

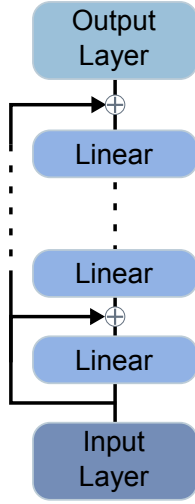
861 **7. Supplementary Material**

Figure 7. MLP Block with Skip-Connections.

862 **7.1. Synthetic Training Data Generation**

863 For the generation we proceed as following: We pick a random
 864 number of objects, sample them in a scene and render
 865 20 PBR images from random camera positions focusing on
 866 a random object in the scene for each camera position. The
 867 scenes vary in objects placed (i) randomly on the ground,
 868 (ii) flying in space, or (iii) in an upright standing position.
 869 We keep generating scenes until all objects have been seen
 870 (*i.e.* at least 40% visibility) at least 64 times in at least 8
 871 different scenes. Afterwards we convert the scene-centric
 872 into a highly efficient object-centric format by mapping the
 873 object id to all images where they occur, giving us at least
 874 64 random views for each object in each scene type. See
 875 Fig. 8 for some example images from the training set.

876 **7.2. Losses**

877 **Encoder loss.** To enforce the encoder to predict towards
 878 ground truth surface points we define a simple regression
 879 loss that minimizes the Euclidean distance between the
 880 estimated surface point $s_i = q_i - \tilde{U}(q_i) \times \frac{dU(q_i)}{dq_i}$ with
 881 $\tilde{U}(q_i) := U(\hat{f}_i^Q(q_i))$ and the ground truth surface point
 882 \bar{s}_i :

$$883 L_{\chi} = \frac{\sum_{i=1}^M \|s_i - \bar{s}_i\|}{M}, \quad (5)$$

884 with $M = |\chi|$. We do not directly enforce any loss on
 885 the NeMO features $\hat{\mathcal{F}}^Q$ so that they can be freely learned
 886 through backpropagation through the decoder. During train-
 887 ing, after Ψ has predicted χ , we randomly rotate, translate
 888 and scale the NeMO points s_i by a random transformation

889 \mathbf{T} to teach the decoder to learn a coordinate system that is
 890 independent of the anchor image I_A .

891 **Dense Prediction Losses.** We define additional losses on
 892 the dense predictions of the decoder Θ . For the modal and
 893 amodal segmentation losses L_{modal} and L_{amodal} we use an
 894 equally weighted dice- and binary cross-entropy loss. The
 895 pointmap is learned using a confidence weighted L1 loss:
 896 We define the set of valid pixels in I_q belonging to the object
 897 as \mathcal{D}_{Obj} . As in [60] we define a regression loss between the
 898 estimated 2D-3D correspondences X and the ground truth
 899 mapping \bar{X} :

$$900 L_{2D3D} = \frac{\sum_{i \in \mathcal{D}_{\text{Obj}}} \tilde{C}_i \|X_i - \bar{X}_i\|}{|\mathcal{D}_{\text{Obj}}|} \quad (6)$$

901 where $\tilde{C}_i = 1 + \exp C_i$ makes sure that a confidence of zero
 902 does not lead to diminishing gradients. Note that during
 903 training we transform the ground truth correspondences \bar{X}
 904 by the same transformation \mathbf{T} as applied to the points s_i
 905 in χ . Different to [60], we explicitly enforce certainty on
 906 pixels belonging to the object and uncertainty on pixels that
 907 do not belong to the object:

$$908 L_{\text{certain}} = \frac{\sum_{i \in \mathcal{D}_{\text{Obj}}} (1 - \tanh(\exp(C_i)))}{|\mathcal{D}_{\text{Obj}}|} \quad (7)$$

$$909 L_{\text{uncertain}} = \frac{\sum_{i \in \mathcal{D}_{\text{Bg}}} (\tanh(\exp(C_i)))}{|\mathcal{D}_{\text{Bg}}|}$$

909 where \mathcal{D}_{Bg} are all pixels not belonging to the object. Based
 910 on the introduced losses, we define the total training loss as:

$$911 L_{\text{total}} = \alpha(L_{\chi} + L_{2D3D})$$

$$912 + \beta(L_{\text{certain}} + L_{\text{uncertain}} + L_{\text{modal}} + L_{\text{amodal}}), \quad (8)$$

913 where α, β are hyper-parameters. During training we set
 914 $\alpha = 1.0, \beta = 0.2$.

915 **7.3. Training Details**

916 **Sample Structure.** We define each sample in a batch as
 917 ten randomly picked images from the training set showing
 918 the same object and 1500 sampled points in the volume
 919 $[-1, 1]^3$. Five of the images are used as template images
 920 for the encoder to generate a NeMO, while all ten images
 921 are used as input to the decoder. The total loss as defined
 922 in Eq. (8) includes the encoder as well as multiple dense
 923 decoder losses, training the encoder and decoder jointly.
 924 While the five images used for the NeMO generation always
 925 show the same object with small variations in the bounding
 926 box, the other five images have a 30% chance of being a
 927 completely random crop from the scene. This allows the

928 decoder to learn to ignore images in which the object is not
929 shown while also improving the decoder’s capability to find
930 the object in cluttered scenes.

931 **Data Augmentation.** In 2/3 of the samples we use 1500
932 uniform sampled 3D points in the volume $[-1, 1]^3$ as input
933 Q while in 1/3 of the samples we use 1500 points sampled
934 close to the surface of the CAD model of the object. Due
935 to the self-attention mechanism in the encoder, this allows
936 us to train the model to use the CAD model information
937 if available, while also being able to work without it. As
938 pixel-level data augmentation to all templates in a sample,
939 we apply:

- 940 • 30% chance of color jitter with 0.5 brightness, 0.5 con-
941 trast, 0.5 saturation and 0.5 hue, and
- 942 • 10% chance of grayscale,
- 943 On template images directly we apply:
- 944 • 50% chance of random affine transformation with up to
945 5° rotation, 0.1 translation and 0.1 scaling, and
- 946 • 50% chance of Gaussian noise with standard deviation
947 between 0 and 0.05.

948 **Initial NeMO Space.** To define the initial coordinate sys-
949 tem based on the anchor image I_A , we use the mean of the
950 ground truth surface points as the center of the NeMO space
951 and the orientation of the anchor image as it is in the camera
952 frame in OpenCV notation. The scale of the NeMO point
953 cloud is defined such that the anchor image’s surface points
954 take 1/3 of the volume $[-1, 1]^3$.

955 7.4. Object Alignment

956 As the BOP challenge compares the predicted object pose
957 against the ground truth pose, we need to align the coord-
958 inate system of our NeMO with the ground truth object
959 pose to be able to evaluate our approach. Note that this is
960 only required for evaluation and generally not a necessity of
961 our method. As we are given the object-to-camera transfor-
962 mation both in the model-based, where we render our tem-
963 plates, and in the model-free setting, where we are given
964 the ground truth pose, we can align our coordinate system
965 by optimizing the scale, rotation and center-offset. We use
966 a simple gradient descent optimization that minimizes the
967 error between the predicted and ground truth pose of the
968 $k = 5$ best rotation predictions of the template images. See
969 Algorithm 1 for the pseudo code of the alignment algorithm.

970 7.5. From Dense Predictions to Detection, Segmen- 971 tation and 6DoF Pose Estimation

972 Since our decoder produces dense predictions, we need ad-
973 ditional steps to get a 6DoF pose estimation and amodal
974 bounding box. These predictions are run in parallel us-
975 ing multiple decoder and different NeMOs. To predict the
976 segmentations and amodal bounding boxes of a scene, we
977 patchify the scene image with two patch sizes (smallest im-
978 age side and half the smallest image side) with a stride of

Algorithm 1 NeMO Alignment Algorithm

Require: Ground truth poses $T_i^{\text{gt}} \in SE(3)$, estimated
poses $T_i^{\text{est}} \in SE(3)$, max iterations N , learning rate
 η , Huber loss threshold δ , number of best poses k

Ensure: Optimal scale s , rotation R , translation t that
aligns T^{est} to T^{gt}

- 1: Select k best estimates with lowest initial rotation error
 - 2: Initialize scale s based on average translation norms
 - 3: Initialize rotation as axis-angle vector $\mathbf{r} \leftarrow 0$, transla-
tion $\mathbf{t} \leftarrow 0$
 - 4: **for** $i = 1$ to N **do**
 - 5: $R \leftarrow \text{AXISANGLETOROTATION}(\mathbf{r})$
 - 6: Construct correction transform $T^{\text{corr}} = \begin{bmatrix} R & \mathbf{t} \\ 0 & 1 \end{bmatrix}$
 - 7: Apply correction: $\hat{T}_i \leftarrow T_i^{\text{est}} \cdot T^{\text{corr}}$
 - 8: Compute predicted \hat{R}_i and scaled \hat{t}_i
 - 9: Compute rotation error $\theta_i = \angle(\hat{R}_i, R_i^{\text{gt}})$
 - 10: Compute translation error $\epsilon_i = \|\hat{t}_i - t_i^{\text{gt}}\|$
 - 11: Compute robust loss using Huber: $L = \text{Huber}(\theta) +$
Huber($\epsilon/10$)
 - 12: Update $s, \mathbf{r}, \mathbf{t}$ via gradient descent with Adam opti-
mizer
 - 13: **if** $L < L_{\text{best}}$ **then**
 - 14: Save current s, R, \mathbf{t} as best parameters
 - 15: **end if**
 - 16: **end for**
 - 17: **return** Optimized s, R , and t
-

1/3 of the patch size. We accumulate the patch predictions
and look for clusters of predictions larger 0 to differenti-
ate between object instances. The amodal bounding box
is retrieved by laying a box around a amodal segmenta-
tion cluster. We ignore small clusters and clusters touching
the image sides. To predict the 6DoF pose estimation, we
crop the region from a bounding box and all pixels with a
minimum confidence value of 0.1 in the 2D-3D correspon-
dences. We then compute the pose by solving the PnP prob-
lem using OpenCV’s solvePnP function with
iterationsCount 500, reprojectionError 6
and min inlier ratio of 30%.

7.6. Additional NeMO Analyses

**Similarity between predicted and ground truth object
surface.** We show the close resemblance of our dense point
cloud by measuring the Chamfer distance to the ground
truth point cloud. The prediction is based on two different
point cloud sampling methods as input to our NeMO. Both
methods achieve strong results in Tab. 7 across the three
datasets with slightly better performance, when the input
point cloud is sampled based on the actual CAD model for
a given object.

| Sampling Method | Chamfer Distance ↓ | | |
|-----------------|--------------------|--------|--------|
| | YCB-V | HOPEv2 | HANDAL |
| Random Sampled | 0.0055 | 0.0029 | 0.009 |
| CAD Surface | 0.0036 | 0.0019 | 0.0022 |

Table 7. **Chamfer distance between ground truth point cloud and dense point cloud prediction based on two different NeMO input sample point clouds.** First, randomly sampled points. Second, points sampled on the CAD surface, which requires the model of a given object. We observe already small errors only in case of randomly sampled point clouds as input to our NeMO representation. We further improve this distance when we sample points based on the provided CAD model.

| ICP Method | Point Sampling | AP ↑ | AP _{MSPD} † | AP _{MSSD} † | Time/Image (s) ↓ |
|------------------|----------------|-------|----------------------|----------------------|------------------|
| - | Random | 0.378 | 0.370 | 0.385 | 0.336 |
| - | CAD Surface | 0.383 | 0.420 | 0.347 | 0.531 |
| Pointmap Surface | Random | 0.475 | 0.394 | 0.557 | 1.720 |
| Pointmap Surface | CAD Surface | 0.507 | 0.447 | 0.567 | 1.479 |
| CAD | Random | 0.671 | 0.639 | 0.703 | 2.049 |
| CAD | CAD Surface | 0.626 | 0.594 | 0.658 | 1.607 |

Table 8. **Precision with different point sampling and ICP.** We report the AP on YCB-V with ground truth detections without refinement.

CAD-Model Conditioned Sampling. We compare the performance of our network using randomly sampled Q and points sampled from the surface of a CAD model. Additionally, we compare different ICP-based refinement strategies: (i) no ICP, (ii) using the pointmap prediction as source point cloud for the ICP, and (iii) using CAD model surface points as source point cloud for the ICP. The results are shown in Tab. 8. We show that our design choice to use a continuous point cloud as input to the encoder can leverage the information from the CAD model to improve the performance of our method. We also show that we can use the pointmap prediction for ICP refinement, although no depth information is given during the encoding step, increasing the average precision compared to purely RGB based (no ICP) pose estimation. An example of estimated object surfaces can be seen in Fig. 10. This experiment demonstrates the versatility of our network, allowing us to adapt to different scenarios.

| Method | T-LESS | TUD-L | YCB-V |
|----------------------------------|--------------|--------------|--------------|
| Co-op (F3DT2D, Coarse, RGBD)[41] | 0.620 | 0.841 | 0.808 |
| FreeZeV2.1[8] | 0.751 | 0.991 | 0.905 |
| FreeZeV2.2[8] | <u>0.777</u> | <u>0.988</u> | <u>0.911</u> |
| FRTPose-WAPP (Default Detection) | 0.826 | 0.978 | 0.918 |
| Ours (CNOS Detection) | 0.291 | 0.560 | 0.592 |
| Ours (NeMO detection) | 0.144 | 0.756 | 0.567 |

Table 9. **Model-Based 6D Detection of Unseen Objects with Refinement.** We report Average Precision (AP) on 3 BOP datasets and compare against other methods published on the public Model-Based Unseen Object 6D Detection leaderboard. [3]. We used a simple ICP as refinement step. † indicates unpublished methods.

| NeMO Generation | Chamfer Distance ↓ | | |
|-----------------|--------------------|--------|--------|
| | YCB-V | HOPEv2 | HANDAL |
| Fused | 0.0080 | 0.0034 | 0.0088 |
| All Templates | 0.0081 | 0.0034 | 0.0084 |

Table 10. **Chamfer distance between ground truth point cloud and dense point cloud prediction based on two different NeMOs.** First, a fused NeMO based on two respective NeMOs generated through 4 images each but each with the same anchor image. Second, predicted dense point cloud based on all 7 template views. This shows the versatility of our NeMO as the point cloud resemblance and error stays constant.

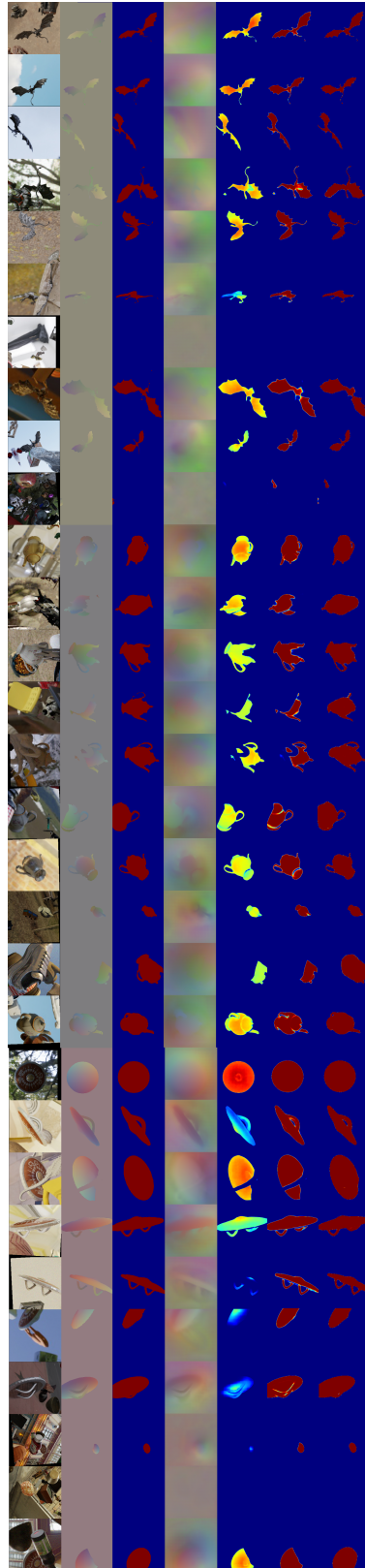


Figure 8. **Training Images.** From left to right, respectively: RGB input, ground truth 2D-3D mapping, ground truth amodal segmentation, estimated 2D-3D mapping and estimated confidence, estimated modal segmentation and estimated amodal segmentation. Shown are three different training samples, each sample consists of 10 images. The first 5 images are used to generate a NeMO, all 10 images are used for dense predictions. Note that in the last 5 templates of a sample, the object crop is stronger augmented and sometimes the object is not visible at all.



Figure 9. **NeMOs from benchmark objects.** We show the NeMOs (left, with PCA feature visualization), dense predictions (middle) and object surface reconstruction (right) from objects selected from different datasets. From top to bottom, the datasets are: HANDAL, HOPEv2, YCB-V and T-LESS. We use 32 template images to generate the NeMOs and use the same 32 images to show the dense prediction results. Note that we show real template images for HANDAL, HOPEv2 and YCB-V and synthetic template images for the T-LESS object. All NeMOs have been aligned with the ground truth pose as described in Sec. 7.4.

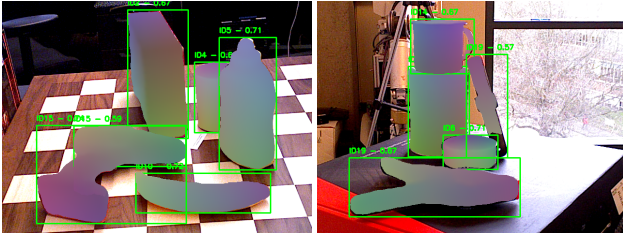


Figure 10. **Pointmap predictions** on the YCB-V dataset. We only show estimations with confidences > 0.1 .

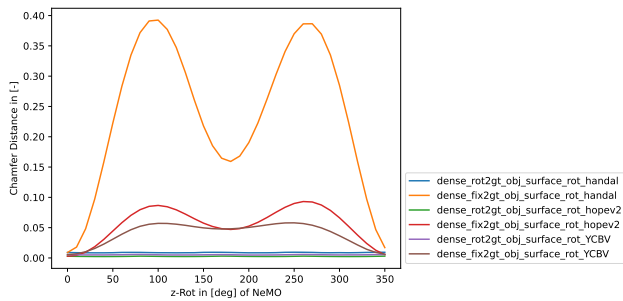


Figure 11. **Chamfer distance between dense predicted and rotated ground truth point cloud** averaged across all objects of respective dataset. In case objects and NeMOs are rotated analogously, we observe small errors *i.e.* the point clouds are very similar. M-shaped error trajectory, with minimum around 180 degrees of rotation around objects' z-axis, in case of fixed 2D-3D dense representation but rotated ground truth object point cloud. This shows that our representation is aware of rotations.