# Expert-LaSTS:
# Expert-Knowledge Guided Latent Space for Traffic Scenarios

Jonas Wurst*, Lakshman Balasubramanian*, Michael Botsch* and Wolfgang Utschick+

*Abstract*— **Clustering traffic scenarios and detecting novel scenario types are required for scenario-based testing of autonomous vehicles. These tasks benefit from either good similarity measures or good representations for the traffic scenarios. In this work, an expert-knowledge aided representation learning for traffic scenarios is presented. The latent space so formed is used for successful clustering and novel scenario type detection. Expert-knowledge is used to define objectives that the latent representations of traffic scenarios shall fulfill. It is presented, how the network architecture and loss is designed from these objectives, thereby incorporating expert-knowledge. An automatic mining strategy for traffic scenarios is presented, such that no manual labeling is required. Results show the performance advantage compared to baseline methods. Additionally, extensive analysis of the latent space is performed.**

*Index Terms*— **Clustering, Novelty Detection, Scenario-Based Testing, Deep Learning**

## I. INTRODUCTION

Scenario-based testing is considered as one possible approach for the validation of *Autonomous Vehicles* (AVs) [1]. Two important tasks to enable the scenario-based approach are the definition of representative scenarios and the identification of potentially unknown and therefore untested scenarios. Representative scenarios can either be defined manually or automatically from collected data. For the latter one, usually clustering is used to define groups and representatives per group. The task of identifying untested scenarios can be realized through novelty detection (e. g. [2], [3]) or by checking if the scenario fits into a group (e. g. [4], [5]). For both tasks, clustering and novelty detection, a good representation or similarity measure is required. When applied directly on the plain data, the result is unsatisfactory (Sec. IV-B). This work proposes a method to design a representation space for traffic scenarios using expert-knowledge constraints. This representation space can be utilized for the tasks of scenario clustering and of detecting novel scenario types.

The method introduced in this work extends the findings of [2], where only the static part of a scenario is considered, i. e. the road infrastructure. The scenery is extended to include the ego dynamics as well. The ego dynamics are considered to be sufficient to represent a scenario in this work.

The methodology of this work can be summarized as follows. First, expert-knowledge based objectives are formulated. Then, it is shown how to design a loss function and network architecture such that those objectives are fulfilled.

*CARISSMA, Technische Hochschule Ingolstadt, 85049 Ingolstadt, Germany {firstname.lastname}@thi.de
+Technical University of Munich, 80333 Munich, Germany utschick@tum.de

To realize the designed loss function, an automatic sample mining process is introduced, that is based on a similarity measure for scenarios. The paper proposes a similarity measure based on the topology graph of the road network and the routes defined by the trajectories. This way, no manual labeling is required. Overall, the objective is to form a latent space, which hierarchically divides samples into groups, based on infrastructure, route, and trajectory.[1]

The resulting latent space shows superior performance with respect to novel type detection, clustering and feature stability compared to alternative approaches. Possible applications of this work are in the field of validation for AVs. It can aid the analysis of existing scenario databases or the detection of novel scenario types.

The contributions of this work can be summarized as:

1) Definition of an expert-knowledge aided loss and architecture to design the latent space as required.
2) Definition of an automatic mining strategy for traffic scenarios.
3) Comprehensive analysis and comparison of various representation spaces.

The remaining is structured as follows. In Sec. II, related work is discussed. Then the method itself is presented in Sec. III. In Sec. IV, various representations are analyzed with respect to novel type detection, clustering and feature stability. The work is concluded, and an outlook is provided in Sec. V.

## II. RELATED WORK

Analyzing traffic scenarios for clustering and novel type detection has been the focus of many works recently. Either the analysis is performed based on an appropriate similarity measure or by generating representations of the scenarios.

*1) Similarity-Based:* In [6] and the successor work [4] a data-adaptive similarity measure is introduced based on the paths through an unsupervised random forest. This similarity measure is used for clustering.

Finding scenario clusters is the objective of [7]. For this, abstract features are defined (e. g. Environment type, street type, curvature, average velocity, etc.) and collected in feature vectors, which are then clustered.

Another approach of using a specific similarity measure for clustering is presented in [8], where the average sum of the differences in the eight-car neighborhood is evaluated.

---

[1] An implementation of the presented method can be found in https://github.com/JWTHI/Expert-LaSTS.

Clustering pairs of trajectories is the objective of [9]. Various approaches were examined, where the result showed that *Dynamic Time Warping* (DTW) [10] in combination with $k$-means performed best in this study.

In [11] trajectories are clustered by hierarchical clustering. Each point of a trajectory is encoded as an area group label. An area group is defined by a Gaussian Mixture Model. The histograms of two encoded trajectories are compared via the Chi-squared distance.

Analyzing traffic scenarios from the ego information and other objects is realized in [12]. A procedure based on DTW and manual thresholds determines if a scenario is known.

Detecting unknown scenarios through novelty detection is the aim of [3]. There, an autoencoder is trained on known data. In test phase, if the reconstruction error is higher than a specific threshold, the scenario is assumed to be unknown.

In [13], novel traffic scenes based on infrastructure images are detected using a novel outlier detection method. The method utilizes local neighborhood similarities.

In contrast to this work the above stated works focus either on a specific similarity measure (e. g. DTW) or features, that are specifically selected to suit the used distance measure. Whereas in this work, the traffic scenarios are projected into a novel representation space.

*2) Representation-Based:* The *Principal Component Analysis* (PCA) is used for dimensionality reduction and hence to generate new representations for clustering in [14]. The PCA is applied to a column-wise normalized feature matrix, which is constructed using DTW on time series data.

In [15] a deep learning network is used to reconstruct the input trajectories through a latent space. In the latent space, cluster analysis is performed.

Another approach utilizing deep learning is presented in [16], where LSTMs [17] are used to encode a trajectory and reconstruct it through a latent representation. The network is trained adversarial using a discriminator network as well. The latent representations are further projected with PCA and t-SNE [18] before clustering. Moreover, the reconstruction error is used to estimate the novelty of a trajectory.

A different setting is presented in [19], where some known classes and also unknown classes are assumed. In a multistep training process, a deep learning network is trained, such that it finds representations suited for classification and clustering of the unknown classes. The steps include self-supervised pre-training, classification and mixed training. For this, a novel representation based on the random forest is introduced. The results in the latent space are clustered. The training input consists of a sequence of images, representing the infrastructure and the objects as well.

Also in [20], an image sequence showing the infrastructure and the objects at each timestamp is used as input. Each frame is fed through an autoencoder which is trained using the reconstruction loss and a triplet loss. For the triplet loss, closer frames (time) shall be closer in the latent space. The latent representations of the image sequence are transformed into a sequence latent representation, which is then used for clustering. This transformation is realized by a recurrent neu-ral network architecture, which aims to reconstruct frames and predict future frames. The sequence representations of the scenarios are clustered.

Instead of using deep learning, in [21] a tool chain of dimensionality reduction techniques and similarity measures is used for clustering. As input, a single trajectory is used.

This work builds on [2], where a method to project infrastructure images for novelty detection by utilizing expert-knowledge about the underlying topologies (c f. Sec. III).

The works summarized in this chapter generate an appropriate representation for traffic scenarios. Especially, [15], [16], [19] and [20] are comparable to this work since all of them use deep learning. However, the only two of them using a comparable input (infrastructure and dynamics) are [19] and [20]. Since [19], assumes known classes, it differs from this work. In contrast to [20], this work focuses on the ego dynamics. Moreover, expert-knowledge about the infrastructure and the trajectory is utilized.

## III. METHOD

The aim of this work is to design a latent space by means of expert-knowledge to represent traffic scenarios. It is shown how this latent space can be utilized to detect unknown traffic scenario types and to cluster traffic scenarios. Here, a traffic scenario is described by the road infrastructure and the dynamic information of the ego. This work extends [2], which is limited to the infrastructure of a traffic scenario.

### A. Preliminaries

A traffic scenario $\mathcal{X} = \{I, \mathcal{T}\}$ consists of the road infrastructure image $I$ and the ego trajectory $\mathcal{T}$. The dataset $\mathcal{D} = \{(\mathcal{X}_0, G_0, R_0), \ldots, (\mathcal{X}_M, G_M, R_M)\}$ consists of $M$ scenarios, its elements are defined in the following.

*1) Infrastructure:* The road infrastructure is represented as a grayscale birds-eye view image $I \in \mathbb{R}^{S \times S}$ and a graph representation $G$. The graph contains $N_G$ lane pieces as vertices $V = \{v_1, \ldots, v_{N_G}\}$ and $N_E$ edges $E = \{e_1, \ldots, e_{N_E}\}$ connecting them. The graph for a scene includes
1) all lanes, which are part of the ego route,
2) all lanes up to and including the next intersection,
3) all lanes of possible intersections in 1) and 2), and
4) all lanes neighboring any lanes in 1) - 3).

This way, the graph contains mainly the relevant lanes, whereas the image contains all lanes within the defined area.

*2) Trajectory:* The trajectory information is represented in two ways. The sequence based representation $\mathcal{T} = \{[x_1, y_1, t_1], \ldots, [x_N, y_N, t_N]\}$ and the route representation $R = \{r_1, \ldots, r_{N_G}\}$. The construction of $R$ is realized as follows. For a trajectory point the corresponding vertices are $v(x_1, y_1)$, which leads to the vertices sequence for the trajectory as $\mathcal{T}_R = \{v(x_1, y_1), \ldots, v(x_N, y_N)\}$. The route representation $R$ is directly linked to $G$ as,

$$r_n = \begin{cases} 2 & \text{if } v_n \in v(x_1, y_1) \\ 1 & \text{if } v_n \in \mathcal{T}_R \setminus v(x_1, y_1) \\ 0 & \text{else} \end{cases} \quad . \quad (1)$$

Hence, the vertex on which the trajectory starts is linked to $r_n = 2$, all other vertices the trajectory passes lead to $r_n = 1$.
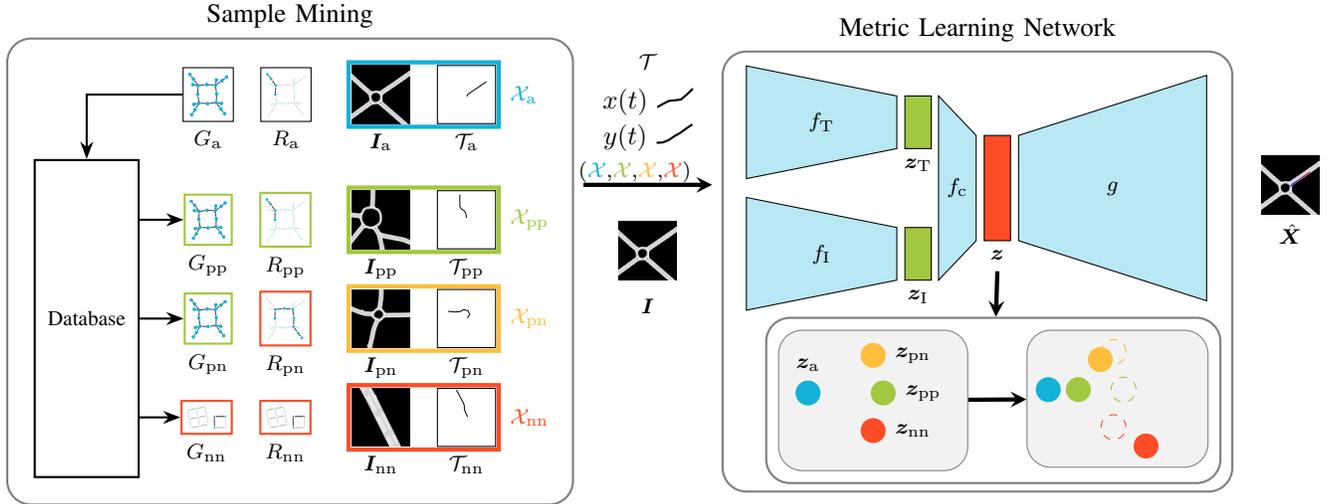
Fig. 1. Metric Learning Network for Traffic Scenarios: The sample mining depicts how the required scenario quadruplet is selected based on the graphs $G$ and routes $R$. Each scenario $\mathcal{X}$ consists of an image $\boldsymbol{I}$ and a trajectory $\mathcal{T}$. Each scenario is processed by the network, leading to the latent representations $\boldsymbol{z}$. Below the network, the quadruplet learning objective is illustrated. The scenario is reconstructed through the decoder $g$ into a merged representation $\hat{\boldsymbol{X}}$

### B. Base Method

In [2], a triplet autoencoder is used to project the infrastructure image $\boldsymbol{I}$ to a latent representation. Using the triplet learning, the latent representations are optimized to represent the topology of the infrastructure as well as their shape. The main components of the base method can be split into the mining of the samples and the triplet learning. This subsection summarizes the base method with only infrastructures [2].

*1) Mining:* For triplet learning it is required to define a data triplet, what is realized in the so-called mining. In this application, a triplet consists of three infrastructure images $(\boldsymbol{I}_\mathrm{a}, \boldsymbol{I}_\mathrm{p}, \boldsymbol{I}_\mathrm{n})$. Given a random anchor image $\boldsymbol{I}_\mathrm{a}$, a positive example $\boldsymbol{I}_\mathrm{p}$ and negative example $\boldsymbol{I}_\mathrm{n}$ must be found. The positive example should be similar to the anchor, while the negative example should be dissimilar.

It is proposed to use the road topologies underlying each image to perform the triplet mining. The possible positive examples are all the images, which have the same graph as the anchor $G_\mathrm{a}$, while all possible negative examples are all other images. An example: if the anchor image shows a four-way roundabout, the positive example would also contain a four-way roundabout. However, the shape of the roundabouts are not necessarily the same. The negative example would show some other topology (e. g. intersection).

The similarity of two graphs is realized through isomorphism. $G_i$ and $G_j$ are isomorphic $G_i \cong G_j$, if there exists a bijection $p : V_i \rightarrow V_j$ such that $(u,v) \in E_i \iff (p(u), p(v)) \in E_j$. This leads to the infrastructure similarity

$$s_\mathrm{i}(G_i, G_j) = \begin{cases} 1 & \text{if } G_i \cong G_j \\ 0 & \text{else} \end{cases} . \qquad (2)$$

The positive example is drawn from the subset $\mathcal{G}_\mathrm{p} = \{G \mid s_\mathrm{i}(G, G_\mathrm{a}) = 1\}$, representing all samples which have the same topology as the anchor sample. Contrary, the negative example is drawn from the remaining samples $\mathcal{G}_\mathrm{n} = \{G \mid s_\mathrm{i}(G, G_\mathrm{a}) = 0\}$.

*2) Triplet Learning:* The network is a triplet autoencoder consisting of an encoder $f : \boldsymbol{I} \mapsto \boldsymbol{z}$ and a decoder $g : \boldsymbol{z} \mapsto \hat{\boldsymbol{I}}$, where $\boldsymbol{z} \in \mathbb{R}^L$ is the latent representation. Training the network is performed simultaneously by two approaches: the autoencoding regime and the triplet learning. This combination enables both, the topology based learning through the triplet strategy, and a low-level image similarity caused by the autoencoding objective.

While training, the samples of the triplet are passed through the encoder $(\boldsymbol{z}_\mathrm{a}, \boldsymbol{z}_\mathrm{p}, \boldsymbol{z}_\mathrm{n})$. These latent representations are used to determine the distances to the anchor $d_\mathrm{ap} = ||f(\boldsymbol{I}_\mathrm{a}), f(\boldsymbol{I}_\mathrm{p})||_2^2$ and $d_\mathrm{an} = ||f(\boldsymbol{I}_\mathrm{a}), f(\boldsymbol{I}_\mathrm{n})||_2^2$. The distances are used in the triplet loss [22], as

$$\mathcal{L}_\mathrm{tri}(\boldsymbol{I}_\mathrm{a}, \boldsymbol{I}_\mathrm{p}, \boldsymbol{I}_\mathrm{n}) = \max(\alpha + d_\mathrm{ap} - d_\mathrm{an}, 0). \qquad (3)$$

Minimizing the triplet loss is achieved by pushing the latent representation of the negative $\boldsymbol{z}_\mathrm{n}$ away and pulling the positive $\boldsymbol{z}_\mathrm{p}$ representation close to $\boldsymbol{z}_a$. This way, the triplet loss realizes the similarity as defined for the mining.

In order to ensure a high visual similarity between neighbors in the latent space, the autoencoder regime is adopted for the anchor sample. Therefore, the reconstruction loss

$$\mathcal{L}_\mathrm{rec}(\boldsymbol{I}_\mathrm{a}) = ||\boldsymbol{I}_\mathrm{a} - g(f(\boldsymbol{I}_\mathrm{a}))||_2^2 \qquad (4)$$

is used for the training as well. The loss to train the network is given as

$$\mathcal{L} = \mathcal{L}_\mathrm{tri} + \mathcal{L}_\mathrm{rec}. \qquad (5)$$

### C. Proposed Method

This work extends the static description from [2] to a scenario by considering the dynamics of the ego vehicle. The network architecture is adjusted and the triplet learning is extended to quadruplet learning, called metric learning in the following. The overall concept is depicted in Fig. 1, which

is divided into the mining process and the metric learning network.

The aim of this work is to design a latent space, which enables the clustering of scenarios and the detection of novel scenario types. Expert-knowledge is used to aid and constrain the training process. For this, the following objectives are formulated:

A) Scenarios with the same infrastructure and similar trajectories shall be close together in the latent space.
B) Scenarios with the same infrastructure but different trajectories shall be close but not as close as A).
C) Scenarios without the same infrastructure shall be farther away than B).
D) The distance of scenarios according to A) shall be adjusted based on the similarity of the underlying actions.
E) Neighbors should have high similarities with respect to trajectory and infrastructure features.

The objectives reflect expert-knowledge based assumptions about the similarity of scenarios in a hierarchical way. Hence, the latent space shall realize these expert-knowledge based hierarchical similarity objectives. In order to achieve the objectives, an automatic mining process, the metric learning as well as the network architecture are presented.

*1) Mining:* In the base method it is shown how the identification of similar infrastructures can be realized through their topology. Hence, distinguishing the cases C) from A) or B) is possible. To realize the required separation between A) and B), and therefore to include the trajectory information to the mining process, the mining definitions are extended.

Given the case that two graphs are isomorphic $G_i \cong G_j$, hence their infrastructure is the same, two trajectories are considered to be similar, if they share the same route within their graphs. The trajectories are transformed into the route representation $R$ which are directly linked to the respective graph (see Sec. III-A). Two scenarios share the same route if there exists a bijection $p$ on $G_i, G_j$ such that $R_i = p(R_j)$, which is formulated as $G_i \cong G_j \mid R_i = p(R_j)$. The route-based similarity measure is defined as

$$s_r(G_i, G_j) = \begin{cases} 1 & \text{if } G_i \cong G_j \mid R_i = p(R_j) \\ 0 & \text{else} \end{cases}. \quad (6)$$

According to D), just defining two trajectories to be similar is not sufficient, instead it shall be adjusted based on the actions of the trajectories. To allow further fine-tuning in the training, for trajectories sharing a similar route an additional similarity measure $s_t$ is defined. Let $\mathcal{A}_i = [\boldsymbol{a}_{\text{lat}}, \boldsymbol{a}_{\text{lon}}, |v|]$ be the accelerations and speed per timestamp for the $i$th trajectory. The dissimilarity between two trajectories is then calculated via $d = d_{\text{DTW}}(\mathcal{A}_i, \mathcal{A}_j) |\text{DTW}_{\text{seq}}|$, where $d_{\text{DTW}}$ denotes the DTW distance and $|\text{DTW}_{\text{seq}}|$ the warping path length divided by maximum sequence length. The intuition is to compare the trajectories which share the same route on an action level, therefore considering the accelerations and speed. The similarity is calculated with respect to maximum dissimilarity ($d_{\text{max}}$) within all trajectories with the same route, as

$$s_t = 1 - \frac{d}{d_{\text{max}}}. \quad (7)$$

The mining of a data quadruplet is realized by randomly sampling an anchor scenario $\mathcal{X}_a$. Based on its graph $G_a$ and route $R_a$, samples for the cases A) - C) are drawn. The example scenario having similar infrastructure and similar route $\mathcal{X}_{\text{pp}}$ is sampled based on $\mathcal{G}_{\text{pp}} = \{G \mid s_i(G, G_a) = 1, s_r(G, G_a) = 1\}$. For sampling the scenario with same infrastructure but different route $\mathcal{X}_{\text{pn}}$, $\mathcal{G}_{\text{pn}} = \{G \mid s_i(G, G_a) = 1, s_r(G, G_a) = 0\}$ is used. Finally, the scenario with a different infrastructure $\mathcal{X}_{\text{nn}}$ is sampled from $\mathcal{G}_{\text{nn}} = \{G \mid s_i(G, G_a) = 0\}$. This mining process is also visualized on the left side of Fig. 1.

*2) Metric Learning:* This section introduces the metric learning and the network, such that the objectives A) - E) are realized.

As shown in Fig. 1, the network consists of two encoders, $f_I : \boldsymbol{I} \mapsto \boldsymbol{z}_I \in \mathbb{R}^{L_I}$ for the image and $f_T : \mathcal{T} \mapsto \boldsymbol{z}_T \in \mathbb{R}^{L_T}$ for the trajectory. The two intermediate representations are concatenated and then passed through the network $f_c : [\boldsymbol{z}_T, \boldsymbol{z}_I] \mapsto \boldsymbol{z}$ to create the final latent representation $\boldsymbol{z} \in \mathbb{R}^L$. Finally, a decoder $g : \boldsymbol{z} \mapsto \hat{\boldsymbol{X}}$ is used to generate a merged representation $\hat{\boldsymbol{X}} \in \mathbb{R}^{S \times S \times 2}$ of infrastructure and trajectory.

The data quadruplet is passed through the encoder, such that the latent representations $\boldsymbol{z}_a, \boldsymbol{z}_{\text{pp}}, \boldsymbol{z}_{\text{pn}}, \boldsymbol{z}_{\text{nn}}$ are generated. The squared distances from the anchor to the examples are defined as $d_{\text{pp}} = \|\boldsymbol{z}_a - \boldsymbol{z}_{\text{pp}}\|_2^2$, $d_{\text{pn}} = \|\boldsymbol{z}_a - \boldsymbol{z}_{\text{pn}}\|_2^2$ and $d_{\text{nn}} = \|\boldsymbol{z}_a - \boldsymbol{z}_{\text{nn}}\|_2^2$. The objectives A) - D) can then be formulated as

$$d_{\text{nn}} \geq d_{\text{pn}} + \alpha_G, \quad (8)$$
$$d_{\text{pn}} \geq \max\{d_{\text{pp}}, \alpha_T\} + \alpha_R, \quad (9)$$
$$d_{\text{pp}} = (1 - s_t)\alpha_T, \quad (10)$$

where $\alpha_{...}$ are margin parameters. Those constraints lead to the following loss formulations

$$\mathcal{L}_G = \max\{\alpha_G + d_{\text{pn}} - d_{\text{nn}}, 0\}, \quad (11)$$
$$\mathcal{L}_R = \max\{\alpha_R + \max\{\alpha_T, d_{\text{pp}}\} - d_{\text{pn}}, 0\}, \quad (12)$$
$$\mathcal{L}_T = |(1 - s)\alpha_T - d_{\text{pp}}|. \quad (13)$$

The losses Eq. 11 and Eq. 12 are basic triplet losses [22], when optimizing both, it leads to the quadruplet loss as presented in [23].

The objective in E) is not directly addressed by the former loss definitions, hence another strategy needs to be adopted. For this purpose, the latent representation of a scenario is used to reconstruct the scenario. In this case, the reconstruction generates an image with two channels, one channel for the infrastructure as in $\boldsymbol{I}$ and one channel for the trajectory information as in $\mathcal{T}$. This way, the network has to connect the image information with the trajectory information. Furthermore, for the decoding to work properly, the neighborhood in the latent space has to share high similarities. To further aid the training, all infrastructure parts, that are not part of the graph are removed from

the target reconstruction image. Since simple reconstruction might fail due to the high sparsity of the generated output, the reconstruction loss is adopted piece-wise for trajectory, infrastructure, and respective background pixels. The weighted sparse reconstruction loss is defined as

$$\mathcal{L}_{\text{Rec}} = \gamma_{\text{I}}\mathcal{R}(\mathcal{I}_{\text{I}}) + \gamma_{\bar{\text{I}}}\mathcal{R}(\mathcal{I}_{\bar{\text{I}}}) + \gamma_{\text{T}}\mathcal{R}(\mathcal{I}_{\text{T}}) + \gamma_{\bar{\text{T}}}\mathcal{R}(\mathcal{I}_{\bar{\text{T}}}), \quad (14)$$

with $\mathcal{R}(\mathcal{I}) = \frac{1}{|\mathcal{I}|}\sum_{i \in \mathcal{I}}(\boldsymbol{X}_{\text{a}}(i) - \hat{\boldsymbol{X}}_{\text{a}}(i))^2$ the reconstruction error for the pixel subset $\mathcal{I}$. The subset $\mathcal{I}_{\text{I}}$ refers to all ground truth infrastructure pixels, $\mathcal{I}_{\bar{\text{I}}}$ all remaining pixels in the infrastructure channel. For the trajectory pixel sets $\mathcal{I}_{\text{T}}$ and $\mathcal{I}_{\bar{\text{T}}}$ the logic applies respectively.

Combining all the loss definitions leads the overall loss as

$$\mathcal{L} = \beta_{\text{M}}\left(\beta_{\text{G}}\mathcal{L}_{\text{G}} + \beta_{\text{R}}\mathcal{L}_{\text{R}} + \beta_{\text{T}}\mathcal{L}_{\text{T}}\right) + \beta_{\text{Rec}}\mathcal{L}_{\text{Rec}}. \quad (15)$$

Training the network with $\mathcal{L}$ aims to realize the expert objectives as defined in the beginning of this section.

## IV. EXPERIMENTS

The quality of the latent space constructed by the proposed method is analyzed through various experiments. The analysis is based on four perspectives: 1) Novel type detection: *Given a base set, are new scenario types detected as novel?* 2) Clustering: *Do the latent representations form meaningful clusters?* 3) Feature stability: *Are scenario features of neighboring latent representations similar?* 4) Visualization: *Can the latent space be analyzed through aided visualizations?* In order to evaluate the impact of the various loss terms and possible network variations, different settings are used for all the experiments.

The section is structured as follows. First, the dataset is explained. Second, the different network settings and training variants are summarized. The novel scenario type detection analysis is shown in the third part, followed by the clustering analysis in fourth and the feature stability in the fifth part. The visual assessment is shown in part six. The different results are summarized in the last part.

### A. Data

The data used to train and analyze the network is generated through simulation. The data generation process is divided in two parts, the infrastructure sampling and the simulation.

*1) Infrastructure Sampling:* Sampling the infrastructures is realized as in [2]. From OpenStreetMap (OSM) [24], random nodes are selected as center for the scenarios. The underlying road infrastructure within the area of $200\,\text{m} \times 200\,\text{m}$ per node is extracted as image ($100\,\text{px} \times 100\,\text{px}$), graph and as map for simulation. To generate the graphs and images, the tools from [2] are used. This way $\approx 70\,000$ infrastructures are extracted.

*2) Simulation:* For each of the infrastructures, simulations in SUMO [25] are performed. One vehicle (ego) is inserted at the center position of the scenario, while other vehicles are randomly spawned, such that the scenarios show a rather high traffic load. The scenario is simulated for a total span of 6 s. The ego information during this timespan is recorded.

The infrastructures as well as the routes are sampled such, that for each scenario, similar infrastructure and route examples are available.

*3) Groups:* For the analysis with respect to groups (clustering and novelty detection), three detail levels are examined. First, a rough level, consisting of the categories: 1) single-lane, 2) multi-lane, 3) intersection, 4) intersection entering, 5) roundabout, 6) roundabout entering, 7) highway and 8) highway entering. This leads to the subscript $\ldots_{\text{C}}$.

The second detail level considers all unique infrastructure graphs in the dataset, leading to 704 groups. Hence, those groups can be used to analyze the performance with respect to the infrastructure. The subscript $\ldots_{\text{G}}$ is used to indicate the usage of the second level.

The third and most detailed level is provided by the 2330 groups, formed by all unique routes combined with their graphs. It provides insight with respect to the complete scenario. Here, the subscript $\ldots_{\text{R}}$ is used.

### B. Comparison

*1) Alternative Approaches:* As mentioned in Sec. II, typical approaches either operate in the input space or perform dimensionality reduction (t.SNE etc.). The proposed method is compared to the following alternatives.

*a) Plain:* The input is used as representation directly. Hence, the image and the trajectory are vectorized and concatenated (10090 dimensions).

*b) UMAP:* The plain input is projected with UMAP [26] to 64-dimensional space.

*c) PCA:* The first 64 principal components (PCA) of the plain input are used.

*d) Classifier:* The latent representations when using the same architecture as the proposed method but replacing the decoder with two classification heads, one for the 704 unique graphs (64-704-704-704) and one for the 2330 unique routes (64-2330-2330-2330). The used representation is the intermediate 64-dimenional latent representation, not the classification output.

*e) Autoencoder:* The latent representations when using the same architecture as the proposed method but without metric learning just operating as autoencoder.

*2) Novel Scenario Type Detection:* Detecting novel scenario types is a crucial task in the validation process of autonomous driving. The latent space designed by the former method suits this need, as shown in this section. In this work, detecting novel scenario types is realized through outlier detection. Therefore, assuming a base data set (the already known scenarios), the task is to identify scenarios which do not fit in the base data set.

The novelty detection is performed as n-vs-1, where one group (Sec. IV-A.3) is excluded from the base dataset. It is tested, how well this left-out group is detected as novel. This procedure is repeated for all groups. The novelty detection performance is measured using the *Area Under Curve* (AUC). *Angle Based Outlier Detection* (ABOD) is used as the novelty detection method.

| Approach | Novelty Detection IV-B.2 | | | Clustering IV-B.3 | | | Feature Stability IV-B.4 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $AUC_C$ | $AUC_G$ | $AUC_R$ | $ACC_C$ | $ACC_G$ | $ACC_R$ | $\bar{d}_I$ | $\bar{d}_T$ | $\bar{d}_v$ | $\bar{d}_{a_{\mathrm{lon}}}$ | $\bar{d}_{a_{\mathrm{lat}}}$ | $\bar{d}_\psi$ |
| **Proposed** | 0.991 | 0.919 | 0.904 | 0.839 | 0.900 | 0.622 | 36.67 | 0.58 | 1.91 | 0.57 | 0.37 | 0.20 |
| `plain` | 0.500 | 0.485 | 0.487 | 0.265 | 0.255 | 0.271 | 28.80 | 1.47 | 5.96 | 0.74 | 0.43 | 0.30 |
| `UMAP` | 0.654 | 0.557 | 0.574 | 0.406 | 0.070 | 0.070 | 36.57 | 1.37 | 5.56 | 0.84 | 0.44 | 0.35 |
| `PCA` | 0.543 | 0.632 | 0.617 | 0.265 | 0.342 | 0.324 | 30.89 | 1.27 | 5.09 | 0.77 | 0.44 | 0.32 |
| `Classifier` | 0.975 | 0.885 | 0.860 | 0.319 | 0.302 | 0.348 | 40.61 | 1.01 | 2.89 | 0.63 | 0.43 | 0.34 |
| `Autoencoder` | 0.786 | 0.628 | 0.641 | 0.503 | 0.218 | 0.167 | 36.53 | 0.54 | 2.02 | 0.65 | 0.37 | 0.18 |

TABLE I

COMPARISON PERFORMANCE SUMMARY: RED INDICATES WORSE THAN THE PROPOSED METHOD, GREEN BETTER AND YELLOW COMPARABLE.

As the results show (Tb. I), detecting novel scenario types is best realized in the latent space formed by the proposed method. The only other method reaching considerable performance is the classifier based approach. The alternative approaches miss out noticeably when compared to the proposed method.

*3) Clustering:* Another task in the field of validating AVs is to cluster scenarios into groups. This way, possible representatives for testing per cluster can be defined. In this subsection, the clustering performance when using the designed latent space is demonstrated.

Because of the highly imbalanced number of samples per group (Sec. IV-A.3), agglomerative clustering suits this task well. As linkage function, average is used. The clustering performance is stated as accuracy $ACC_{...}$ [27]. For this, the best mapping between the ground truth labels and the predicted labels is determined. Given this mapping the accuracy can be determined.

For the clustering, none of the alternative approaches reached comparable performance to the proposed method (Tb. I). Even the classifier based model is not able to provide sufficient clustering results. This clearly shows the necessity for the expert-knowledge designed latent space for clustering traffic scenarios.

*4) Feature Stability:* As stated in the design requirements E), one of the objectives is that neighbors in the latent space share high similarities with respect to various features. An analysis accessing this is realized in this section.

To analyze the stability within a neighborhood, for each data point, the 15 nearest neighbors in the latent space are considered for the further analysis. The average differences from the data points in focus to their neighbors are determined. For calculating the differences $d_{...}$ various features are used: 1) $d_I$ image difference (like in [2]), 2) $d_T$ trajectory difference (average displacement), 3) $d_v$ average velocity difference, 4) $d_{a_{\mathrm{lon}}}$ average longitudinal acceleration difference, 5) $d_{a_{\mathrm{lat}}}$ lateral acceleration difference and 6) $d_\psi$ average orientation difference. Those values are averaged over the complete data set, leading to $\bar{d}_{...}$. The smaller those average values, the more similar the features within the neighborhood, hence the better the objective is fulfilled. It is important to note, that the features 3 - 6 are not part of the input.

The results are listed in Tb. I column feature stability. The performance for the difference of images is better for the most of the alternative approaches except for the classifier. When comparing the trajectory features, however, the proposed model outperforms all alternative approaches except the autoencoder. If neighboring data points shall share high similarities with respect to features from both domains (infrastructure and trajectory), the latent spaces provided by the proposed method or the autoencoder are the best choices.

### C. Ablation

*1) Model Variants:* To assess the impact of the various possible settings of the network and the learning, they are varied and compared.

**Proposed Setting:** The setting which shows overall good perfomance is as follows: $f_I$: ResNet-18 [28], $f_T$: Transformer-Encoder [29], $L_I = 64$, $L_T = 16$, $L = 64$, $\beta_M = \beta_G = \beta_R = \beta_T = 1$, $\beta_{\mathrm{Rec}} = 10$, $\gamma_I = \gamma_T = 5$, $\gamma_{\bar{I}} = 10$, $\gamma_{\bar{T}} = 20$, $\alpha_G = \alpha_R = \alpha_T = 1$ and `random` negative sampling.

In the Transformer-Encoder ($f_T$), an embedding token is used like in [2]. As alternative $f_T$, a LSTM [17] is used. And as alternative image encoder $f_I$, a ViT [30] with patch-size of 10, dimensionality of 256, MLP dimensionality of 128, 16 layers and 16 heads is used.

All the other variants used in the ablation, adjust few parameters from the proposed setting. For example, in the variant $\beta_M = 0$ just the according value is changed, all other values are as stated above.

For the negative sampling, the following strategies are examined. `random`: from all graphs that are different, `group`: from all graphs that are different but only inside the same category (highway, etc.) and `random-excl`: from all graphs that are different excluding the same category.

*2) Novel Scenario Type Detection:* In Tb. II, the results for various model variations are shown. The description in the left column, indicates what parameters are changed compared to the proposed setting.

As one can see, detecting novel scenario types is realized best either with the proposed setting, $\alpha_G = 10$ and $\alpha_R = 5$, $L_{...} = L_{...}*2$ or $f_T$ : LSTM settings. Hence, the choice of the margin parameters has neglectable effect on the performance. Also, the size of the latent space size is rather irrelevant in terms of novelty detection. The same holds for the selected trajectory encoder.

The importance of each loss terms can be seen from Tb. II (rows 2-5). The metric learning related losses ($\mathcal{L}_M$, $\mathcal{L}_G$, $\mathcal{L}_R$ and $\mathcal{L}_T$) are required for good performance.

| Setting | Novelty Detection IV-C.2 | | | Clustering IV-C.3 | | | Feature Stability IV-C.4 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $AUC_C$ | $AUC_G$ | $AUC_R$ | $ACC_C$ | $ACC_G$ | $ACC_R$ | $\bar{d}_I$ | $\bar{d}_T$ | $\bar{d}_v$ | $\bar{d}_{a_{lon}}$ | $\bar{d}_{a_{lat}}$ | $\bar{d}_\psi$ |
| **Proposed** | 0.991 | 0.919 | 0.904 | 0.839 | 0.900 | 0.622 | 36.67 | 0.58 | 1.91 | 0.57 | 0.37 | 0.20 |
| $\beta_T = 0$ | 0.986 | 0.904 | 0.879 | 0.948 | 0.880 | 0.507 | 36.59 | 0.63 | 2.30 | 0.61 | 0.37 | 0.20 |
| $\beta_T = 0, \beta_R = 0$ | 0.987 | 0.883 | 0.847 | 0.757 | 0.753 | 0.379 | 36.99 | 0.57 | 2.05 | 0.59 | 0.37 | 0.19 |
| $\beta_M = 0$ | 0.786 | 0.628 | 0.641 | 0.503 | 0.218 | 0.167 | 36.53 | 0.54 | 2.02 | 0.65 | 0.37 | 0.18 |
| $\beta_{Rec} = 0$ | 0.980 | 0.943 | 0.912 | 0.712 | 0.917 | 0.541 | 39.91 | 0.77 | 1.43 | 0.51 | 0.41 | 0.31 |
| $\alpha_G = 10, \alpha_R = 5$ | 0.990 | 0.919 | 0.898 | 0.705 | 0.857 | 0.383 | 37.66 | 0.64 | 2.19 | 0.57 | 0.38 | 0.21 |
| $L_{...} = L_{...} * 2$ | 0.991 | 0.915 | 0.904 | 0.611 | 0.910 | 0.595 | 36.21 | 0.55 | 1.80 | 0.55 | 0.36 | 0.19 |
| $f_I$ : ViT | 0.965 | 0.883 | 0.880 | 0.768 | 0.722 | 0.620 | 36.70 | 0.59 | 2.05 | 0.56 | 0.37 | 0.20 |
| $f_T$ : LSTM | 0.991 | 0.919 | 0.904 | 0.848 | 0.906 | 0.605 | 36.54 | 0.62 | 2.23 | 0.61 | 0.37 | 0.21 |
| random-excl | 0.999 | 0.866 | 0.855 | 0.996 | 0.611 | 0.460 | 36.12 | 0.56 | 1.94 | 0.59 | 0.36 | 0.19 |
| group | 0.821 | 0.873 | 0.871 | 0.344 | 0.542 | 0.520 | 36.34 | 0.58 | 1.94 | 0.61 | 0.37 | 0.20 |

TABLE II

ABLATION PERFORMANCE SUMMARY: RED INDICATES WORSE THAN THE PROPOSED METHOD, GREEN BETTER AND YELLOW COMPARABLE.

*3) Clustering:* The clustering accuracies for the model variants are shown in the corresponding columns in Tb. II. Only the model variant when using LSTM encoder is achieving comparable results to the proposed setting.

As for the novelty detection, the metric learning related losses ($\mathcal{L}_M$, $\mathcal{L}_G$, $\mathcal{L}_R$ and $\mathcal{L}_T$) are important.

*4) Feature Stability:* The double sized latent space setting $L_{...} = L_{...} * 2$ achieves better results than the proposed setting. Changing the negative sampling strategy to `random-excl` does only slightly affect the performance in terms of feature stability.

Using only the reconstruction loss ($\beta_M = 0$) does not outperform the proposed setting in terms of feature stability. Hence, the expert-knowledge aided losses ($\mathcal{L}_G$, $\mathcal{L}_R$, $\mathcal{L}_T$) help in structuring the latent space also for the stability criterion. Not using the reconstruction loss ($\beta_{Rec} = 0$) has a negative effect for the most features. This supports the designed intuition to use the autoencoder regime to achieve feature stability.
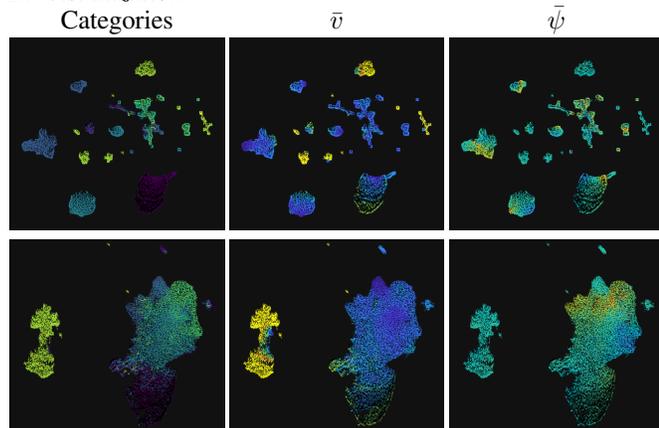
## D. Visualization



Fig. 2. UMAP visualizations of the latent spaces of the proposed setting (top) and with $\alpha_M = 0$ (bottom). *Categories*: single-lane , multi-lane , intersection , intersection-enter , roundabout , roundabout-enter , highway and highway-enter . $\bar{v}$: 0 m/s ... 42 m/s . $\psi$: $-\pi$ ... $+\pi$ .

The resulting latent representations can be assessed by visualizing them. Therefore, UMAP is used to project the representation into two-dimensional space. In Fig. 2, the projections of two latent representations are depicted. The upper row shows the projection for the proposed setting, and the lower row shows the projection for the setting $\alpha_M = 0$ (turning off the metric learning). In the columns different color codings are used. In the first, the categories as in Sec IV-A.3 are used. The second shows the average velocity of the trajectory and the last show values related to the orientation of the trajectory. The two latent representations were picked to demonstrate, how the latent representations differ, and how they can be analyzed using the visualization. It becomes clear, that the latent representation of the proposed setting provides well structure behavior in terms of categories, since the various infrastructure types are clearly separated. The model without the metric loss fails in separating the categories, instead two big clusters can be seen, one for the highway scenarios and one with all other scenarios. There is a strong relationship between the internal cluster structure and the shown features when using the proposed method. Therefore, the features ($\bar{v}$, $\bar{\psi}$) show smooth course within the clusters (e. g. in the lower right cluster, the average speed increases from top to bottom). This is also true for the model without metric loss, but here in a more global scale. The analysis can support in understanding and validating the latent representations. Readers interested in exploring the projections in more detail may refer to the website published alongside the paper `https://jwthi.github.io/Expert-LaSTS/`. There, also the projections for the other settings as well as some alternative approaches are shown.

## E. Summary

The different analysis perspectives highlight the performance with respect to specific tasks. Here, the overall performance is summarized and best model variants are discussed.

The only alternative approach able to perform considerably well in one of the perspectives is the classifier. However, when considering the other perspectives, the classifier does not seem to be a good choice either. All the other alternative approaches perform worse in all the perspectives, and hence are not appropriate for the presented problem setting.

Over all perspectives, three model variants should be highlighted. First, the proposed setting performs well throughout the perspectives. It seems to be the best selection when solving all tasks considerably well. The double size latent space setting is the second which performs well on different perspectives. With respect to feature stability it even outperforms the proposed setting. But, in terms of clustering it is worse. Therefore, if the focus is towards feature stability and less towards clustering this might be a good selection. The third and last model setting to be highlighted is using the LSTM encoder. With respect to detecting novel scenario types and clustering, it performs equally well as the proposed setting. However, it misses out slightly on the feature stability.

## V. CONCLUSION

In this work, a method to design a latent space for traffic scenarios by means of expert-knowledge is presented. An automated mining strategy for traffic scenarios is introduced and used to find similar infrastructures and routes. This way, relative similarities as defined by expert objectives can be realized. The resulting latent space outperforms alternative approaches on various analysis perspectives, namely detecting novel scenario types, clustering and feature stability. The ablation study provides deep insight to the impact of various model parameters on the performance.

The method presented in this work can be used in the validation process for AVs. More precisely, it can support the analysis of scenarios as well as the detection of representative and novel scenarios.

Including further objects can be one possible direction for further research on the proposed method. Also, the performance when using real-world data can be analyzed in a next step.

## 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, pp. 491–496.

[2] J. Wurst *et al.*, "Novelty detection and analysis of traffic scenario infrastructures in the latent space of a vision transformer-based triplet autoencoder," in *2022 IEEE Intelligent Vehicles Symposium (IV)*, 2021.

[3] 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, jun 2018.

[4] F. Kruber *et al.*, "Unsupervised and supervised learning with the random forest algorithm for traffic scenario clustering and classification," in *2019 IEEE Intelligent Vehicles Symposium (IV)*, 2019, pp. 2463–2470.

[5] L. Balasubramanian *et al.*, "Open-set recognition based on the combination of deep learning and ensemble method for detecting unknown traffic scenarios," in *2021 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, jul 2021.

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

[7] 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.

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

[9] W. Wang *et al.*, "Clustering driving encounter scenarios using connected vehicle trajectories," *IEEE Transactions on Intelligent Vehicles*, pp. 1–1, 2020.

[10] H. Sakoe and S. Chiba, "Dynamic programming algorithm optimization for spoken word recognition," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 26, pp. 43–49, 1978.

[11] J. Bernhard, M. Schutera, and E. Sax, "Optimizing test-set diversity: Trajectory clustering for scenario-based testing of automated driving systems." Indianapolis, IN, USA: IEEE, 2021, pp. 1371–1378.

[12] L. Ries *et al.*, "Trajectory-based clustering of real-world urban driving sequences with multiple traffic objects." Indianapolis, IN, USA: IEEE, 2021, pp. 1251–1258.

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

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

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

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

[17] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, p. 1735–1780, Nov. 1997.

[18] L. v. d. Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, pp. 2579–2605, 2008.

[19] L. Balasubramanian *et al.*, "Traffic scenario clustering by iterative optimisation of self-supervised networks using a random forest activation pattern similarity," in *2021 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, jul 2021.

[20] J. Zhao *et al.*, "Large scale autonomous driving scenarios clustering with self-supervised feature extraction," Mar. 2021.

[21] F. Hoseini, S. Rahrovani, and M. H. Chehreghani, "Vehicle motion trajectories clustering via embedding transitive relations." Indianapolis, IN, USA: IEEE, 2021, pp. 1314–1321.

[22] 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.

[23] X. Zhang *et al.*, "Embedding label structures for fine-grained feature representation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

[24] OpenStreetMap contributors, "Data from March 2021 via Geofabrik," https://www.openstreetmap.org, 2020.

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

[26] 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.

[27] Y. Yang *et al.*, "Image clustering using local discriminant models and global integration," *IEEE Transactions on Image Processing*, vol. 19, pp. 2761–2773, oct 2010.

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

[29] 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.

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