

# Point Cloud Hand/Object Segmentation using Multi-Modal Imaging with Thermal and Color Data for Safe Robotic Object Handover

Yan Zhang <sup>1</sup>, Steffen Müller<sup>2</sup>, Benedict Stephan<sup>2</sup>, Horst-Michael Gross<sup>2</sup> and Gunther Notni <sup>1,3</sup>

<sup>1</sup> Group for Quality Assurance and Industrial Image Processing Technische Universität Ilmenau, Germany

<sup>2</sup> Neuroinformatics and Cognitive Robotics Lab Technische Universität Ilmenau, Germany

<sup>3</sup> Fraunhofer Institute for Applied Optics and Precision Engineering IOF Jena, Germany

\* Correspondence: yan.zhang@tu-ilmenau.de

**Abstract:** This paper presents an application of neural networks operating on multi-modal 3D data (3D point cloud, RGB, thermal) to effectively and precisely segment human hands and objects held in hand to realize a safe human-robot object hand over. We discuss the problems encountered for building a multi-modal sensor system, while the focus is on the calibration and alignment of a set of cameras including RGB, thermal, and NIR cameras. We propose the use of a copper-plastic chessboard calibration target with an internal active light source (near-infrared and visible light). By brief heating, the calibration target could be simultaneously and legibly captured by all cameras. Based on our multi-modal dataset captured by our sensor system, PointNet [1], PointNet++ [2] and RandLA-Net [3] are utilized to verify the effectiveness of applying multi-modal point cloud data for hand/object segmentation. These networks were trained on various data modes (XYZ, XYZ-T, XYZ-RGB and XYZ-RGB-T). The experimental results show a significant improvement in the segmentation performance of XYZ-RGB-T (mean Intersection over Union: 82.8% by RandLA-Net) compared to the other three modes (77.3% by XYZ-RGB, 35.7% by XYZ-T, 35.7% by XYZ), in which, it is worth mentioning that the Intersection over Union for the single class of hand achieves 92.6%.

**Citation:** Zhang Y.; Müller, St.; Stephan, B.; Gross, H.-M.; Notni, G. Point Cloud Hand/Object Segmentation using Multi-Modal Imaging with Thermal and Color Data for Safe Robotic Object Handover. *Sensors* **2021**, *1*, 0. <https://doi.org/>

Received:

Accepted:

Published:

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Copyright:** © 2021 by the authors. Submitted to *Sensors* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

**Keywords:** multi-modal imaging, thermal, deep neural network, hand segmentation, point cloud segmentation

## 1. Introduction

Nowadays, robot vision plays an important role in the robotics industry. To enable a robot to navigate or grasp objects as intelligently and safely as a human, a correct understanding of its working environment is a necessary prerequisite. For this task, currently there are many state-of-the-art solutions based on object detection, such as YOLO [4]. However, our work focuses on the vision system of an assistant robot which is used to transport objects between humans. In order to pick up the object from a human hand without injuring the person, the challenge is exact and efficient pixel-level segmentation and 3D representation of the object and obstacles in the interaction area. In this regard, it is not sufficient to separate hand and object with only a bounding box. Therefore, the discussion in this article will focus on hand/object segmentation.

To solve the segmentation problem, the current mainstream approaches can be classified into two categories.

The first one is color image segmentation based on texture information on the surface of objects. Extensive research has been done on this subject and some of the achievements are impressive, such as the MASK R-CNN network [5] or the PointRend network [6]. However, there are a number of difficulties in hand segmentation, such as the effect of lighting conditions, confusion with objects whose color resembles human skin, and the variety of skin tones.

The second category is 3D point cloud segmentation based on geometric features of objects. In this respect, some challenges such as the articulated nature of the human body, changes in appearance and partial occlusions [7] make hand segmentation in point clouds more difficult than in RGB images.

Although, the deep learning technology has repeatedly surprised in the field of image processing, the above mentioned particular difficulties for hand segmentation can never be solved completely. For example, in [8], the authors explicitly mention that their VGG-16 [9] based hand segmentation network (2D RGB segmentation) can achieve a 91.0% mean IoU (Intersection over Union) on their dataset. *If the hand has a complex interaction with other objects, such as holding a complex-shaped object in the hand, it is hard using their approach to detect the hand in the contact areas.* Nevertheless, the segmentation of real-world data seen in interactions with humans is just the core challenge for an assistant robot aiming to grasp objects from a human hand.

Since humans are warm-blooded, our body temperature stays almost constant while skin color, light conditions, and hand posture are varied. Therefore, body temperature is a more stable and robust feature for hand recognition or segmentation compared to RGB data alone. We propose to apply an additional LWIR camera (thermal camera) (LWIR: longwave infrared) to mitigate the problems for hand segmentation mentioned above. However, there are also some difficulties with thermal image segmentation. As mentioned in [10], for an outdoor intelligent surveillance system with a thermal camera, in summer or on a hot day, the contrast of human and background becomes very low and makes it difficult to distinguish human areas from the background in the thermal image. *This low contrast problem holds also for a couple of indoor scenarios e.g. in industrial facilities where there are differently tempered objects in the background.* In addition, an object that is held in hand for a longer time will become similar in temperature to the hand, lowering the contrast to the fingers as well. In this case, hand and object segmentation will also become tough and additional features are needed. In the research of Kim et al. [11], 2D multi-modal imaging with fusing LWIR and RGB-D images was used for first-person view hand segmentation, and their results of using a DeepLabV3+ [12] network showed that using LWIR there were 5% better hand IoU performance than using just RGB-D frames.

Therefore, in this work, we will explain a multi-modal 3D sensor system composed of a 3D sensor, an RGB camera, and a thermal camera, which is able to capture point cloud data with 7 channels (XYZ-RGB-T). None of these channels are all-purpose, but in combination the information of each channel compensate for their respective weaknesses. It is reasonable to expect that the multi-modal 3D data carries more potential features compared to 3D data alone, and that these complex features can be learned by a neural network as well. Besides that, the calibration and registration approach for the sensors will be described. By using this sensor system, a multi-modal dataset was captured in order to evaluate the performance of applying the multi-modal 3D data for hand and object segmentation. The state-of-the-art methods PointNet, PointNet++, and RandLA-Net were trained and compared on that dataset.

## 2. Related Work

In this section previous studies on the application of thermal imaging in the field of human recognition will be reviewed. Wang et al. [13] presented a thermal pedestrian detector, in which an edge feature (Shape Context Descriptor) and an Adaboost cascade classifier were adopted. Jeon et al. [10] showed for an outdoor surveillance thermal camera, that it is hard to segment the human body from the background, if the ambient temperature is similar to or higher than the human body temperature (e.g. in summer). To solve this problem, they attempt to do background subtraction using the sequence of thermal images and a pre-recorded background thermal image. However, both of the two studies require a fixed background as a prerequisite, which is not possible for a mobile robot application.



In the research of Setjo et al. [14], Haar cascade classifiers have been applied to detect human faces in thermal images, and a comprehensive evaluation was conducted with a thermal image dataset that was acquired with variation of human poses and environmental conditions. They showed, that precision and recall of human detection decreases with greater distance to the camera. In addition, the detection results were also affected by the orientation of the face. For such problems in [15] an integrated analysis for RGB-T (thermal) fusion was proposed to detect human skin using a skin segmentation algorithm (Skindiff) [16]. Their results indicated that the use of the fusion sensor system allows to work in environments with many warm objects. An RGB-T dataset and a discussion of the advantages of RGB-T fusion over single RGB or T, such as when objects of interest may not have easily discernible thermal signatures but have strong cues from RGB, can be found in [17]. In addition, there are a number of articles on this topic, such as [18] [19] [20]. However, all these articles are based on traditional image processing methods. That means that a few parameters or thresholds in the system need to be adjusted manually, and they are usually dependent on the varying camera environment or the state of the camera. For example, in [15] it is mentioned, that the response of a thermal camera depends on the up-time of their specific device, which has an effect on the human recognition rate as images become more saturated.

Palermo et al. [7] introduced a RGB-D-T dataset and used HOG (Histogram of oriented gradient) and random forests for human segmentation in an indoor environment. In [21], transfer learning of YOLO [4], a deep learning model, was performed on thermal images for human detection in a night environment. In [22], YOLO was applied on 6-channel 2D images containing RGB color and various geometric features (point density, difference of normal, and curvature). A further CNN (convolutional neural network) model named MCNet for 2D thermal image semantic segmentation of nighttime driving scenes has been published in [23]. For point cloud data, in [24] a PointNet [1] based hand segmentation network is explained but they are not using multi-modal data at all.

In summary, the above mentioned studies can be roughly categorized into three groups:

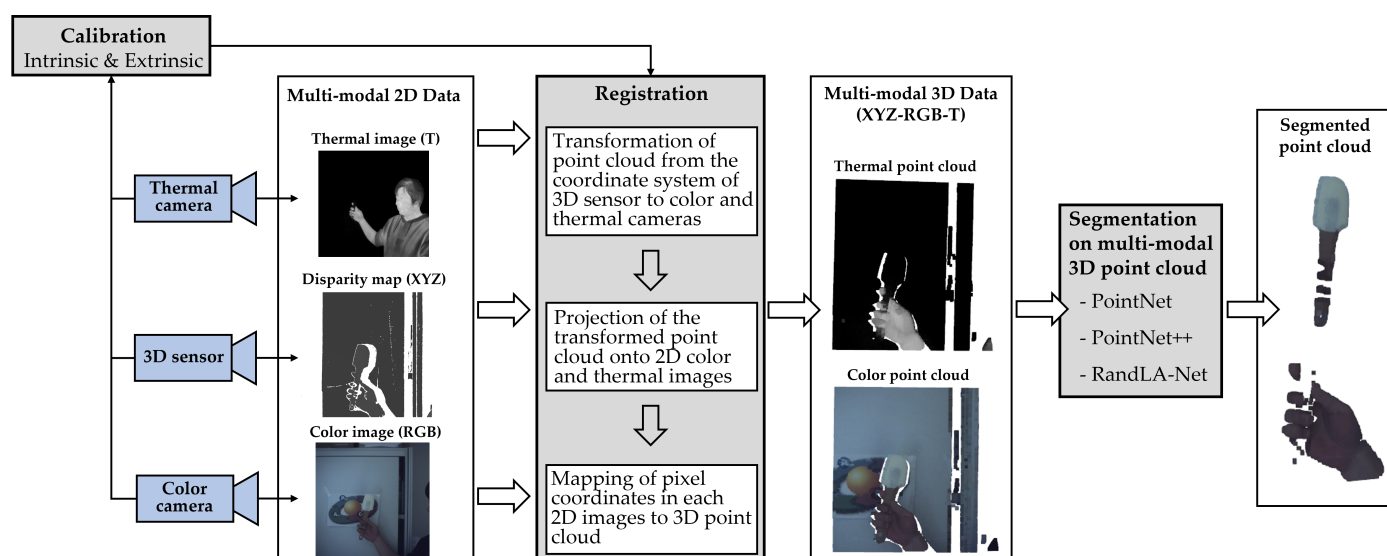
1. Using traditional methods to analyze multi-modal data for human recognition, such as [13] [10] [17].
2. Using a neural network approach to detect humans but not on multi-modal data, such as [24] [8].
3. Using a neural network approach to process multi-modal data for human detection, however almost all of them are regarding the task of autonomous driving in urban scenarios, such as [23] [21].

To our knowledge, there has not been a comprehensive study using deep learning technology and multi-modal 3D data to specifically address the problem of indoor human hand segmentation in point clouds for an assistant robot. Therefore, in this work we will provide a detailed discussion on this issue.

### 3. Method Overview

As shown in Figure 1, the entire pipeline for hand and object segmentation based on multi-modal 3D data is divided into 3 steps: Calibration, Registration and Segmentation.

- **Calibration:** For a multi-modal sensor system containing a 3D sensor, a color camera and a thermal camera, the intrinsic parameters of each sensor and the extrinsic parameters of color and thermal camera in respect to the 3D sensor should be calibrated. It will be explained in section 4.2.
- **Registration:** By using the intrinsic and extrinsic parameters, the color and thermal pixel values in 2D images should be mapped onto the 3D point cloud, in order to build a multi-modal 3D point cloud, in which each point integrates multi-modal information of color (RGB), temperature (T) and coordinates in 3D space (XYZ). It will be described in section 4.3.



**Figure 1.** Workflow for a hand/object segmentation approach using a multi-modal 3D sensor system containing a 3D sensor, a RGB camera and a thermal camera.

- Segmentation:** With the help of neural networks, potential multi-modal features hidden in the point cloud can be learned for secure and robust hand/object segmentation. In this work, PointNet [1], PointNet++ [2], and RandLA-Net [3] have been used as segmentation approaches. These approaches will be briefly discussed in section 5. In section 7, comparative experiments (training on various data modes of XYZ, XYZ-T, XYZ-RGB and XYZ-RGB-T) will be provided to evaluate the application of multi-modal 3D data for hand/object segmentation.

#### 4. Sensor System

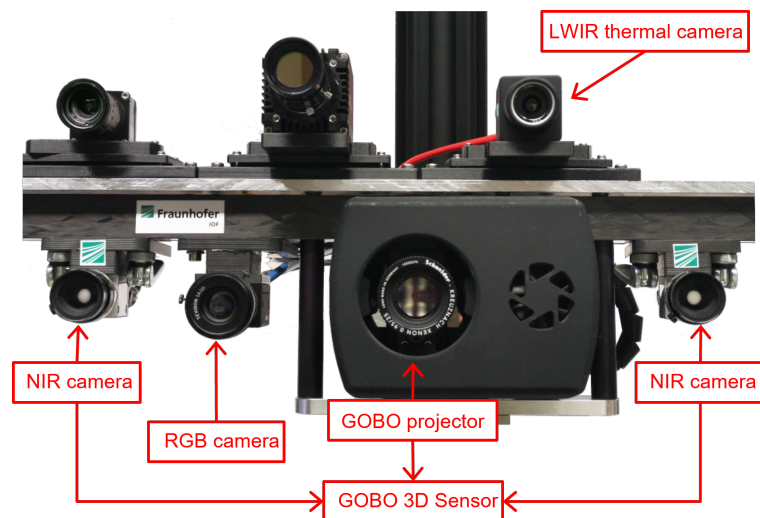
##### 4.1. Multi-modal Sensors

The Figure 2 shows our multi-modal 3D imaging system, which consists of an active stereo-vision 3D sensor based on GOBO (Goes Before Optics) projection [25], a color camera (FLIR Grasshopper3 [26]) and a thermal camera (FLIR A35 [27]). It has been used to record a multi-modal dataset (XYZ-RGB-T) of humans holding objects as described in Sec. 6.

The 3D sensor utilizes two NIR (near-infrared) (850 nm) cameras and a NIR (850 nm) GOBO projector [28] for projecting a temporally varying aperiodic sinusoidal pattern into the scene. By means of that pattern, corresponding 3D points can be identified in an image sequence, enabling a robust reconstruction of pixel disparities and therefore depth of point cloud points. The GOBO system yields point clouds with 0.32–1.18 mm resolution and roughly 0.15 mm measurement error in a relatively small field of view of  $48^\circ \times 44^\circ$  in a limited range of 0.4–1.5 m at 36Hz. The FLIR Grasshopper3 provides color images with a resolution of 2048x2048 pixels in the field of view of  $50^\circ \times 50^\circ$  at 90Hz, and the FLIR A35 captures thermal images of 320x256 pixels in a range of  $-25^\circ\text{C}$  to  $135^\circ\text{C}$  at 60Hz. It has a field of view of  $63^\circ \times 50^\circ$  and therefore covers the whole point cloud area as does the rgb camera. Table 1 shows the technical data of the GOBO 3D sensor and the additional cameras.

##### 4.2. Calibration Target

In order to fuse the image data of each camera, precise mapping of pixel coordinates in each image to the 3D point cloud is required and the camera system needs to be calibrated. Normally, a printed checkerboard pattern can be used as a calibration target. It works fine for RGB and NIR cameras, but for thermal camera it is challenging. Because black and white grids of the printed pattern have almost the same emittance in LWIR,



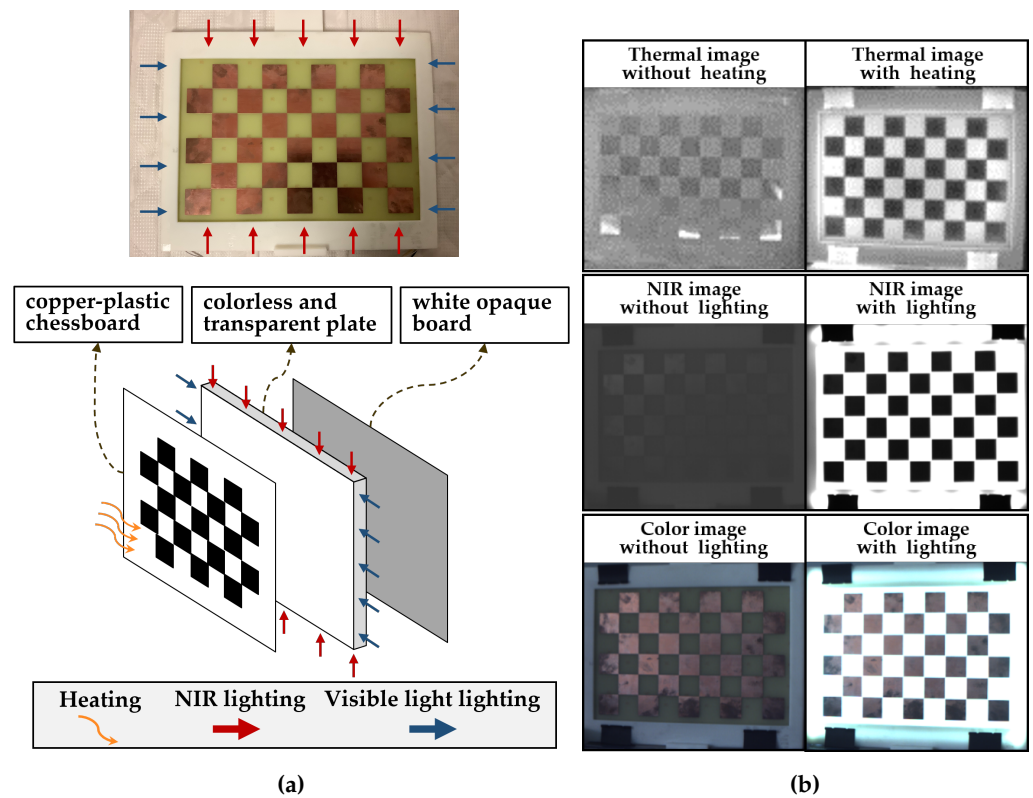
**Figure 2.** A multi-modal 3D sensor system consisting of an active stereo-vision 3D sensor based on GOBO projection, a RGB camera (FLIR Grasshopper3), and a thermal camera (FLIR A35).

**Table 1.** Technical data of the used camera systems

	GOBO 3D Sensor	FLIR A35 thermal camera	FLIR Grasshopper3 color camera
Resolution	1024 x 1024	320 x 256	2048 x 2048
Image frequency	36 Hz	60 Hz	90 Hz
Field of view	48° x 44°	63° x 50°	50° x 50°
Mean depth error	0.15-0.5 mm [29]	-	-
Range	0.4-2 m	-	-
Wavelength band	850nm	8 – 14μm	R: 640 G: 525 B: 470 (nm)
Thermal sensitivity	-	< 0.05°C	-
Temperature range	-	-25 to 135°C	-

the chessboard pattern cannot be captured by thermal cameras. Hence, as shown in Figure 3a, inspired by [30], a copper-plastic chessboard calibration target has been manufactured to solve this problem. Before the actual calibration the target needs to be heated for example by means of a hair dryer. After a few seconds of cooling down, the copper grids and plastic grids will have different temperatures and different grey values (copper dark and plastic bright) in the thermal image, because they have different emissivity coefficients, as shown in Figure 3b.

This copper-plastic chessboard works fine for thermal camera calibration, but it brings another problem for RGB and NIR camera. The surface of copper plating is always smooth, resulting in an overexposure problem because of specular reflections with external and passive light sources. Furthermore, the low contrast of texture on the chessboard surface in wavelength of visible light and NIR leads to the fact that the grid in the calibration images is not sharp enough for the corners to be detected. Therefore, as shown in Figure 3a, a calibration target with an internal active light source is proposed. A colorless and transparent plate is mounted behind the chessboard and visible light and NIR LEDs are mounted at the edges of the plate as active light sources. A white opaque board is set behind the plate as a diffusor. Figure 3b shows the comparison of calibration images with passive lighting and active lighting. With the help of the active light source, images with sufficient contrast can be captured to calibrate the intrinsic and extrinsic parameters for our multi-modal cameras.



**Figure 3.** (a) A copper-plastic chessboard calibration target (upper) and its principle (bottom); (b) Comparison of calibration images with and without active lighting for color image, NIR image and thermal image.

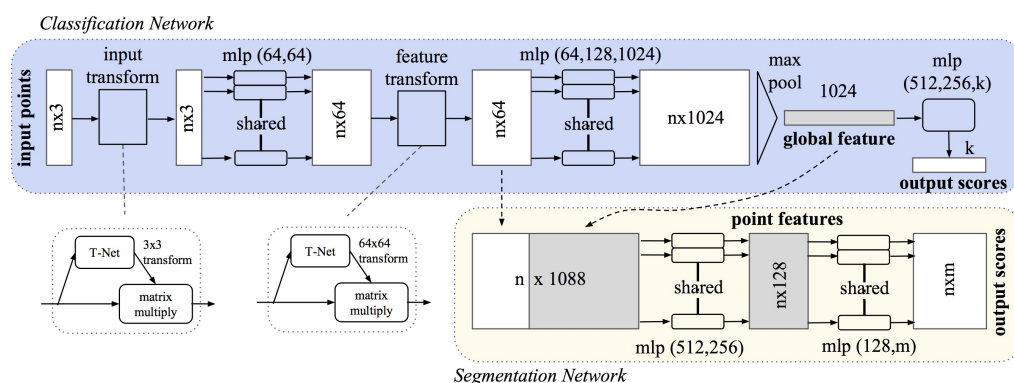
#### 4.3. Calibration and Registration

The intrinsic parameters  $\mathbf{K}$  of each camera can be simply calculated by using the Zhang's calibration algorithm [31]. The 3D sensor is used as a reference camera for calibration of extrinsic parameters. That means that rotation  $\mathbf{R}$  and translation  $\mathbf{T}$  of each camera (except 3D sensor) with respect to the coordinate system of the 3D sensor are calculated from a series of image tuples showing the calibration target. By using the parameters, alignment of the multi-modal point cloud can be performed with the following method:

For generating the multi-modal point cloud, each point of the original point cloud from the 3D sensor, is projected onto the image plane of the color and thermal camera. To that end, the point cloud can be transformed from the coordinate system of the 3D sensor to the coordinate system of the target camera with the extrinsic parameters  $(R_t, T_t)$ . For each 3D point of the transformed point cloud, a 2D projection pixel  $(u_t, v_t)$  on the target sensor plane can be calculated with intrinsic parameter  $\mathbf{K}_t$  of the thermal camera and RGB camera respectively. If the projection pixel is located on the sensor, i.e.  $0 \leq u_t < width, 0 \leq v_t < height$ , it will be determined as a corresponding pixel of this 3D point, as shown in equation 1. Once 2D corresponding pixels of all 3D points are determined, thermal and RGB values can be mapped onto the 3D point cloud.

$$s \cdot \begin{bmatrix} u_t \\ v_t \\ 1 \end{bmatrix} = \mathbf{K}_t \cdot \left( \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ z \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix} \right), \quad \mathbf{K}_t = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \quad (1)$$

where  $\mathbf{K}_t$  is the matrix of intrinsic parameters of target camera,  $c_x$  and  $c_y$  are the principal point coordinates,  $f_x$  and  $f_y$  are the focal lengths of thermal or RGB camera's lens,  $r_{ij}$  represent the rotation matrix, and  $[t_x, t_y, t_z]$  is the translation vector defining the extrinsic



**Figure 4.** The principle of PointNet. Point positions get transformed into spatial features by two stages of MLPs, before they are pooled into a global feature vector describing the whole object. Afterwards, a combination of local and global features can be used for segmentation purposes.[1]

calibration parameters.  $x$ ,  $y$  and  $z$  are the point coordinates in the coordinate system of the 3D sensor.

## 5. Point Cloud Segmentation Networks

### 5.1. PointNet

In the field of 3D point cloud segmentation, PointNet [1] is a milestone study. The article proposed the idea of using shared multi-layer perceptrons (MLP) to extract global features from a point cloud. By using a novel T-Net (transformation-network), a reference frame for the point cloud can be learned and utilized to keep features rotationally invariant. Usually in traditional methods, principal component analysis (PCA) was used to solve this problem instead. In addition to that, max-pooling was recommended as a symmetric aggregation function to solve the problem that usually a point cloud is an unordered set. As shown in Figure 4, global features of a point cloud can be efficiently extracted using PointNet. Finally, the global features and the output features of the last feature transformation unit will be concatenated to be used as input for another network to achieve pixel level segmentation.

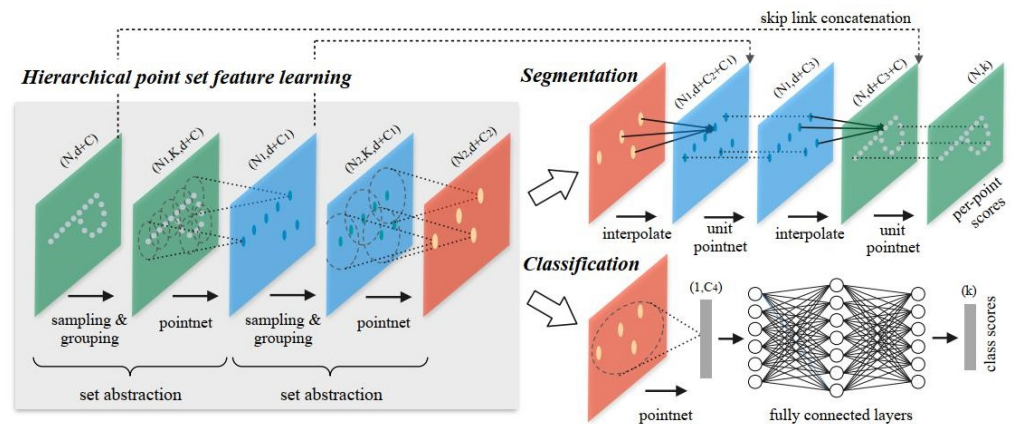
### 5.2. PointNet++

Obviously, with only global features, PointNet has insufficient ability to represent the semantic information for a local region. PointNet++ [2] describes a multi-level architecture, as shown in Figure 5. By using a farthest point sampling (FPS) algorithm, in each level the input point cloud is progressively downsampled and the point density decreases. Each point in the sampled sparse point cloud is used as a centroid for a neighborhood search in the dense point cloud. Then a mini-PointNet is utilized to extract the global features of this neighborhood that will be used as the local feature of this centroid point. A hierarchical propagation strategy with distance based interpolation and across level skip links is adopted to upsample the enriched point clouds to the original size.

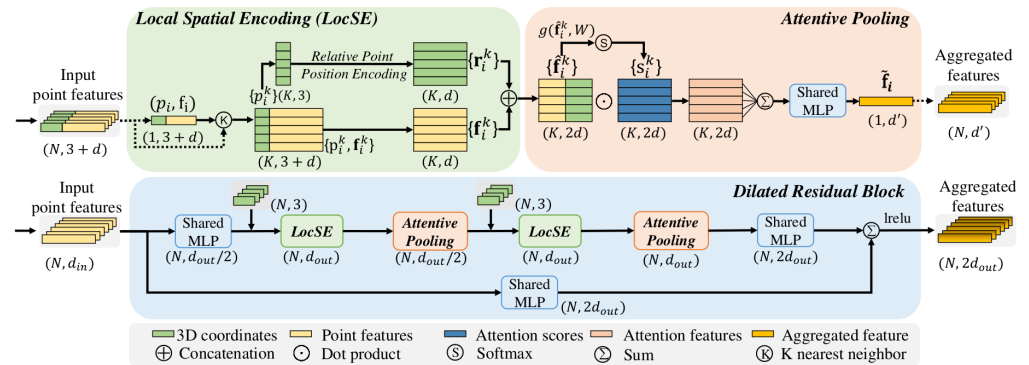
### 5.3. RandLA-Net

RandLA-Net [3] is a state-of-the-art neural network designed for large-scale 3D point cloud semantic segmentation. Similar to PointNet++, RandLA-Net is also a multi-level architecture, which in contrast uses random downsampling instead of FPS in order to reduce memory requirements and speed up computation. However, random sampling has a drawback of missing some useful point features occasionally. To overcome that issue, a powerful local feature aggregation module was designed in that approach, as shown in Figure 6. By using a local spatial encoding module (LocSE) in each neighborhood various spatial information are explicitly concatenated and encoded. Therefore, XYZ-coordinates of all points as well as euclidean distances and XYZ-differences be-





**Figure 5.** The multi-level architecture of PointNet++. Explicit neighborhood search in the point clouds is used to extract local features by means of a locally applied PointNet in multiple stages. These strong local feature can be used for object classification (lower branch) or for segmentation (upper branch) [2]

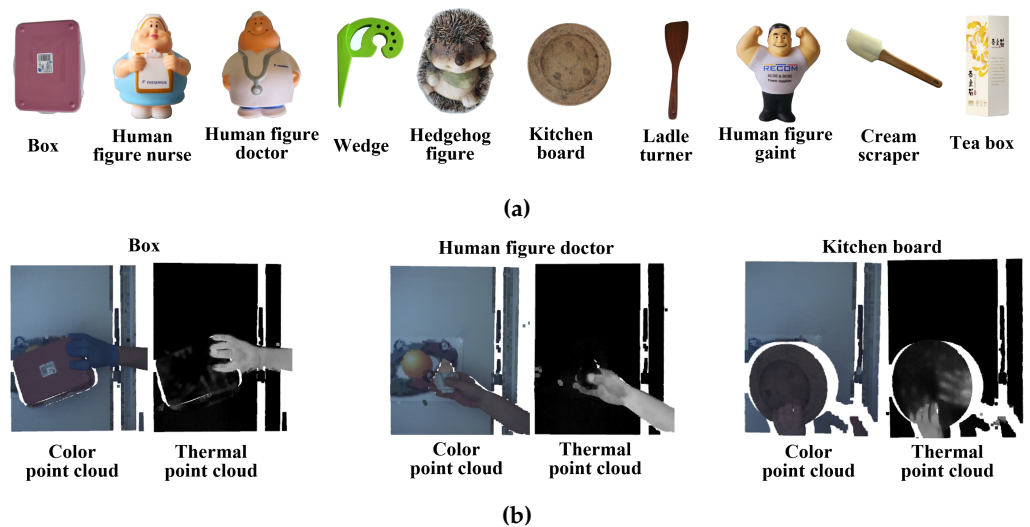


**Figure 6.** The architecture of the local feature aggregation module of RandLA-Net, which consists of multiple Local Spatial Encoding layers (LocSE) and Attentive Pooling layers (AP). In LocSE, the geometric information of a local area in the point cloud is encoded and then concatenated with the point features for local features extraction. The local features are further aggregated by an AP layer. [3]

between the centroid point and all neighboring points are calculated. Then the spatial information and point features are concatenated and local features can be extracted using a shared MLP. Additionally, in between two adjacent levels, an attentive pooling is utilized to aggregate the features. Then, multiple LocSE and attentive pooling units with a skip connection are stacked as a dilated residual block, which is repeatedly used in the RandLA-Net. Overall, RandLA-Net is built by stacking multiple dilated residual blocks to aggregate local features and an upsampling method identical to PointNet++ is used to interpolate the downsampled point clouds.

## 6. Datasets for Hand/Object Segmentation

In this paper, a dataset captured by our sensor system and named GOBO-Dataset is used to evaluate the performance of the multi-modal 3D data hand/object segmentation, as shown in Figure 7. The hand with one of the objects was placed roughly one meter in front of the sensors. In half of the data the human hand is recorded with opaque rubber gloves, in the other half without. In some samples of the dataset, the objects have taken the temperature of the holding hand caused by the long time holding them (see the right thermal point cloud). We used our own semi-autonomous annotation tool for labeling these multi-modal point clouds. The tool takes advantage of the simple separation of the background in the point cloud and uses region growing on the thermal



**Figure 7.** Overview of the GOBO-Dataset with 12 classes (10 objects, background, and hand): (a) all objects; (b) examples of multi-modal 3D data (Box, Human figure doctor and Kitchen board)

or color channel for an initial separation of hand and held object, which afterwards can be refined manually.

The GOBO-Dataset provides 600 multi-modal point clouds labeled with 12 classes (10 objects, background, and human hand). The samples have been split into a training set with 420 point clouds, a validation set with 60 point clouds and a test set with another 120 point clouds. In the dataset we have multi-modal point clouds with 7-channels containing spatial data (XYZ), color data (RGB) and thermal data (T). In comparative experiments, the networks mentioned in section 5 were trained on different modalities of the dataset in order to understand the influence of the individual parts (XYZ, XYZ-RGB, XYZ-T and XYZ-RGB-T).

## 7. Segmentation Experiment

As mentioned in section 5, PointNet has the simplest architecture with only plain global feature extraction. Local features can be extracted by PointNet++ and RandLA-Net, and especially RandLA-Net has a more powerful and complex local feature aggregation with respect to PointNet++. Therefore, in this experiment we have chosen these three networks with different performance for training on the GOBO-Dataset to evaluate the influence of multi-modal 3D data on hand/object segmentation in general. Thus, the findings should generalize to future architectures.

### 7.1. Evaluation Approach

A measuring method is required for performance evaluation of the point cloud segmentation. The Intersection over Union (IoU) was applied to intuitively reflect segmentation performance. For each class, the IoU could be calculated by using equation 2. Moreover, the mean IoU of all the classes (hand and objects) was calculated to present the overall performance, and the mean IoU of all the objects is also provided in the experiment section. In the scene of grabbing an object from a human hand by an assistant robot, we did not measure the IoU for background, because the segmentation of hand and held object is much more important than background, while the background class otherwise dominates the results.

$$IoU = \frac{TP}{TP + FP + FN} \quad (2)$$

where  $TP$  is the number of true positive predicted point classes,  $FP$  is the number of false positive predictions and  $FN$  is the number of false negative predictions.



**Figure 8.** The convergence curves in training phase of RandLA-Net, in which RandLA-Net was trained for 400 epochs on data XYZ, XYZ-T, XYZ-RGB and XYZ-RGB-T. Upper: the training curve, bottom: the validation curve.

## 7.2. Training Details

PointNet, PointNet++ and RandLA-Net were trained on XYZ, XYZ-RGB, XYZ-T and XYZ-RGB-T for 400 epochs without any pre-training. The learning rate setting is as shown in table 2. We used a framework of PointNet and PointNet++ available from [32]. For PointNet++, for sake of efficiency, we replaced the farthest point sampling by a uniform random sampling similar to RandLA-Net, and multi-scale grouping (MSG) was also adopted. The used implementation of RandLA-Net was available on [33]. For PointNet++ and RandLA-Net, a 4-level architecture was used, and in each level the point cloud was progressively downsampled with a factor of 1/3.

**Table 2.** Learning rate schedule for the experiments

Epochs	0-100	100-200	200-400
learning rate	0.01	0.001	0.0001

Due to removal of invalid points, 3D sensors will inevitably produce point clouds of different sizes and depending on the scene, the number of valid points varies. However, neural networks require batches of data with the same size for training. Hence, before training, the multi-modal point clouds in the training dataset have been uniformly downsampled to standardized size (10,000 points in our case).

The training of the networks used the cross-entropy loss function and the Adam optimizer [34]. The point clouds in our dataset are imbalanced in the number of points per class (the ratio of background, hand, and object is approximately 3309:267:1). Therefore, different weights for each class were used for weighting the loss function, as shown in

**Table 3.** The quantitative segmentation results on test split of the GOBO-Dataset (IoU %)

		mIoU overall	IoU Hand	mIoU Object	Box	Nurse Figure	Doctor Figure	Wedge	Hedge- hog	Kitchen Board	Spatula	Human Figure	Ice Scraper	Tea Box
RandLA-Net	XYZ-RGB-T	<b>82.8</b>	<b>92.6</b>	<b>81.9</b>	91.7	77.5	65.9	80.3	71.6	85.2	70.7	95.3	88.8	91.5
	XYZ-RGB	77.3	88.9	76.2	77.0	59.9	33.0	80.0	96.2	83.8	68.8	88.1	91.5	83.4
	XYZ-T	35.7	84.1	30.9	44.2	15.0	5.0	28.4	57.7	47.9	27.0	25.3	26.5	32.1
	XYZ	35.7	76.7	31.6	34.3	2.5	36.1	29.4	49.6	65.4	17.8	28.2	22.5	30.6
PointNet++	XYZ-RGB-T	<b>55.0</b>	<b>79.8</b>	<b>52.5</b>	58.0	40.2	33.7	65.9	87.1	71.7	25.5	43.1	24.1	75.2
	XYZ-RGB	45.5	66.6	43.4	66.2	30.2	23.0	58.0	69.7	57.0	23.6	41.1	19.3	46.1
	XYZ-T	25.9	52.0	23.2	45.2	13.3	1.8	20.5	37.6	38.3	15.1	5.5	12.0	43.1
	XYZ	25.8	60.6	22.4	55.1	15.6	12.7	7.8	27.8	43.3	24.3	6.4	15.8	14.7
PointNet	XYZ-RGB-T	<b>45.9</b>	<b>79.9</b>	<b>42.5</b>	62.2	23.5	26.7	62.6	36.2	52.7	28.8	39.3	42.2	50.9
	XYZ-RGB	43.5	78.4	40.1	60.1	20.5	22.6	63.5	34.9	51.8	27.0	38.6	42.9	39.2
	XYZ-T	24.8	72.1	20.1	18.1	31.4	13.9	16.1	25.3	23.6	9.1	26.3	18.9	18.2
	XYZ	22.0	52.9	18.9	36.4	15.8	13.9	10.4	24.0	31.0	12.9	18.1	11.6	15.2

equation 3. The normalized weight  $w_i$  depends on the probability  $p_i$  of a point to belong to the  $i^{th}$  of  $K$  classes in the entire dataset.

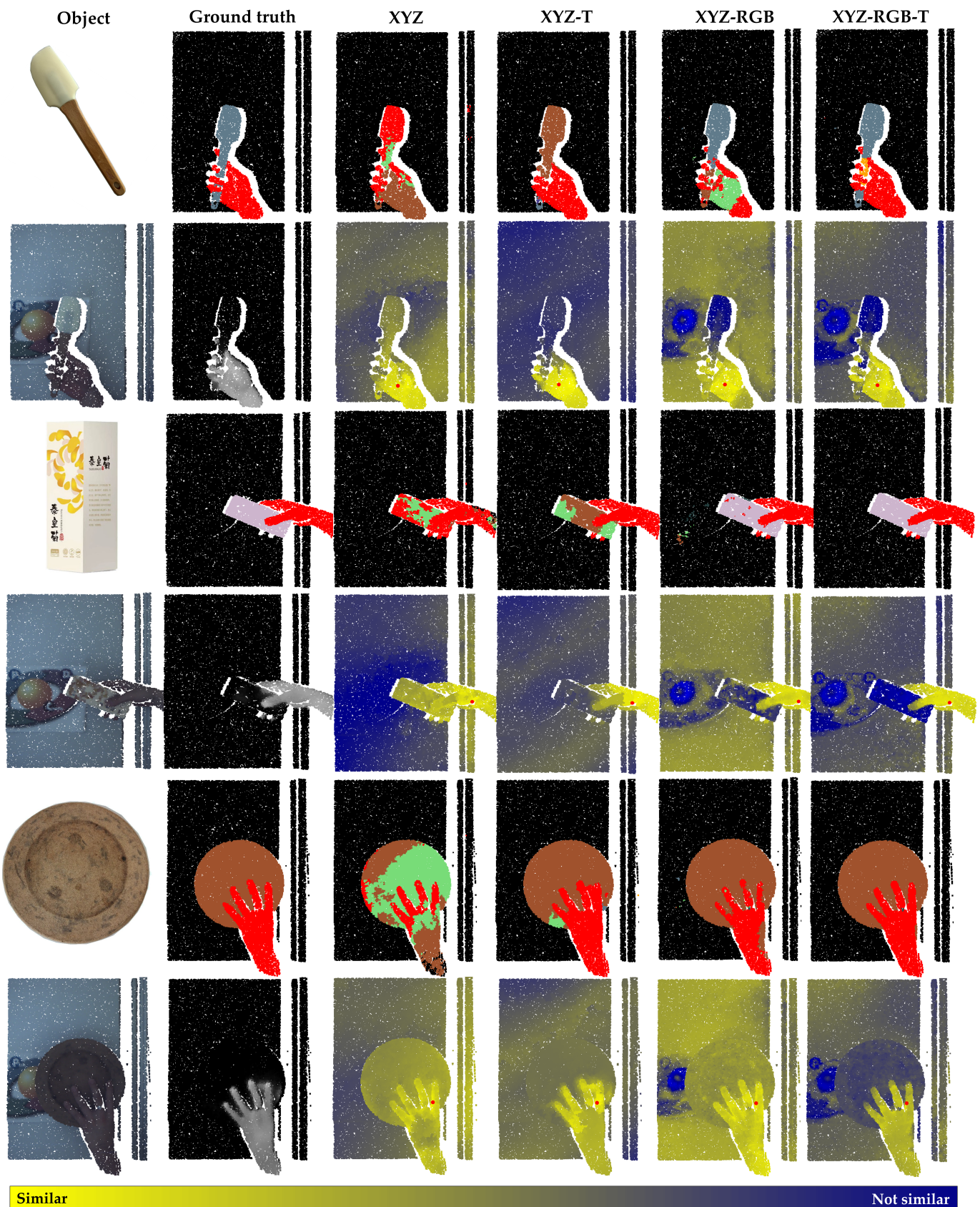
$$w_i = \frac{\log(p_i)}{\sum_{j=0}^K \log(p_j)} \quad (3)$$

### 7.3. Segmentation Results

Figure 8 shows convergence curves of RandLA-Net in the training phase on training dataset and validation dataset. At epoch 400, although the training curve still shows an improving trend, the validation curve indicates that the results are no longer improving. So we interrupted the training at 400th epoch. Obviously, with the help of the strong feature of color, the RandLA-Net has a significant superiority by XYZ-RGB-T and XYZ-RGB over by XYZ and XYZ-T. With the use of thermal, XYZ-RGB-T has further improved over XYZ-RGB. This is in line with our expectation. Meanwhile, trends of the convergence curves show that the convergence rate of the four modes were almost the same. This indicates that the multi-modal point cloud does not lead to a longer training time due to more channels. The training phases of PointNet and PointNet++ feature almost the same tendency.

Table 3 shows the detailed quantitative segmentation results based on the test dataset. For all the three networks, the overall mIoU shows similar relations for the individual input channels used, while the absolute performance of the three networks differs reflecting their individual abilities. However, XYZ-T has almost no improvement over XYZ in the test dataset independent on the used architecture. For this, the second and third columns of the table provide the explanation. For example, for RandLA-Net, the mIoU of objects by XYZ-T actually decreased by 0.7% compared to XYZ. It is possible that this is because in some samples of our dataset, the object took on hand temperature at some parts of the surface. The points in these areas could be confused for hand without any additional color information. Compared to XYZ-RGB, although XYZ-RGB-T should not have a dominant advantage for predicting the object points, it has a significantly better object mIoU. This indirect improvement is due to reduction of false positive points in the interaction area of hand and object, which can be better predicted as hand, as shown in Figure 9. In comparison, the object mIoU, as well as the mIoU of the hand class has an obvious improvement from XYZ to XYZ-RGB-T, proving that multi-modal data significantly supports a more robust segmentation independent of the actual method used.





**Figure 9.** Visualization of experimental results for RandLA-Net on individual samples of the test dataset; The first row shows the ground truth and segmentation by XYZ, XYZ-T, XYZ-RGB and XYZ-RGB-T, while the hand class is labeled in red. The second row shows the color point cloud, thermal point cloud, and the feature point cloud generated by XYZ, XYZ-T, XYZ-RGB and XYZ-RGB-T. In the feature point cloud, the Euclidean distances between a reference point (red point) and all the other points are calculated and normalized in features space. The distances are color coded (light yellow: similar points, dark blue: dissimilar points).



#### 338 7.4. Visualization of Segmentation Results

339 Figure 9 shows a visualization of the segmentation results. For each object, the first  
 340 row shows the ground truth and RandLA-Net predictions by XYZ, XYZ-T, XYZ-RGB  
 341 and XYZ-RGB-T. The second row shows the color point cloud, thermal point cloud and  
 342 feature point clouds. The features extracted by the last feature layer of RandLA-Net  
 343 were used to generate these feature point clouds. Inspired by [35], we used the following  
 344 method to generate the feature point cloud:

345 First, we choose a reference point (red point) that is located on the hand. In the  
 346 corresponding feature space, the euclidean distances between this reference point and  
 347 all the other points of this point cloud were calculated. The 3D sensor inevitably will  
 348 generate some outlier points (incorrectly reconstructed points). In feature space, the  
 349 distances between these points and other points may be exceptionally large. Therefore,  
 350 the distances for visualization was normalized to the 97% quantile and were presented  
 351 with gradient colors (light yellow to dark blue). Hence, in the feature point cloud, the  
 352 greater color contrast between two points indicates that they have greater dissimilarity.

353 It is clearly visible, that for all the objects, the feature point cloud of XYZ-RGB-T  
 354 has higher contrast than any other, i.e. the points of the hand have greater distances  
 355 in the feature space to the object and background. Although the final segmentation  
 356 result is still dependent on the classifier, these distances make it easier to cluster points  
 357 in the feature space, and implies that the segmentation will be better. Figure 9 shows  
 358 XYZ-RGB-T has the best segmentation results in the interaction area of hand and object.  
 359 For example, in the pictures of the first object, the segmentation of the fingers and the  
 360 object is refined when using XYZ-RGB-T. For the second object, some areas of the surface  
 361 possess a similarity to the hand in the feature point cloud by XYZ-RGB because of the  
 362 color texture. The segmentation results by XYZ-RGB show that some points on these  
 363 areas are indeed predicted as hand. In comparison, the segmentation of XYZ-RGB-T  
 364 is much more precise. For the third object, as shown in the thermal point cloud, the  
 365 boundary area of the kitchen board has a similar temperature as the hand, causing  
 366 the points in this area to be predicted as hand when using XYZ-T. In contrast, this  
 367 similar temperature does not affect the prediction when using XYZ-RGB-T. However,  
 368 the pictures of the third object show, the middle finger of the hand with a ring has points  
 369 that were mistakenly predicted as object by XYZ-RGB and XYZ-RGB-T classifiers.

#### 370 7.5. Time Consumption Analysis

371 The experiment was conducted on the computing platform of Intel Core i9-9960x  
 372 (CPU) and GeForce RTX 2080 Ti (GPU). We recorded the time consumption for processing  
 373 a multi-modal point cloud with 10k points. By using a parallel computing by OpenMP  
 374 [36], multi-modal data fusion consumes 14 milliseconds (ms) approximately, and the  
 375 inference time consumption by the three networks PointNet, PointNet++ and RandLA-  
 376 Net are approximately 7 ms, 124 ms and 102 ms. As we can see, with respect to the  
 377 inference by PointNet++ and RandLA-Net, data fusion occupies only a fraction of the  
 378 time consumption for the entire process. PointNet++ and RandLA-Net have the multi-  
 379 level architecture, leading that multiple k-nearest-neighbors (KNN) based neighborhood  
 380 searches are required for each of two adjacent levels. As a result, these two approaches  
 381 are not as efficient as PointNet. The neighborhood searches for the 4-level architecture  
 382 has a time consumption of 81 ms and hence is the major part. The additional effort  
 383 to achieve an improvement through sensor fusion seems to be justified in view of the  
 384 runtimes of increasingly complex networks necessary to improve the results otherwise.

#### 385 8. Discussion

386 To enable a precise segmentation of hand and object for an assistant robot to grasp  
 387 objects from a human hand safely, in this work, we presented a multi-modal 3D sensor  
 388 system. We also focused on the challenges for calibration and alignment of a multi-modal  
 389 sensor system with a thermal camera. The successful experiments showed, that applying

a copper-plastic chessboard calibration target with internal and active light source (NIR and visible light) effectively solves the calibration problem. As it can be captured by each camera with sufficient contrast simultaneously, the use of such a calibration target makes the calibration and alignment of the multi-modal camera systems no longer tedious.

The segmentation experiments using PointNet, PointNet++ and RandLA-Net on our dataset could confirm our hypothesis, that multi-modal data significantly supports point wise segmentation. RandLA-Net as the strongest state-of-the-art network has achieved the remarkable results on XYZ-RGB-T (overall mIoU: 82.8%). In contrast, the mIoU for XYZ, XYZ-T and XYZ-RGB were 35.7%, 35.7%, and 77.3% respectively. Surprisingly, XYZ-T has almost no improvement over XYZ, this is partly because some objects have the similar temperature as human hands, which confuses the prediction on XYZ-T without any additional cues. In addition, a visualization of feature point cloud extracted by RandLA-Net intuitively demonstrates the feasibility of using a neural network to extract the potential features of multi-modal data. As mentioned in the section 1, in the multi-modal data, none of the channels are all-purpose, but the information of all channels can be integrated to make up for their respective weaknesses.

In recent years, the computing performance of computers has improved tremendously. Therefore, deep learning technology has started to be widely studied and applied. Under this condition, the computing performance required for efficient multi-modal sensor fusion is also achievable. On our computing platform of Intel Core i9-9960x (CPU) and GeForce RTX 2080 Ti (GPU), data fusion consumes approximately 14 ms for a point cloud with 10k points. Therefore, we propose the application of multi-modal data to reduce the complexity of image processing tasks. It is a matter of data and improvements in the segmentation methods which in future will allow to raise the limits of the hand/object segmentation results further. Nevertheless, for safety critical applications, the IoU results alone will not be a sufficient criterion. In order to rely on machine learning based safety critical features, other questions like explainability and robustness in case of adversarial or out of distribution data have to be considered. We are sure, that multi-modal data helps to reach an acceptable level of robustness more easily either way.

In the future, in order to make the multi-modal sensor system usable in a real-world environment, we would like to expand our dataset and further evaluate it in practical scenarios. Currently, with training of RandLA-Net on our dataset, we can precisely segment hand and objects in real-time point clouds from the sensor. Nevertheless, some points on the background will be identified as hand or object occasionally, which may be caused by the different point density compared to the training data. Therefore, we address the portability of our models to other sensor setups and unseen objects in future work.

**Author Contributions:** Conceptualization, Gunther Notni; Data curation, Yan Zhang; Funding acquisition, Horst-Michael Gross and Gunther Notni; Investigation, Yan Zhang; Methodology, Yan Zhang; Project administration, Horst-Michael Gross and Gunther Notni; Software, Yan Zhang and Steffen Mueller; Supervision, Horst-Michael Gross and Gunther Notni; Validation, Steffen Mueller and Benedict Stephan; Visualization, Yan Zhang; Writing – original draft, Yan Zhang; Writing – review & editing, Yan Zhang, Steffen Mueller, Benedict Stephan and Gunther Notni.

**Funding:** This research was funded by Thueringer Aufbaubank, Thueringer Zentrum fuer Maschinenbau and Freistaat Thüringen aus Mitteln des Europäischen Sozialfonds (Project: "SONARO", Project No.: "2018 FGR 0097").

**Acknowledgments:** We thank Chen Zhang for assistance for building the camera system. We thank Richard Fütterer, Xiaojiang Han, Jinxin Zhu and Yujian Yuan for assistance of production of the calibration target. We thank Yang Li, Xiao Ling, Jingyu Wang for building the datasets. We acknowledge support for the publication costs by the Open Access Publication Fund of the Technische Universität Ilmenau.

**Conflicts of Interest:** This article has no conflict of interest with any organization.

## References

1. Qi, C.R.; Su, H.; Mo, K.; Guibas, L.J. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
2. Qi, C.R.; Yi, L.; Su, H.; Guibas, L.J. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413* **2017**.
3. Hu, Q.; Yang, B.; Xie, L.; Rosa, S.; Guo, Y.; Wang, Z.; Trigoni, N.; Markham, A. Randla-net: Efficient semantic segmentation of large-scale point clouds. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11108–11117.
4. Redmon, J.; Divvala, S.; Girshick, R.; Farhadi, A. You only look once: Unified, real-time object detection. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
5. He, K.; Gkioxari, G.; Dollár, P.; Girshick, R. Mask r-cnn. *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
6. Kirillov, A.; Wu, Y.; He, K.; Girshick, R. Pointrend: Image segmentation as rendering. *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9799–9808.
7. Palmero, C.; Clapés, A.; Bahnsen, C.; Møgelmoose, A.; Moeslund, T.B.; Escalera, S. Multi-modal rgb–depth–thermal human body segmentation. *International Journal of Computer Vision* **2016**, *118*, 217–239.
8. Zhao, S.; Yang, W.; Wang, Y. A new hand segmentation method based on fully convolutional network. 2018 Chinese Control And Decision Conference (CCDC). IEEE, 2018, pp. 5966–5970.
9. Simonyan, K.; Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* **2014**.
10. Jeon, E.S.; Kim, J.H.; Hong, H.G.; Batchuluun, G.; Park, K.R. Human detection based on the generation of a background image and fuzzy system by using a thermal camera. *Sensors* **2016**, *16*, 453.
11. Kim, S.; Chi, H.G.; Hu, X.; Vegesana, A.; Ramani, K. First-Person View Hand Segmentation of Multi-Modal Hand Activity Video Dataset. *BMVC*, 2020.
12. Chen, L.C.; Zhu, Y.; Papandreou, G.; Schroff, F.; Adam, H. Encoder-decoder with atrous separable convolution for semantic image segmentation. *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.
13. Wang, W.; Zhang, J.; Shen, C. Improved human detection and classification in thermal images. 2010 IEEE International Conference on Image Processing. IEEE, 2010, pp. 2313–2316.
14. Setjo, C.H.; Achmad, B.; others. Thermal image human detection using Haar-cascade classifier. 2017 7th International Annual Engineering Seminar (InAES). IEEE, 2017, pp. 1–6.
15. Correa, M.; Hermosilla, G.; Verschae, R.; Ruiz-del Solar, J. Human detection and identification by robots using thermal and visual information in domestic environments. *Journal of Intelligent & Robotic Systems* **2012**, *66*, 223–243.
16. Ruiz-del Solar, J.; Verschae, R. Robust skin segmentation using neighborhood information. 2004 International Conference on Image Processing, 2004. ICIP'04. IEEE, 2004, Vol. 1, pp. 207–210.
17. Shivakumar, S.S.; Rodrigues, N.; Zhou, A.; Miller, I.D.; Kumar, V.; Taylor, C.J. Pst900: Rgb-thermal calibration, dataset and segmentation network. 2020 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2020, pp. 9441–9447.
18. Nishi, K.; Demura, M.; Miura, J.; Oishi, S. Use of thermal point cloud for thermal comfort measurement and human pose estimation in robotic monitoring. *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 1416–1423.
19. Rosenberger, M.; Zhang, C.; Zhang, Y.; Notni, G. 3D high-resolution multimodal imaging system for real-time applications. *Dimensional Optical Metrology and Inspection for Practical Applications IX*. International Society for Optics and Photonics, 2020, Vol. 11397, p. 1139704.
20. Zhang, C.; Gebhart, I.; Kühmstedt, P.; Rosenberger, M.; Notni, G. Enhanced Contactless Vital Sign Estimation from Real-Time Multimodal 3D Image Data. *Journal of Imaging* **2020**, *6*, 123.
21. Ivašić-Kos, M.; Krišto, M.; Pobar, M. Human detection in thermal imaging using YOLO. *Proceedings of the 2019 5th International Conference on Computer and Technology Applications*, 2019, pp. 20–24.
22. Zhang, Y.; Zhang, C.; Rosenberger, M.; Notni, G. 6D Object Pose Estimation Algorithm Using Preprocessing of Segmentation and Keypoint Extraction. 2020 IEEE International Instrumentation and Measurement Technology Conference (I2MTC). IEEE, 2020, pp. 1–6.
23. Xiong, H.; Cai, W.; Liu, Q. MCNet: Multi-level Correction Network for thermal image semantic segmentation of nighttime driving scene. *Infrared Physics & Technology* **2021**, *113*, 103628.
24. Ge, L.; Cai, Y.; Weng, J.; Yuan, J. Hand pointnet: 3d hand pose estimation using point sets. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8417–8426.
25. Heist, S.; Zhang, C.; Reichwald, K.; Kühmstedt, P.; Notni, G.; Tünnermann, A. 5D hyperspectral imaging: fast and accurate measurement of surface shape and spectral characteristics using structured light. *Optics express* **2018**, *26*, 23366–23379.
26. FLIR Grasshopper 3 overview. Available online: <https://www.edmundoptics.com/p/gs3-u3-41c6c-c-1-grasshopper-usb-30-color-camera/30772/> (accessed 17 August 2021).
27. FLIR A5 product overview. Available online: <https://www.flir.com/products/a35/> (accessed on 17 August 2021).

28. Heist, S.; Lutzke, P.; Schmidt, I.; Dietrich, P.; Kühmstedt, P.; Tünnermann, A.; Notni, G. High-speed three-dimensional shape measurement using GOBO projection. *Optics and Lasers in Engineering* **2016**, *87*, 90–96.
29. Heist, S.; Dietrich, P.; Landmann, M.; Kühmstedt, P.; Notni, G.; Tünnermann, A. GOBO projection for 3D measurements at highest frame rates: a performance analysis. *Light: Science & Applications* **2018**, *7*, 1–13.
30. Landmann, M.; Heist, S.; Dietrich, P.; Lutzke, P.; Gebhart, I.; Kühmstedt, P.; Notni, G. Multimodal sensor: high-speed 3D and thermal measurement. *Photonics and Education in Measurement Science 2019*. International Society for Optics and Photonics, 2019, Vol. 11144, p. 1114403.
31. Zhang, Z. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence* **2000**, *22*, 1330–1334.
32. Benny. Pointnet-Pointnet2-pytorch. Available online: [https://github.com/yanx27/Pointnet\\_Pointnet2\\_pytorch/](https://github.com/yanx27/Pointnet_Pointnet2_pytorch/) (accessed on 17 August 2021).
33. Qiqihaer. RandLA-Net. Available online: <https://github.com/qiqihaer/RandLA-Net-pytorch/> (accessed 17 August 2021).
34. Kingma, D.P.; Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* **2014**.
35. Wang, Y.; Sun, Y.; Liu, Z.; Sarma, S.E.; Bronstein, M.M.; Solomon, J.M. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)* **2019**, *38*, 1–12.
36. OpenMP overview. Available online: <https://www.openmp.org/> (accessed on 17 August 2021).