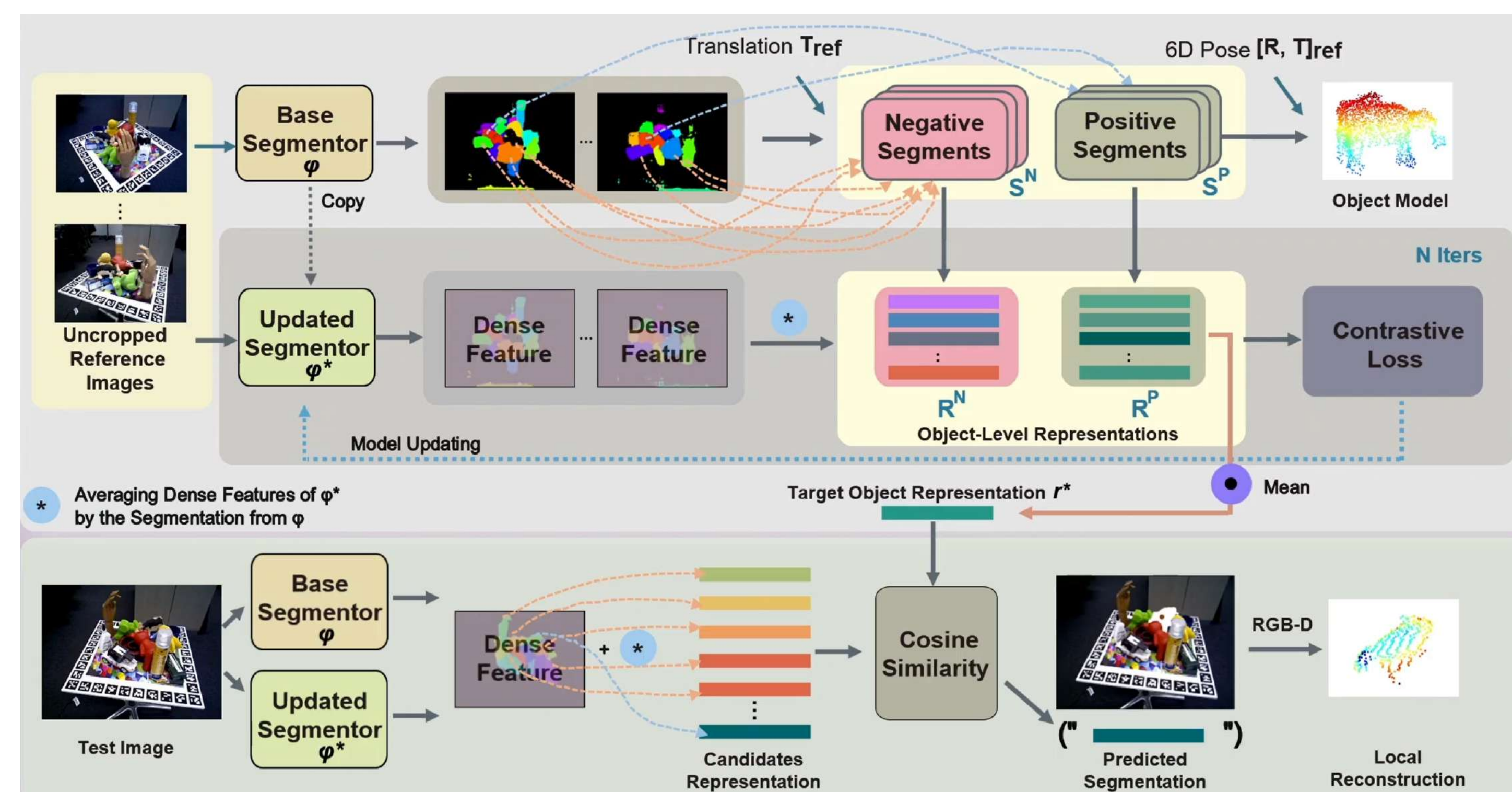
## OVERVIEW



We present a generalizable and category-agnostic few-shot 6D object pose estimator using a small number of posed RGB-D images as references. Compared to existing methods, our approach provides robust and accurate predictions of novel objects against occlusions without requiring retraining or any object information.

## PIPELINE



**Overview.** SA6D includes three modules: i) The *online self-adaptation module* discovers and segments the target object (*milk cow*) from a cluttered scene giving a few posed RGB-D images as reference. ii) The *region proposal module* outputs a robust region of interest (ROI) of the target object against occlusion by incorporating visual and geometric features. A coarse 6D pose is then estimated by Gen6D and iii) further fine-tuned using ICP.

**Onlie Adaptation Segmentation:**



**Online self-adaptation module**. A pretrained segmentor $\varphi$ is applied on reference images to predict segmentations. With the ground-truth translation of the target object in the reference images $T_{ref}$, the object center can be reprojected to the image. For each reference image, one segment is chosen as a positive sample if it includes the reprojected object center while the remaining segments are considered as negative samples. Subsequently, an object-level representation of each segment is computed by averaging the pixel-wise dense features from $\varphi*$. A contrastive loss is then applied over the positive and negative object representations and updates $\varphi*$ iteratively. After adaptation, $\varphi*$ generates the target object representation $r*$ by averaging over all positive representations from reference images. Given a test image, we compute the cosine similarity between each candidate and $r*$ and the most similar candidate is chosen as the segment of the target object.
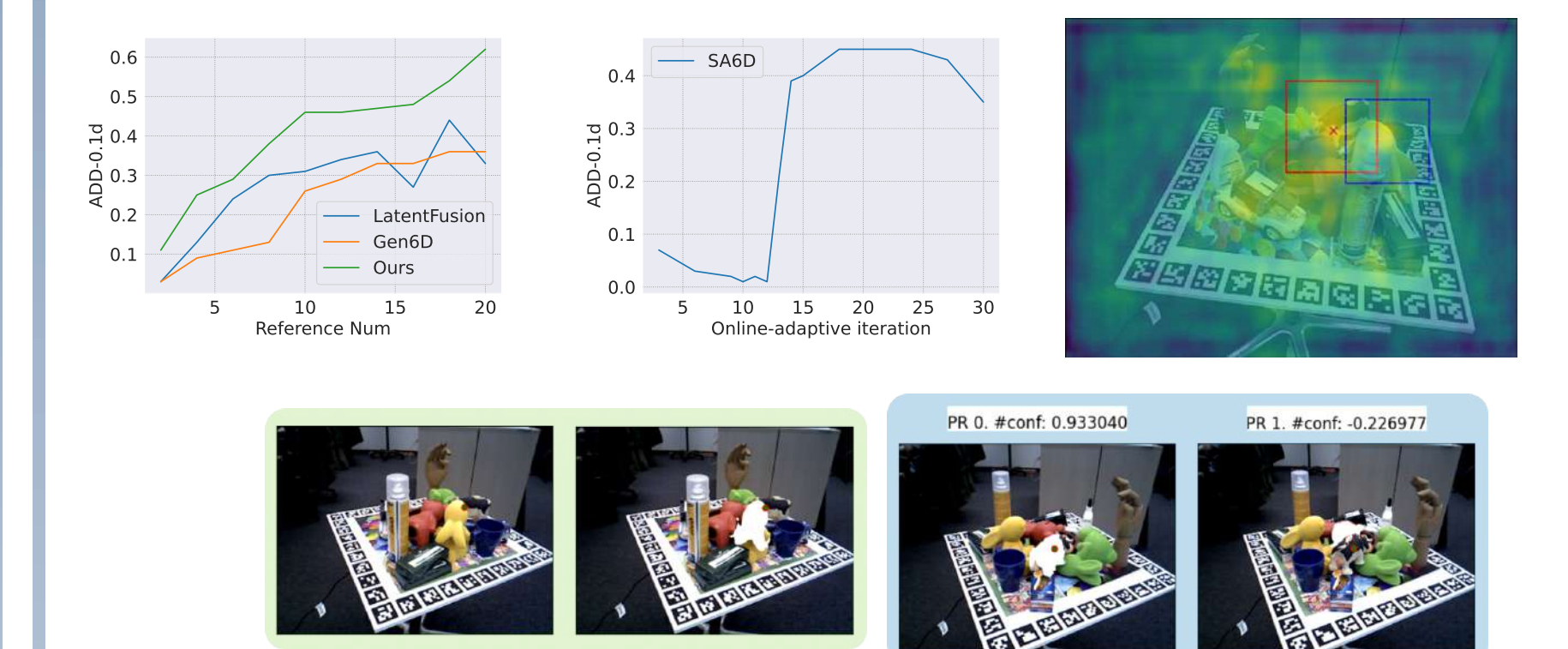
## EXPERIMENTS

**Contrastive loss:**

$$l_{ij} = -\log \frac{\exp(\text{sim}(r_i^P, r_j^P)/\tau)}{\sum_{r' \in R^N \cup \{r_j^P\}} \exp(\text{sim}(r_i^P, r')/\tau)}, \quad (1)$$
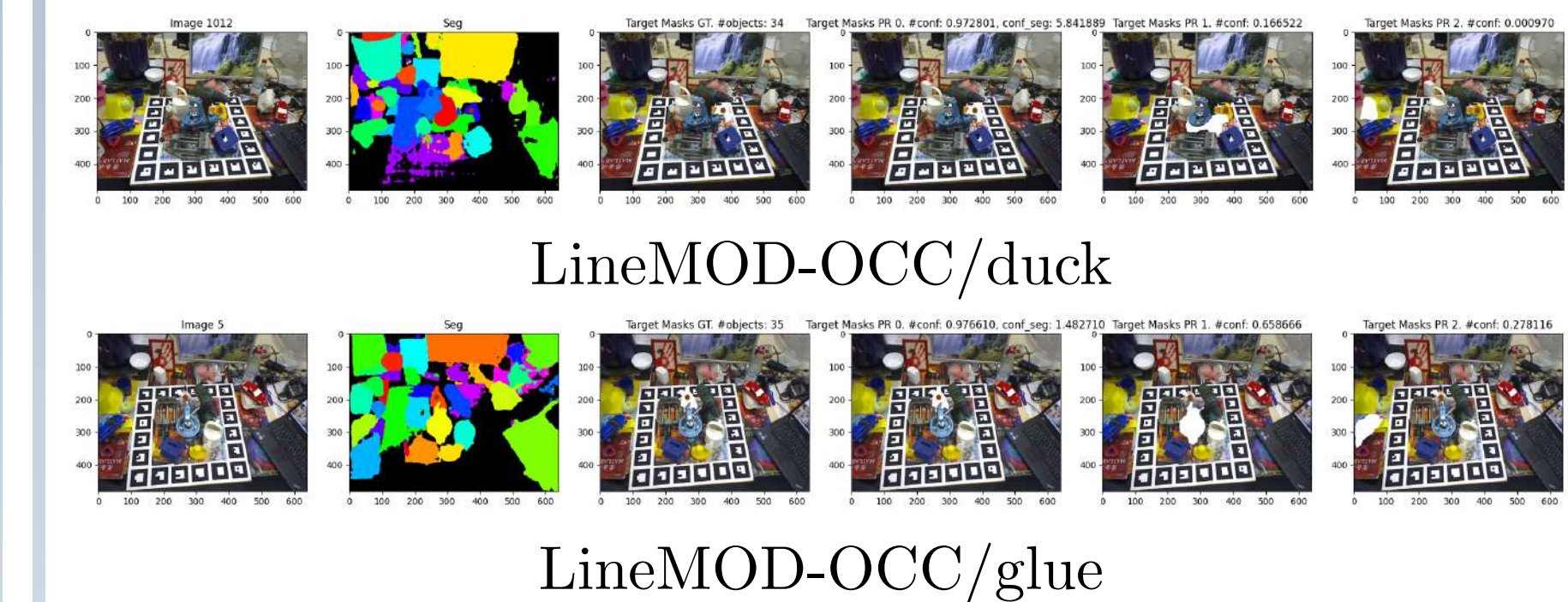
**Qualitative results:**



**More analysis:**



**Robust prediction of target segmentation:**



LineMOD-OCC/duck



LineMOD-OCC/glue

**Quantitative results:**



Evaluation of ADD-0.1d

| Method | ADD-0.1d | ADD-0.3d | ADDs-0.1d | ADDs-0.3d |
|---|---|---|---|---|
| LFLatentFusion | 0.1162 | 0.1738 | 0.1299 | 0.1907 |
| Gen6DGen6D | 0.3571 | 0.6399 | 0.6399 | 0.7530 |
| SA6D (wo/ RFM) | 0.4018 | 0.7292 | 0.6964 | **0.8780** |
| SA6D | **0.5595** | **0.7887** | **0.8393** | **0.8780** |

| Method | IOU$_{0.5}$ | 5°2cm | 5°5cm | 10°5cm |
|---|---|---|---|---|
| CASSess | 0.01 | 0.0 | 0.0 | 0.0 |
| Shape-Priorshapeprior | 0.33 | 0.03 | 0.04 | 0.14 |
| DualPoseNetdualposenet | 0.70 | 0.18 | 0.23 | 0.37 |
| RePoNetWild6D | **0.71** | 0.30 | 0.34 | **0.43** |
| SA6D | 0.65 | **0.37** | **0.40** | 0.42 |

Evaluation on FewSOL          Evaluation on Wild6D