

Traffic Scenario Clustering by Iterative Optimisation of Self-Supervised Networks Using a Random Forest Activation Pattern Similarity

Lakshman Balasubramanian¹, Jonas Wurst¹, Michael Botsch¹ and Ke Deng²

Abstract—Traffic scenario categorisation is an essential component of automated driving, for e.g., in motion planning algorithms and their validation. Finding new relevant scenarios without handcrafted steps reduce the required resources for the development of autonomous driving dramatically. In this work, a method is proposed to address this challenge by introducing a clustering technique based on a novel data-adaptive similarity measure, called Random Forest Activation Pattern (RFAP) similarity. The RFAP similarity is generated using a tree encoding scheme in a Random Forest algorithm. The clustering method proposed in this work takes into account that there are labelled scenarios available and the information from the labelled scenarios can help to guide the clustering of unlabelled scenarios. It consists of three steps. First, a self-supervised Convolutional Neural Network (CNN) is trained on all available traffic scenarios using a defined self-supervised objective. Second, the CNN is fine-tuned for classification of the labelled scenarios. Third, using the labelled and unlabelled scenarios an iterative optimisation procedure is performed for clustering. In the third step at each epoch of the iterative optimisation, the CNN is used as a feature generator for an unsupervised Random Forest. The trained forest, in turn, provides the RFAP similarity to adapt iteratively the feature generation process implemented by the CNN. Extensive experiments and ablation studies have been done on the highD dataset. The proposed method shows superior performance compared to baseline clustering techniques.

I. INTRODUCTION

In recent years, there have been rapid developments in the field of autonomous driving and driver assistance systems [1]. As new and improved autonomous driving functions are introduced, the autonomous system must be capable of handling various driving scenarios. Traffic scenario categorisation is a key component for downstream tasks like path planning [2], behaviour planning and scenario-based validation methods [3], [4]. Hence, representative traffic scenarios (e.g. overtake, cut-in, etc.) are necessary for developing and testing the behaviour of autonomous vehicles [5]. The representative scenarios can be defined by means of expert knowledge, can be generated from simulations or can be identified from real-world driving data.

The approaches using expert knowledge and simulations have important constraints that limit the generation of a list of representative scenarios, such as the limited knowledge of the experts and the ability of the simulation environment to model the complex interactions between traffic participants.

¹Technische Hochschule Ingolstadt, Research Center CARISSMA, Esplanade 10, 85049 Ingolstadt, Germany, {firstname.lastname}@thi.de

²Royal Melbourne Institute of Technology, Melbourne, Australia, {firstname.lastname}@rmit.edu.au

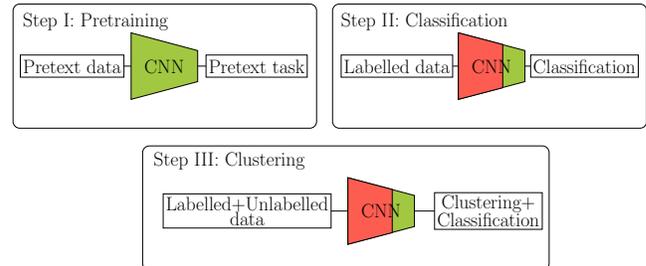


Fig. 1: Overview of the three step clustering process. The red portion in the figure shows frozen layers of the CNN and the green portion represents the layers that are trained.

Thus, a promising way to obtain scenario categories is to identify them automatically from real-driving data [3], [6], [7]. In this work, the problem of identifying new traffic scenario categories from real-driving data is studied. Unsupervised scenario clustering methods have been proposed in previous works. In [3], [8], an Unsupervised Random Forest (URF) algorithm is proposed to cluster traffic scenarios. The clustering is based on custom features that are selected based on expert knowledge. In [7], an automated way for deriving traffic scenarios is presented. Even though the dependence on handcrafted features are reduced, still some features and distance measures have to be selected by a human.

In comparison to the aforementioned methods, this work does not rely on any handcrafted features and the proposed method can be applied to all kinds of scenarios irrespective of the number of vehicles present. More importantly, unlike the works mentioned above, instead of considering this problem as completely unsupervised, this work explores how the knowledge from the available labelled scenarios can help in guiding the clustering of the unlabelled scenarios. This idea is based on the assumption that representations learned for the labelled scenarios can also provide a good representation for the unlabelled scenarios. So, in this work it is assumed that both labelled and unlabelled traffic scenarios are available.

The clustering process consists of three steps as shown in Fig. 1. The steps are as follows:

- 1) Model initialisation using a self-supervised objective - *pre-training*.
- 2) Fine tuning the model with labelled scenarios - *classification*.
- 3) Optimising the model for clustering - *clustering*.

The first step is to pre-train a 3D CNN [9] using self-supervision objective i.e., a *pretext* task. The pretext task here means defining a supervised task to train the network without the need for actual ground truth labels, i.e., the network is trained using labels that are generated in an automated fashion without any human input. In computer vision, predicting the angle of rotation, colourising images are used as pretext tasks. Hence, both the labelled and unlabelled scenarios are used to pre-train the 3D CNN. It is important to note that the actual ground truth labels from the labelled scenarios are not used in this step.

As a second step, a classification head is added to the pre-trained 3D CNN. Only the last layers of the 3D CNN and the classification head are tuned on the labelled data. This ensures the scenario classification along with maintaining the general feature extraction capability from the pre-trained 3D CNN. This is shown in Fig. 1, Step II. The red portion in the figure represents the frozen layers and the green portion represents the fine-tuned layers of the 3D CNN.

As a third and final step, a clustering head is added to the fine-tuned 3D CNN. The last layers of the 3D CNN along with the classification and clustering heads are iteratively optimised using both the labelled and unlabelled scenarios. The classification head is optimised using the ground truth labels from the labelled scenarios. The clustering head is optimised using a data-adaptive similarity generated from an URF algorithm called Random Forest Activation Pattern (RFAP) similarity. The implementation of the architecture is made publicly available¹.

The contributions of this work are the following:

- Introduction of a self-supervised pretext task for traffic scenarios.
- Introduction of novel data-adaptive features called RFAPs. The RFAPs provide a data-adaptive similarity measure for the unlabelled data.
- Presenting a method for unsupervised clustering of unlabelled traffic scenarios given few labelled scenarios.
- Comprehensive analysis and ablation studies on the proposed methods using the highD [10] dataset.

The remainder of the paper is organised as follows: Section II presents the related work. Section III discusses the proposed method. Section IV illustrates the experiments and analysis on the highD dataset. Section V discusses the ablation studies. Finally, the paper is concluded in Section VI.

II. RELATED WORKS

A. Traffic Scenario Clustering

Unsupervised traffic scenario clustering methods have been studied in previous works [3], [6], [7], [11]. There are also methods [12], [13] that focuses on trajectory clustering which compares and clusters trajectories from a single vehicle. Since this work focus on clustering scenarios with multiple traffic participants, publications about clustering trajectories are not discussed further. In [8], an URF algorithm is proposed to identify clusters from simulated driving data.

The authors suggest using the path based proximity from an URF algorithm along with Hierarchical Clustering (HC) to solve the task. But, the method relies on features selected by experts and the clusters are selected visually, which limits the amount of data that can be clustered at a single time. In [7], an automatic scenario clustering method on real-world driving data is proposed. The automatic clustering uses dynamic time warping to compare distances between trajectories based on handcrafted features and the generated distance measure is used to construct scenario clusters. In [4], a spatial filter to determine the relevant target objects around the ego and a custom distance metric along with HC is proposed to cluster scenarios.

There are also approaches like [14] which cluster traffic scenarios based on the interaction between two vehicles using Long Short Term Memory (LSTM)+CNN. A deep learning based method for clustering traffic scenarios is presented in [6]. Spatio-temporal autoencoders and recurrent neural networks are used for solving the task of traffic scenario clustering.

In comparison to the above mentioned works except [6], which does feature engineering, the scenarios in this work are described as a sequence of occupancy grids. Such a representation can be generated from most of the common autonomous vehicle sensor suites. Also, in this work the scenarios considered are not limited by the number of traffic participants around the ego. As opposed to all the clustering methods discussed above, the problem setting addressed in the proposed method is different, this work utilises labelled traffic scenarios in guiding the clustering of unlabelled traffic scenarios.

B. Method Comparison

The method presented in this work is based on [15], but extends the architecture in the following ways: (1). The traffic scenarios used in this work are described as time series data as opposed to images in [15], (2). A new self-supervised objective is introduced for training a 3D CNN for traffic scenarios, (3). A data-adaptive similarity measure based on novel features generated from a RF algorithm [16] is introduced for clustering purposes. Also, intended application of the proposed method is different.

III. METHODOLOGY

This section details the clustering methodology. A dataset with unlabelled traffic scenarios $\mathcal{D}_u = \{\mathbf{G}_1^u, \dots, \mathbf{G}_{M_u}^u\}$, where the traffic scenarios are represented as a sequence of occupancy grids \mathbf{G} , is available. Also, a dataset with labelled traffic scenarios $\mathcal{D}_l = \left\{(\mathbf{G}_1^l, y_1), \dots, (\mathbf{G}_{M_l}^l, y_{M_l})\right\}$, where y_{m_l} is the scenario label. M_l and M_u are the number of labelled and unlabelled data respectively. The clustering is realised in a three step process. In step I, a 3D CNN is pre-trained on a defined pretext task using both \mathcal{D}_l and \mathcal{D}_u to learn robust feature representations (see Sec. III-B). In step II, a classification head is added to the pre-trained 3D CNN. The classification head and the last layers of the 3D CNN are fine-tuned only using \mathcal{D}_l (see Sec. III-C). In step

¹<https://github.com/lab176344/TrafficScenarios-RFAPsClustering>

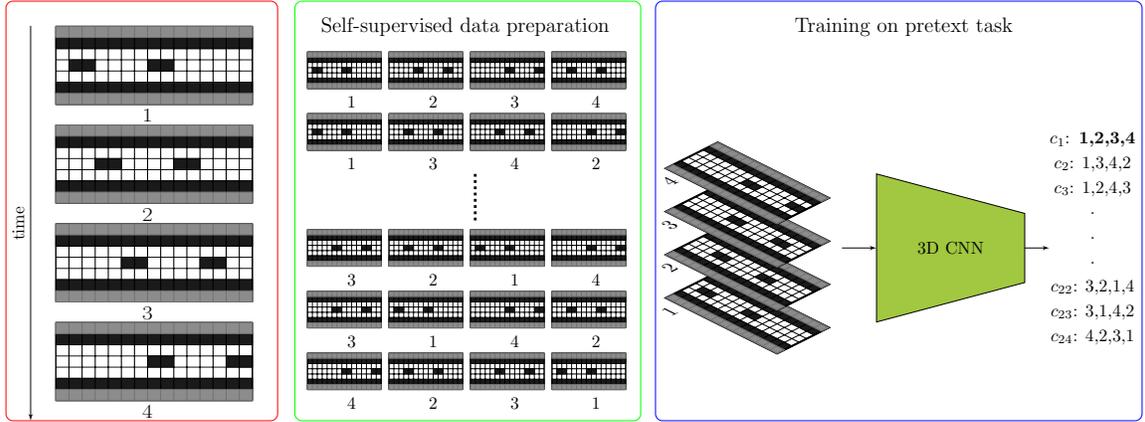


Fig. 2: Self-supervised learning for traffic scenarios. The red box shows the original scenario in the correct temporal order, the green box shows the self-supervised data preparation, and the blue box shows the pretext task of classifying the input to one of the 24 classes.

III, a clustering head is added to the fine tuned 3D CNN. An iterative optimisation procedure is performed on the last layers of the 3D CNN, the clustering and the classification heads. For solving the clustering task using both \mathcal{D}_1 and \mathcal{D}_u are used (see Sec. III-D). The traffic scenario representation used in this work is described first before explaining the methodology in detail.

A. Traffic Scenario Representation

Traffic scenarios in this work are described as introduced in [2], a discretised space-time representation of the environment around the ego, from the time t_{-3} to t_0 . The t_{-3} refers to a time before t_0 , i.e., a time before the traffic situation becomes interesting. The traffic scenario at each time instance is represented as a 2D occupancy grid $\mathbf{G}_t \in \mathbb{R}^{I \times J}$. The occupancy probability of each cell $\mathbf{G}_t(i, j)$ is assigned either as 1 for occupied space, 0 for free space or 0.5 for an unknown region. For a time span of $t_0 - t_{-3} = 1.5s$ and Δt of 0.5s, a traffic scenario is represented as $\mathbf{G} = [\mathbf{G}_{t_{-3}}, \mathbf{G}_{t_{-2}}, \mathbf{G}_{t_{-1}}, \mathbf{G}_{t_0}]$, where $\mathbf{G} \in \mathbb{R}^{I \times J \times N_{ts}}$ with I rows, J columns and a depth of N_{ts} . The number of time steps is $N_{ts} = 1 + \frac{t_0 - t_{-3}}{\Delta t}$. An exemplary representation is shown in Fig. 2. The trigger used to determine the time t_0 at which the traffic situation becomes interesting is determined based on Time-Headway (THW).

B. Step I: Self-Supervised Learning for Traffic Scenarios

The main aim in this step is to train a 3D CNN network with a pretext task to provide robust low-level features for the traffic scenarios. To do this, the most naive way is to train the 3D CNN only on the labelled data first. The trained model can be used as a feature extractor for the unlabelled data and clustering is done on the extracted features. But, this does not guarantee that a model learned from the labelled data will generate good features for the unlabelled data. The model might be biased towards the classification of the labelled data. Hence, this work proposes a self-supervised

pre-training similar to [15], [17] to generate robust low-level features which later can be fine tuned for the clustering.

The self-supervised objective does not require any data annotations, so it can be applied for both labelled and unlabelled data. The idea behind self-supervision is that in the course of solving the self-supervised objective the model will learn some semantic structure in the data and learn robust low-level features. The self-supervised objective used in this work is a pretext task.

The pretext task is to predict the correct temporal order given a sequence of occupancy grids. Consider a sequence of four occupancy grids \mathbf{G} that describe a traffic scenario. The sequence of four occupancy grids in the correct temporal order can be shuffled and 24 different combinations of sequences can be produced i.e., say the correct temporal order is (1, 2, 3, 4) after shuffling one can have 24 different orders like $\{c_1 = (1, 2, 3, 4), \dots, c_{23} = (3, 1, 4, 2), c_{24} = (4, 2, 3, 1)\}$. This can be seen in the green box of Fig. 2. So, each of the shuffled 24 sequence of grids can be assigned to one of the order from $\{c_1, \dots, c_{24}\}$. The problem this way is converted to a 24-class supervised classification task for assigning an input grid \mathbf{G} to one of the 24 orders $\{c_1, \dots, c_{24}\}$. This can be seen in the blue box of Fig. 2, where \mathbf{G} with the order (1, 2, 3, 4) is given as input and the network chooses c_1 as the output. The labels for this pretext task is generated without using any ground truth labels from \mathcal{D}_1 are used here. Temporal order shuffling has been explored as pretext task for video classification in [18]. The rationale behind temporal shuffling is to make the model understand the temporal structure in the data by reasoning out how the vehicles moves in the scenario snippets. The classification here in this work is done with a 3D CNN $f(\mathbf{G})$, where the convolution happens both in spatial and temporal dimensions. The 3D CNN is trained with categorical cross entropy to classify a shuffled $f(\mathbf{G})$ to one of the 24 classes as shown in the blue box in Fig. 2.

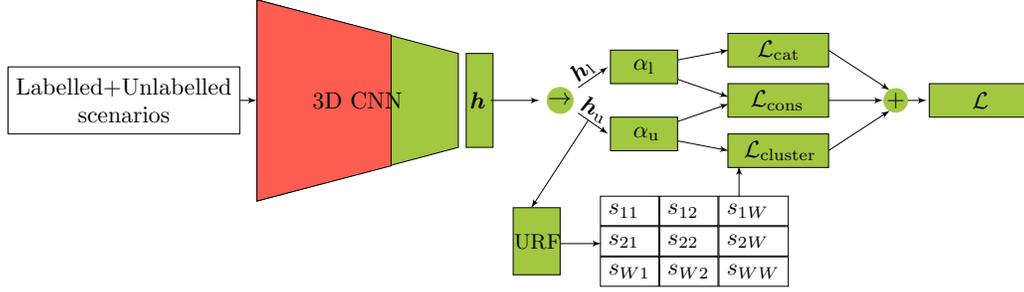


Fig. 3: Iterative optimisation procedure - single epoch. The operator \rightarrow separates the labelled and unlabelled data.

C. Step II: Fine-Tuning with Labelled Data

The goal in this step is to fine-tune $f(\mathbf{G})$ which is pre-trained in step I using the ground truth labels from \mathcal{D}_l . The mapping $\mathbf{G} \mapsto \mathbf{y}$ is realised, where $\mathbf{y} \in \mathbb{R}^K$ is the one-hot representation of the scenario label with K labelled classes. Before fine-tuning, the fully connected layers from $f(\mathbf{G})$, which were trained for the 24-class problem, are removed. A new classification head α_l is added to $f(\mathbf{G})$ to classify the input to one of K classes. The classification head here refers to a fully connected layer followed by a softmax activation layer. Only the last layers of $f(\mathbf{G})$ and α_l are tuned to retain the robust low-level features from pre-training and to learn only high level features for scenario classification. The network $f(\mathbf{G}) : \mathbb{R}^{I \times J \times N_s} \rightarrow \mathbb{R}^F$ takes a traffic scenario $\mathbf{G} \in \mathbb{R}^{I \times J \times N_s}$ as input and produces the vectorised representation $\mathbf{h} \in \mathbb{R}^F$, where F is the dimensionality of the vectorised representation. Given \mathcal{D}_l , the classification head $\alpha_l(\mathbf{h}^l) : \mathbb{R}^F \rightarrow \mathbb{R}^K$ takes the representation \mathbf{h}^l as input and produces the target vector $\hat{\mathbf{y}} \in \mathbb{R}^K$, where K is the number of classes in the labelled data. The last layers of the network $f(\mathbf{G})$ and the classification head $\alpha_l(\mathbf{h}^l)$ are trained using the categorical cross entropy,

$$\mathcal{L}_{\text{cat}} = -\frac{1}{M} \sum_{m=1}^M \sum_{k=1}^K \mathbf{y}_{m,k} \log(\hat{\mathbf{y}}_{m,k}). \quad (1)$$

D. Step III: Iterative Optimisation with Labelled and Unlabelled Data

The goal of clustering is to divide the given unlabelled data into Q number of groups. To realise this as a final step $f(\mathbf{G})$ is iteratively optimised to learn feature representations for the unlabelled data in \mathcal{D}_u . In this section the following contents are discussed: the iterative optimisation procedure to learn feature representations for the scenarios in \mathcal{D}_u , and similarity calculation based on RFAPs.

1) *Iterative Optimisation*: The intent of doing iterative optimisation task is to fine-tune the last layers and the heads of $f(\mathbf{G})$ using the losses $\mathcal{L}_{\text{cluster}}$, \mathcal{L}_{cat} and $\mathcal{L}_{\text{cons}}$ to cluster the unlabelled scenarios. $\mathcal{L}_{\text{cluster}}$ is the clustering loss to learn feature representations for unlabelled scenarios. \mathcal{L}_{cat} is the categorical cross entropy loss which retains the knowledge from labelled scenarios. $\mathcal{L}_{\text{cons}}$ is used to ensure stability when training with unlabelled scenarios.

The clustering is realised by iteratively-tuning $f(\mathbf{G})$ trained on both \mathcal{D}_u and \mathcal{D}_l , i.e. tuning the network epoch by epoch. In a single epoch the following procedures are performed. First, the representation set $\mathcal{D}_u^h = \{\mathbf{h}_1^u, \mathbf{h}_2^u, \dots, \mathbf{h}_N^u\}$ for the scenarios in \mathcal{D}_u is extracted. Followed by that, an URF algorithm $g(\mathbf{h}^u)$ is trained on the extracted features \mathcal{D}_u^h . Given a mini-batch, the trained URF algorithm can then be used to compute the proximity/similarity matrix \mathbf{S} . The similarity in this work is based on RFAPs which is explained in Section III-D.2. Using \mathbf{S} and extracted features \mathbf{h}_u and \mathbf{h}_l of the labelled and unlabelled scenarios, the three losses mentioned above are constructed and used for optimisation in a single epoch. These procedures are repeated for a selected number of epochs. At each epoch a new URF model $g(\mathbf{h}^u)$ is trained using the extracted features \mathcal{D}_u^h at that epoch. The complete schematic procedure for a single epoch is shown in Fig. 3. The procedures at a single epoch are discussed in detail.

a) *Clustering Loss*: Given an URF $g(\mathbf{h}^u)$ trained on \mathcal{D}_u^h , the similarity matrix \mathbf{S} is computed for a mini-batch of unlabelled traffic scenarios using the RFAPs. The computation of the similarity with RFAPs is explained in detail in Section III-D.2. The value at $S_{ij} \in (0, 1]$, gives the pairwise-similarity between the i -th and j -th scenario from the mini-batch denoting how similar the scenarios are. The optimisation objective for clustering is constructed using the matrix \mathbf{S} . A new clustering head $\alpha_u(\mathbf{h}^u) : \mathbb{R}^F \rightarrow \mathbb{R}^Q$ is added to $f(\mathbf{G})$ parallel to the head $\alpha_l(\mathbf{h}^l)$. Here, Q is the number of clusters in the dataset \mathcal{D}_u . Following [19], [20], [21], the parameter Q is assumed to be known but in Section V a study is performed, where the parameter Q is kept as unknown and Q is estimated as suggested in [22]. The clustering loss to fine-tune the $f(\mathbf{G})$ for the unlabelled dataset \mathcal{D}_u is based on binary cross entropy. It is given by

$$\mathcal{L}_{\text{cluster}} = \frac{1}{W^2} \sum_{i=1}^W \sum_{j=1}^W S_{ij} \log P(i=j) + (1 - S_{ij}) \log P(i \neq j), \quad (2)$$

where W is the mini batch size. $P(i=j)$ denotes the probability that the traffic scenario i and j are in the same cluster. A similar loss was used in [15], [19], but \mathbf{S} in

this work is given by the data-adaptive similarity from the RFAPs. As shown in [19], if the number of clusters Q is fixed and i, j are independent, $P(i = j)$ can be modelled as the inner product between the vectors $\alpha_u(\mathbf{h}_i^u)$ and $\alpha_u(\mathbf{h}_j^u)$. The final clustering loss is given as

$$\mathcal{L}_{\text{cluster}} = \frac{1}{W^2} \sum_{i=1}^W \sum_{j=1}^W S_{ij} \log(\alpha_u(\mathbf{h}_i^u)^\top \alpha_u(\mathbf{h}_j^u)) + (1 - S_{ij}) \log(1 - (\alpha_u(\mathbf{h}_i^u)^\top \alpha_u(\mathbf{h}_j^u))). \quad (3)$$

Now as the clustering loss is defined, the last layers of $f(\mathbf{G})$ and the new clustering head can be fine-tuned for the clustering task. An argmax on the vector $\alpha_u(\mathbf{h}^u)$ will provide the cluster number to which an input scenario belongs.

But, training only on the unlabelled data will destroy the representation learned for the labelled dataset, leading to catastrophic forgetting [23]. Hence, both heads α_l and α_u are trained in parallel on categorically cross entropy using the dataset \mathcal{D}_l and the clustering loss using the dataset \mathcal{D}_u respectively. Additionally, the head α_l is extended to classify the clusters identified by the clustering head α_u . So, the head α_l does the mapping $\mathbb{R}^F \rightarrow \mathbb{R}^{K+Q}$ and it is trained on the categorical cross entropy. This way the knowledge of the labelled traffic scenarios is also preserved when optimising for clustering the unlabelled traffic scenarios. The knowledge from the labelled traffic scenarios is shown to improve clustering accuracy from in Section V.

b) Consistency Loss: The similarity S_{ij} between two given traffic scenarios is updated every epoch as the representation is optimised by the categorical cross entropy and the clustering losses. This might lead to instability or inconsistency in the training for the unlabelled dataset. Following [15], [19], a consistency constrain, which is used in semi-supervised learning setting [24], is also introduced. The consistency loss constrains the model to produce the same output for a given traffic scenario $\mathbf{G}_i \mapsto \mathbf{h}_i$ and an augmented version of the traffic scenario $\hat{\mathbf{G}}_i \mapsto \hat{\mathbf{h}}_i$ (e.g. random erasing, adding random noise). The consistency loss is given by

$$\mathcal{L}_{\text{cons}} = \frac{1}{M_l} \sum_{i=1}^{M_l} (\alpha_l(\mathbf{h}_i^l) - \alpha_l(\hat{\mathbf{h}}_i^l)) + \frac{1}{M_u} \sum_{i=1}^{M_u} (\alpha_u(\mathbf{h}_i^u) - \alpha_u(\hat{\mathbf{h}}_i^u)). \quad (4)$$

c) Total Loss: The total loss used to optimise $f(\mathbf{G})$ for the clustering task while keeping the knowledge from the labelled data is given by,

$$\mathcal{L} = \mathcal{L}_{\text{cat}} + \mathcal{L}_{\text{cluster}} + \omega(\beta) \mathcal{L}_{\text{cons}}. \quad (5)$$

Following [25], $\omega(\beta) = \lambda \exp(-5(1 - \frac{\beta}{T})^2)$ is the ramp-up function where β being the epoch, T is the ramp-up length and $\lambda \in \mathbb{R}_+$. Using \mathcal{L} , the last layers of the $f(\mathbf{G})$ and the heads α_l and α_u are tuned. This procedure is repeated a given number of epochs.

2) Data-adaptive similarity with RFAPs: In this work, the similarity matrix \mathbf{S} is defined by a data-adaptive similarity measure based on features called RFAPs. In the following, the generation of RFAPs from a URF algorithm and calculation of \mathbf{S} with RFAPs are explained.

a) Generation of RFAPs: The RFAPs are generated by an URF algorithm $g(\mathbf{h}^u)$ trained on the extracted features \mathbf{h}^u of the scenarios from \mathcal{D}_u . The generation of RFAPs and the calculation of similarity will be discussed in this section.

The generation of RFAPs is based on a novel encoding scheme for indexing the nodes of the URF trees. The RFAPs in turn will be used to calculate the data-adaptive similarity. Let the B trees in the URF $\{T_1, \dots, T_B\}$ be fully grown. The b -th tree T_b divides the input space \mathbb{R}^F into many small hypercubes. So, if two data points \mathbf{h}_i and \mathbf{h}_j end up in the same hypercube they are similar to each other. This is termed as proximity in [16]. Similarly, the paths taken by the data points to reach the terminal nodes can also provide a similarity measure by comparing the common paths taken [8].

The paths taken by the data point \mathbf{h}_i in all the trees can be represented by the vector $\mathbf{r}_i \in \mathbb{N}^B$. Each element of \mathbf{r}_i is an ordinal number representing a path through a tree. By using a special node encoding, a single number is sufficient to identify a specific path in a tree. In the following, the encoding of a data point's path in one tree is presented. This encoding is the basis for RFAPs. In a first step, all nodes of each tree have to be indexed. Looking at one specific node in the tree T_b the following quantities are used for indexing: the maximum depth of the tree d_b , the depth of the considered node $k_b \in \{1, \dots, d_b\}$, and the number of nodes N_b in T_b .

Algorithm 1 Indexing for a tree in the RF

Input: T_b, d_b, N_b

Output: Indexed T_b

```

1:  $id^1 \leftarrow 0$ 
2: for  $n = 2$  to  $N_b$  do
3:    $id^{\text{pr}} = \text{getParentNodeid}(id^n)$ 
4:   if  $n$ .isleft then
5:      $id^n \leftarrow id^{\text{pr}} + 10^{d_b - k_b}$ 
6:   else
7:      $id^n \leftarrow id^{\text{pr}} + 2 * 10^{d_b - k_b}$ 
8:   end if
9: end for

```

The algorithm for indexing a tree is depicted in Alg. 1. The index of the root node id^1 at $k_b = 1$ in a tree is set to 0. Each level in a tree is indexed as follows: if the node is a left child node then the index for the child node is assigned as the sum of the index of its parent node id^{pr} and $10^{d_b - k_b}$, while the right child node is indexed as the sum of its id^{pr} and $2 * 10^{d_b - k_b}$. With this indexing of nodes, the path of data point \mathbf{h}_i in the tree T_b is encoded in the ordinal number assigned to the terminal node reached by \mathbf{h}_i . An example indexed tree T_b is shown in Fig. 4 with $N_b = 9$ and $d_b = 4$. Assume a data point \mathbf{h}_i , which is reaching the terminal

node 211. The index 211 is encoding the path shown in red: $0 \rightarrow 200 \rightarrow 210 \rightarrow 211$. Considering all B trees in a RF, the RFAP-representation for the data point \mathbf{h}_i is represented as $\mathbf{r}_i = [id_1^i, id_2^i, \dots, id_B^i]^T$, where id_b^i denotes the index of the terminal node reached by \mathbf{h}_i in the tree T_b .

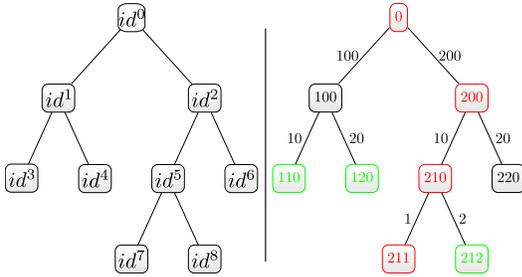


Fig. 4: RF tree before indexing (left) and indexed tree (right).

b) Calculating Similarity with RFAPs: The advantage of RFAPs is that the complete path taken by two data points can be used to calculate similarity. The digits of the RFAP indices encode the complete path information. The more similar the digits are between two RFAP indices, the deeper the paths are shared to reach the terminal nodes. Hence, the similarity with RFAPs is calculated using the hamming distance. The indices of the nodes are of equal length, and if the indices are treated as sequences of digits, the similarity S_{ij} can be given as,

$$S_{ij} = 1 - \frac{1}{B} \sum_{b=1}^B \frac{|\{o \in \{1, \dots, |\mathbf{r}_i^b|\} \mid \mathbf{r}_i^b[o] \neq \mathbf{r}_j^b[o]\}|}{|\mathbf{r}_j^b|}. \quad (6)$$

Here, \mathbf{r}_i^b refers to the RFAP index at the b -th position in \mathbf{r}_i . The numerator checks how many digits are not equal when comparing \mathbf{r}_i to \mathbf{r}_j . The $|\cdot|$ returns the digits length of elements in the vector \mathbf{r}_i^b . For e. g., in Fig. 4, with the terminal node indices 211 and 212, the similarity in a single tree is $1 - (1/|211|)$, where $|211| = 3$. A S_{ij} value of 1 means the two data points i and j took the same path in all the trees.

In summary, the clustering is done in a three step process: As a first step, in order to get robust low-level features a 3D CNN is trained on a pretext task. As a second step, the trained 3D CNN is fine-tuned on labelled classes. As the final step, the fine-tuned network is iteratively optimised using the three losses combined for clustering while preserving the knowledge from the labelled data.

IV. EXPERIMENTS AND RESULTS

This section reports the experiments and results of the proposed method applied on the highD dataset.

A. Dataset

The highD dataset is a naturalistic vehicle trajectory dataset recorded using drones on German motorways. The dataset consists of 16.5 hours of drone video records in six different locations and contains around 110 000 vehicles. The

data analysed in this work is restricted to traffic scenarios where the ego vehicle has a leading vehicle. Since this work focuses on traffic scenarios to validate the clustering method, 7 common highway scenarios are extracted from the dataset:

- 1) Ego lane change to the right lane,
- 2) Ego lane change to the left lane,
- 3) Leader cutting into ego's lane from the left lane,
- 4) Leader lane change from ego's lane to the left lane,
- 5) Leader cutting into ego's lane from the right lane,
- 6) Ego following the leader in a lane,
- 7) Leader lane change from ego's lane to the right lane.

As described in III-A, $\text{THW} < 4\text{s}$ is used as the criterion for finding interesting scenarios. The environment at each time instance is represented with \mathbf{G}_t of span $15\text{m} \times 200\text{m}$ and a resolution of $0.5\text{m} \times 1\text{m}$. The interval $t_{\text{fb}} - t_0$ and Δt used in this work are 2s and 0.5s respectively. So, a traffic scenario is represented with \mathbf{G} of size $30 \times 200 \times 4$. The grids are generated in an ego-centric fashion fixed at t_0 . In total 4480 scenarios are extracted and is split as 70% for training, 10% for validation and 20% for testing.

B. Implementation Details

The $f(\mathbf{G})$ is a 3D-Resnet [26], a 3D version of the normal Resnet [27] with spatio-temporal operations. The URF is trained with $B = 500$ trees. The parameters T and λ are set to 100 and 5 respectively. The $f(\mathbf{G})$ is trained using SGD optimiser [28] with moment and decay and with a batch size W of 32. Only the last Residual block of the 3D-ResNet is trained in the Step II and III.

C. Evaluation Metric

The evaluation metric used for measuring the clustering accuracy (ACC) is the unsupervised clustering accuracy [29]. The best mapping between the labels obtained from clustering and the ground truth is computed by the Hungarian algorithm. The ACC is defined as follows

$$\text{ACC} = \max_m \frac{\sum_{j=1}^{M_u} 1(y_j = m(c_j^u))}{M_u}, \quad (7)$$

where y_j is the ground truth and c_j^u is the predicted label for the unlabelled sample \mathbf{G}_j . The range of the ACC is $[0, 1]$, with 1 referring perfect clustering.

D. Baselines

The proposed method is compared with the following methods: (a) K -means clustering [30] directly on the dataset \mathcal{D}_u . (b) A spatio Temporal Autoencoder + Heirarchical Clustering (STAE+HC) [6], a spatio-temporal autoencoder is trained on the dataset \mathcal{D}_u and HC is performed on the latent space of the autoencoder. (c) Autonovel [15], a spatio-temporal extension of the proposed method in [15] is done and the method uses rank statistics as similarity measure for clustering. (d) Comparison with other similarity measures like cosine, L2, rank statistics [15] and the similarity from a K -Nearest Neighbour (KNN) algorithm for determining \mathcal{S} .

E. Results

1) *Clustering*: In this experiment, the following problem setting is used. The first 4 classes out of the available 7 scenarios are treated as labelled classes and the remaining 3 classes are treated as unknown. In the ablation study, an analysis is also conducted with a different combination of labelled and unlabelled classes. The aim of this experiment is to measure the clustering accuracy on the 3 unlabelled classes. As suggested in [19], [20], the experiment is repeated 5 times and an average of the accuracy (ACC) is reported in the Table I. There, it can be seen that the proposed method outperforms the baselines and provides good clustering accuracy.

TABLE I: highD - Clustering accuracy.

Method	ACC (\uparrow)
<i>K</i> -means [30]	0.391
STAE+HC [6]	0.52
Autonovel [15]	0.794
Proposed method (RFAPs)	0.810

2) *Comparison with Other Similarities*: In this experiment, the proposed RFAP based pairwise similarity used in the mini-batch for clustering is compared with standard similarity measures like cosine, L2, rank statistics [15] and similarity from KNN. From the results shown in Table II, it can be seen the RFAP similarity provides superior performance when compared to other standard similarity measures.

TABLE II: Clustering accuracy with different similarities.

Similarity	cosine	L2	KNN	rank [15]	RFAPs
ACC (\uparrow)	0.707	0.703	0.793	0.794	0.810

V. ABLATION STUDY

A. Importance of Step I

To underline the importance of self-supervised initialisation, another experiment is setup. For this, as a first step all the layers of the 3D CNN is trained on the labelled classes without any Self-Supervised Learning (SSL). As a second step, the iterative optimisation as described in the Section III. As seen in Table III, the ACC improves with SSL initialisation.

TABLE III: ACC with and without SSL initialization.

Method	ACC (\uparrow)
with SSL	0.810
without SSL	0.537

B. Estimating Q

Until now the number of classes Q in the unlabelled traffic scenario dataset is assumed to be given. In this study, Q in the unlabelled traffic scenarios is estimated as suggested in [22]. In [22], the number of clusters in the unlabelled samples are found using the Silhouette index [31]. The experiment setting

is identical to the one used for clustering with 4 labelled and 3 unlabelled data. The Q is estimated to be 4 compared to the ground truth (GT) 3 as seen in the Table IV.

TABLE IV: Number of classes.

Method	GT # of classes	Predicted # of classes	Error
RFAPs	3	4	1

C. Randomly Chosen Labelled Classes

This study is conducted to analyse the influence of the chosen labelled and unlabelled classes. Until now, out of the 7 scenarios, the first 4 were considered labelled and the remaining 3 were considered unlabelled. Here, three possible combinations of labelled and unlabelled classes are chosen randomly and the clustering performance is analysed. As shown in Table V, the results remain consistent.

TABLE V: ACC for randomly chosen unlabelled classes.

Labelled classes	Unlabelled classes	ACC (\uparrow)
1,2,3,4	5,6,7	0.810
2,3,6,7	5,1,4	0.82
1,2,4,5	3,6,7	0.807

D. Importance of Step II

The aim in this study is to understand the influence of using the labelled data in the clustering step (step III) of the method. To perform this, step II is skipped in the proposed method where 3D CNN is fine tuned with labelled data. In the step III, the labelled data is not used and the categorical cross-entropy loss is also dropped. The experiment set-up has 4 labelled classes and the remaining 3 classes are considered as unlabelled data. The clustering accuracy with and without the labelled data is shown in Table. VI. The results shows experimentally that the labelled data indeed helps in guiding the clustering process.

TABLE VI: ACC with and without labelled data.

Method	ACC (\uparrow)
with labelled data	0.810
without labelled data	0.565

E. Importance of Step III

The intent of the iterative optimisation in Section III is to provide good feature representations for the unlabelled samples. The experiment setup is similar to the one used in Section IV. The feature representations for the considered 3 unlabelled classes are extracted. The feature representations before and after the iterative optimisation are projected on to a 2D space by UMAP [32] as shown in the Fig. 5. It can be seen that there is a good separation of clusters. The clusters, the leader cut-out from ego lane to the right lane (class 7) and the ego following a leader (class 6) are close and few datapoints are mixed up. This is because in both of these scenarios there is a leader vehicle in front of the ego in most

of the time for a scenario. The clustering accuracy of the 3 classes before iterative optimisation is 48.65% and after iterative optimisation is 81.0%. This shows that the iterative optimisation procedure improves the clustering accuracy.

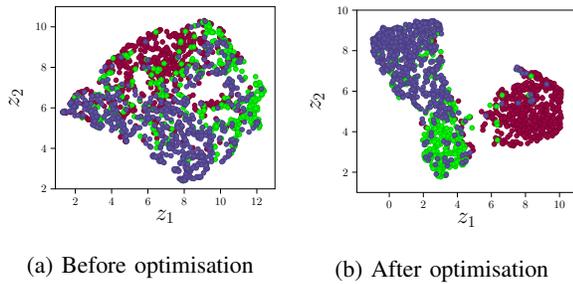


Fig. 5: Feature representation study, ■ Ego following, ■ Cut-in from right, ■ Leader cut-out to right.

VI. CONCLUSION

Traffic scenario clustering is an important problem to be solved for identifying relevant and new traffic scenarios. The new traffic scenario categories are important for the development of motion planning algorithms and the validation of autonomous functionalities. This work proposes a method to cluster traffic scenarios automatically without any handcrafted features using self-supervised learning and a data-adaptive similarity based on novel features called RFAPs. The problem set-up in this work uses labelled scenarios and retains the knowledge about labelled scenarios to aid the clustering task. The 3D CNN is trained robust feature representation on a defined pretext task followed by fine-tuning using labelled traffic scenarios. The clustering is addressed by an iterative optimisation procedure using the labelled and unlabelled traffic scenarios. Experiments on the highD dataset have verified the advantages of the solution proposed when compared to considered baseline methods.

ACKNOWLEDGEMENT

This work is supported by Bavarian State Ministry for Science and Art under the funding code VIII.2-F1116.IN/18/2.

REFERENCES

- [1] H. Winner et al., *Handbook of Driver Assistance Systems: Basic Information, Components and Systems for Active Safety and Comfort*, 1st ed. Springer Publishing Company, Incorporated, 2015.
- [2] A. Chaulwar et al., “A machine learning based biased-sampling approach for planning safe trajectories in complex, dynamic traffic-scenarios,” in *2017 IEEE Intelligent Vehicles Symposium (IV)*, June 2017, pp. 297–303.
- [3] F. Kruber et al., “An unsupervised random forest clustering technique for automatic traffic scenario categorization,” in *2018 IEEE International Conference on Intelligent Transportation Systems*, 11 2018, pp. 2811–2818.
- [4] J. Kerber et al., “Clustering of the scenario space for the assessment of automated driving,” in *Intelligent Vehicles Symposium 2020*, IEEE, Ed., 2020, (accepted).
- [5] T. Menzel et al., “Scenarios for development, test and validation of automated vehicles,” *2018 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1821–1827, 2018.
- [6] N. Harmening et al., “Deep representation learning and clustering of traffic scenarios,” *ArXiv*, vol. abs/2007.07740, 2020.

- [7] F. Hauer et al., “Clustering traffic scenarios using mental models as little as possible,” in *PrePrint for the proceedings of IEEE Intelligent Vehicles Symposium 2020*, 2020.
- [8] 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.
- [9] S. Ji et al., “3d convolutional neural networks for human action recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, pp. 221–231, 2010.
- [10] R. Krajewski et al., “The highD dataset: A drone dataset of naturalistic vehicle trajectories on german highways for validation of highly automated driving systems,” in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 2118–2125.
- [11] “Gauging investment in self-driving cars,” <https://www.brookings.edu/research/gauging-investment-in-self-driving-cars/>, (Accessed on 14/10/2018).
- [12] A. Demetriou et al., “A deep learning framework for generation and analysis of driving scenario trajectories,” *ArXiv*, vol. abs/2007.14524, 2020.
- [13] J. Langner et al., “Logical scenario derivation by clustering dynamic-length-segments extracted from real-world-driving-data,” in *VEHITS*, 2019.
- [14] W. Wang et al., “Clustering of driving scenarios using connected vehicle datasets,” *ArXiv*, vol. abs/1807.08415, 2018.
- [15] K. Han et al., “Automatically discovering and learning new visual categories with ranking statistics,” in *International Conference on Learning Representations (ICLR)*, 2020.
- [16] L. Breiman, “Random forests,” *Machine Learning*, vol. 45, pp. 5–32, Oct 2001.
- [17] T. Chen et al., “A simple framework for contrastive learning of visual representations,” *ArXiv*, vol. abs/2002.05709, 2020.
- [18] I. Misra et al., “Unsupervised learning using sequential verification for action recognition,” *ArXiv*, vol. abs/1603.08561, 2016.
- [19] S.-A. Rebuffi et al., “Lsd-c: Linearly separable deep clusters,” *ArXiv*, vol. abs/2006.10039, 2020.
- [20] Y. Wada et al., “Spectral embedded deep clustering,” *Entropy*, vol. 21, 2019.
- [21] X. Guo et al., “Improved deep embedded clustering with local structure preservation,” in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, ser. IJCAI’17. AAAI Press, 2017, p. 1753–1759.
- [22] K. Han et al., “Learning to discover novel visual categories via deep transfer clustering,” in *International Conference on Computer Vision (ICCV)*, 2019.
- [23] S. Rebuffi et al., “icarl: Incremental classifier and representation learning,” *CoRR*, vol. abs/1611.07725, 2016.
- [24] T. Miyato et al., “Virtual adversarial training: A regularization method for supervised and semi-supervised learning,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, pp. 1979–1993, 2019.
- [25] A. Tarvainen and H. Valpola, “Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results,” in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS’17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 1195–1204.
- [26] K. Hara et al., “Learning spatio-temporal features with 3d residual networks for action recognition,” *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pp. 3154–3160, 2017.
- [27] K. He et al., “Deep residual learning for image recognition,” *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- [28] I. Sutskever et al., “On the importance of initialization and momentum in deep learning,” in *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ser. ICML’13. JMLR.org, 2013, p. III–1139–III–1147.
- [29] 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.
- [30] J. Macqueen, “Some methods for classification and analysis of multivariate observations,” in *In 5-th Berkeley Symposium on Mathematical Statistics and Probability*, 1967, pp. 281–297.
- [31] P. Rousseeuw, “Silhouettes: A graphical aid to the interpretation and validation of cluster analysis,” *J. Comput. Appl. Math.*, vol. 20, p. 53–65, Nov. 1987.

- [32] L. McInnes and J. Healy, "Umap: Uniform manifold approximation and projection for dimension reduction," ArXiv, vol. abs/1802.03426, 2018.