

Cal or No Cal? - Real-Time Miscalibration Detection of LiDAR and Camera Sensors

Iilir Tahiraj*, Jeremialie Swadiryus, Felix Fent, Markus Lienkamp

Abstract—The goal of extrinsic calibration is the alignment of sensor data to ensure an accurate representation of the surroundings and enable sensor fusion applications. From a safety perspective, sensor calibration is a key enabler of autonomous driving. In the current state of the art, a trend from target-based offline calibration towards targetless online calibration can be observed. However, online calibration is subject to strict real-time and resource constraints which are not met by state-of-the-art methods. This is mainly due to the high number of parameters to estimate, the reliance on geometric features, or the dependence on specific vehicle maneuvers. To meet these requirements and ensure the vehicle’s safety at any time, we propose a miscalibration detection framework that shifts the focus from the direct regression of calibration parameters to a binary classification of the calibration state, i.e., calibrated or miscalibrated. Therefore, we propose a contrastive learning approach that compares embedded features in a latent space to classify the calibration state of two different sensor modalities. Moreover, we provide a comprehensive analysis of the feature embeddings and challenging calibration errors that highlight the performance of our approach. As a result, our method outperforms the current state-of-the-art in terms of detection performance, inference time, and resource demand. The code is open source and available on <https://github.com/TUMFTM/MiscalibrationDetection>.

I. INTRODUCTION

Ensuring safety is the primary challenge in the development of autonomous driving systems. One key aspect of it is the development of an accurate environment model, which helps autonomous vehicles (AVs) to understand and predict their surroundings by fusing data from camera and LiDAR sensors [1], [2], [3]. Active research continues to improve sensor fusion and 3D object detection, with efforts focused not only on enhancing the accuracy but also on making sensors more reliable in challenging conditions such as adverse weather conditions, sensor failures, spatial and temporal misalignment, and feature sparsity [4], [5], [6], [7].

While sensor fusion is an enabler for reliable scene understanding, it also introduces new challenges. In particular, sensor synchronization and calibration. Sensor fusion requires aligned sensor data to ensure the accuracy and reliability of object detection algorithms. Therefore, accurate sensor calibration is critical to the safety of AVs.

Sensor calibration, and more specifically extrinsic sensor calibration, is the process of obtaining the rigid transformations between sensor coordinate systems. To address this problem, two strategies are used in autonomous driving:

* Corresponding author. iilir.tahiraj@tum.de.

Authors are with the TUM School of Engineering and Design, Chair of Automotive Technology, Technical University of Munich.

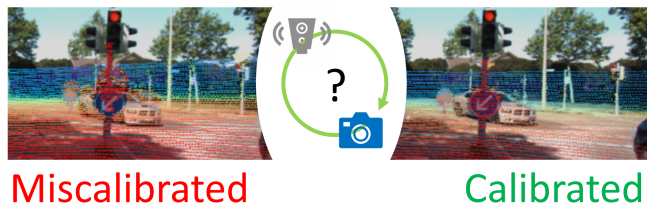


Fig. 1: Our miscalibration detection framework takes RGB images and LiDAR point clouds as input. The point clouds are projected onto the image and used to detect miscalibration between the sensors.

Calibration before or after the deployment of the autonomous vehicle. The former concept can be performed using specific calibration targets in a dedicated environment such as a vehicle factory. The latter involves the use of online calibration methods that aim to (re-)calibrate the sensors in an unstructured environment during vehicle operation. The second approach is mainly motivated by the fact that sensor miscalibration is caused by temperature changes and/or vehicle vibrations that occur after vehicle deployment.

Despite considerable efforts to develop online calibration methods, several challenges remain [8]. First, identifying common features in multimodal sensor data and open scenes remains a complex task, preventing online calibration methods from achieving the accuracy of offline target-based approaches. Second, online calibration — particularly learning-based approaches (see Section II-A) — still suffers from limited generalization in real-world scenarios.

Online sensor calibration traditionally involves the regression of at least 12 parameters for each individual extrinsic sensor-to-sensor calibration, namely the components of the rotation matrix and the translation vector. The number of parameters to regress increases even further when estimating more than two sensor calibrations or the intrinsics simultaneously. Therefore, regression-based calibration algorithms are resource intensive [9], [10] and difficult to verify for real-world, safety-critical systems [8]. Furthermore, many online calibration algorithms require certain environmental conditions to be met, such as the presence of distinctive geometric features or the execution of specific driving maneuvers [11], [12]. These requirements effectively limit the capacity of online calibration algorithms to ensure anytime-safety of the system. In addition, online recalibration algorithms are mainly trained to regress the extrinsic parameters. Intrinsic calibration errors, such as focal length or principal point offsets, can lead to incorrect recalibration of the extrinsic parameters because such intrinsic errors cannot be easily

distinguished from extrinsic errors.

In response to these challenges, we provide a continuous sensor monitoring approach and propose shifting the focus from a regression to a classification-based framework. By formulating the online calibration problem as a classification task, we can simplify the model while improving both efficiency and reliability. Our classification-based framework allows the identification of binary calibration states rather than the regression of extrinsic parameters, enabling the system to detect and respond to miscalibrations in real time with less computational effort compared to a full recalibration. Our main contributions are as follows:

- We introduce an open-source framework for miscalibration detection using a self-supervised learning architecture and feature-based classification.
- This is the first two-stage learning approach for miscalibration detection, where one stage learns robust representations of the miscalibrated sensor data. The second stage uses these representations to train the classification task.
- We are the first to provide a framework that shows that it can robustly detect a miscalibration in the presence of intrinsic errors.
- Using Centered Kernel Alignment (CKA), we analyze feature representations, enabling a simple architecture with faster inference and reduced computational needs suited for real-time use.
- Our approach achieves state-of-the-art detection, with inference time $6\times$ faster and a model size 42% smaller than existing methods.

II. RELATED WORK

Current state-of-the-art LiDAR-to-camera sensor calibration methods estimate transformation parameters by identifying correspondences and optimizing the alignment between LiDAR points and camera images. These methods often use appearance-based optimizations or deep learning techniques to achieve the desired calibration accuracy. We categorize these methods as regression-based approaches, focusing on those that perform online calibration to address on-site sensor misalignment. In contrast, we categorize miscalibration detection methods as classification-based tasks that aim to identify whether two or more sensors are properly aligned.

A. Regression-based Concepts

LiDAR-to-camera calibration methods mainly build correspondences from LiDAR point clouds to images or vice versa. These correspondence-based methods can be further divided into explicit and learned correspondence search [8].

Explicit Correspondences: Examples of approaches that match explicit correspondences using geometric features are presented in [11], [12], [13]. Yuan *et al.* [11] rely on finding edge correspondences to perform LiDAR-to-camera calibration. They first extract edges from 3D point clouds and estimate the transformation matrix by aligning them with edges from 2D images. In [12], distinctive planar features are captured from both the LiDAR and the camera

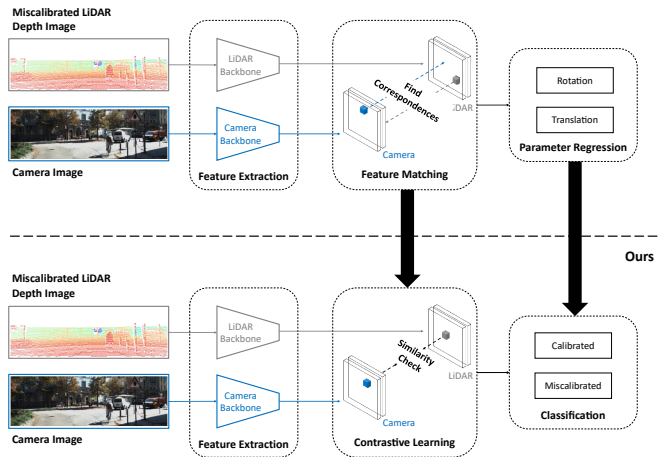


Fig. 2: **Above:** The general approach for online estimation of calibration parameters. The feature matching layer serves the correspondence finding process. The final layer regresses the translation and rotation of the extrinsic calibration transforms. **Below:** Our approach aims to learn representations that describe the similarity or dissimilarity between the multimodal features. The final layer classifies whether the inputs are calibrated or miscalibrated.

data to estimate the extrinsic calibration matrix via plane-constrained bundle adjustment. CRLF [13] performs a two-step calibration process by using line features extracted from the image and LiDAR point clouds. First, an initial calibration matrix is estimated by line fitting in both LiDAR and camera space, before a refinement step estimates the final transformation.

Recent approaches aim to perform online calibration using non-geometric features. The approach in [14] takes advantage of the fact that the projected intensities (LiDAR-to-camera) have a high variance when there are errors in the calibration parameters and formulates an optimization problem based on the intensity variance. The calibration algorithms presented in [15], [16], [17] rely on semantic information extracted from the sensor modalities, i.e. correspondences exist between point clouds and images based on the same semantic information.

Learned Correspondences: Methods with learned correspondences have proven to be powerful for use in extrinsic sensor calibration. These methods use deep learning techniques to find correspondences between the LiDAR and camera domain and estimate the calibration matrix in an end-to-end fashion [18], [19], [20], [21]. Schneider *et al.* [18] introduced RegNet, a LiDAR-to-camera calibration framework with an end-to-end approach: a convolutional neural network that performs feature extraction, correspondence matching, and the regression of a 6 DoF calibration matrix. CalibNet [19] proposes a more general approach, which is independent of camera characteristics and sensor configurations. In general, these methods are based on the architectures shown in Fig. 2.

Wang *et al.* [21], similar to traditional techniques in extrinsic sensor calibration, incorporate the coarse-to-fine

calibration procedure into a learning-based framework. Hu *et al.* [22] follow the same idea of iterative refinement and introduce HIFMNet, which constructs different cost volumes, one to initialize with a coarse extrinsic calibration to find large deviations between the depth maps and the RGB images and another for an iterative refinement of the extrinsics.

SOAC [23] propose to use Neural Radiance Fields (NeRF) for extrinsic and temporal sensor calibration. It is trained with observations from different sensors to generate a 3D scene representation and optimize the extrinsic calibration parameters to fit this representation. The work in [10] assumes that matching correspondences within the same data representation is most effective. Therefore, images and point clouds are transformed into depth maps, which are fed into a feature extraction, feature matching, and parameter regression module.

While methods based on explicit correspondences generalize well, they require a good initial estimate of the calibration matrix. Learning-based approaches, on the other hand, achieve good calibration performance but suffer from generalization capabilities [8] and complex architectures [9], [23], [21]. In addition, both explicit and learned-based approaches often depend on the presence of features such as rich texture, edges, or lines in the environment, which effectively limits the real-time capabilities of such models. More importantly, none of these methods consider or discuss the detection of miscalibrations. This is an important requirement in real-world applications and for both safety and robustness, miscalibration must be detected before online calibration can be performed.

B. Classification-based Concepts

Analogous to Section II-A, we will categorize the classification-based concepts into methods that require explicit or learned correspondences to detect miscalibrations. However, finding correspondences differs between regression-based and classification-based approaches. The regression-based approaches aim to find correspondences, i.e. matching features to align them and compute the calibration parameters. Classification-based approaches check for spatial similarity of features to discriminate between different modalities.

Explicit Correspondences: Levinson and Thrun [24] are the first to specifically consider miscalibration detection before estimating calibration parameters. They use depth discontinuities and image edges to detect misalignment between sensors. This is incorporated into an objective function that indicates whether small adjustments to the current calibration will result in a decrease in the cost function. This objective function measures "edginess" based on LiDAR and image data.

Learned Correspondences: Recent work in the field of miscalibration detection also tackles this problem by using deep learning techniques. These methods are mainly applied in sensor monitoring techniques, not only in LiDAR-to-camera setups but also in a wider range of systems beyond the field of autonomous driving [25], [26]. In the field

of autonomous driving, Chen *et al.* [27] presents a sensor monitoring framework that detects the spatial and temporal misalignment of sensor data as well as single sensor error sources such as image blur, point cloud noise, and perturbation. They train a Siamese network using a contrastive loss to detect sensor and cross-sensor inconsistencies. Wei *et al.* [28] introduce a self-checking framework that specifically detects sensor miscalibrations. For feature extraction, they use a patch transformer in the LiDAR and Unet in the camera domain to check for cosine similarity. Their framework implements the LCCNet [9] recalibration algorithm. Since they are the first and currently the only ones to perform learning-based miscalibration detection, we use this method as a baseline.

The role of classification-based concepts in sensor calibration is an important consideration that has not yet received much attention in the literature. It can be considered as a continuous monitoring system and a trigger for the recalibration process or any safety-related modification in the fusion strategy. As introduced in Section I, for safety reasons, these methods must be designed with resource efficiency and fast inference time in mind. Addressing this research gap, we present a two-stage learning approach and analyze the architecture (see Fig. 2 below) with a focus on monitoring for extrinsic calibration errors. We design the network to meet the requirements of model size and inference time.

III. METHOD

A. Fault Injection

To generate miscalibrated data, we introduce perturbations during the projection of a 3D point cloud into a 2D image plane from the KITTI dataset [29]. First, a 3D point is transformed from the LiDAR coordinate frame to the camera coordinate frame using the transformation matrix \mathbf{T}_{lid}^{cam} , which consists of the rotation matrix \mathbf{R}_{lid}^{cam} and translation vector \mathbf{t}_{lid}^{cam} . The transformed point x is then projected onto the rectified, rotated image plane of the i -th camera using the camera projection matrix $\mathbf{P}_{rect}^{(i)}$ and the rectifying rotation matrix of the reference camera $\mathbf{R}_{rect}^{(0)}$. $\mathbf{P}_{rect}^{(i)}$ is the intrinsic calibration matrix with the parameters focal length f_u, f_y , the principal point offsets c_x, c_y , and the axis skew γ . The full transformation is shown in Eq. 1.

$$\begin{aligned} \mathbf{y} &= \mathbf{P}_{rect}^{(i)} \cdot \mathbf{R}_{rect}^{(0)} \cdot \mathbf{T}_{lid}^{cam} \cdot \mathbf{x} \\ &= \mathbf{P}_{rect}^{(i)} \cdot \begin{bmatrix} \tilde{\mathbf{R}}_{rect}^{(0)} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{R}_{lid}^{cam} & \mathbf{t}_{lid}^{cam} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \cdot \mathbf{x}. \end{aligned} \quad (1)$$

To simulate the extrinsic calibration error, we augmented the dataset by introducing perturbations in both the rotational and translational components of this transformation, as shown in Eq. 2. For the rotational component, we applied errors to the roll θ_1 , pitch θ_2 , and yaw θ_3 angles by multiplying the transformation matrix by the corresponding rotation matrices $\mathbf{R}_{i,err}$. For the translational component, the error is introduced by adding the translation vector \mathbf{t}_{err} to the relative positions between the LiDAR and camera sensors:

$$\tilde{\mathbf{T}}_{lid}^{cam} = \mathbf{T}_{lid}^{cam} \cdot \mathbf{R}_{1,err} \cdot \mathbf{R}_{2,err} \cdot \mathbf{R}_{3,err} + \mathbf{t}_{err}. \quad (2)$$

B. Miscalibration Detection

Our approach consists of two steps and is shown in Figure 3. First, we utilize multimodal, pixel-wise contrastive learning to learn the distinct input representations between correctly calibrated and miscalibrated inputs. The embeddings of the frozen encoders are then used for the detection task. This two-step approach using pretext and downstream tasks [30], [31], [32] generally showed promising results compared to end-to-end supervised learning. As described by He *et al.* [30], in the pretext task, the model solely learns input feature representations. The downstream task is the primary task of the model’s application, in our case, miscalibration detection.

The generation of negative samples presents a significant challenge in contrastive learning. For our problem, however, the constraints imposed by the six degrees of freedom in the relative pose between LiDAR and camera enable the efficient generation of negative samples for contrastive learning. These constraints allow the systematic generation of negative pairs by introducing plausible calibration errors into the dataset, as described in the previous section. Because of its ability to generate discriminative representations for similar and dissimilar data points, contrastive learning is particularly interesting for miscalibration detection.

1) *Pixel-wise contrastive learning*: For those reasons, we chose the model architecture based on contrastive learning with two-stream encoders, similar to the architecture proposed by Hadsell *et al.* [33]. However, in contrast to the more common contrastive learning models that utilize augmentations of RGB images, such as [33], [34], [30], our model processes both RGB images and LiDAR data, which could be seen as different views of the same scene or pair, analogous to how contrastive learning typically uses multiple views of an image.

As with other contrastive methods, the model must learn to distinguish between positive and negative pairs. In this framework, positive (correct) pairs consist of RGB images and corresponding LiDAR data. Negative (incorrect) pairs, on the other hand, are the data samples with misalignments between images and LiDAR points due to miscalibrations.

The overall pipeline of the model is shown in Figure 3. First, the 3D LiDAR points are projected into the 2D image plane. In this step, miscalibration errors are introduced through the miscalibrated transformation matrix as explained in Section III-A. Each modality is processed by separate encoders: the RGB images are processed by an image encoder, and the LiDAR data by a LiDAR encoder. These encoders share a common architecture, based on the ResNet18 model, but differ in the input channels: three channels for RGB images and one channel for LiDAR data, representing depth information. Due to different modalities and input channels, unlike unimodal contrastive learning, weight sharing of the encoders cannot be applied.

To learn a feature representation, we only use the first two residual blocks of the ResNet18 architecture for the encoders. This compact model reduces computational load, making it suitable for real-time applications, such as autonomous

driving, where resources are often constrained. Furthermore, smaller models may preserve more spatial details, which can be advantageous when comparing the embeddings [27], which will be discussed in more detail in Section VI. Additionally, given the inherent sparsity of LiDAR data, using a compact model ensures that meaningful features are extracted even when depth information is limited. While Wang *et al.* [27] use a sparsity mask to account for the inherent sparsity of LiDAR data, our framework achieves strong results without the additional step indicating masking has negligible effect.

After feature extraction, the resulting embeddings from the image and LiDAR data samples are compared using a contrastive loss function to measure the distance of the embeddings. The key difference in our approach compared to common contrastive learning methods is the employment of a pixel-wise comparison, which ensures that even calibration errors at the pixel level are considered. The pixel-wise contrastive loss function is adapted from the standard contrastive loss function [33], but instead of projecting the image and LiDAR embeddings into a one-dimensional vector for comparison, the spatial resolution of the embeddings is maintained. The contrastive loss function is defined as follows:

$$L = \frac{1}{NHW} \sum_{n,h,w} [(1-y)D_{nhw}^2 + y \cdot \max(0, m - D_{nhw})^2], \quad (3)$$

where N is the batch size, H and W are the dimension of the embeddings, and m is the margin. The pixel-wise distance is defined as $D_{nhw} = \sqrt{\sum_{c=0}^C (E_{I,nchw} - E_{L,nchw})^2}$. The label $y = 1$ for the incorrect pairs (miscalibrated inputs) and $y = 0$ for the correct pairs (calibrated inputs).

Unlike previous work [28], which uses a cost volume layer to compare and constrain the similarity between the LiDAR and RGB features in the embedding layer before the decoder, we directly optimize both similarity and dissimilarity as the training objective. Additionally, our approach efficiently learns (dis)similarity and alignment of the input representations by omitting the decoder.

2) *Detection task*: After completing the pretext task of training the encoders using our multimodal contrastive learning framework, the learned representations are applied to the subsequent miscalibration detection task. This task determines whether the given sensor inputs are correctly calibrated or miscalibrated based on the embeddings learned from the encoders and is therefore designed as a binary classification task.

As illustrated in Fig. 3, the embeddings from both the LiDAR and image encoders are concatenated along the channel dimension, forming a spatially aligned feature representation that integrates information from both sensor modalities. This concatenation captures the relationships between the LiDAR and corresponding image embeddings, which is crucial for accurately detecting miscalibrations.

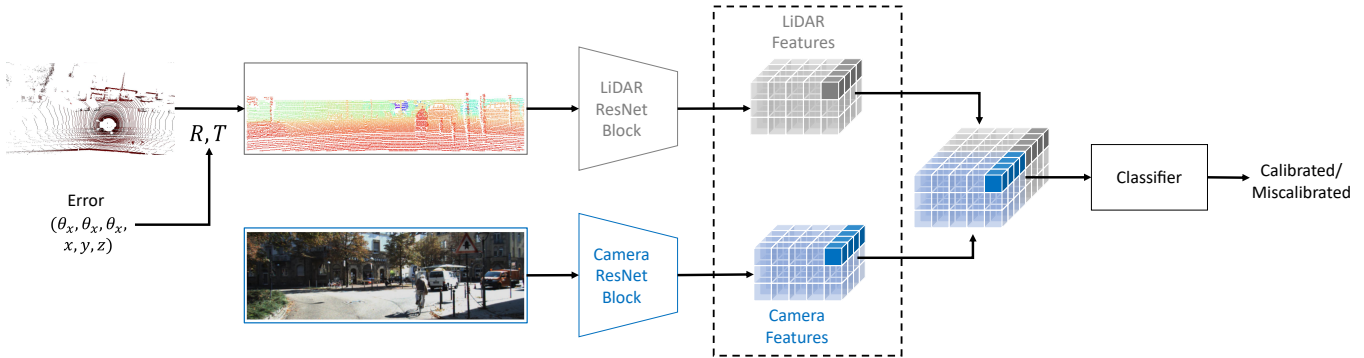


Fig. 3: Pipeline for the classification task within the sensor miscalibration detection. Both image and LiDAR inputs are processed through ResNet blocks, with feature extraction performed for each modality. The ResNet blocks are frozen, and the extracted features are then concatenated and passed to a classifier, which predicts whether miscalibration is detected or not. The feature embeddings inside the dotted box are trained in the first stage using contrastive loss.

The classifier head consists of three 3x3 convolutional layers to further reduce the dimension, followed by global average pooling and three fully connected (MLP) layers with 512 and two 216 neurons. The final layer uses a sigmoid activation function to produce a probability score that indicates whether the sensor inputs are miscalibrated (class 1) or calibrated (class 0). Binary cross-entropy is used as the loss function.

IV. IMPLEMENTATION DETAILS

In this section, we outline the implementation details of our approach, starting with the dataset perturbations used to simulate various calibration errors. We then describe the training configuration for our architecture.

A. Datasets

To train and evaluate our model, we used the KITTI Odometry dataset, which is also used by Wei *et al.* [28]. For comparability, we used the same sequences as in the previous works – the sequences from '01' to '20' comprising 39,011 frames for training and '00' with 4541 frames for the test set. Table I shows the range of errors included in the calibrated and miscalibrated datasets. The calibrated dataset includes small allowed shifts as tolerance to simulate environmental noise and relative deviations between different calibrations in the KITTI dataset. The tolerances are based on benchmarks that investigate the robustness of miscalibration on object detection performance [2], [35]. Accordingly, the ranges are chosen with the expectation that our model will detect even more challenging misalignments. The miscalibrated dataset contains larger perturbations in both translational and rotational components, which reflect the errors we want to detect. In addition, we also introduced a small margin between the calibrated and miscalibrated datasets. Without this margin, the miscalibrated and calibrated data points at the interval boundaries would be similar or even identical. This similarity leads the model to encounter nearly identical or the same input labeled as miscalibrated and calibrated, which introduces ambiguity and ultimately degrades data quality. During the evaluation, we first introduce small noise to the calibrated dataset to monitor whether negligible errors will

Datasets	Translation Error [m]	Rotation Error [°]
Calibrated	[0, 0.02]	[0, 0.3]
Miscalibrated	[0.04, 0.1]	[0.5, 5]

TABLE I: Absolute error ranges of calibrated and miscalibrated dataset for training and validation.

trigger false alarms of miscalibration. We evaluate the model for various miscalibration errors for the test as described in Table II. The goal of using these datasets is to evaluate the model's performance in detecting different error modes with different magnitudes and the generalization ability to detect unseen errors. Lastly, we exclude small intervals between the calibrated and miscalibrated dataset, assuming that the model is allowed to overlook errors within the range.

Datasets	Translation Error [m]	Rotation Error [°]
Noise	[0, 0.005]	[0, 0.1]
Miscalibrated	[0.04, 0.1]	[0.5, 5]
Unseen	[0.1, 0.2]	[5, 10]
All Errors	[0.1, 0.2]	[0.5, 1]
Rot hard	0	[0.5, 1]
Rot easy	0	[1, 5]
Trans hard	[0.04, 0.1]	0
Trans easy	[0.1, 0.2]	0

TABLE II: Absolute error ranges for different dataset configurations for the evaluation and they cover different error modes and magnitudes.

B. Training Configurations

Each phase of training was conducted on an A100 GPU with a batch size of 64, evenly split into 32 calibrated and 32 miscalibrated samples to maintain class balance. The AdamW optimizer was used, with an initial learning rate of 0.001 and a weight decay of 0.05. The learning rate was reduced to 1e-4 after 30 epochs for fine-tuning, and training continued for a total of 50 epochs. For contrastive loss, a margin $m = 4$ was selected.

V. RESULTS AND DISCUSSION

This chapter evaluates the model’s performance across various miscalibration conditions, real-time capability, and the learned representations. For the model performance, we use the evaluation metrics accuracy, precision, and recall to measure the prediction outputs. In terms of miscalibration detection, a high precision score reduces the rate of false positives, which would trigger false alerts for system recalibrations. Higher recall values ensure that most miscalibrations are detected, minimizing the risk of missing actual calibration errors. This is especially important for safety-critical systems like autonomous driving, where undetected miscalibrations could degrade sensor performance and compromise safety. Furthermore, we utilized Centered Kernel Alignment analysis to evaluate the resulting embeddings across layers. CKA analysis is a technique used to compare the similarity of embeddings (i.e., representations) from different layers of the same or different neural networks [36], [37], [38]. This analysis provides insights into how the model processes and represents input data at various stages, allowing us to better understand the internal representation learned by each layer.

A. Miscalibration Detection Performance

The results on the miscalibration performance are shown in Table III and V. We first compare the performance of our model to LCCNet [9] and Wei *et al.* [28].

Metrics	Methods	All Errors	Rot Hard	Trans Easy
Accuracy	LCCNet	90.91%	86.48%	90.44%
	Wei	95.13%	86.28%	92.05%
	Ours	99.08%	99.00%	99.97%
Precision	LCCNet	88.79%	78.69%	85.63%
	Wei	92.02%	78.24%	86.59%
	Ours	100.00%	100.00%	100.00%
Recall	LCCNet	94.04%	99.51%	97.38%
	Wei	99.05%	99.96%	99.65%
	Ours	98.17%	97.99%	99.93%

TABLE III: Results on KITTI Odometry and comparison with existing methods.

Compared to the existing methods, our model shows better accuracy and precision performance with an accuracy of around 99% and a precision of 100% highlighting the robustness of our approach. A robust miscalibration detection refers to a system that does not trigger false positive recalibration for example. Considering that current online calibration algorithms do not achieve target-based calibration accuracies, avoiding false positive recalibrations also contributes to the overall safety of the system. In terms of recall, our method shows a slightly lower performance in the range of around 1% for combined and 2% for rotational errors. The translational errors, however, are more accurately detected by our approach. Looking at the overall performance presented in Table III, the existing methods show over-sensitivity to miscalibration errors that caused reduced precision and poor trade-off between the precision and recall metrics. Note that LCCNet is not specifically provided with a miscalibration

detection, but was extended by [28] with a classification module to be able to benchmark against existing methods. Table IV shows that our model has significantly smaller model size and inference time. This highlights the real-time capability of our model, while achieving good performance in the evaluation metrics.

Methods	Inference Time [ms]	Model Size [M]
LCCNet	97.1	210
Wei	160.7	49
Ours	26.5	28

TABLE IV: Comparison of the inference time and model size evaluated on an Nvidia A100 GPU. The inference time results were obtained according to the evaluation method presented in [28]. Model size refers to the number of parameters of the models in millions.

Table V extends the error modes and magnitudes (referring to the results of the first three rows). Compared to Table III, additional modes such as unseen errors as well as larger rotational errors and more challenging translational errors are evaluated. We observe an increase in performance for larger rotational calibration errors and maintain high detection performance for various sets of miscalibrated datasets including unseen calibration errors. This highlights the strength of our framework to capture very challenging translational errors and to distinguish between rotational and translational errors for the combination of errors.

B. CKA Analysis

We utilized CKA analysis to evaluate the similarity of feature representations learned by the image and LiDAR encoder in different layers. As depicted in Figure 4a, the CKA values across the layers of the two encoders when taking the calibrated dataset as inputs shows that the deeper layers of the network exhibit a higher degree of similarity. When analyzing miscalibrated samples (Figure 4b), the CKA values remain generally low across most layers, indicating a clear difference between feature embeddings generated from miscalibrated data. This shows that the model can effectively

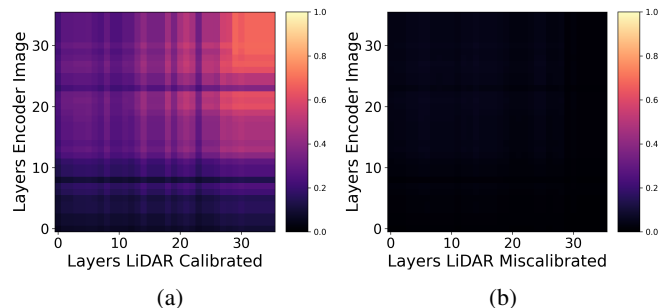


Fig. 4: CKA analysis between LiDAR and image encoders, evaluated on calibrated 4a and miscalibrated input 4b. A CKA value of 1 indicates identical feature embeddings between the encoders in a given layer, whereas a value close to 0 means minimal similarity in their representations.

Metrics	Encoders	Miscalibrated	Unseen	Rot Easy	Rot Hard	Trans Easy	Trans Hard
Accuracy	ResNet18-Small	99.42%	99.97%	99.99%	99.00%	99.97%	99.22%
Precision		100.00%	100.00%	100.00%	100.00%	100.00%	100.00%
Recall		98.84%	99.93%	99.98%	97.99%	99.93%	98.44%
Accuracy	ResNet18-All	98.21%	99.36%	99.11%	94.02%	99.33%	98.56%
Precision		99.36%	99.37%	99.37%	99.30%	99.37%	97.75%
Recall		97.05%	99.35%	98.84%	88.66%	99.29%	99.36%

TABLE V: Performance metrics of different encoders for different test conditions, as shown in Table II.

differentiate between calibrated and miscalibrated inputs, as miscalibration distorts the feature alignment between the two encoders. This behavior is particularly crucial for detecting sensor misalignment.

C. Intrinsic Miscalibration

We evaluate the performance of our algorithm in the presence of intrinsic calibration errors in $\mathbf{P}_{rect}^{(i)}$. Specifically, we introduce relative (percentage) errors in the focal lengths (f_u, f_v) and principal point coordinates (c_u, c_v) by $n\%$. For the skew γ we introduce an error of $n\%$ of f_u , where $n = 10 - 20$, $n = 5 - 10$, $n = 3 - 5$ are denoted as easy, medium and hard, respectively. The results are shown in Table VI. The performance of our approach decreases only slightly for the accuracy and precision metrics when intrinsic calibration errors are introduced. However, the results of the recall metric indicate that not all intrinsic calibration errors can be detected. Nevertheless, we still achieve significantly high scores for all metrics, considering that the framework is trained only on extrinsic errors and therefore showing good generalization capability for unseen error patterns.

Metrics	easy	medium	hard
Accuracy	96%	94%	94%
Precision	99%	99%	99%
Recall	93%	88%	88%

TABLE VI: Results of the miscalibration detection for intrinsic calibration errors.

We now want to discuss the practical implications of this work. It is clear that relying solely on miscalibration detection is not sufficient for autonomous driving tasks. The main practical implication lies in the continuous monitoring of the calibration state of the system, which can be achieved with the computational resources and inference time provided by this method. It should be noted, however, that after recalibrating sensors using regression-based methods, zero calibration errors cannot be guaranteed. Especially, when intrinsic errors are compensated using extrinsic calibration methods. This becomes more severe with axis skew errors, which cannot be compensated using extrinsic recalibration methods. In such cases, our approach can quickly trigger recalibration when initial miscalibration occurs due to real-world effects such as sensor drift or temperature changes, or when corrections are needed after an erroneous recalibration.

VI. ABLATION STUDY

This chapter explores the ablation of encoders by using the entire ResNet18 model (in the following denoted as

ResNet18-All). For this analysis we exclude its average pooling and fully connected layers to account for the deeper encoders, which result in reduced feature dimensions. We also modified the classifier head by simplifying it to a single 3x3 convolutional layer.

A. CKA Analysis

The results in Fig. 5 show that the learned embeddings between the image and the LiDAR encoder trained with ResNet18-All have very low similarity. The reasons are twofold: In the loss function, we only imposed the distance constraint and not the similarity. It is important to note that

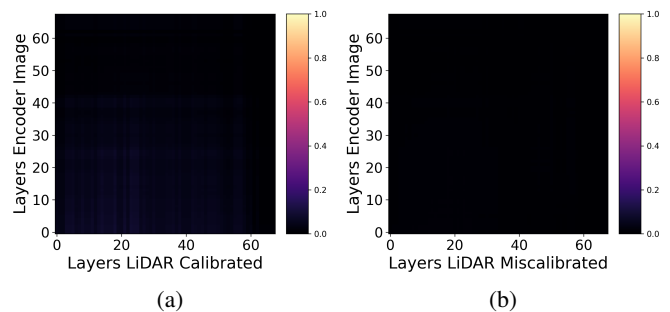


Fig. 5: CKA analysis between LiDAR and image encoders of whole ResNet18 architecture evaluated on calibrated 5a and miscalibrated input 5b.

a small distance between embeddings, as measured by the contrastive loss, does not necessarily imply that the embeddings are similar in terms of their feature representations. In addition, since we could not apply weight sharing, the features of two encoders across different layers are not aligned and most likely learn modality-specific features.

B. Miscalibration Detection Performance

Table V describes the influence of using ResNet18-All as encoders on the performance of the miscalibration detection. In contrast to ResNet18-Small, a lower overall performance in almost all calibration conditions and metrics can be observed. Both the CKA analysis as well as the evaluation metrics validates that our smaller model indeed preserves relevant spatial details. The results indicate that the smaller ResNet18-Small encoder not only achieves better accuracy, precision, and recall but also allows us to employ a more compact model architecture.

VII. CONCLUSION

We introduce a novel framework for safe and real-time miscalibration detection in LiDAR-camera sensor setups,

focusing on binary classification to distinguish calibrated from miscalibrated states. Our approach leverages contrastive learning and a simple and efficient architecture to achieve state-of-the-art performance. Through our comprehensive analysis, we identify the depth at which latent features should be compared, enabling us to simplify the feature extraction and matching layers without compromising accuracy. Our experiments demonstrate that our method not only achieves state-of-the-art performance in various miscalibration scenarios, but also requires less computational resources and inference time to detect miscalibration, making it well suited for real-time applications.

REFERENCES

- [1] L. Wang *et al.*, “Multi-modal 3d object detection in autonomous driving: A survey and taxonomy,” *IEEE Transactions on Intelligent Vehicles*, vol. 8, pp. 3781–3798, 7 2023.
- [2] K. Yu *et al.*, “Benchmarking the robustness of lidar-camera fusion for 3d object detection,” in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2023.
- [3] Y. Dong *et al.*, “Benchmarking robustness of 3d object detection to common corruptions in autonomous driving,” in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 8, 3 2023. [Online]. Available: <http://arxiv.org/abs/2303.11040>
- [4] X. Bai *et al.*, “Transfusion: Robust lidar-camera fusion for 3d object detection with transformers,” in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [5] Y. Xie *et al.*, “Sparsefusion: Fusing multi-modal sparse representations for multi-sensor 3d object detection,” in *ICCV*, 2023. [Online]. Available: <https://github.com/yichen928/SparseFusion>.
- [6] Z. Song *et al.*, “Graphalign: Enhancing accurate feature alignment by graph matching for multi-modal 3d object detection,” in *ICCV*, 2023.
- [7] Z. Liu *et al.*, “Befv-fusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation,” in *IEEE International Conference on Robotics and Automation*, 5 2023. [Online]. Available: <http://arxiv.org/abs/2205.13542>
- [8] P. An *et al.*, “Survey of extrinsic calibration on lidar-camera system for intelligent vehicle: Challenges, approaches, and trends,” *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–25, 7 2024.
- [9] X. Lv *et al.*, “Lccnet: Lidar and camera self-calibration using cost volume network,” in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2021. [Online]. Available: <https://github.com/LvXudong-HIT/LCCNet>
- [10] J. Zhu, J. Xue, and P. Zhang, “Calibdepth: Unifying depth map representation for iterative lidar-camera online calibration,” in *IEEE International Conference on Robotics and Automation*, vol. 2023-May. Institute of Electrical and Electronics Engineers Inc., 2023, pp. 726–733.
- [11] C. Yuan *et al.*, “Pixel-level extrinsic self calibration of high resolution lidar and camera in targetless environments,” *IEEE Robotics and Automation Letters*, vol. 6, pp. 7517–7524, 10 2021.
- [12] F. Chen *et al.*, “Pbcalib: Targetless extrinsic calibration for high-resolution lidar-camera system based on plane-constrained bundle adjustment,” *IEEE Robotics and Automation Letters*, vol. 8, pp. 304–311, 1 2023.
- [13] T. Ma, Z. Liu, G. Yan, and Y. Li, “Crlf: Automatic calibration and refinement based on line feature for lidar and camera in road scenes,” 3 2021. [Online]. Available: <http://arxiv.org/abs/2103.04558>
- [14] R. Ishikawa *et al.*, “Lidar-camera calibration using intensity variance cost,” in *IEEE International Conference on Robotics and Automation (ICRA)*. Institute of Electrical and Electronics Engineers Inc., 2024, pp. 10 688–10 694.
- [15] Z. Liu *et al.*, “Semalign: Annotation-free camera-lidar calibration with semantic alignment loss,” in *IEEE International Conference on Intelligent Robots and Systems*. Institute of Electrical and Electronics Engineers Inc., 2021, pp. 8845–8851.
- [16] Z. Luo *et al.*, “Zero-training lidar-camera extrinsic calibration method using segment anything model,” in *IEEE International Conference on Robotics and Automation*. Institute of Electrical and Electronics Engineers Inc., 2024, pp. 14 472–14 478.
- [17] Z. Lin *et al.*, “Sgcalib: A two-stage camera-lidar calibration method using semantic information and geometric features,” in *IEEE International Conference on Robotics and Automation (ICRA)*. Institute of Electrical and Electronics Engineers Inc., 2024, pp. 14 527–14 533.
- [18] N. Schneider *et al.*, “Regnet: Multimodal sensor registration using deep neural networks,” in *2017 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2017, p. 1927.
- [19] I. Ganesh, R. K. Ram, K. Murthy, and K. M. Kirshna, “Calibnet: Geometrically supervised extrinsic calibration using 3d spatial transformer networks,” in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018.
- [20] K. Yuan, Z. Guo, and Z. J. Wang, “Rggnet: Tolerance aware lidar-camera online calibration with geometric deep learning and generative model,” *IEEE Robotics and Automation Letters*, vol. 5, pp. 6956–6963, 10 2020.
- [21] G. Wang, J. Qiu, and Y. Guo, “Fusionnet: Coarse-to-fine extrinsic calibration network of lidar and camera with hierarchical point-pixel fusion,” in *IEEE International Conference on Robotics and Automation (ICRA)*, 2022.
- [22] X. Hu *et al.*, “Lidar-camera extrinsic calibration with hierarchical and iterative feature matching,” in *IEEE International Conference on Robotics and Automation (ICRA)*. Institute of Electrical and Electronics Engineers Inc., 2024, pp. 16 691–16 697.
- [23] Q. Herau, N. Piasco, M. Bennehar, L. Roldão, D. Tsishkou, C. Migniot, P. Vasseur, and C. Demonceaux, “Soac: Spatio-temporal overlap-aware multi-sensor calibration using neural radiance fields,” in *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Institute of Electrical and Electronics Engineers Inc., 2024.
- [24] J. Levinson and S. Thrun, “Automatic online calibration of cameras and lasers,” in *Robotics: Science and Systems*, 2013.
- [25] M. Ma *et al.*, “Deep coupling autoencoder for fault diagnosis with multimodal sensory data,” *IEEE Transactions on Industrial Informatics*, vol. 14, pp. 1137–1145, 3 2018.
- [26] J. Qian *et al.*, “A review on autoencoder based representation learning for fault detection and diagnosis in industrial processes,” *Chemometrics and Intelligent Laboratory Systems*, vol. 231, 2022.
- [27] Y. Chen *et al.*, “Cross-modal matching cnn for autonomous driving sensor data monitoring,” in *Proceedings of the IEEE International Conference on Computer Vision Workshop*, vol. 2021-October. Institute of Electrical and Electronics Engineers Inc., 2021, pp. 3103–3112.
- [28] P. Wei *et al.*, “Online lidar-camera extrinsic parameters self-checking and recalibration,” *Measurement Science and Technology*, vol. 35, 10 2024.
- [29] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, “Vision meets robotics: The kitti dataset,” *International Journal of Robotics Research (IJRR)*, vol. 32, 2013. [Online]. Available: <http://www.cvlibs.net/datasets/kitti>.
- [30] K. He *et al.*, “Momentum contrast for unsupervised visual representation learning,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 9729–9738.
- [31] I. Misra and L. v. d. Maaten, “Self-supervised learning of pretext-invariant representations,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 6707–6717.
- [32] X. Chen, S. Xie, and K. He, “An empirical study of training self-supervised vision transformers,” in *ICCV*, 2021, pp. 9640–9649.
- [33] R. Hadsell, S. Chopra, and Y. LeCun, “Dimensionality reduction by learning an invariant mapping,” in *2006 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2. IEEE, 2006, pp. 1735–1742.
- [34] T. Chen *et al.*, “A simple framework for contrastive learning of visual representations,” in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [35] M. Fürst *et al.*, “Learned fusion: 3d object detection using calibration-free transformer feature fusion,” in *International Conference on Pattern Recognition Applications and Methods (ICPRAM-2024)*, 2024. [Online]. Available: <https://orcid.org/0009-0000-0711-229X>
- [36] S. Kornblith *et al.*, “Similarity of neural network representations revisited,” in *International conference on machine learning*. PMLR, 2019, pp. 3519–3529.
- [37] C. Cortes, M. Mohri, and A. Rostamizadeh, “Algorithms for learning kernels based on centered alignment,” *The Journal of Machine Learning Research*, vol. 13, pp. 795–828, 2012.
- [38] M. Raghu *et al.*, “Do vision transformers see like convolutional neural networks?” *Advances in neural information processing systems*, vol. 34, pp. 12 116–12 128, 2021.