# Novelty Detection and Analysis of Traffic Scenario Infrastructures in the Latent Space of a Vision Transformer-Based Triplet Autoencoder

Jonas Wurst[1], Lakshman Balasubramanian[1], Michael Botsch[1] and Wolfgang Utschick[2]

*Abstract*— Detecting unknown and untested scenarios is crucial for scenario-based testing. Scenario-based testing is considered to be a possible approach to validate autonomous vehicles. A traffic scenario consists of multiple components, with infrastructure being one of it. In this work, a method to detect novel traffic scenarios based on their infrastructure images is presented. An autoencoder triplet network provides latent representations for infrastructure images which are used for outlier detection. The triplet training of the network is based on the connectivity graphs of the infrastructure. By using the proposed architecture, expert-knowledge is used to shape the latent space such that it incorporates a pre-defined similarity in the neighborhood relationships of an autoencoder. An ablation study on the architecture is highlighting the importance of the triplet autoencoder combination. The best performing architecture is based on vision transformers, a convolution-free attention-based network. The presented method outperforms other state-of-the-art outlier detection approaches.

## I. INTRODUCTION

A simple statistical proof of an *Autonomous Vehicle's* (AV's) safety is infeasible. Approaches like scenario-based testing are used for validation. There, the testing of an AV is focused on relevant scenarios, instead of exclusively driving randomly millions of kilometers. Identifying representative scenarios is required for this approach [1]. The scenarios can be constructed by expert-knowledge or from real world driving data. A key aspect of the latter strategy is to identify new scenarios which have not been tested yet. The two commonly used approaches to detect unknown scenarios are cluster assignment quality or outlier detection. The latter is realized in this work.

A traffic scenario is described by multiple aspects, for example the dynamics and the environment [2]. Publications often focus on the dynamics to identify representative and novel scenarios (e. g. [3], [4], [5], [6]). Besides dynamics, another crucial component of a scenario is the infrastructure. Here, birds-eye view images, representing the infrastructure, are used to detect novel traffic scenarios.

In this work, a method to detect unknown and potentially untested infrastructures is presented. For this purpose, the infrastructure images are projected into a latent space using a deep learning pipeline. As will be shown, a high performance increase is achieved when using the latent space instead of the input space for novelty detection. In the latent space, simple outlier detection methods can be used to identify novel infrastructures, which indicates that the method is able to

generate strong representations. In order to create this latent space, an autoencoder architecture utilizing metric learning via triplet loss is used. The triplet mining is based on the connectivity of the infrastructures. Through the combination of the autoencoder scheme and the triplet learning, expert-knowledge is used for shaping the latent space.

Extensive evaluation of the presented method is performed. For the encoder, state-of-the-art networks such as *Vision Transformers* (ViTs) [7] and ResNet-18 [8] are evaluated. Experiments demonstrate the influence of the triplet loss as well as the autoencoder scheme. The resulting architecture combinations are outperforming the alternative methods. An implementation of the architecture is made publicly available[1].

The contributions of this work can be summarized as:

1) A new method for novelty detection of infrastructure images in the latent space of triplet loss-based autoencoder networks is presented.
2) Automated triplet mining of road infrastructures without manual labeling is introduced.
3) The performance is evaluated against various state-of-the-art methods and shows significant improvements.

This work is organized as follows. In Sec. II the related work in the field of novel traffic scenario detection and clustering are discussed and compared to this paper. The method, consisting of the data generation, triplet mining, triplet network and the outlier detection is introduced in Sec. III. Sec. IV presents the experimental results, when applying the proposed method to a road infrastructure data set. Finally, the work is concluded in Sec. V.

## II. RELATED WORK

### A. Traffic Scenario Identification

The validation of AVs through scenario-based testing is considered to be one of the possible approaches to prove their safety [1]. In the survey [9], various works utilizing the scenario-based validation approach are summarized. The survey also lists works trying to identify and define relevant scenarios, of which some are based on machine learning.

Several works use clustering for novelty detection. The underlying assumption is that scenarios which do not belong to a certain cluster, can be treated as novel scenarios. Since in this work a learned representation for the infrastructural part of a scenario is introduced, first the works are grouped with respect to their used scenario similarity or the used scenario representation. In [10] and [11] a similarity measure based on

[1]CARISSMA, Technische Hochschule Ingolstadt, 85049 Ingolstadt, Germany {firstname.lastname}@thi.de
[2]Technical University of Munich, 80333 Munich, Germany utschick@tum.de

[1]https://github.com/JWTHI/ViTAL-SCENE

an unsupervised random forest is used. More similar to this work are [3] (LSTM+CNN), [4] (SeqDSPN), [5] (RC-GAN), where the latent representations of deep neural networks are used to cluster scenarios. Also a lot of non-machine learning approaches were used to determine representations or similarities of scenarios, such as in [3] (DTW), [6] (DTW+PCA), [12] (Dynamic-Length-Segmentation), [13] (custom similarity measure). In this work, an expert-knowledge aided latent representation is introduced, and hence differs from the aforementioned works. To the best of the authors knowledge, this is the first work utilizing a triplet-based autoencoder scheme for the novelty detection of traffic scenarios.

Another viewpoint to compare this work to others is the information used from the scenario. In this work, only the static information is utilized. The most infrastructural information is considered in [12], where categorical and continuos variables are used to roughly describe the environment. In the works [10], [11] only little static information are used. Whereas in the works [3], [4], [5], [6] and [13] only spatial and dynamic information is used. This work focuses on the static environment of a traffic scenario. Contrary to the former works images of the infrastructure are used here.

In this work, novelty detection is used instead of clustering. While some works focus on the detection of anomalies or corner cases in videos or images (e. g. [14],[15]), there are only a few works addressing the detection of novel traffic scenarios through outlier detection on a scenario level. In [16] the novelty of traffic scenarios is estimated through an autoencoder. The autoencoder is trained on a data set containing known scenarios. If a scenario is passed through the network, it is assumed to produce a high reconstruction error if it is novel. This reconstruction assumption is a widely used mechanism when performing outlier detection with deep learning. In this work, it is refrained from using the reconstruction paradigm, rather to learn latent representations which can be used with simpler outlier detection mechanisms. In [17] the novel outlier method ULEF based on directed $k$-nearest neighbor graphs is presented. The method is applied to road infrastructure images. In this work, the road infrastructure image generation as in [17] is used. Moreover, the performance of ULEF will be compared to this work.

### B. Triplet Learning and Outlier Detection

This work is using triplet learning to form the latent space. Triplet networks [18] are used to perform deep metric learning through ranking loss. In tile2vec [19], triplet networks are used for geo-spatial analysis. The sampling is determined based on the distances of the tiles. The architecture used in this work consists of an autoencoder structure. In [20], an autoencoder network is combined with triplet learning. The difference of this work to [20] lies in the application, the specific network architecture and the triplet mining.

To aid the detection of outliers, some works are using metric learning. An example is [21], where out-of-distribution data is used for the training.

In the field of deep learning various approaches to detect outliers exist. The most common approach is to use the

reconstruction paradigm. For example in [22], a GAN is used as generator network, where an additional encoder is trained to learn the mapping from an input image to the GAN's latent representation. An extension of the reconstruction paradigm is using the hidden reconstructions additionally [23].

In this paper, standard outlier detection methods are applied in the latent space, namely the *Local Outlier Factor* (LOF) [24], *Angle-Based Outlier Detection* (ABOD) [25], *Isolation Forest* (IF) [26] and the *One-Class Support Vector Machine* (OCSVM) [27].

### III. METHOD

In this work, an autoencoder network utilizing triplet loss is used to project road infrastructure images into the latent space, where the outlier detection is performed. By using the proposed architecture, a latent representation for road infrastructure images is learned. During training, the infrastructure is assessed in two ways. First, the visual appearance of the images, hence the shape of roads etc., is considered via the reconstruction loss. Second, the similarity of infrastructure topologies is taken into account via the similarity of their connectivity graphs, allowing data triplets to be identified. A data triplet consists of three data points: the anchor, a sample similar to the anchor and a sample dissimilar to the anchor, which are required for the triplet loss.

This section is split into the following parts. First, the data generation and used similarity measure is explained. Second, the realized triplet autoencoder network as well as the used triplet mining is summarized. Last, the traffic scenario detection in the latent space is described.

### A. Data Generation and Similarity Measure

In this paper, the static part of a traffic scenario is described through its road infrastructure. For this purpose, a black and white image of the infrastructure in top-view is generated. Furthermore, a connectivity graph is generated, which will be required to generate the similarity between infrastructure images.

The used tool-chain consists of the following steps: a) position selection, b) map data collection, c) image generation and d) connectivity graph generation. *OpenStreetMap* (OSM) [28] is used for the map data. Also the position selection is realized via OSM.

*a) Position Selection:* Within a search area, all OSM nodes can be considered as positions. For this work the valid positions are restricted to public and car-drivable roads. Identifying the nodes and positions can be realized through the API of OSM for example. Each selected position leads to a single data set entry, consisting of an image and the corresponding connectivity graph.

*b) Map Data Collection:* For each selected position, it is required to get the associated map data. In this work, the OSM map with a parameterizable bounding box around the location is converted into an openDRIVE [29] map, such that the image generation tool as introduced in [17] can be used. For the conversion from OSM to openDRIVE, the netconvert tool of SUMO [30] is used.
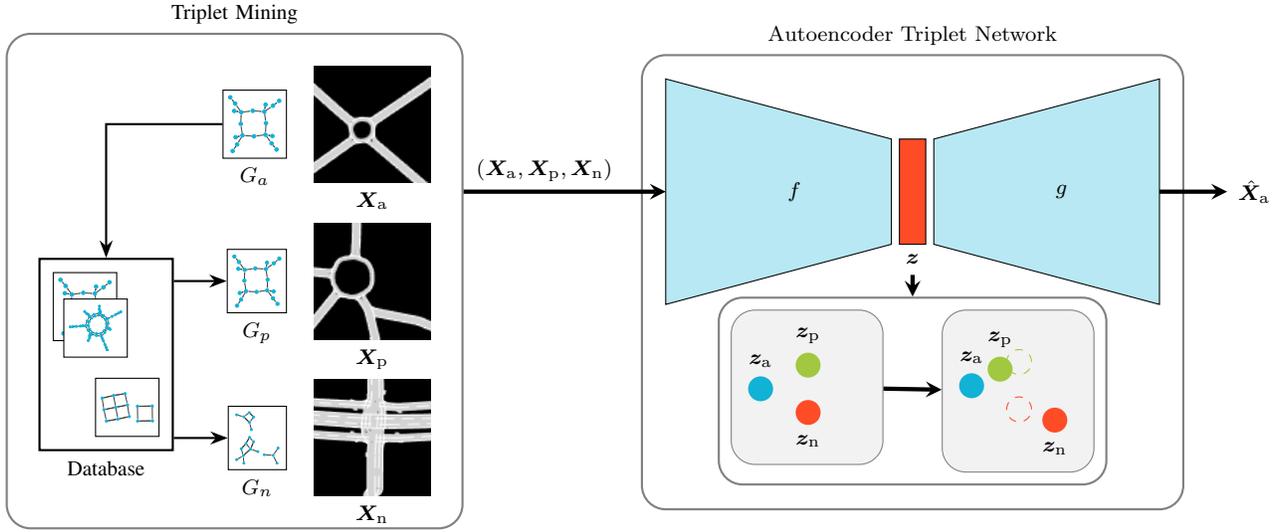
Fig. 1. Triplet Mining and Autoencoder Network

*c) Image Generation:* The image generation is realized as proposed in [17]. Hence, the roads are colored gray, lane markings white and the background black. Furthermore, all non-public and non-car accessible roads are not rendered. For each position an image is generated.

*d) Graph Generation:* The graph generation is new with respect to [17] and is realized as follows. First, the complete openDRIVE map is converted into a network graph, using the connectivity and neighbor information. The result is somewhat comparable to the routing graph realized in Lanelet2 [31]. In the next step, the selected location is assigned to one of the graph's node. Hence, the position on the road is identified. Then, the graph is cropped and simplified using the following rules:

- Use all nodes which can be reached within the time $t_{\max}$, given the allowed speed, up to and including all nodes of the first junction (e. g. crossing, roundabout).
- Use all nodes which are neighboring lanes.

The above steps are used to generate the data set $\mathcal{D} = \{(\boldsymbol{X}_1, G_1), \ldots, (\boldsymbol{X}_M, G_M)\}$, where $M$ is the number of selected positions, and hence the number of resulting images and graphs. The image for the $m$-th position is given by the matrix $\boldsymbol{X}_m \in \mathbb{R}^{S \times S}$ with $S$ the selected resolution of the image. Analogous, the graph for the $m$-th position is given by $G_m = (V_m, E_m)$, where $V_m$ are the vertices and $E_m$ the edges of the graph. The actual selected position is not used after the corresponding image and graph are extracted.

*e) Similarity Measure:* For the triplet-based learning, a notation of being similar or dissimilar is required. In this work, the graphs are used for this purpose. Two cropped areas are considered to be similar if their connectivity graphs are similar. The graphs need to be permuted versions of each other to be similar. This is the case if there exists an isomorphism between the two graphs. Given the graphs $G_i$ and $G_j$, they are similar if there exists a bijection $p : V_i \rightarrow V_j$ such that $(u, v) \in E_i \iff (p(u), p(v)) \in E_j$. Using

the notation $G_i \cong G_j$ for two graphs being isomorph, the similarity function is defined as

$$s(G_i, G_j) = \begin{cases} 1 & \text{if } G_i \cong G_j \\ 0 & \text{else} \end{cases}. \tag{1}$$

This similarity measure considers all data points as same, when their connectivity is the same. Within the triplet mining block of Fig. 1 examples of extracted graphs alongside their images are shown.

*B. Triplet Autoencoder*

The main part of this work is the triplet-based autoencoder network. It is used to produce latent representations for given input images of road infrastructures. This section will address the architectural details, the used loss, the triplet mining and details on the used networks.

Triplet learning realizes metric learning via a ranking loss. The objective is to enforce similarity in the latent space based on data triplets. As explained before, each triplet consists of an anchor, a positive example and a negative example. The anchor and the positive example are similar, according to the used similarity measure. Consequently, the negative example is dissimilar to the anchor. For example, in [18], the anchor is an image of a face, the positive example is another image of the same person and the negative is an image of another person. The objective of the training is to push the latent representation of the negative example away from the latent representation of the anchor while pulling the latent representation of the positive example closer. This way, the metric learning is realized.

In this work, a triplet learning scheme is used, where the road infrastructure images are used as input, with the anchor $\boldsymbol{X}_a$, the positive example $\boldsymbol{X}_p$ and the negative example $\boldsymbol{X}_n$. Let $f$ be a trainable network, realizing the mapping from the input representation to the latent representation $f : \boldsymbol{X} \mapsto \boldsymbol{z}$ with $\boldsymbol{z} \in \mathbb{R}^L$ being the latent representation with dimensionality $L$. During the training, each sample of

the triplet $(\boldsymbol{X}_\mathrm{a}, \boldsymbol{X}_\mathrm{p}, \boldsymbol{X}_\mathrm{n})$ is passed through $f$ separately, leading to the latent triplet $(\boldsymbol{z}_\mathrm{a}, \boldsymbol{z}_\mathrm{p}, \boldsymbol{z}_\mathrm{n})$. During inference, only single samples will be passed through $f$. The training of $f$ is partially realized by using the triplet loss

$$\mathcal{L}_\mathrm{tri}(\boldsymbol{X}_\mathrm{a}, \boldsymbol{X}_\mathrm{p}, \boldsymbol{X}_\mathrm{n}) = \max\left(\alpha + d_\mathrm{ap} - d_\mathrm{an}, 0\right), \quad (2)$$

with the squared distance between the anchor representation and the positive example representation $d_\mathrm{ap} = ||f(\boldsymbol{X}_\mathrm{a}), f(\boldsymbol{X}_\mathrm{p})||_2^2$, the squared distance between the anchor representation and the negative example representation $d_\mathrm{an} = ||f(\boldsymbol{X}_\mathrm{a}), f(\boldsymbol{X}_\mathrm{n})||_2^2$ and the margin $\alpha$. The triplet loss' objective is twofold, to lower the distance between the positive example and the anchor $d_\mathrm{ap}$ while simultaneously increase the distance between the negative example and the anchor $d_\mathrm{an}$ until $d_\mathrm{ap} + \alpha$. This way the latent representation is forced to follow the definition of similarity as provided by the triplets.

One key point of the triplet scheme is the triplet mining. While sampling $\boldsymbol{X}_\mathrm{a}$ is realized by randomly picking an image from the data set $\mathcal{D}$, determining the positive and negative examples is based on the anchor. The similarity of two road infrastructure images is defined through Eq. 1. Hence, the positive example is randomly picked out of all samples having the same connectivity as the anchor. Given the graph of the anchor as $G_\mathrm{a}$ the graph of the positive example $G_\mathrm{p}$ is picked from $\mathcal{G}_\mathrm{p} = \{G|s(G, G_\mathrm{a}) = 1\}$, given $G \in \mathcal{G}$, where $\mathcal{G}$ is the set of all graphs of $\mathcal{D}$. Using the same notation, the graph of the negative example $G_\mathrm{n}$ is picked from $\mathcal{G}_\mathrm{n} = \{G|s(G, G_\mathrm{a}) = 0\}$.

The selection of especially the negative sample has a significant influence on the training. Considering the case where the negative is already too far away and hence there is no training contribution for this triplet, since $d_\mathrm{an} \geq d_\mathrm{ap} + \alpha$ would lead to $\mathcal{L}_\mathrm{tri} = 0$. Such samples are called easy negatives. On the other hand-side having hard negatives, i.e. $d_\mathrm{an} < d_\mathrm{ap}$, might lead to bad local minima early in the training [18]. Both types should be prevented. Therefore, the objective is to sample data points that are called semi-hard negative samples. They are further away than the positive example but still within the margin $\alpha$, such that $d_\mathrm{ap} < d_\mathrm{an} < d_\mathrm{ap} + \alpha$ holds.

The triplet loss conditioned training might be sufficient to separate different connectivity types, for example roundabout with 4 versus 3 exit roads. However, another objective in this work is to ensure that neighboring points in the latent space have visually very similar input images. The reasoning for this is twofold. First, it is assumed that the shape of the road also has an influence on the novelty of a scenario. Furthermore, it is assumed, the more similar the neighborhood in the latent space is, the more reliable the projection is for unknown data. By this, even within the same connectivity group further refinement is achieved. For this purpose, a decoder is introduced into the overall architecture. The decoder is used to regularize the latent space, such that it can be used to reconstruct the anchor $\boldsymbol{X}_\mathrm{a}$. The underlying assumptions is, that for the reconstruction to work, the latent space has to be formed in such away, that the neighbors

are not sharing only connectivity-based similarities but also visual similarities. The trainable decoder network $g$ is used to generate the reconstructed anchor image $\hat{\boldsymbol{X}}_\mathrm{a}$ by $g : \boldsymbol{z} \mapsto \hat{\boldsymbol{X}}$. The reconstruction loss of the anchor,

$$\mathcal{L}_\mathrm{rec}(\boldsymbol{X}_\mathrm{a}) = ||\boldsymbol{X}_\mathrm{a} - g\left(f\left(\boldsymbol{X}_\mathrm{a}\right)\right)||_2^2 \quad (3)$$

is used to enforce the above described objective.

The complete architecture is trained using the loss

$$\mathcal{L} = \mathcal{L}_\mathrm{tri} + \mathcal{L}_\mathrm{rec}, \quad (4)$$

such that both objectives are fulfilled: connectivity-based structure refined by visual similarity. The overall pipeline of the triplet autoencoder network including the triplet mining can be seen in Fig. 1. The exemplary graphs are the results for the shown images. This highlights the need for the local refinement by the decoder, since the shown anchor and its positive example are sharing the same connectivity but have quite different visual appearance. Below the network, the learning process indicated by the triplet loss is symbolized.

The triplet autoencoder network can be trained on a huge data set such as OSM, ensuring a strong projection method. After training, the network can be used during inference phase to project new data. This way, it is possible to train the network only once. Furthermore, it allows one to store the data in a compressed format, i.e. the latent representations.

The decoder network is kept fairly simple in order to limit the complexity of the latent space. The decoder network structure is kept the same for all experiments, only the encoder is varied. For the encoder network the ResNet and ViT are tested and are therefore briefly explained in the following.

*a) ResNet:* In [8] the widely used ResNets 18, 50 etc. are proposed. Residual networks consist of multiple residual blocks with multiple convolutional layers each. The input to each residual block is added to the output again. Since those networks are widely known, further details will not be given here but can be found in [8].

*b) Vision Transformer:* Most recently, ViTs have been introduced in [7], using the attention mechanism-based transformers as in [32]. ViT is a convolution-free network for image classification, showing state-of-the-art performance.

The transformer usually uses an encoder and a decoder for sequence to sequence modelling, but the ViT is using only the encoder part, as shown in Fig. 2. First, an image is split into flattened patches which are linearly projected, leading to the input embeddings (marked blue). Then, an additional embedding token (orange) is concatenated to the input embeddings before added to the learnable positional embedding (green). The sum is fed as an input to the transformer. Inside the transformer, a multi-layer multi-head-attention network is realized.

The output token, corresponding to the additional embedding token, is fed through a *Multi Layer Perceptron* (MLP). Here, in difference to the original ViT, the output vector is the latent representation $\boldsymbol{z}$ (red), whereas in ViT it is used to predict a class label.
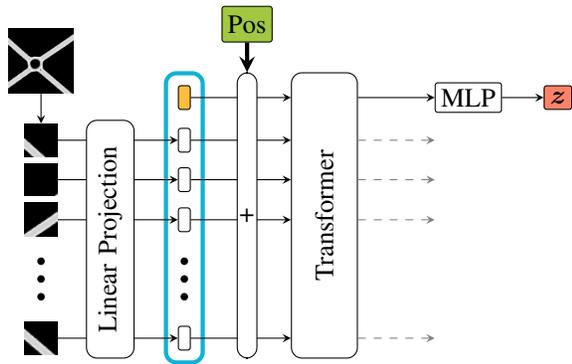
Fig. 2. Vision Transformer as used in this work. Inspired from [7].

## C. Novelty Detection

The overall objective in this work is to detect unknown road infrastructures. For this purpose, outlier detection methods can be applied. Let the base data set $\mathcal{D}_{\mathrm{base}}$ be the data already known, for example, the scenarios an autonomous driving function has already been tested on. If a data point is tested with respect to its outliership it is basically investigated if the point fits into the known data $\mathcal{D}_{\mathrm{base}}$. The training data $\mathcal{D}$, used to train the projection $f$, is not required to be the same as the base data $\mathcal{D}_{\mathrm{base}}$.

Provided with the latent infrastructure representation $z$, simple outlier mechanisms can be applied directly on the latent space instead of the input space. This way, the enforced similarity utilizing connectivity and shape will mainly be responsible for detecting outliers. In the following, let $o$ be the outlier ranking function fitted on $\mathcal{D}_{\mathrm{base}}$, mapping each input to an estimated outlier score $v = o(z)$.

Within this work various outlier detection methods are used. Methodological details are not covered in this work. Interested readers may refer to the corresponding works. In the following the used methods are briefly summarized.

*a) Local Outlier Factor:* The LOF [24] is analyzing how isolated a data point is with respect to its neighbors. For this the local densities are used.

*b) Isolation Forest:* The IF [26] is estimating the outlierness through the number of splits required to isolate a data point, given randomly grown trees. It holds, that the fewer splits the more outlying.

*c) Angle-Based Outlier Detection:* Another neighborhood-based approach is ABOD [25]. The variance of the direction vector's angle from the point under investigation to all other data points is investigated. The higher the variance the lower the outlierness. In this work, the fast version of ABOD is used, considering only the $k$ nearest neighbors of the investigated point.

*d) One-Class Support Vector Machine:* The OCSVM [27] is the extension of the normal SVM to the problem of outlier detection. The data is mapped to the feature space as by the kernel. In the feature space, the data is separated from the origin via a hyperplane. This way, the decision boundary is meant to encapsulate the data in the input space.

*e) UMAP-based Local Entropy Factor:* In [17] the neighborhood of data point's is used to define local similarity measures. A points outliership is evaluated based on how well it fits into its neighbors' neighborhoods using the entropy and the point-wise similarity.

## IV. EXPERIMENTS

The introduced method is evaluated in this section. Various architecture realizations and outlier detection methods are compared. This section is split into the following parts. First, the details about the data used for the outlier detection and the data used for visual analysis is explained. Second, the analyzed architectures are briefly summarized. The outlier detection performance is discussed in the third part. The various resulting latent space visualizations are shown and discussed in the fourth part. In the fifth subsection, the local visual similarity quality is assessed, proposing another possibility to evaluate the performance of the shown architectures. The results are summarized in the last part.

### A. Data sets

The training data set $\mathcal{D}$ is decoupled from the base data set used for the outlier detection $\mathcal{D}_{\mathrm{base}}$. Both are briefly described in the following. For both the images are showing a region of size $100\,\mathrm{m} \times 100\,\mathrm{m}$, while the reachable time for generating the graphs is selected as $t_{\max} = 5\,\mathrm{s}$.

The training data set $\mathcal{D}$ consists of $\approx 70\,000$ pairs of images and graphs. To provide an insight to the data set, it has been analyzed with respect to rough groups, but those groups are not used in the training. In total, approx. $13\,400$ highway, $16\,900$ roundabouts, $18\,000$ crossings, $19\,900$ single lane and $1\,700$ multiple lane non-highway pairs are used. The extraction region for the highway and roundabout pairs is selected to be the complete district of Upper Bavaria in Germany. The extraction region for the remaining types is the city of Ingolstadt in Germany with its adjacent counties.

The base data set $\mathcal{D}_{\mathrm{base}}$ to fit the models of the outlier detection methods is taken from [17]. Therefore, the data considered to be known are highway images only. The outlying test data is taken from [17] as well. It consists of rural images and inner-city images, being gathered in the anomaly data set $\mathcal{D}_{\mathrm{ano}}$. It holds that $\mathcal{D}_{\mathrm{base}} \in \mathcal{D}$ and $\mathcal{D}_{\mathrm{ano}} \in \mathcal{D}$.

### B. Architectures

In order to investigate the influence of various architecture realizations, different versions of the network were used for the experiments. The architecture of the decoder is kept the same for all the experiments, enabling a fair comparison. The encoder is realized through different networks. Furthermore, the loss function was changed for some experiments, highlighting the importance of the overall pipeline. Here, the various used options will be briefly summarized.

For all architectures, the following parameters hold: image size $64 \times 64$, epochs 200, latent dimensionality 50.

*a) Res-S:* A small ResNet, as ResNet-18 but consisting of fewer parameters. Learnable parameters $\approx 0.3 \cdot 10^{6}$.

*b) Res-18:* A larger ResNet. Here the basic implementation of ResNet-18 [8] is adjusted for single channel inputs. Learnable parameters $\approx 11.0 \cdot 10^{6}$.

TABLE I
OUTLIER DETECTION PERFORMANCE – AUC

| Architecture | | | OD Input | LOF | ULEF | IF | OCSVM | ABOD | CAE-ORD | CAE-SAP | f-AnoGAN |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| $f$ | $g$ | $\mathcal{L}_{\mathrm{tri}}$ | | | | | | | | | |
| Res-S | ✗ | ✓ | $z$ | 0.913 | 0.775 | 0.935 | **0.959** | 0.892 | – | – | – |
| Res-S | ✓ | ✓ | $z$ | **0.954** | 0.696 | 0.933 | 0.950 | **0.954** | – | – | – |
| Res-18 | ✓ | ✗ | $z$ | 0.696 | 0.776 | 0.770 | 0.752 | **0.784** | – | – | – |
| Res-18 | ✓ | ✓ | $z$ | 0.949 | 0.781 | 0.917 | 0.912 | **0.956** | – | – | – |
| ViT-S | ✓ | ✓ | $z$ | 0.730 | 0.689 | 0.698 | 0.910 | **0.920** | – | – | – |
| ViT-L | ✓ | ✓ | $z$ | 0.900 | 0.793 | 0.707 | 0.937 | **0.956** | – | – | – |
| – | | | $X$ | 0.446 | 0.612 | 0.196 | 0.247 | 0.700 | 0.855 | 0.845 | 0.758 |

*c) ViT-S:* A small version of the ViT. Learnable parameters $\approx 0.2 \cdot 10^6$. (patch size $8 \times 8$, layers 6, dim. of input embedding $n_{\mathrm{patches}} \times 64$, internal MLP dim. 128).

*d) ViT-S:* A large version of the ViT. Learnable parameters $\approx 6.7 \cdot 10^6$. (patch size $8 \times 8$, layers 20, dim. of input embedding $n_{\mathrm{patches}} \times 256$, internal MLP dim. 128).

Furthermore, ✗$g$ is highlighting architectures where the decoder network is deactivated for the experiments. Therefore, the loss is only based on the triplet part. The triplet loss is not used for the architectures marked with ✗$\mathcal{L}_{\mathrm{tri}}$. Hence, only the reconstruction loss is used.

### C. Novelty Detection

The various architectures are evaluated with respect to their novelty detection performance. For this, the outlier detection methods LOF, ULEF, OCSVM and ABOD are applied in the latent space of the resulting network. The outlier detection methods are also applied in the input space to highlight the performance gain when using the latent representation (last row in Tb. I). The networks are trained using $\mathcal{D}$, while the outlier detection is performed using $\mathcal{D}_{\mathrm{base}}$ as known and considering the remaining $\mathcal{D}_{\mathrm{ano}}$ as unknown. Therefore, the novelty detection is applied to the example where only highway images $\mathcal{D}_{\mathrm{base}}$ are known. Then unknown data points, such as inner-city images, $\mathcal{D}_{\mathrm{ano}}$ are investigated with respect to their outliership. If the outlier detection is able to identify that for example inner-city images are outlying with respect to highway images, the novelty detection task is fulfilled successfully. This evaluation shows the outlier detection capabilities with respect to connectivity classes.

Additionally, the reconstruction-based outlier detection methods f-AnoGAN [22] and RaPP [23] are evaluated. For the RaPP, a convolutional autoencoder is trained using $\mathcal{D}_{\mathrm{base}}$ then the normal reconstruction error is used for outlier detection (CAE-ORD) and the simple aggregation along pathway (CAE-SAP) as introduced in [23]. For the f-AnoGAN the network is also trained only on $\mathcal{D}_{\mathrm{base}}$.

In Tb. I, the resulting *Area Under Curve* (AUC) values are shown for the various combinations. The AUC will be 1 if all outliers are detected correctly. As a result of the experiments, it is shown that when the latent representations are used for the outlier detection, the performance of each method is improving (LOF, ULEF, IF, OCSVM, ABOD). The networks are providing a powerful latent representation, where simple outlier detection methods can perform well. In fact, most combinations (network with a basic outlier

detection method) are outperforming other baseline methods such as CAE-SAP and f-AnoGAN. Using one of the architectures in combination with ABOD is the preferable solution, since it yields the highest performance for all triplet autoencoding-based schemes (✓$g$ and ✓$\mathcal{L}_{\mathrm{tri}}$) and provides the highest performance overall. The results show that the proposed approach to include domain-knowledge in shaping the latent space improves the outlier detection performance significantly in this application.

The relevance of the triplet loss becomes clear when comparing the simple autoencoder architecture (Res-18 ✗$\mathcal{L}_{\mathrm{tri}}$) against the one including the triplet loss (Res-18 ✓$\mathcal{L}_{\mathrm{tri}}$). The use of the triplet loss increases the performance for all architectures remarkably. The influence of the decoder can be identified from Res-S ✗$g$ versus ✓$g$. Its contribution to the outlier detection is not as clear as for the triplet loss. Indeed, for some methods the outlier detection is getting slightly worse, however, for some it is getting better. The reason for introducing the decoder is the local visual similarity, which is not represented by the outlier detection analysis, since this is only covering the class oriented scale. For this purpose, another analysis will be carried out in the following section. Further comparison can be drawn from the different encoder types when using the triplet loss and a decoder. Because of its stable and high performance only ABOD is considered for further discussions. The Res-18 is only slightly better than the smaller Res-S. The performance difference for the two ViT version is more significant. Here, the ViT-L would be the architecture of choice. In conclusion, using the triplet loss as well as the decoder is clearly beneficial, while either the Res-S, Res-18 or ViT-L can be used from the outlier detection performance point of view.

### D. Latent Space Visualization

Another approach to access the quality of the latent space is provided through visual investigation. For this purpose, the latent representations $\mathcal{Z} = \{z_1, \ldots, z_M\}$ are projected into a two-dimensional space using UMAP [33]. In Fig. 3, the projections for the various architectures are shown.

On this scale, the difference of Res-S ✓$g$ versus ✗$g$ is hard to figure out. That difference will be further investigated via the local similarity analysis. For the difference produced by the triplet loss (middle column in Fig 3), this scale is sufficient. It is visible, that training the autoencoder without the triplet loss leads to a less separable latent representation. The difference between the Res-S and Res-18 (upper) is mainly in the more diffuse distribution of the roundabouts
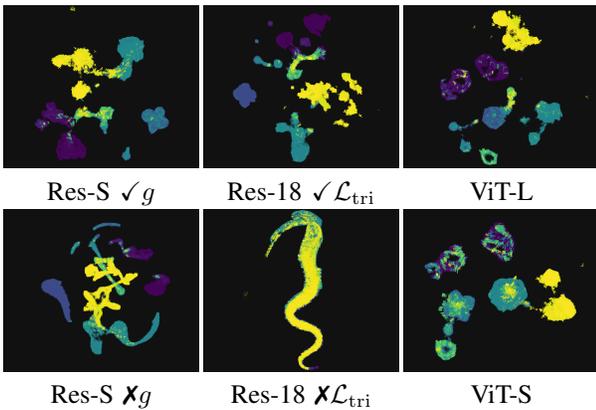
Fig. 3. UMAP projections of $\mathcal{Z}$. ■ highway, ■ crossing, ■ single lane, ■ multiple lane, ■ roundabout.
Interactive visualization tool SCENATLAS: `https://jwthi.github.io/SCENATLAS/`



Fig. 5. Average distance of the data points to their nearest neighbors corresponding to the respective latent representations. Smaller distance values are better. ─·─·: Res-18 ✗$\mathcal{L}_{\mathrm{tri}}$, ────: Res-18, ────: Res-S, ─ ─ ─: Res-S ✗$g$, ────: ViT-S and ────: ViT-L

when using Res-18. When comparing the ViT-L against ViT-S, the larger one is showing clear advantage, since in contrary to the smaller version, it is able to distinguish between highway and multiple lane. Comparing the architectures Res-S (up) Res-18 (up) and ViT-L, all of them show clear grouping of the analysis classes, and hence, from this perspective are equally well suited. Interested readers may refer to `https://jwthi.github.io/SCENATLAS/`, where the interactive visualization tool SCENATLAS is provided, which allows to discover the latent spaces more intuitively.
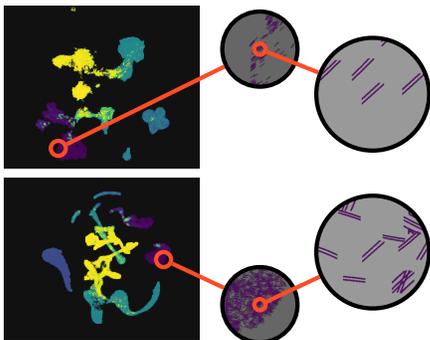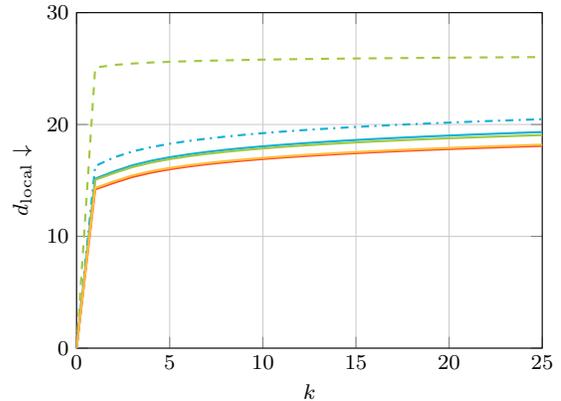
*E. Local Similarity Analysis*



Fig. 4. Local visual similarity motivation. Res-S Up: ✓$g$, Down: ✗$g$

As stated above, the main motivation of introducing the decoder is the local visual similarity of the latent space. This property has not yet been analyzed. So far, only the more class oriented outlier performance was evaluated, i.e. differentiate highway versus non-highway. Therefore, another analysis is performed for the local scale. In Fig. 4 the problem is visualized. It is showing two zoom levels of the embeddings with (up) and without (down) the decoder. Therefore, the embedding with the decoder appears to be more visually similar. In the following this quantity is assessed through numerical evaluation.

Given the latent representations $\mathcal{Z}$, the $k$ nearest neighbors for the $i$-th sample in the latent space are determined, leading

to the nearest neighbor set $\mathcal{K}_i$. Then, the average distance in the input space between a point and its neighbors is determined. By using

$$d_{\mathrm{local}}(k) = \frac{1}{M} \sum_{m=1}^{M} \frac{1}{k} \sum_{j \in \mathcal{K}_i} \|\boldsymbol{X}_i - \boldsymbol{X}_j\|_2 \qquad (5)$$

the average distance between the latent space neighbors is evaluated. Since a low value states a high average visual similarity between the $i$-th point and its neighbors, it is considered that lower values are better. In Fig. 5 the resulting values are shown for all the architecture variants. Therefore, the reasoning for the usage of the decoder is clearly provided. Moreover, the usage of the triplet loss is increasing the performance as well. The Res-S and Res-18 are again on a comparable level, but the ViT based architectures outperform all others. This indicates, that the ViT based method is preferable to identify local, shape based outliers.

*F. Summary*

TABLE II
ARCHITECTURE OVERVIEW

| Architecture | | | AUC ↑ | $d_{\mathrm{local}}(5)$ ↓ | Nparams |
|---|---|---|---|---|---|
| $f$ | $g$ | $\mathcal{L}_{\mathrm{tri}}$ | | | |
| Res-S | ✗ | ✓ | 0.892 | 25.61 | $0.3 \cdot 10^6$ |
| Res-S | ✓ | ✓ | 0.954 | 16.90 | $0.3 \cdot 10^6$ |
| Res-18 | ✓ | ✗ | 0.784 | 18.28 | $11.0 \cdot 10^6$ |
| Res-18 | ✓ | ✓ | **0.956** | 17.07 | $11.0 \cdot 10^6$ |
| ViT-S | ✓ | ✓ | 0.920 | **16.00** | $0.2 \cdot 10^6$ |
| ViT-L | ✓ | ✓ | **0.956** | 16.12 | $6.7 \cdot 10^6$ |

To sum up this analysis, the important values are gathered in Tb. II. Here, the results using ABOD are used. The overall best performance is provided by the ViT-L. However, also the small networks Res-S and ViT-S perform remarkably well. The final suggestion is to use ViT-L with ABOD to achieve high outlier detection performance as well as high visual similarity in the latent space.

## V. CONCLUSION

A method to identify novel traffic scenarios based on their static component is presented in this work. The introduced pipeline is outperforming existing outlier detection methods. It has the additional advantage that it can be trained on a huge data set (e. g. OSM), such that no retraining is necessary. In contrary, methods relying on the reconstruction paradigm would require retraining when detecting unknown scenarios. This work presents a possibility to incorporate expert-knowledge of a scenario's static environment for shaping the latent space of a triplet autoencoder. The presented results show that methods like outlier detection can significantly benefit from a latent space shaped in this way.

Another approach to detect unknown infrastructure based on their topology, could be realized through a simple categorization logic using the graphs, leading to a mixture of experts. However, the presented method provides an insight to the relationship between infrastructure types and additionally the local shape similarity.

Future work might focus on the inclusion of further parts of a scenario. Given this latent space, a possible improvement in clustering performance can also be investigated. Another possible extension is to include more infrastructural information into the graphs.

In conclusion, the suggested pipeline consists of a connectivity graph-based similarity definition, an autoencoder triplet network with a ViT as encoder and the ABOD method performing the novelty detection in the latent space. It shows superior performance with respect to its novelty detection capabilities and with respect to the neighborhood similarity evaluation. The interactive visualization SCENATLAS of the latent spaces is provided (`https://jwthi.github.io/SCENATLAS/`) and the code implementing the method is made publicly available.

## VI. ACKNOWLEDGEMENT

## REFERENCES

[1] P. Junietz *et al.*, "Evaluation of different approaches to address safety validation of automated driving," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, Nov 2018.

[2] "Research project pegasus," https://www.pegasusprojekt.de/en/, accessed: 2020-03-31.

[3] W. Wang *et al.*, "Clustering driving encounter scenarios using connected vehicle trajectories," *IEEE Transactions on Intelligent Vehicles*, 2020.

[4] N. Harmening, M. Biloš, and S. Günnemann, "Deep representation learning and clustering of traffic scenarios," 2020.

[5] A. Demetriou *et al.*, "A deep learning framework for generation and analysis of driving scenario trajectories," *cs.CV*, 2020.

[6] F. Hauer *et al.*, "Clustering traffic scenarios using mental models as little as possible," in *IEEE Intelligent Vehicles Symposium (IV)*, 2020.

[7] A. Dosovitskiy *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021.

[8] K. He *et al.*, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016.

[9] S. Riedmaier *et al.*, "Survey on scenario-based safety assessment of automated vehicles," *IEEE Access*, vol. 8, pp. 87 456–87 477, 2020.

[10] F. Kruber, J. Wurst, and M. Botsch, "An unsupervised random forest clustering technique for automatic traffic scenario categorization," in *21st International Conference on Intelligent Transportation Systems (ITSC)*, 2018.

[11] F. Kruber *et al.*, "Unsupervised and supervised learning with the random forest algorithm for traffic scenario clustering and classification," in *IEEE Intelligent Vehicles Symposium (IV)*, 2019.

[12] J. Langner *et al.*, "Logical scenario derivation by clustering dynamic-length-segments extracted from real-world-driving-data," in *Proceedings of the 5th International Conference on Vehicle Technology and Intelligent Transport Systems*. SCITEPRESS - Science and Technology Publications, 2019.

[13] J. Kerber *et al.*, "Clustering of the scenario space for the assessment of automated driving," in *IEEE Intelligent Vehicles Symposium (IV)*, 2020.

[14] M. Hasan *et al.*, "Learning temporal regularity in video sequences," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016.

[15] J.-A. Bolte *et al.*, "Towards corner case detection for autonomous driving," in *IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019.

[16] J. Langner *et al.*, "Estimating the uniqueness of test scenarios derived from recorded real-world-driving-data using autoencoders," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018.

[17] J. Wurst *et al.*, "An entropy based outlier score and its application to novelty detection for road infrastructure images," in *IEEE Intelligent Vehicles Symposium (IV)*, 2020.

[18] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015.

[19] N. Jean *et al.*, "Tile2vec: Unsupervised representation learning for spatially distributed data," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 3967–3974, jul 2019.

[20] Y. Yang, H. Chen, and J. Shao, "Triplet enhanced AutoEncoder: Model-free discriminative network embedding," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence Organization, aug 2019.

[21] M. Masana *et al.*, "Metric learning for novelty and anomaly detection," in *British Machine Vision Conference (BMVC)*, 2018.

[22] T. Schlegl *et al.*, "f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks," *Medical Image Analysis*, vol. 54, pp. 30–44, may 2019.

[23] K. H. Kim *et al.*, "Rapp: Novelty detection with reconstruction along projection pathway," in *International Conference on Learning Representations*, 2020.

[24] M. M. Breunig *et al.*, "LOF: identifying density-based local outliers," in *Proceedings of the 2000 ACM SIGMOD international conference on Management of data - SIGMOD 00*. ACM Press, 2000.

[25] H.-P. Kriegel, M. S. hubert, and A. Zimek, "Angle-based outlier detection in high-dimensional data," in *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD 08*. ACM Press, 2008.

[26] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in *8th IEEE International Conference on Data Mining*. IEEE, dec 2008.

[27] B. Schölkopf *et al.*, "Support vector method for novelty detection," in *Advances in neural information processing systems*, 2000, pp. 582–588.

[28] OpenStreetMap contributors, "Data from october 2020 via Overpass API," https://www.openstreetmap.org, 2020.

[29] M. Dupuis *et al.*, *OpenDRIVE – Format Specification Rev1.4*, 2015.

[30] P. A. Lopez *et al.*, "Microscopic traffic simulation using sumo," in *The 21st IEEE International Conference on Intelligent Transportation Systems*. IEEE, 2018.

[31] F. Poggenhans *et al.*, "Lanelet2: A high-definition map framework for the future of automated driving," in *Proc. IEEE Intell. Trans. Syst. Conf.*, Hawaii, USA, November 2018.

[32] A. Vaswani *et al.*, "Attention is all you need," in *Advances in Neural Information Processing Systems 30*, I. Guyon *et al.*, Eds. Curran Associates, Inc., 2017, pp. 5998–6008.

[33] L. McInnes, J. Healy, and J. Melville, "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction," *arXiv:1802.03426 [cs, stat]*, Feb. 2018, arXiv: 1802.03426.