

Real-time texturing for 6D object instance detection from RGB Images

Pavel Rojtberg *
Fraunhofer IGD, Darmstadt
TU Darmstadt

Arjan Kuijper †
Fraunhofer IGD, Darmstadt
TU Darmstadt

ABSTRACT

For object detection, the availability of color cues strongly influences detection rates and is even a prerequisite for many methods. However, when training on synthetic CAD data, this information is not available. We therefore present a method for generating a texture-map from image sequences in real-time. The method relies on 6 degree-of-freedom poses and a 3D-model being available. In contrast to previous works this allows interleaving detection and texturing for upgrading the detector on-the-fly. Our evaluation shows that the acquired texture-map significantly improves detection rates using the LINEMOD [5] detector on RGB images only. Additionally, we use the texture-map to differentiate instances of the same object by surface color.

Index Terms: I.2.10 [ARTIFICIAL INTELLIGENCE]: Vision and Scene Understanding—Modeling and recovery of physical attributes; I.5.5 [PATTERN RECOGNITION]: Implementation—Interactive systems;

1 INTRODUCTION

In recent years there has been great progress on the task of object detection and 6D pose estimation by means of convolutional neural networks (CNN) [15, 7]. Combined with a successive local refinement method [10, 14], it is now possible to obtain a precise object pose from a single RGB image only.

However, these methods take advantage of 3D scans of the target objects to generate training data. The 3D scans not only provide geometry but also surface information. Although scans are reasonably easy to acquire [11], the need of a 3D scan restricts the availability of the methods.

For many use-cases only a model created with computer aided design (CAD) software is available. Such models are coloured by semantics (e.g. engine, wheel) instead of material appearance (e.g. metal, rubber). Therefore, when using CAD reference geometry for object detection, one can not rely on surface colour or texture. This scenario is common in industrial environments or with 3D printing, where a manifold of materials can be used to represent the same CAD geometry.

In these cases depth-only variants of some algorithms [15, 7] can be used — however at the cost of degraded performance.

This work therefore focuses on the real-time acquisition of surface texture data and dynamic detection-model augmentation. This allows capturing and using the surface color information on-the-fly.

In the context of real-time surface reconstruction there is notably the work by Whelan et al. [17], who extends the Kinect Fusion [11] algorithm to colors. For this, an additional 3D color volume of the same size as the geometry voxel grid is used. This means that the surface resolution is tied to the geometry resolution and the memory consumption has cubic complexity. Furthermore, RGB-D data is required. The mentioned 3D scans are typically acquired by variations of [11] and consequently store surface information as

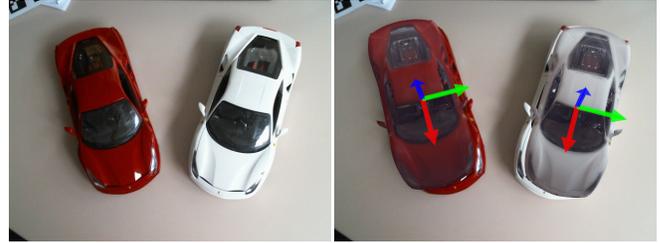


Figure 1: Our extension of LINEMOD [5] is able to correctly detect multiple instances of the same object based on surface color. The corresponding 6D poses are visualized using the textured mesh.

vertex colors. Here, [18] improve surface resolution by subdividing the scanned geometry. This is similar to using texture maps but results in inhomogeneous sampling and inefficient storage.

On the other hand, in the context of 3D scanning [3], 2D textures are often used which only have quadratic storage complexity and allow decoupling surface resolution from geometric resolution. More specifically, structure-from-motion based methods [16] only require RGB images to reconstruct both surface and geometry.

However, these methods operate on a-priori recorded data-sets which prevents real-time operation. Notably, the global optimization step alone, as employed by those methods, takes up to several hours. Our method in contrast operates in real-time while only requiring RGB frames to incrementally generate a 2D texture. To this end, we assume all geometry as fixed and given and only optimize locally for color consistency. The closest method to our work is by Magnenat et al. [9], who also map a 2D camera image to texture in real-time. However, their work specifically only addresses a single view and focuses on the in-painting aspect.

To employ our texturing method for object detection, we build upon the LINEMOD detection framework by [5]. They employ a two-stage, handcrafted feature descriptor, specifically tuned for texture-less object detection. In the first step gradient templates (DOT) are matched to the input image in a sliding window fashion, which capture the contour of an object. In a successive outlier-rejection step, surface color is used to filter implausible matches, based on the interior color.

Even though this no longer provides state-of-the-art detection performance [15, 7], the internal separation allows computing the DOT features on CAD geometry only and add the surface color at run-time. This is generally not possible with deep learning based approaches, which rely on surface color being available during training. Efforts to train on an abstract representation [12] to allow for different object appearances, typically result in a degraded performance compared to training on real images. In contrast, our extension of [5] improves its performance, while allowing to differentiate several instances of the same geometric object by their surface properties.

Based on the above, our key contributions are;

1. an incremental, real-time texture-map extraction pipeline and
2. efficient integration of texture-maps for object instance recognition.

*pavel.rojtberg@igd.fraunhofer.de

†arjan.kuijper@igd.fraunhofer.de

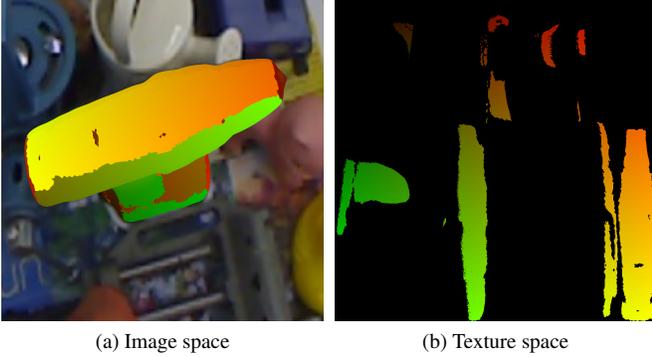


Figure 2: We are mapping from image to texture space, which is the reverse direction compared to rendering. Texture coordinates are encoded as red-green.

This paper is structured as follows: in Section 2 we present our texture extraction algorithm in detail. Section 3 then describes the application of the texture-maps for object instance recognition, while in Section 4 the use of texture-maps for object detection is evaluated using the LINEMOD dataset [5].

We conclude with Section 5 giving a summary of our results and discussing the limitations and future work.

2 TEXTURE EXTRACTION

In rendering, the process of "texture mapping" consist of the following two steps

1. creating a mapping from a texture to the surface of a 3D model and
2. projecting the model and simultaneously mapping the texture into a 2D image.

The first step is also called "texture atlas creation" and is typically performed by an artist during mesh creation. In contrast, our focus is the reverse direction, namely mapping from a 2D image of a projected 3D model back to the surface image as specified by the texture atlas (see Figure 2).

Generally texture atlas coordinates are not included in CAD data and therefore have to be generated. However, automatic texture atlas generation is still an active area of research [8] and outside of the scope of this work. Here, we just use the angle-based "Smart UV Project" algorithm implemented in the Blender toolset (v2.79b) to generate the texture atlas and instead focus on the second step of texture mapping.

In the remainder of this section we first discuss a simple exposure normalization scheme, before we present our texture extraction method in detail and finally turn to merging multiple views into one texture. The full pipeline is illustrated in Figure 3.

2.1 Exposure normalization

As our method does not explicitly compensate for different exposure times we pre-process the image stream to homogenize the brightness. For this we use the first captured frame as reference and modify the successive frames to match its brightness and contrast levels.

Here we follow the idea of Reinhard et al. [13] of adapting an input image \mathbf{I} to match a reference image as

$$\mathbf{I}_n = \frac{\sigma_{ref}}{\sigma_I} \cdot (\mathbf{I} - \mu_I) + \mu_{ref} \quad (1)$$

where μ_I, σ_I and μ_{ref}, σ_{ref} are the mean and variance of the input image and the reference image, respectively.

However, whereas [13] apply the transfer for all channels in the Lab color space, we only apply it to the luma component Y in the YUV color space as we explicitly want to preserve the chrominance information.

This step is omitted if the exposure can be fixed during capturing.

2.2 Texture-space to image-space mapping

Texture mapping can be formalized as follows: given a triangulated mesh, each vertex $\mathbf{v}_i = [X, Y, Z, 1]$ with an associated texture coordinate $\mathbf{t}_i = [u, v]$ is projected into the current view by a world-to-image transform \mathbf{P} as $\mathbf{p}_i = \mathbf{P} \cdot \mathbf{v}_i$. Here $\mathbf{p} = [x, y, 1]$ is a normalized pixel location in the image \mathbf{I} .

On the interior of the triangle formed by $(\mathbf{t}_i, \mathbf{t}_j, \mathbf{t}_k)$, a texture coordinate $\hat{\mathbf{t}}$ is interpolated and used for lookup in texture \mathbf{T} as

$$\mathbf{I}(\mathbf{p}) = \mathbf{T}(\hat{\mathbf{t}}). \quad (2)$$

This mapping is continuous in texture space and therefore allows for bi-linear interpolation to avoid aliasing artifacts.

For texture extraction however we are interested in the reverse mapping, namely

$$\mathbf{T}(\mathbf{t}) = \mathbf{I}(\hat{\mathbf{p}}). \quad (3)$$

Instead of iterating over the mesh topology as defined by \mathbf{v}_i in 3D, we now iterate over \mathbf{t}_i as defined by the texture-atlas in 2D. Conversely, we now require a continuous value of $\hat{\mathbf{p}}$ in image space for lookup. This is computed by interpolating in the triangle formed by $(\mathbf{p}_i, \mathbf{p}_j, \mathbf{p}_k)$, of which each point is obtained as above by $\mathbf{p}_i = \mathbf{P} \cdot \mathbf{v}_i$.

Here, visibility must be explicitly computed; with equation (2), we implicitly assumed overlapping points to be resolved by a depth-test, only retaining the points closest to the camera. This can no longer be exploited, as points do not overlap in the texture space.

To handle visibility we therefore introduce an additional depth buffer and render depth from the camera view. This allows comparing the depth of an interpolated coordinate $\hat{\mathbf{p}}$ to the actually visible depth value. However, this leads to aliasing; with non-planar objects the view resolution cannot be adapted to match the texture space resolution.

To remedy the aliasing artifacts we apply techniques from the shadow mapping domain [2] where the same problem occurs when a scene is rendered from a shadow camera and an observer camera view. Particularly, we

1. focus the camera on the object bounding box to increase the sampling rate in image space and
2. apply a slope-scale depth-bias to account for the remaining differences in sampling rates during visibility testing.

The latter is especially important; as the texture atlas has a higher sampling rate than the depth buffer, several points $\hat{\mathbf{p}}$, interpolated in the texture space, map to the same point \mathbf{p} in the image depth-buffer. At steep angles $\hat{\mathbf{p}}$ has a strong depth variation and thus neighboring points alternatively fail and pass the visibility test when compared to a single reference value (see Figure 4b).

To account for this we store a biased depth that allows for a sampling offset of 1px in image space. The bias b depends on the depth slope dz per pixel dx and the minimal depth buffer resolution r as:

$$b = \frac{dz}{dx} + r. \quad (4)$$

The bias is large in steep regions while minimal for faces parallel to the camera. This computation can be implemented efficiently on the GPU by using e.g. *glPolygonOffset*. The effect can be observed

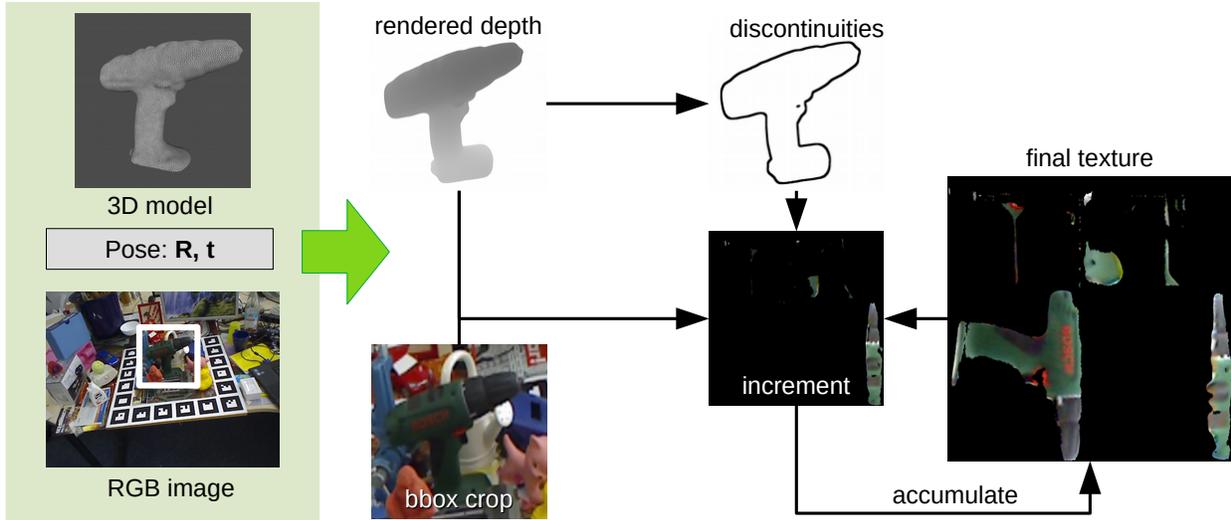


Figure 3: Our texture extraction pipeline. Given a 3D model and its pose in a RGB frame, we first render the depth to determine visibility. Image regions around depth discontinuities are discarded as they are unreliable. Next an texture-increment is extracted and a per-pixel score is computed to decide whether to merge the visible pixels into the final texture. Only the following buffers are required on the GPU; "final texture", "increment" and "discontinuities".

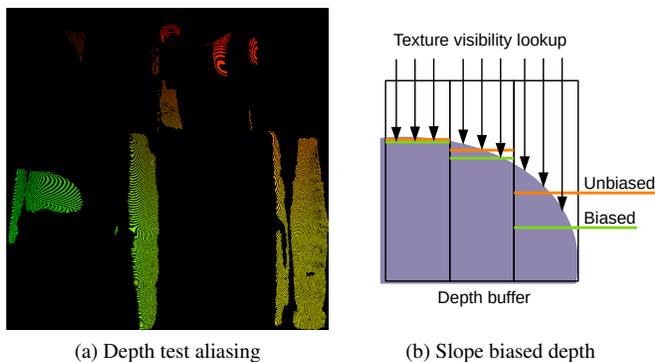


Figure 4: We store slope-scaled biased depth values to avoid aliasing errors during the visibility test.

by comparing Figure 4a and Figure 2b. This allows us to map each visible pixel from a single image into the texture to record the object surface. The resulting reconstruction can already be applied for detecting the object in similar views (see Section 3).

2.3 Merging multiple views

Generally the object surface is only partially visible from a single view and therefore multiple images are needed to reconstruct the full texture.

Assuming that the same texture point \mathbf{t} will be observed in different images as $\mathbf{c}_0, \dots, \mathbf{c}_N$ where $\mathbf{c}_i = \mathbf{I}_i(\hat{\mathbf{p}}_i)$, we discard edge-pixels at object boundaries or strong depth discontinuities. These measurements are unreliable as they might come from different surfaces due to pose imprecisions and limited camera resolution. Instead, we aim for a view where \mathbf{t} is not at an observed edge. A pixel is considered to be part of an edge if the depth change is larger than 10% of the object diameter. All points in a 5px neighborhood of an edge-pixel are discarded as well (see Figure 6a).

To combine multiple valid observations \mathbf{c}_i of \mathbf{t} , we define score

s that, inspired by [3], weighs each observation by the distance to camera d and the angle α between surface normal and view direction as

$$s = \cos \alpha \cdot (1 - d), \quad (5)$$

where d is assumed in normalized device coordinates ranged $[0; 1]$ and α is computed based on the interpolated surface normal, which can be defined per vertex \mathbf{v}_i (e.g. for a sphere) and therefore is not required to be constant for a single face.

Using s we implemented two merging strategies; a weighted arithmetic mean

$$\mathbf{t} = \frac{s_0 \cdot \mathbf{c}_0 + \dots + s_n \cdot \mathbf{c}_n}{s_0 + \dots + s_n}, \quad (6)$$

and only retaining the best view

$$\mathbf{t} = \arg \max_{\mathbf{c}_i} \{s_0, \dots, s_n\}. \quad (7)$$

Both equations can be efficiently implemented on the GPU using a single RGBA buffer for accumulation, as $RGBA = [\mathbf{c}, s]$.

Figure 5 shows exemplary results. Eq. (6) produces a smooth surface, while retaining more detail than vertex coloring. However, the averaging over slightly inaccurate object poses results in a loss of fine detail when compared to Eq. (7).

Using Eq. (7) on the other hand retains all details, but emphasizes inconsistencies in exposure or object pose as seams between neighboring increment texture-patches.

To alleviate this problem we blend increment-patches at their boundaries into the existing texture during accumulation. Instead of simply overwriting the texture content with the new maximum, we compute the distance transform to the patch boundaries over a 5x5px support using the L2 norm. Using the distance we then linearly interpolate between the old and the new color value \mathbf{c} and pixel score s .

Figure 6b shows an increment-patch for Figure 6a, projected onto the object. Note the gradient at the edges, which is linear in texture space.

The blending not only produces visually more pleasing results (compare Figures 5c and 5d), but is crucial for computing the LINEMOD descriptor which relies on local gradient orientation.

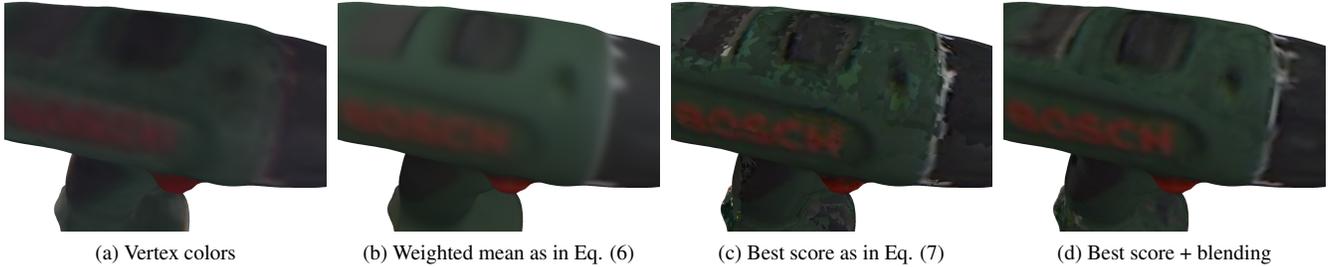


Figure 5: Exemplary surface color reconstructions of the "Driller" object Texture merging strategies using (a) KinectFusion and (b, c, d) variations of our algorithm.

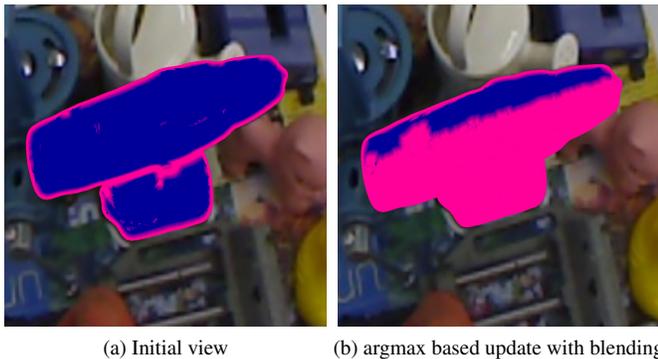


Figure 6: Merge-maps of two successive frames when using Eq. (7). Valid pixels are colored blue.

3 OBJECT INSTANCE DETECTION

In this section we describe how to employ the extracted textures for object instance detection i.e. differentiating multiple instances of the same object. Here, we extend the color based outlier rejection of [5] to multiple color hypotheses to simultaneously perform classification.

The idea of color based outlier rejection in [5] is to store the expected color of the object projection alongside the LINEMOD template and at run-time count how many pixels in the camera frame have the expected color.

To make the check robust against lighting variations, they convert the images to the HSV colour space and compare only the hue component. However, hue does not cover the colors black ($V = 0$) and white ($V = 1, S = 0$). Therefore, these are mapped to blue and yellow respectively, which completes the color based descriptor (see Figure 7a).

To extend this scheme for object instance detection as well as for on-the-fly recorded textures, we separate the expected color from the expected surface visibility. To this end, we store the texture coordinates of the object projection (compare Figure 2a) instead of storing the expected color directly. The template surface-texture is stored separately. At runtime we now use the texture coordinates to perform a lookup into the template-texture to retrieve the expected color, which gives us the same information as in [5].

However, it is now possible to easily swap the surface-texture to globally change the expected colors. Here a live-reconstructed texture can provide more accurate template colors and notably multiple template-textures can be used for object instance detection (see Figure 7b).

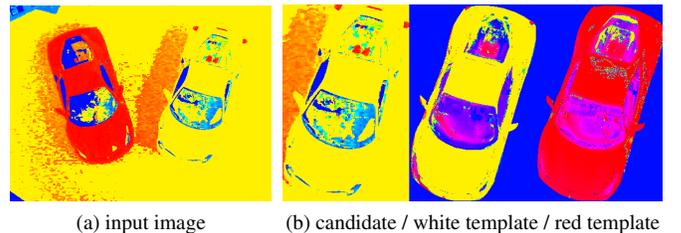


Figure 7: Hue based instance detection. The input image (Figure 1) is cropped based on the template bounding box and compared to a set of hue templates.

Finally, the outlier rejection scheme needs a slight modification for classification. Instead of returning the first inlier based on the expected color, it needs to allow multiple matches without repetition. For this, after finding an inlier, only the corresponding texture-template is removed and the remaining candidates are checked until all template-textures are found or all candidates are rejected.

While this is an integral part of the LINEMOD pipeline, it can be optionally integrated as a post-processing step to an CNN based architecture that is capable to abstract the object appearance to some degree. E.g. it can be executed after non-maximum-suppression in [15] to compute agreement with the color template.

4 EVALUATION

The presented method is evaluated in the context of object detection. To this end we train the LINE2D variant of the LINEMOD detector on the corresponding dataset [5]. The dataset does not contain views specifically for surface reconstruction and thus represents reconstruction during detection well. We use the publicly available LINEMOD implementation in OpenCV.

There are 15 sequences for different objects, consisting of RGB-D frames with ground-truth poses and recorded at distances of 65cm-115cm. We select a subset of 8 objects for which a 3D mesh is available and that are large enough to provide a reasonable texture resolution. The meshes included in the dataset were recorded using a variation of KinectFusion [11] and thus encode surface information as vertex-colors.

We apply our texturing algorithm on each sequence using the ground-truth poses to merely simulate a tracking algorithm for better reproducibility. Then we train LINEMOD on synthetic renderings using the generated textures as well as included vertex-colors as a baseline. We parametrize training and testing as [5], particularly;

- We use 89 views on the upper hemisphere around the object,

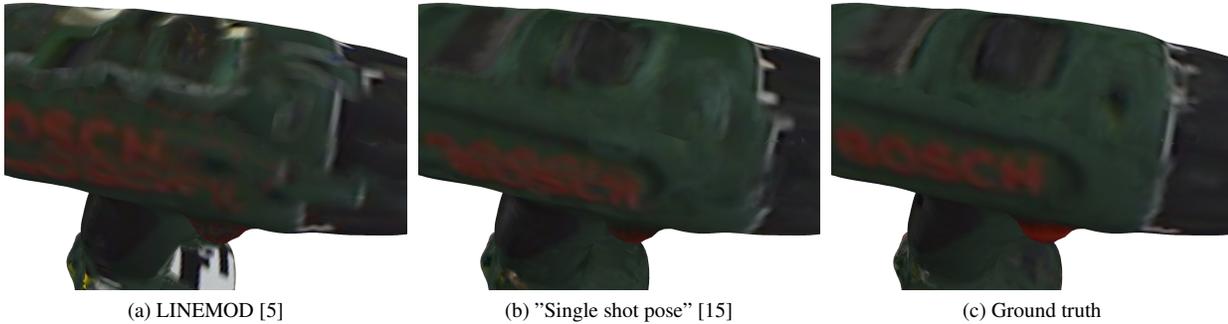


Figure 8: Qualitative results for on-the-fly surface color reconstructions of the "driller" object in relation to different pose detection methods.

Table 1: True positive rates on the linemod dataset with different training data. The ambient occlusion (AO) variant does not include any outlier rejection.

Object	AO	vertexcolor	texture (7)	texture (6)
benchvise	0.54	0.75	0.82	0.82
driller	0.15	0.43	0.63	0.54
iron	0.53	0.71	0.67	0.68
can	0.38	0.67	0.78	0.83
glue	0.07	0.21	0.17	0.17
cam	0.1	0.28	0.62	0.55
eggbox	0.47	0.6	0.79	0.79
holepuncher	0.2	0.62	0.59	0.65
average	0.3	0.53	0.64	0.63

derived by subdividing an icosahedron twice recursively.

- For each view there are 7 in-plane rotations with roll angles between -45° and 45° .
- Furthermore 6 distances, with 10cm increments, between 65cm and 115cm are used.
- During color based outlier rejection we discard candidates where less than 70% of the pixels have the expected color. The threshold on per-pixel hue difference is set to 54° .

This results in a total of 3738 templates per object for training. However, in contrast to [5], we are only using RGB data without depth — therefore we do not restrict the color gradient features to the contour, but compute them on the interior as well.

For testing, we measure the true positive rate on the sequences. As in [4] we consider an object successfully detected when it is within a fixed radius r around the ground truth position. We globally set $r = 11\text{cm}$ in our experiments to allow for depth misclassification by one step.

For keeping interactive performance we only consider the first 30 LINEMOD candidates for matching and outlier rejection.

To simulate the CAD data use-case without any surface information available, we additionally perform training using a white diffuse material for all objects. For generating gradient features on the interior of the object, we use ambient occlusion (AO) [1] as a lighting approximation. Ambient occlusion is a purely geometrical method that is independent of actual light and surface properties. We skip the outlier-rejection step as no color information is available.

Table 1 shows the true positive rates for the variants mentioned above — as can be seen the texture based variants outperform the

vertex-color baseline of [5] by a margin of 10% on average. However, there are strong variations between the individual objects, therefore it remains inconclusive whether variant (6) or variant (7) of our algorithm is preferable.

Notably the AO variant cannot reach the performance of the other methods. With some objects where it even becomes unusable (e.g. driller, cam). This emphasizes the need of surface information for object detection.

4.1 Using noisy pose data

To evaluate the applicability of our method for on-the-fly texturing with noisy pose data, we additionally used the state-of-the-art "single shot pose" (SSP) detector [15] instead of relying on ground-truth poses.

Figure 8 shows qualitative results of texturing using ground-truth, SSP and LINEMOD poses. While the LINEMOD results only allow for for a rough color based outlier rejection, the results using SSP poses are very similar to using the ground truth. To further quantify this, we repeated the training of SSP using synthetic renderings of the "driller" object instead of using cross-validation as in the original paper. At this, we measured the true positive rate (TPR) using the 5cm, 5deg metric. Here, training with textured renderings (Fig. 8c) resulted in a TPR of 0.37. Using the imperfectly textured objects (Fig. 8b) resulted in a TPR of 0.34, which supports the qualitative impression. When training with vertex colored renderings only, the performance was significantly degraded, resulting in a TPR of 0.14.

4.2 Speed

The evaluation was performed on a notebook with an Intel i7-7700HQ CPU at 2.80GHz and an Intel HD 630 iGPU. The average time to accumulate one video frame into a 1024x1024 px sized texture is 2.69 ms. This allows running the texturing algorithm in parallel to tracking to reconstruct a texture on-the-fly.

The average time to perform a texture lookup as described in Section 3 is 0.82 ms using the software remap implementation in OpenCV. This step can be therefore applied generally without requiring GPU usage.

4.3 Multi instance detection

For the multi-instance detection we performed a qualitative analysis using a separate sequence where two toy cars are alternately and simultaneously visible. The surface colors are white and red which are adjacent in HSV space (white is mapped to yellow as described in section 3). Furthermore, the surface exhibits specular reflection which is not filtered during texturing.

Nevertheless, our method was able to robustly discriminate both objects (see Figure 1 and supplemental material¹).

5 CONCLUSION AND FUTURE WORK

We have presented a method for real-time texturing that can be used to improve detection on-the-fly. At this we have shown that texturing itself is crucial for detection of CAD data where no surface information is available. However, even for meshes where vertex-colors were previously available, our approach improves detection performance significantly. Furthermore, we successfully applied the resulting textures to extend LINEMOD for object-instance recognition. By interleaving detection and texture extraction it now becomes possible to extend detection algorithms by color cues on-the-fly.

Our method currently requires the camera exposure to be fixed or relies on a global exposure compensation approach which is error-prone. Here reading the actual camera exposure could be used for accurate exposure fusion of the images. The surface specularities could be explicitly considered during merging [6]. Currently we assume diffuse reflection, which systematically over-brightens specular surfaces. At this a plausibility test during merging could be used to reject implausible colors as caused by e.g. occlusion. As the LINEMOD detector is no longer state-of-the-art and further investigation is needed to similarly integrate our approach into an existing CNN based method. This will require to breaking up the end-to-end trained "black-box" to make the color information explicit.

REFERENCES

- [1] L. Bavoil and M. Sainz. Screen space ambient occlusion. *NVIDIA developer information: <http://developers.nvidia.com>*, 6, 2008.
- [2] S. Brabec, T. Annen, and H.-P. Seidel. Practical shadow mapping. *Journal of Graphics Tools*, 7(4):9–18, 2002.
- [3] M. Callieri, P. Cignoni, M. Corsini, and R. Scopigno. Masked photo blending: Mapping dense photographic data set on high-resolution sampled 3d models. *Computers & Graphics*, 32(4):464–473, 2008.
- [4] S. Hinterstoisser, S. Holzer, C. Cagniard, S. Ilic, K. Konolige, N. Navab, and V. Lepetit. Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 858–865. IEEE, 2011.
- [5] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab. Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes. In *Asian conference on computer vision*, pages 548–562. Springer, 2012.
- [6] J. Jachnik, R. A. Newcombe, and A. J. Davison. Real-time surface light-field capture for augmentation of planar specular surfaces. In *2012 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 91–97. IEEE, 2012.
- [7] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab. Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again. In *Proceedings of the International Conference on Computer Vision (ICCV 2017), Venice, Italy*, pages 22–29, 2017.
- [8] M. Li, D. M. Kaufman, V. G. Kim, J. Solomon, and A. Sheffer. Optcuts: joint optimization of surface cuts and parameterization. In *SIGGRAPH Asia 2018 Technical Papers*, page 247. ACM, 2018.
- [9] S. Magnenat, D. T. Ngo, R. W. Sumner, F. Zund, M. Ryffel, G. Noris, G. Rothlin, A. Marra, M. Nitti, P. Fua, et al. Live texturing of augmented reality characters from colored drawings. *IEEE Transactions on Visualization & Computer Graphics*, (11):1201–1210, 2015.
- [10] F. Manhardt, W. Kehl, N. Navab, and F. Tombari. Deep model-based 6d pose refinement in rgb. *arXiv preprint [arXiv:1810.03065](https://arxiv.org/abs/1810.03065)*, 2018.
- [11] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*, pages 127–136. IEEE, 2011.
- [12] J. Rambach, C. Deng, A. Pagani, and D. Stricker. Learning 6dof object poses from synthetic single channel images. In *2018 IEEE International Symposium on Mixed and Augmented Reality Adjunct (ISMAR-Adjunct)*, pages 164–169. IEEE, 2018.
- [13] E. Reinhard, M. Adhikhmin, B. Gooch, and P. Shirley. Color transfer between images. *IEEE Computer graphics and applications*, 21(5):34–41, 2001.
- [14] B. Seo, H. Park, J. Park, S. Hinterstoisser, and S. Ilic. Optimal Local Searching for Fast and Robust Textureless 3D Object Tracking in Highly Cluttered Backgrounds. *IEEE Transactions on visualization and computer graphics (TVCG)*, 2013.
- [15] B. Tekin, S. N. Sinha, and P. Fua. Real-time seamless single shot 6d object pose prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 292–301, 2018.
- [16] M. Waechter, N. Moehrle, and M. Goesele. Let there be color! large-scale texturing of 3d reconstructions. In *European Conference on Computer Vision*, pages 836–850. Springer, 2014.
- [17] T. Whelan, H. Johannsson, M. Kaess, J. J. Leonard, and J. McDonald. Robust real-time visual odometry for dense rgb-d mapping. In *Robotics and Automation (ICRA), 2013 IEEE International Conference on*, pages 5724–5731. IEEE, 2013.
- [18] Q.-Y. Zhou and V. Koltun. Color map optimization for 3d reconstruction with consumer depth cameras. *ACM Transactions on Graphics (TOG)*, 33(4):155, 2014.

¹<https://youtu.be/IB19rTXUot8>