

# Self-supervised Vision Transformers for 3D Pose Estimation of Novel Objects

Stefan Thalhammer<sup>1</sup><sup>[0000-0002-0008-430X]</sup>, Jean-Baptiste Weibel<sup>1</sup><sup>[0000-0003-0201-4740]</sup>, Markus Vincze<sup>1</sup><sup>[0000-0002-2799-491X]</sup>, and Jose Garcia-Rodriguez<sup>2</sup><sup>[0000-0002-7798-3055]</sup>

- <sup>1</sup> Automation and Control Institute, TU Wien, Vienna, Austria {`thalhammer, weibel, vincze`}@acin.tuwien.ac.at  
<sup>2</sup> Department of Computer Technology, University of Alicante, San Vicente del Raspeig, Spain  
`jgarcia@dtic.ua.es`

**Abstract.** Object pose estimation is important for object manipulation and scene understanding. In order to improve the general applicability of pose estimators, recent research focuses on providing estimates for novel objects, that is objects unseen during training. Such works use deep template matching strategies to retrieve the closest template connected to a query image. This template retrieval implicitly provides object class and pose. Despite the recent success and improvements of Vision Transformers over CNNs for many vision tasks, the state of the art uses CNN-based approaches for novel object pose estimation. This work evaluates and demonstrates the differences between self-supervised CNNs and Vision Transformers for deep template matching. In detail, both types of approaches are trained using contrastive learning to match training images against rendered templates of isolated objects. At test time, such templates are matched against query images of known and novel objects under challenging settings, such as clutter, occlusion and object symmetries, using masked cosine similarity. The presented results not only demonstrate that Vision Transformers improve in matching accuracy over CNNs, but also that for some cases pre-trained Vision Transformers do not need fine-tuning to do so. Furthermore, we highlight the differences in optimization and network architecture when comparing these two types of network for deep template matching.

**Keywords:** Object pose estimation · Template matching · Vision transformer · Self-supervised learning

## 1 Introduction

Object pose estimation is an important yet difficult vision problem. Many downstream tasks, such as grasping [37], augmented reality [25] and reconstruction [35] benefit from the availability of object poses. Classical object pose estimation

---

source code: <https://github.com/sThalham/TraM3D>

approaches encode latent representations of multiple object views per object, during training. During run-time these are matched against an observation to retrieve a coarse object pose [20, 24, 12]. After retrieving the pose of the closest template, poses are refined using Iterative-Closest-Points [17] algorithm or other algorithms to optimize the rigid transformations between two corresponding sets of points. In contrast, learning-based solutions using Convolutional Neural Networks (CNN) learn a feature representation to infer object class and geometric correspondences during testing [34, 38, 47, 26, 48, 50, 1, 43, 9]. Yet, training pose estimators for each object instance [34, 38], or each set of object instances [48, 50] is insufficient to be usable in real world scenarios where object instances are manifold and constantly changing. As a consequence research shifts towards category-level [51, 39] and novel object pose estimation [32, 42, 30]. These recent novel object pose estimation approaches are similar to classical ones in the sense that queries are matched against templates.

The approach of [32] employs a CNN backbone to learn occlusion-aware template matching for novel object pose estimation. Real observations are matched against rendered templates and tested for 3D pose estimation. While they show that such strategies are expedient for novel object pose estimation it has been shown that Vision Transformers (ViT) [11, 49, 5] learn more discriminative feature spaces than CNNs when trained in such unsupervised manners. This advantage of ViTs over CNNs, however, has primarily been empirically demonstrated by matching to distinct object classes and not by matching views of the same object class for more complex reasoning, such as 3D object pose estimation [5, 8].

In this work we empirically demonstrate that ViTs excel over CNNs when used for novel object pose estimation. Modifying the approach of [32] for comparing two similarly sized feature extractors, ResNet50 [18] with 23M and ViT-s [49] with 21M parameters, we show that these improvements are manifold. Training self-supervised ViTs for 3D object pose estimation not only improves the template matching accuracy, but also reduces the training time. Depending on the dataset and metric, template matching accuracy for seen objects ranges from 1% on Linemod [20], over 4% on Linemod-Occlusion [4], to 19% on T-LESS [22]. For unseen objects, the respective improvements are 3%, 5% and 18%. Achieving these improvements using ViT-s takes one fourth of the training time and iterations on LM and LM-O, and only one twenty-fifth of it on T-LESS. More remarkably, testing ViT-s on T-LESS in a zero-shot fashion, thus without fine-tuning, already improves over using fine-tuned ResNet50 by 7% and 9%, for seen and unseen objects respectively. Finally, works such as [5, 8] train self-supervised ViTs to retrieve the object class of seen objects assuming the availability of templates in the same domain. These assumptions are impractical for novel object pose estimation. Uniform coverage of the pose space is crucial and thus rendering templates is expedient. Furthermore, handling unseen objects is desired to further generalize real-world deployment of pose estimators. As a consequence, this work provides ablations on the matter of network architecture used for matching. While the aforementioned works [5, 8] benefit from using high-dimensional, multi-layered projection heads, we empir-

ically show that these increase the template matching error on unseen objects when matched against rendered templates. In summary we:

- Show that Vision Transformers not only exhibit reduced template matching errors compared to CNNs for matching synthetic templates to known objects, but also to novel objects. The relative improvements for novel object pose estimation range from 3% to 18%, depending on the dataset and metric used.
- Demonstrate that pre-trained Vision Transformers exhibit excellent matching performance for zero-shot matching. On the T-LESS dataset, non fine-tuned Vision Transformers exhibit a relative improvement over fine-tuned CNNs of 7% and 9%, on known and novel objects respectively. Fine-tuning further improves to 19% and 18% respectively.
- Highlight the differences in matching procedure and optimization of fine-tuning Vision Transformers for template matching. Our results indicate that Vision Transformers encode relevant features over a broad range of descriptor sizes for seen and novel objects. As compared to CNNs, where there is a trade-off when choosing the descriptor size for either seen or novel objects. Our results additionally indicate that high-dimensional, multi-layered projection heads increase the template matching error for the problem at hand.

The remainder of the manuscript is organised in the sections Related Work, Method, Experiments and Conclusion. The next section presents the state of the art for object pose estimation, focusing on deep template matching for deriving poses of novel objects, and self-supervised vision transformers.

## 2 Related Work

This sections presents the state of the art for object pose estimation with the focus on novel object pose estimation. Subsequently, ViTs and self-supervised training for them is presented.

Learning-based object pose estimation research focuses on multi-staged pipelines [28, 34, 50, 43] that often train separate networks for instance-level pose estimation [34, 50], in order to improve the estimated pose’s accuracy. Different streams of research improve on the scalability of instance-level pose estimation, presenting solutions for improved multi-object handling [1, 47, 56] and reducing the number of stages needed for providing reliable pose estimates [48, 9, 55]. Yet, re-training pose estimators every time novel objects or object sets are encountered is cumbersome and delays the deployment in the real world. As a consequence, recent works overcome these shortcomings by training for category-level pose estimation [51, 39] or by training deep template matching for novel object pose estimation [32, 42, 30].

**Deep Template Matching** Matching observations against predefined templates is a long-standing concept of object pose estimation [20, 24, 12]. Recent learning-based solutions adopt this strategy, since it has two major advantages [32, 42, 30]. First, training time is low since encoding templates does not

require learning a representation of each object individually. Creating a latent representations for each relevant template only requires one network forward pass. Thus, template encoding is done in the magnitude of seconds for an object of interest, as compared to training instance-level pose estimators, which takes hours to days, depending on the number of objects and the hardware [34, 50, 48]. Second, training instance-level pose estimators encodes a latent representation of the object, respectively objects, of interest. This representation does not generalize to novel objects. This shortcoming has to be addressed by either category-level object pose estimation, or by deep template matching.

The approach of [52] introduces deep descriptors for matching query objects against templates for retrieving the 3D pose using nearest neighbor search. In [3] the authors improve over [52] by guiding learning in pose space, also accounting for object symmetries in the process. Recently, [32] proposed further improvements. They replace the triplet loss-based training with an InfoNCE-based one and improve occlusion handling by masking the feature embedding using the template’s mask and an occlusion threshold. We adopt and improve over their approach for deep template matching by using ViTs for descriptor extraction, which have not yet been adopted by the community. As such, we demonstrate their advantage with respect to their generality as deep template matcher and show empirical evaluations highlighting their advantages for the problem of novel object pose estimation.

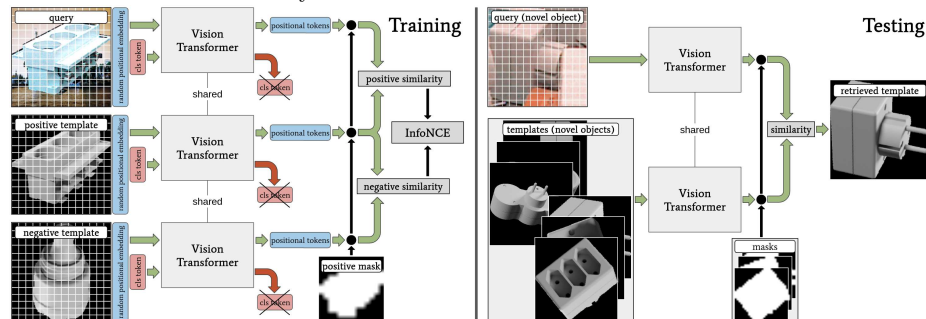
**Vision Transformer** It has recently been shown that ViTs [11, 36] learn superior features when trained in a self-supervised fashion [8, 5, 49]. These mainstream works focus on training object classifiers from scratch, and using large datasets with little domain shift between query images and templates. Such large datasets are difficult to obtain for object pose estimation due to the complexity of generation accurate 6D pose annotations. Additionally, it is relevant for pose estimation to effectively cover the viewing sphere around objects of interest [45]. This implies training on comparably small datasets and preferably using synthetically creating templates, i.e. using rendering for template creation [10]. As such, in this work, ViTs are assumed to be pre-trained, and templates are rendered. We thus show the potential of self-supervised ViTs under that shifted perspective and also highlight the differences in network design as compared to the mainstream research direction.

### 3 Method

This section presents our self-supervised learning framework for matching real observations to synthetic templates for novel object 3D pose estimation. Figure 1 provides an abstract visualization of the presented method.

Self-supervised training is done using contrastive learning. Contrastive learning aims at maximizing the similarity of semantically close training samples, referred to as positive pairs, while minimizing the similarity for samples that are semantically dissimilar, that is negative pairs. More precisely, one training sample consists of a tuple of a query crop ( $I_q$ ), a positive example ( $I_{pos}$ ),

**Fig. 1. Method overview** During training, a query image, a positive and a negative template is processed by a Vision Transformer to encode a feature embedding. The number of the positional tokens is retained for the feature map. InfoNCE [33] is used in a Triplet loss-like fashion with the input feature map being masked with the positive template. During testing, novel query objects are matched against templates to retrieve object class and 3D pose from the matched template. Template retrieval is guided using the masked cosine similarity.



and a negative one ( $I_{neg}$ ). The positive and negative template are rendered using physically-based rendering (pbr) [10]. Where the positive sample correlates with respect to object class and rotation with the query image. The negative sample deviates with respect to both properties. Crops are tokenized using random patch embedding and a shared pre-trained ViT-s [49] is used for extracting features of the query and the template images. In contrast to self-supervised ViT-frameworks for classification [5, 8] we discard the class token and employ the positional tokens for similarity calculation. Using such spatial output enables dropping tokens based on the positive template’s mask. Optimization is guided using InfoNCE-loss [33] with the positive and negative similarities as input. During testing, similarities are computed between real object observations and pbr-templates of seen and novel objects. Thus, in contrast to contemporary ViT-research, similarities have to bridge the synthetic-to-real gap, since templates are created using rendering [5, 8]. The real observations are compared against templates that represent uniformly distributed object views of the potentially new objects. Ultimately, the class and the 3D rotation of the matched template are retrieved.

### 3.1 Feature Embedding

The aim of this work is novel object pose estimation. Recent works shows that deep contrastively-learned template matching strategies are well suited for this task [32, 42, 30]. In order to exhibit high similarities between similar view points of the same object in different domains, the learned feature embedding has to represent the object view as accurately as possible. It has been shown that Vision

Transformers [49, 11, 36], trained in an unsupervised way, learn to accurately model long-range image relationships, improving over CNNs [5].

This work adopts the ViT-s network, presented in [49] as feature extractor. The weights are pre-trained on ImageNet [29] in a self-supervised manner [5]. ViT-s is used by [5] only retaining the class token for training and testing. In this work, the class token is discarded and the positional tokens are retained in order to benefit from the spatial nature of the output. Diverse works indicate that augmenting feature extractors with deep multi-layered heads, for projecting embeddings to higher dimensions, improves performance when training on ImageNet [8, 5, 7, 15]. The presented results in Section 4 indicate this finding does not apply to pose estimation. A single linearly-activated fully-connected layer projects the feature embedding, coming from the pre-trained backbone, to a lower dimensionality. It has to be noted that this different behavior is connected to the difference in problem; a) the backbone is initialized with pre-trained weights, b) the problem at hand matches real observations against rendered templates and c) testing is partially done on novel objects, thus data unseen during training. We hypothesize that using deeper heads overfit to the training data characteristics.

The authors of [8] note that randomly initialized patch embedding stabilizes training on ImageNet and thus improves classification accuracy. Accordingly, the patch embedding layer is not updated during fine-tuning. Results are provided in Section 4.

### 3.2 Contrastive Learning Framework

The feature embeddings extracted using ViT-s are processed by a contrastive learning framework for learning to increase similarity between object crops of the same class and a similar viewpoint. As similarity measure, the cosine similarity is employed:

$$\text{sim}(emb_{I_q,t}, emb_{*,t}) = \frac{emb_{I_q,t} \cdot emb_{*,t}}{\|emb_{I_q,t}\|_2, \|emb_{*,t}\|_2} \quad (1)$$

Where  $*$  is either  $I_{pos}$  or  $I_{neg}$ . The similarity is computed locally and aggregated for locations indicated by the mask image:

$$\text{sim}_{pos/neg} = \sum_{t=1}^T \text{sim}(I_q, *) \times M_t \begin{cases} \text{sim} & \text{if } M_t == 1, \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Where  $T$  refers to the number of feature map locations, i.e. the number of positional tokens. The negative similarity is summed over all embedded tokens inside the template’s object mask, while the positive similarity is computed globally with  $M = 1^{\text{size of } I_q}$ . Both similarities are used in a triplet loss fashion [6] using InfoNCE loss [16, 33]. Each positive sample is compared against all negative samples in a batch, resulting in  $B = (b \cdot b) - b$  negative samples per iteration.

$$L = - \sum_{i=1}^b \log \frac{\exp \frac{sim_{pos,i}}{\tau}}{\sum_{k=1}^B \frac{sim_{neg,k}}{\tau} \forall i \neq k} \quad (3)$$

Where  $\tau$  is a temperature parameter set to 0.1. For more details consult [32].

### 3.3 Template Matching

During testing templates of seen and novel objects, are matched against the query image. Embeddings are created for the query crop and all templates. The cosine similarity in Equation 1 is reused, yet modified to:

$$sim_q = \sum_{t=1}^T sim(I_q, *) \times M \begin{cases} sim \text{ if } M_t == 1, \\ \text{and } sim_t > \delta, \\ 0 \text{ otherwise} \end{cases} \quad (4)$$

Where  $\delta$  is a hyperparameter set to 0.2, which is meant to increase robustness against occluded image regions, as introduced by [32]. The class and 3D rotation of the template leading to the highest cumulative cosine similarity are retrieved.

## 4 Experiments

Presented results compare the CNN-based baseline methods, like [32], to our approach that uses ViT-s as feature extractor. Additional results evaluate the generality of the self-supervised pre-trained ViT-s without fine-tuning, showing that even without fine-tuning the template matching error is low and even improves over the baseline method on T-LESS. Ultimately, we present diverse ablations that highlight the differences between ViT- and CNN-architectures for 3D pose estimation. The experiments section is concluded by providing an ablation with respect to the projection head used for our approach, highlighting the fundamental difference that for the addressed problem shallow heads are beneficial, as compared to approaches used for classification on ImageNet [5, 8, 7, 15].

### 4.1 Experimental Setup

In the following paragraphs data retrieval and processing is detailed. Following that, template creation for matching is explained. In order to evaluate the proposed approach, standard metrics from concurrent, conceptually similar approaches and presented.

**Datasets** Results are provided on three standard datasets for object pose estimation, Linemod [20] (LM), Linemod-Occlusion [4] (LM-O), and T-LESS [22]. These datasets are processed to provide crop-level data in order to evaluate template matching accuracy and compare against the baseline method.

**LM and LM-O** These are two of the most-used datasets for evaluating object pose estimation approaches. LM features 13 objects. For each object a set of  $\approx 1200$  scene-level images is available. Annotations are only provided for the respective object, though each set contains multiple objects of the dataset in the cluttered background. The main characteristics of the dataset are texture-poor objects of different geometry, sizes and colors. Annotated object views exhibit virtually no occlusion. As a consequence, [4] created annotations for all 8 dataset objects in the Benchvise’s set, thus introducing LM-O as a test set specifically for strongly occluded object views.

With respect to training and test we follow [32], in order to provide a fair comparison. For evaluation on seen and unseen objects the LM-objects are partitioned into three sets, see Table 1. As training data, 90% of LM images’ per object set are used, and the remaining 10% are used for testing. As a consequence training images are without occlusion. The images of LM-O are exclusively used for testing, yet for evaluation also split accordingly, into seen and unseen objects. In order to evaluate on all objects, one split is used for testing on unseen objects, while the other two are used training.

**Table 1. LM/LM-O object splits.** Two of the sets are used for training and testing on seen objects, while the third is used for testing on unseen objects, as done by [32].

Split	Objects
1	Ape, Benchvise, Camera and Can
2	Cat, Driller, Duck and Eggbox
3	Glue, Holepuncher, Iron, Lamp and Phone

**T-LESS** On T-LESS we follow the protocol of [44]. Isolated object views of the object 1 – 18 are used for training and are pasted on a randomly chosen image of SUN397 [53], using the cut-paste strategy [13]. These 18 objects are considered as seen objects. The remaining objects, 19 – 30, are used as novel ones. Test images are cropped from the primesense test set.

**Template Generation** In contrast to works that train self-supervised ViTs for image classification [5, 8], this work considers matching the closest template for viewpoint classification, thus for 3D pose retrieval. The major difference is that templates uniformly distributed in the viewing sphere, respectively hemisphere, are required. Which is not relevant to the workings of aforementioned works. Consequently, templates to match against are created using physically-based rendering for the task at hand [10].

**LM and LM-O** The training and test dataset for LM and LM-O are processed as done by [52] and [32]. These works crop the images from the real dataset by omitting in-plane rotations. Thus, effectively only considering azimuth and elevation as degrees of freedom. Objects are cropped in a way that the image



space at object distance projects 0.4 by 0.4 meters. Thus, all objects appear at the same distance to the camera, independent of their size. Furthermore, neither the LM nor LM-O training and test images show objects from the lower viewing hemisphere. Due to these constraints 301 templates are sufficient for training and testing on LM and LM-O.

**T-LESS** For T-LESS objects are cropped in a way to tightly encapsulate the objects. Additionally, objects appear in arbitrary views in the test set. As a consequence 92,232 templates are used for training and testing on T-LESS, as done by [32, 44].

**Evaluation** This section presents the metrics used in this work. The approach of [32] introduces *Acc15* for evaluating template matching accuracy and classification. The *VSD*-score, as proposed by [23] is a standard metric for evaluating 6D object pose estimation accuracy. The following paragraphs provide detailed explanations how these metrics are used in this work.

**Acc15** This metric is introduced by [32]. It represents the accumulated true positive rate for matched templates that are below 15 deg rotational error with respect to the object class and ground truth rotation of the query crop:

$$Acc15 = \sum_{n=1}^n \begin{cases} 1 & \text{if } \arccos \frac{R_q \times R_t}{\|R_q\|_2 \cdot \|R_t\|_2} < 15 \text{ deg} \\ & \text{and } C_q == C_t, \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

Where  $n$  refers to the number of query crops,  $R_q$  and  $R_t$  to the three-dimensional rotation vectors, and  $C_q$  and  $C_t$  to the object class of the queries' ground truth and the template, respectively. Thus, matched templates with a rotation deviation of more than 15 deg from the ground truth, or which have a different class than the query image, are considered as false positives.

**VSD** This metric has been proposed by [23]. For each query object crop the deviation of the estimated pose  $\hat{P}$  to the ground truth  $P$  is projected to a scalar value using:

$$e_{VSD} = \underset{p \in \hat{V} \cup V}{avg} \begin{cases} 0 & \text{if } p \in \hat{V} \cap V \wedge |\hat{D}(p) - D(p)| < \tau, \\ 1 & \text{otherwise} \end{cases} \quad (6)$$

where  $\hat{V}$  and  $V$  are sets of image pixels;  $\hat{D}$  and  $D$  are distance maps and  $\tau$  is a misalignment tolerance with the standard value of  $20mm$ . Distance maps are rendered and compared to the distance map of the test image to derive  $\hat{V}$  and  $V$ . Since  $\hat{P}$  and  $P$  need to represent 6D poses, including the 3D translation, we need to raise estimates to 6D, the strategy of [46, 32] is adopted. Using the bounding box of the observation  $box_{obs}$ , and that of the template  $box_{tmp}$ , the corresponding intrinsics  $f_{obs}$  and  $f_{tmp}$ , and the template distance to the camera  $z_{tmp}$ , enables deriving the observed object's distance  $\hat{z}_{obs}$ :

$$\hat{z}_{obs} = z_{tmp} \cdot \frac{\|box_{tmp,x}^2 \cdot box_{tmp,y}^2\|_2}{\|box_{obs,x}^2 \cdot box_{obs,y}^2\|_2} \cdot \frac{f_{obs}}{f_{tmp}} \quad (7)$$

Using  $\hat{z}_{obs}$ , the relative translation between the observation and template of the other two translation parameters are derived. Where  $\bullet$  is a placeholder for  $x$  and  $y$ :

$$\Delta_{\bullet obs} = \frac{(box_{obs,\bullet} - c_{obs,\bullet}) \cdot \hat{z}_{obs}}{f_{obs,\bullet}} - \frac{(box_{tmp,\bullet} - c_{tmp,\bullet}) \cdot \hat{z}_{tmp}}{f_{tmp,\bullet}} \quad (8)$$

The 3D translation vector is ultimately composed as  $t_{obs} = \{x_{tmp} + \Delta x_{obs}, y_{tmp} + \Delta y_{obs}, \hat{z}_{obs}\}$ .

The *VSD*-score is then defined as:

$$VSD = \sum_{n=1}^n \frac{1}{n} \begin{cases} 1 & e_{VSD,n} < 0.3, \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

where  $n$  again refers to the number of the query sample in an evaluated test set.

## 4.2 Implementation Details

This sections outlines the base method for comparing ViT to CNN-based template matching. Following that the training procedure and the network architecture are detailed.

**Baseline method** For demonstrating the difference of CNNs and ViTs for self-supervised matching of real query crops to synthetic templates the baseline method of [32] is modified. In order to provide a fair comparison all results are generated comparing backbones with a similar number of trainable parameters, ResNet50 [18] with 23M and ViT-s [49] with 21M parameters, pre-trained in a self-supervised manner [5] on [40]. The following paragraph details training procedure and optimization settings.

**Optimizer Setting** As optimizer AdamW [54] is used. The batch size is set to 16, which is also the case for the reference method [32]. The ViT networks are only trained for five epochs, as compared to the baseline, which is trained for 20 epochs. The linear scaling rule  $lr = lr_b \cdot batch\_size/256$  [14] is adopted for choosing the learning rate. A grid search was used to determine the base learning rate ( $lr_b$ ) of  $2.5 \cdot 10^{-5}$ . No learning rate scheduling is used. Cosine weight decay scheduling, starting at 0.04 and ending at 0.4 after two epochs, is employed.

The input image size is  $224^2$  and the template’s mask size  $14^2$ . A patch size of 16 is used for input image tokenization. A single linear layer is used to project the backbone feature size of 384 to 32. This stands in contrast to works like [5, 8, 7], where multi-layered high-dimensional projectors are used. The input to the projection head is normalized using batch normalization [27]. The output of the projector is normalized using [2]. Section 4.5 ablates mask and descriptor size, as well as the choice for the projection head.

### 4.3 Main Results

This section presents experiments comparing ResNet50 [18] as feature extractor to ViT-s [49]. Evaluations are provided comparing to the state of the art for 3D template matching to the presented approach.

**Results on LM/LM-O** Table 2 compares the presented approach to those of [52], [3] and [32] for template matching on LM and LM-O. Reported are the true positive rates of matched templates with respect to object class and rotational error below 15 deg ( $Acc15$ ), as defined in [52]. We follow the paradigm of [32] and report the results of the best-performing epoch during fine-tuning. The results show that using ViTs as feature extractor consistently outperforms the CNN approach for objects seen and unseen during training. Both, conceptually similar approaches, use backbones with a comparable amount of parameters, ResNet50 [18] with 23M and ViT-s [49] with 21M. It has to be mentioned that the method of [32] is fine-tuned for 20 epochs while the ViTs are fine-tuned for only 5.

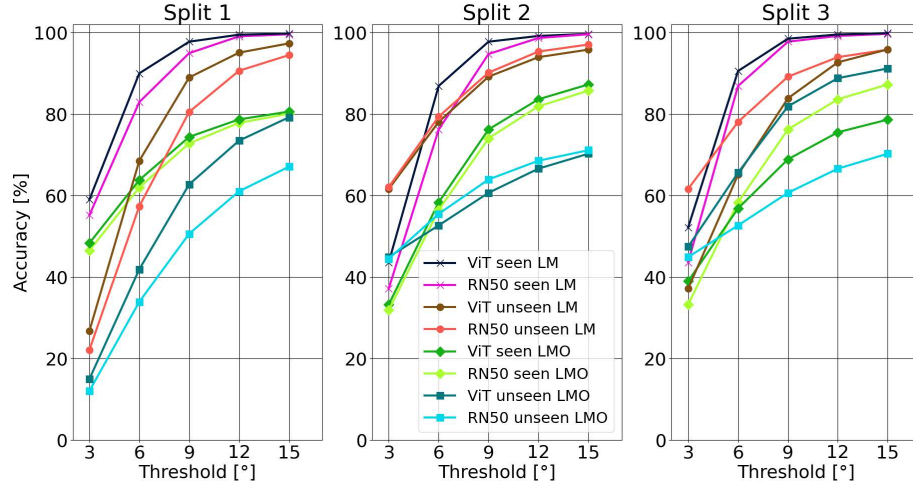
**Table 2. Comparison on LM/LM-O.** Amount of true poses for a rotational error threshold of 15 deg ( $Acc15$  [52]) for objects seen and unseen during training, see Table 1. The compared backbones have similar parameters, 23M for ResNet50 [18] and 21M for ViT-s [49]. Results for the methods indicated with † are taken from [32].

Method	Backbone	seen		unseen	
		LM	LM-O	LM	LM-O
[52]†	RN50[18]	98.1	67.5	45.1	29.9
[3]†	RN50[18]	96.1	64.7	44.3	29.1
[32]	RN50[18]	99.1	79.4	93.5	76.3
Ours	ViT-s[49]	<b>99.8</b>	<b>82.2</b>	<b>96.4</b>	<b>80.2</b>

Figure 2 shows a detailed comparison for the individual data splits of LM and LM-O, using ResNet50 [18] and ViT-small [49] as feature extractors for template matching. Tendentiously, ViT-s improves in pose estimation with respect to all rotational error thresholds on all the splits. The only exceptions are the seen LM split 3, unseen LM-O split 2 and seen LM-O split 3.

**Results on T-LESS** Table 3 compares the proposed approach to the approaches of [32], [44] and [46]. We follow the evaluation paradigm of [44] and report the VSD-score [23] using the standard thresholds, and the ground truth bounding box as basis for translation estimation. We report the performance after one epoch of fine-tuning, as compared to the 25 epochs for [32]. The results show that our approach, using ViT-small [49] as feature extractor, consistently outperforms the competing approaches for objects seen and unseen during training. Especially relevant is the comparison to the conceptually similar approach

**Fig. 2. Results on LM and LM-O splits in detail.** Reported is the percentage of true poses for different rotational error thresholds, of the CNN- and ViT-backbone for the seen and unseen object splits.



of [32], which again uses ResNet50 [18] as backbone. These results show that ViTs work well for industrial objects of T-LESS, resulting in similar pose estimation accuracy for seen and unseen objects. The following section presents pose estimation results using ViTs without fine-tuning.

**Table 3. Comparison on T-LESS.** Results are presented using the *VSD*-score with the standard thresholds presented in [21].

Method	seen: Objects 1-18	unseen: Objects 19-30	Average
[46]	35.60	42.45	38.34
[44]	35.25	33.17	34.42
[32]	59.62	57.75	58.87
Ours	<b>70.65</b>	<b>68.03</b>	<b>69.71</b>

#### 4.4 Feature Extractor Fine-Tuning

This section discusses and presents results using only ImageNet-pretrained ViTs as feature extractor. In order to use the pre-trained backbone without fine-tuning, the last linear projection layer is discarded. The output dimensionality per feature map location is 384. Table 4 compares the presented approach with

and without fine-tuning (indicated with "f.t." in the table) on LM, LM-O and T-LESS. The pre-trained ViT-s demonstrate tremendous generality with respect to feature embedding. On the LM and LM-O datasets the matching accuracy using *Acc15* is higher than that of [52] and [3], see Table 2. Yet, fine-tuning improves for all test cases. The matching accuracy on both, seen and unseen, T-LESS sets, evaluated using the *VSD*-metric, is higher than for all methods compared against in Table 3, even without fine-tuning. Fine-tuning further improves performance. The presented evaluation shows that ViTs pre-trained in a self-supervised fashion learn features that translate well to new tasks with a large shift in object categories, even without fine-tuning.

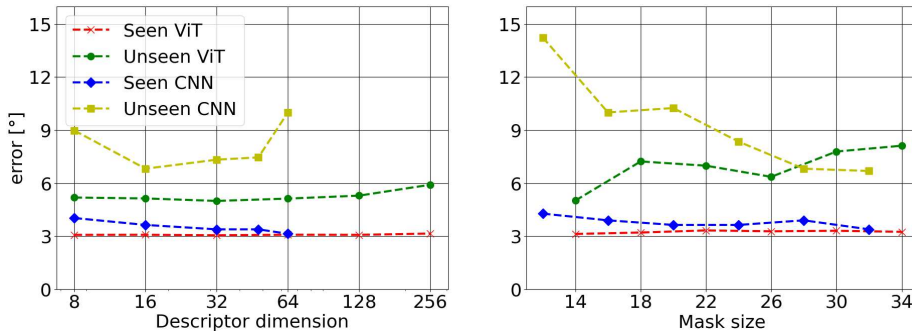
**Table 4. Influence of fine-tuning.** Result comparison for fine-tuning (f.t.) the ViT-s backbone versus only using the pre-trained feature extractor without fine-tuning.

Metric	f.t.	seen		unseen	
<i>Acc15</i> [52]		LM	LM-O	LM	LM-O
	X	81.3	56.3	85.1	63.6
	✓	<b>99.8</b>	<b>82.2</b>	<b>96.4</b>	<b>80.2</b>
<i>VSD</i> [21]		T-LESS			
	X	63.93		62.93	
	✓	<b>70.65</b>		<b>68.03</b>	

#### 4.5 Ablation Study

This sections discusses the difference in output space size and descriptor size for CNNs and ViTs. ViT and CNN approaches benefit from multi-layered, high-dimensional projection heads [8, 5, 7, 15]. Ultimately, we present experiments on the influence of projection head on our approach and additional architecture choices.

**Descriptor Size** The left plot of Figure 3 evaluates the influence of the descriptor size on the presented approach, and the ResNet50 baseline one on the seen and unseen sets of LM. The cumulative rotational error on LM decreases steadily with increasing descriptor size when using ResNet50. Yet, the optimal dimensionality is 16 for minimizing the rotational error for the unseen LM objects. While the descriptor size has a large influence on the seen LM set and even more on the unseen one, the behaviour using ViT-s is vastly different. For ViT-s the descriptor dimensionality has little influence and leads to low errors over a broad range of dimensions for seen and unseen objects. While for ResNet50 the error progression is different for both sets, the dimensionality that minimizes the error on both sets is 32 when using ViT-s.



**Fig. 3. Influence of the descriptor and mask size on LM seen and unseen** The left plot shows the influence on the rotational error of the retrieved templates when using ResNet50 and ViT-s with different descriptor sizes. The mask size is set to 32 for ResNet50 and to 14 for ViT-s. The right plot show the same comparison for different mask sizes. The descriptor size is set to 16 for ResNet50 and to 32 for ViT-s.

**Mask Size** The matching accuracy of the baseline method [32] increases when using spatially higher-dimensional feature maps since occlusion handling improves. In order to use larger feature maps for computing the template similarities we adopt the projection head of the baseline. Instead of using two convolutional layers for downsampling, we employ two transposed convolutional layers for upsampling. Both are ReLU [31]-activated. The first projecting the 384 dimensional feature vectors output by the backbone to 256, the second one to 32. Both convolutions apply no feature map padding, slide with a stride of one over the feature map and use the same kernel size, which is set depending on the desired mask size to either 3, 5, 7, 9, or 11. This projector replaces the projection head detailed in Section 4.2.

The right plot of Figure 3 evaluates the influence of the mask size on the rotational error of the matched templates. For the presented comparison the ResNet50 baseline approach [32] uses a descriptor size of 16, and ViT-s is used with a descriptor size of 32. With the ResNet50 backbone, for both the the seen and unseen objects the rotational error reduces with increasing mask size. Using the proposed ViT-s approach the behaviour is again vastly different. While the influence of the mask size is negligible for the seen objects, the rotational error for the novel objects increases significantly when increasing the mask size. This indicates that ViT-s learns relevant features for the seen objects during fine-tuning with projection heads with larger spatial output. As such, the template matching accuracy remains constant. However, increasing the feature map size used for matching is detrimental for novel objects. This correlates with the results presented in Section 4.4, which indicate that ViTs already generalize well without fine-tuning. The feature projection learned by a projection head with increased spatial output is less general and thus increases template matching error for novel objects.

**Table 5. Network architecture.** Reported is the average rotational error on LM and LM-O. The projection heads output a feature dimensionality of 32. When no head is used the standard ViT-s dimensionality of 384 is output. The column patch embedding (p.e.) indicates if the patch embedding layer is updated (l) or frozen (r) during fine-tuning.

			seen		unseen	
Head	p.e.	act.	LM	LM-O	LM	LM-O
none	l		3.14	10.95	7.80	15.44
	r		3.27	10.96	5.87	13.05
linear	l		3.07	11.05	5.39	12.83
	r		3.14	10.69	5.02	<b>11.78</b>
	r	ReLU	3.04	10.56	5.37	12.85
	r	GELU	3.12	10.28	4.98	12.75
[7]	r		3.11	10.47	<b>4.67</b>	12.20
[7]	r	ReLU	3.04	10.53	4.92	12.06
[7]	r	GELU	<b>3.02</b>	10.69	5.52	13.59
[15]	r		3.12	10.66	5.11	12.56
[15]	r	ReLU	<u>3.17</u>	<b>10.20</b>	5.14	<u>14.49</u>
[15]	r	GELU	3.04	10.70	5.17	13.56
[8]	r		3.07	10.92	5.28	12.67
[8]	r	ReLU	3.05	<u>11.22</u>	<u>5.69</u>	14.45
[8]	r	GELU	3.14	10.87	5.01	12.98

**Network Architecture Design** This section ablates different aspects of network design choices when using self-supervised learning frameworks. We investigate patch embedding and projection head design. Table 5 reports the average rotational error on LM and LM-O for the investigated aspects.

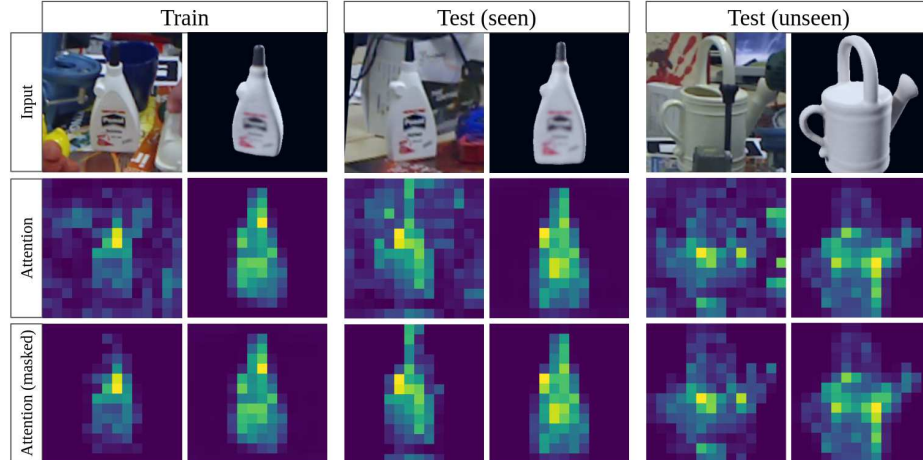
**Projection Head** The works of [7, 15, 8] use high-dimensional, multi-layered projection heads to project the feature output of the backbone to the desired dimensionality. The work of [7] uses a two-layered MLP, with the first ReLU [31] and the second layer linearly activated. The work of [15] and [5] both use three-layered MLPs, yet different versions. The latter using GELU [19]-activated hidden layers and weight normalization [41]. In [8], the projection head of [15] and the prediction head of [7] are combined. Features are normalized using batch normalization [27], and hidden layers are ReLU-activated. We compare using these projection heads to using no head or a single linear layer as head. Since using no head requires using the backbone’s output as it is, the descriptor dimensionality per feature map location is 384. For all the evaluated projection heads a hidden dimension of four times the output dimension of the previous stage and batch normalization are used. We have tested with and without using weight normalization as used by [5]. Compared to batch normalization both consistently lead to increased rotational error of the matched templates.

Table 5 compares the average rotational errors of different projection heads on LM and LM-O. The lowest error per set is indicated in bold, the highest is indicated with an underline. The lowest errors on seen and unseen LM, and un-

seen LM-O occur with heads with less layers. Using no head leads to comparably high errors. When using projection heads, the highest errors over all sets occur using higher dimensional heads. In general, for the seen objects the results are similar for all heads. Yet, projection heads with a smaller number of layers lead to less rotational error on unseen objects. This evaluation stands in contrast to self-supervised ViTs for classification that use projection heads with  $\geq 3$  layers and high dimensional hidden and last layers [8, 5]. The choice of activation appears to have little influence. Yet, heads with a lower number of layers shows reduced error on unseen objects when using no activation function.

**Patch embedding** The authors of [8] propose to use random patch embedding to increase stability during training. We experiment with the initialization of the convolution layer used for patch embedding. The second column in Table 5 (p.e.) ablates the influence. Updating the pre-trained patch embedding layer during fine-tuning is referred to as learned (l). With a slight abuse of denotation we refer to not updating the patch embedding layer during fine-tuning as random (r). We observe a similar effect as in [8]. While the error difference for the seen objects is insignificant, using random patch embedding leads to significantly less error on the unseen objects.

#### 4.6 Self-Attention



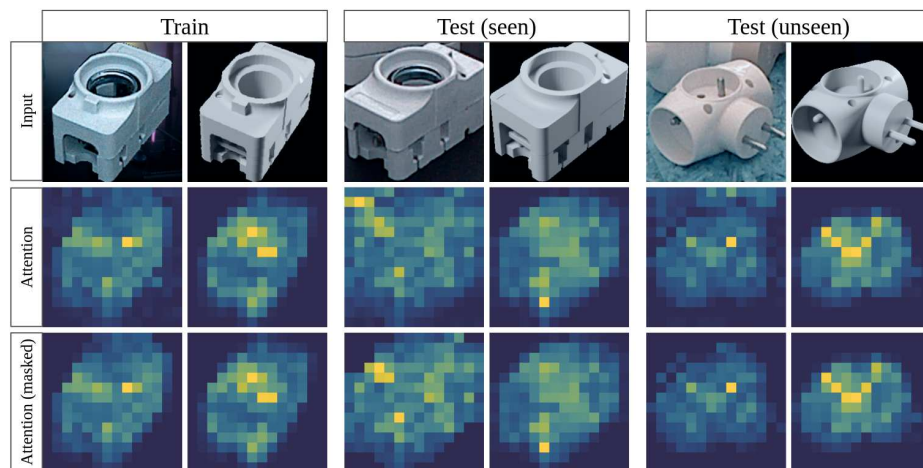
**Fig. 4. Self-Attention on LM/LM-O** Visualized is the self-attention of the first head of the last self-attention layer using the positional tokens as input.

Figures 4 and 5 visualize self-attentions maps on the training and test sets of LM/LM-O and T-LESS, respectively. The same projection mechanism as in [5] is used.



On LM/LM-O, Figure 4, ViT-s effectively learns to encode relevant features of the seen objects. The unseen test case shows that the learned self attentions not only transfer the concept of objectness to unseen objects, but also manages to distinguish relevant from irrelevant feature map locations.

On T-Less, Figure 5, object crops often show dataset objects in front or behind the query object, as is visualized in the seen and unseen test images. Cropping the feature map using the template’s mask is important in order to improve matching accuracy.



**Fig. 5. Self-Attention on T-LESS** Visualized is the self-attention of the first head of the last self-attention layer using the positional tokens as input.

## 5 Conclusion

This work presents diverse empirical analyses for using ViTs for self-supervised template matching for 3D pose retrieval. The presented findings are threefold. Using ViTs for deep template matching improves matching accuracy for seen and novel objects, in comparison to CNNs. Using pre-trained ViTs in a zero-shot fashion, that is without fine-tuning, already exhibits strong matching accuracy. Depending on the object set and metric used for evaluation, even improving over using a similar, fine-tuned CNN-based approach. For the problem of self-supervised synthetic template to real query object matching the network architecture is different to a comparable CNN approach and to self-supervised ViTs for image classification. In comparison to CNNs, ViTs benefit more from pre-training due to their feature extraction being more general. And in comparison to self-supervised ViTs for image classification, large, multi-layered projector heads are detrimental to the matching accuracy on novel objects. We hypothesize that

this occurs due to the stronger overfitting of deeper heads on the seen examples during fine-tuning, in turn harming the generality of the features learned during pre-training. Future work will investigate how to effectively exploit the features learned during ViT pre-training.

## References

1. Aing, L., Lie, W.N., Lin, G.S.: Faster and finer pose estimation for multiple instance objects in a single rgb image. *Image and Vision Computing* **130**, 104618 (2023). <https://doi.org/https://doi.org/10.1016/j.imavis.2022.104618>, <https://www.sciencedirect.com/science/article/pii/S0262885622002475>
2. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. *arXiv preprint arXiv:1607.06450* (2016)
3. Balntas, V., Doumanoglou, A., Sahin, C., Sock, J., Kouskouridas, R., Kim, T.K.: Pose guided rgb-d feature learning for 3d object pose estimation. In: *Proceedings of the IEEE international conference on computer vision*. pp. 3856–3864 (2017)
4. Brachmann, E., Krull, A., Michel, F., Gumhold, S., Shotton, J., Rother, C.: Learning 6D object pose estimation using 3D object coordinates. In: *Proceedings of the European Conference on Computer Vision*. pp. 536–551 (2014)
5. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 9650–9660 (2021)
6. Chechik, G., Sharma, V., Shalit, U., Bengio, S.: Large scale online learning of image similarity through ranking. *Journal of Machine Learning Research* **11**(3) (2010)
7. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: *International conference on machine learning*. pp. 1597–1607. PMLR (2020)
8. Chen, X., Xie, S., He, K.: An empirical study of training self-supervised vision transformers. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 9640–9649 (2021)
9. Dede, M.A., Genc, Y.: Object aspect classification and 6dof pose estimation. *Image and Vision Computing* **124**, 104495 (2022). <https://doi.org/https://doi.org/10.1016/j.imavis.2022.104495>, <https://www.sciencedirect.com/science/article/pii/S026288562200124X>
10. Denninger, M., Sundermeyer, M., Winkelbauer, D., Zidan, Y., Olefir, D., Elbadrawy, M., Lodhi, A., Katam, H.: Blenderproc. *CoRR* **abs/1911.01911** (2019), <http://arxiv.org/abs/1911.01911>
11. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
12. Drost, B., Ulrich, M., Navab, N., Ilic, S.: Model globally, match locally: Efficient and robust 3d object recognition. In: *2010 IEEE computer society conference on computer vision and pattern recognition*. pp. 998–1005. Ieee (2010)
13. Dwibedi, D., Misra, I., Hebert, M.: Cut, paste and learn: Surprisingly easy synthesis for instance detection. In: *Proceedings of the IEEE international conference on computer vision*. pp. 1301–1310 (2017)
14. Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., He, K.: Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677* (2017)

15. Grill, J.B., Strub, F., Alché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al.: Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems* **33**, 21271–21284 (2020)
16. Gutmann, M., Hyvärinen, A.: Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. pp. 297–304. *JMLR Workshop and Conference Proceedings* (2010)
17. Hartley, R., Zisserman, A.: *Multiple view geometry in computer vision*. Cambridge university press (2003)
18. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
19. Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415* (2016)
20. Hinterstoisser, S., Lepetit, V., Ilic, S., Holzer, S., Bradski, G., Konolige, K., Navab, N.: Model based training, detection and pose estimation of texture-less 3D objects in heavily cluttered scenes. In: *Proceedings of the Asian Conference on Computer Vision*. pp. 548–562 (2012)
21. Hodan, T., Barath, D., Matas, J.: Epos: Estimating 6d pose of objects with symmetries. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 11703–11712 (2020)
22. Hodan, T., Haluza, P., Obdržálek, Š., Matas, J., Lourakis, M., Zabulis, X.: T-less: An rgb-d dataset for 6d pose estimation of texture-less objects. In: *2017 IEEE Winter Conference on Applications of Computer Vision*. pp. 880–888. IEEE (2017)
23. Hodan, T., Michel, F., Brachmann, E., Kehl, W., GlentBuch, A., Kraft, D., Drost, B., Vidal, J., Ihrke, S., Zabulis, X., et al.: Bop: Benchmark for 6d object pose estimation. In: *Proceedings of the European conference on computer vision (ECCV)*. pp. 19–34 (2018)
24. Hodaň, T., Zabulis, X., Lourakis, M., Obdržálek, Š., Matas, J.: Detection and fine 3d pose estimation of texture-less objects in rgb-d images. In: *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. pp. 4421–4428. IEEE (2015)
25. Hou, T., Ahmadyan, A., Zhang, L., Wei, J., Grundmann, M.: Mobilepose: Real-time pose estimation for unseen objects with weak shape supervision. *arXiv preprint arXiv:2003.03522* (2020)
26. Huang, L., Hodan, T., Ma, L., Zhang, L., Tran, L., Twigg, C., Wu, P.C., Yuan, J., Keskin, C., Wang, R.: Neural correspondence field for object pose estimation. In: *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part X*. pp. 585–603. Springer (2022)
27. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: *International conference on machine learning*. pp. 448–456. *pmlr* (2015)
28. Jiang, Z., Wang, X., Huang, X., Li, H.: Triangulate geometric constraint combined with visual-flow fusion network for accurate 6dof pose estimation. *Image and Vision Computing* **108**, 104127 (2021). <https://doi.org/https://doi.org/10.1016/j.imavis.2021.104127>, <https://www.sciencedirect.com/science/article/pii/S0262885621000329>
29. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. *Communications of the ACM* **60**(6), 84–90 (2017)

30. Labbé, Y., Manuelli, L., Mousavian, A., Tyree, S., Birchfield, S., Tremblay, J., Carpentier, J., Aubry, M., Fox, D., Sivic, J.: Megapose: 6d pose estimation of novel objects via render & compare. In: 6th Annual Conference on Robot Learning (2022)
31. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: Proceedings of the 27th international conference on machine learning (ICML-10). pp. 807–814 (2010)
32. Nguyen, V.N., Hu, Y., Xiao, Y., Salzmann, M., Lepetit, V.: Templates for 3d object pose estimation revisited: Generalization to new objects and robustness to occlusions. In: Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2022)
33. Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018)
34. Park, K., Patten, T., Vincze, M.: Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (Oct 2019)
35. Park, K., Patten, T., Vincze, M.: Neural object learning for 6d pose estimation using a few cluttered images. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16. pp. 656–673. Springer (2020)
36. Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N., Ku, A., Tran, D.: Image transformer. In: International conference on machine learning. pp. 4055–4064. PMLR (2018)
37. Patten, T., Park, K., Vincze, M.: Dgcm-net: dense geometrical correspondence matching network for incremental experience-based robotic grasping. *Frontiers in Robotics and AI* **7**, 120 (2020)
38. Peng, S., Liu, Y., Huang, Q., Zhou, X., Bao, H.: Pvnet: Pixel-wise voting network for 6dof pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2019)
39. Remus, A., D’Avella, S., Felice, F.D., Tripicchio, P., Avizzano, C.A.: i2c-net: Using instance-level neural networks for monocular category-level 6d pose estimation. *IEEE Robotics and Automation Letters* **8**(3), 1515–1522 (2023). <https://doi.org/10.1109/LRA.2023.3240362>
40. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International journal of computer vision* **115**(3), 211–252 (2015)
41. Salimans, T., Kingma, D.P.: Weight normalization: A simple reparameterization to accelerate training of deep neural networks. *Advances in neural information processing systems* **29** (2016)
42. Shugurov, I., Li, F., Busam, B., Ilic, S.: Osop: A multi-stage one shot object pose estimation framework. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6835–6844 (2022)
43. Sun, H., Wang, T., Yu, E.: A dynamic keypoint selection network for 6dof pose estimation. *Image and Vision Computing* **118**, 104372 (2022). <https://doi.org/https://doi.org/10.1016/j.imavis.2022.104372>, <https://www.sciencedirect.com/science/article/pii/S0262885622000014>
44. Sundermeyer, M., Durner, M., Puang, E.Y., Marton, Z.C., Vaskevicius, N., Arras, K.O., Triebel, R.: Multi-path learning for object pose estimation across domains. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 13916–13925 (2020)

45. Sundermeyer, M., Hodan, T., Labbe, Y., Wang, G., Brachmann, E., Drost, B., Rother, C., Matas, J.: Bop challenge 2022 on detection, segmentation and pose estimation of specific rigid objects. arXiv preprint arXiv:2302.13075 (2023)
46. Sundermeyer, M., Marton, Z.C., Durner, M., Brucker, M., Triebel, R.: Implicit 3d orientation learning for 6d object detection from rgb images. In: Proceedings of the european conference on computer vision (ECCV). pp. 699–715 (2018)
47. Thalhammer, S., Leitner, M., Patten, T., Vincze, M.: Pyrapose: Feature pyramids for fast and accurate object pose estimation under domain shift. In: 2021 IEEE International Conference on Robotics and Automation (ICRA). pp. 13909–13915. IEEE (2021)
48. Thalhammer, S., Patten, T., Vincze, M.: Cope: End-to-end trainable constant runtime object pose estimation. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 2860–2870 (2023)
49. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International conference on machine learning. pp. 10347–10357. PMLR (2021)
50. Wang, G., Manhardt, F., Tombari, F., Ji, X.: Gdr-net: Geometry-guided direct regression network for monocular 6d object pose estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16611–16621 (2021)
51. Wang, H., Sridhar, S., Huang, J., Valentin, J., Song, S., Guibas, L.J.: Normalized object coordinate space for category-level 6d object pose and size estimation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2642–2651 (2019)
52. Wohlhart, P., Lepetit, V.: Learning descriptors for object recognition and 3d pose estimation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3109–3118 (2015)
53. Xiao, J., Hays, J., Ehinger, K.A., Oliva, A., Torralba, A.: Sun database: Large-scale scene recognition from abbey to zoo. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 3485–3492 (June 2010). <https://doi.org/10.1109/CVPR.2010.5539970>
54. You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., Hsieh, C.J.: Large batch optimization for deep learning: Training bert in 76 minutes. arXiv preprint arXiv:1904.00962 (2019)
55. Zhang, X., Jiang, Z., Zhang, H.: Real-time 6d pose estimation from a single rgb image. *Image and Vision Computing* **89**, 1–11 (2019). <https://doi.org/https://doi.org/10.1016/j.imavis.2019.06.013>, <https://www.sciencedirect.com/science/article/pii/S0262885619300964>
56. Zhang, X., Jiang, Z., Zhang, H.: Out-of-region keypoint localization for 6d pose estimation. *Image and Vision Computing* **93**, 103854 (2020). <https://doi.org/https://doi.org/10.1016/j.imavis.2019.103854>, <https://www.sciencedirect.com/science/article/pii/S0262885619304470>