

DeepStay: Stay Region Extraction from Location Trajectories using Weak Supervision

Christian Löwens¹, Daniela Thyssens¹, Emma Andersson², Christina Jenkins², and Lars Schmidt-Thieme¹

Abstract—Nowadays, mobile devices enable constant tracking of the user’s position and location trajectories can be used to infer personal points of interest (POIs) like homes, workplaces, or stores. A common way to extract POIs is to first identify spatio-temporal regions where a user spends a significant amount of time, known as stay regions (SRs).

Common approaches to SR extraction are evaluated either solely unsupervised or on a small-scale private dataset, as popular public datasets are unlabeled. Most of these methods rely on hand-crafted features or thresholds and do not learn beyond hyperparameter optimization. Therefore, we propose a weakly and self-supervised transformer-based model called DeepStay, which is trained on location trajectories to predict stay regions. To the best of our knowledge, this is the first approach based on deep learning and the first approach that is evaluated on a public, labeled dataset. Our SR extraction method outperforms state-of-the-art methods. In addition, we conducted a limited experiment on the task of transportation mode detection from GPS trajectories using the same architecture and achieved significantly higher scores than the state-of-the-art. Our code is available at <https://github.com/christianl19/deepstay>.

I. INTRODUCTION

Extracting stay regions (SR) from location trajectories identifies segments where a subject stays in the same place. It supports fine-grained spatio-temporal analysis of human and animal behavior and is often an intermediate step in point of interest (POI) mapping or POI extraction.

Common SR extraction approaches apply unsupervised clustering algorithms and use thresholds for time, distance, and velocity, among others. These thresholds are determined by a qualitative analysis or a quantitative hyperparameter optimization. Typically, all experiments are performed either on small private datasets, in some cases with manually annotated labels, or on large datasets without any labels. This makes it difficult to compare different approaches and makes the problem less suited for supervised learning that requires a large amount of labeled data.

Even though most trajectories do not contain ground truth SR labels, it is still possible to derive so-called “weak labels” from OpenStreetMap (OSM). For example, we can classify any location point lying within a building as part of a stay and a point near a road as part of a “non-stay” (see Figure 1). Given the large number of weak labels available, we assume

¹Christian Löwens, Daniela Thyssens, and Lars Schmidt-Thieme are with the Information Systems and Machine Learning Lab (ISMLL), University of Hildesheim, 31141 Hildesheim, Germany. loewensc@uni-hildesheim.de, {thyssens,schmidt-thieme}@isml1.uni-hildesheim.de

²Emma Andersson and Christina Jenkins are with Devoteam Creative Tech, 211 20 Malmö, Sweden. {emma.andersson,christina.jenkins}@devoteam.com

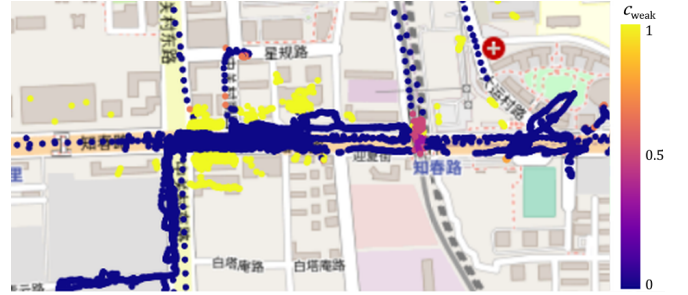


Fig. 1: Weak supervision of trajectories using OSM data and additional heuristics. Blue points indicate weakly labeled non-stays, yellow indicate stays.

that this data contains enough signal to learn useful latent representations. To this end, we apply a transformer model [1] that takes a trajectory as a time series of location points and classifies each point as either part of a stay or a non-stay.

To our knowledge, this is the first approach to extract SRs from trajectory data using deep learning. Furthermore, we use publicly available data for training and evaluation to ensure reproducibility. We derived a ground truth dataset from the field of activity recognition and use it to compare our model with baselines from related work.

II. PROBLEM STATEMENT

We define a location trajectory $\mathcal{X} = \{g_1, g_2, \dots, g_{|\mathcal{X}|}\}$ as a time series of consecutive location points $g_i = (t_i, x_i, y_i)$, where $x, y \in \mathbb{R}$ denote the 2D coordinates and $t \in \mathbb{R}^{\geq 0}$ the ascending timestamp. The sample rate $\Delta t_i = t_i - t_{i-1}$ is defined as the time difference between two consecutive points and is either constant or fluctuating, depending on the dataset.

SR extraction can be viewed as a time series segmentation task, where the trajectory \mathcal{X} is split in a set of segments $\mathcal{TS} = \{ts_1, \dots, ts_q\}$. Each segment $ts_j = (t_{\text{start}_j}, t_{\text{end}_j}, c_j)$ is defined by its start t_{start} and end time t_{end} and the binary class $c \in \{0, 1\}$ indicating whether the user is staying at one place ($c = 1$) or is moving around ($c = 0$) within the time window $t_{\text{start}} \leq t < t_{\text{end}}$. Moreover, we define

$$\begin{aligned} t_{\text{start}_1} &= t_1, \\ t_{\text{end}_q} &= \infty, \\ t_{\text{start}_j} &< t_{\text{end}_j} && \forall j \in \{1, \dots, q\}, \\ t_{\text{end}_j} &= t_{\text{start}_{j+1}} && \forall j \in \{1, \dots, q-1\}, \\ c_j &\neq c_{j+1} && \forall j \in \{1, \dots, q-1\}. \end{aligned}$$

The set of stay regions \mathcal{SR} is a subset of all segments, where

$$\mathcal{SR} = \{ts_j | ts_j \in \mathcal{TS} \wedge c_j = 1\}. \quad (1)$$

The task of SR extraction is now to predict \mathcal{SR} (and therefore \mathcal{TS}) solely from the trajectory data \mathcal{X} .

III. RELATED WORK

Trajectory segmentation is an important research topic with many examples such as activity recognition, transportation mode detection (TMD) and SR extraction. In TMD, each segment is assigned to a mode, e.g. walking, car, bus, etc. [2]. A special binary case of this task is SR extraction with only two possible modes: stay and non-stay. In most cases, it functions as a preprocessing step for tasks such as POI mapping/extraction/prediction. In POI mapping, each SR is assigned to a visit to one of several POIs [3].

SR extraction identifies segments of a user's trajectory where the subject remains at the same place. A virtual location, usually the centroid of an SR, is called a stay point. So this task is also called stay point extraction/recognition/identification/detection.

A. Threshold-based Clustering

The vast majority of published work uses threshold-based spatio-temporal clustering methods, where the clusters represent stay segments of the trajectory. Commonly used thresholds are a minimum duration T_{\min} and a maximum distance D_{\max} [4], [5], [6], [7], [8]. Here, the task is to find the maximum sets of consecutive location points $\mathcal{P} = \{g_m, g_{m+1}, \dots, g_n\}$ in the trajectory \mathcal{X} , such that:

$$t_n - t_m \geq T_{\min} \quad (2)$$

$$\text{dist}(g_i, g_j) \leq D_{\max} \quad \forall g_i, g_j \in \mathcal{P} \quad (3)$$

Others apply additional thresholds for velocity, acceleration, and heading change [9], [10], [11], [12], [13].

B. Adapted Density-based Clustering

Other approaches adapt density-based clustering methods such as DBSCAN [14] and OPTICS [15]. If the trajectory is sampled at a constant rate, prolonged stays will result in dense spatial data and thus can be detected. Unlike k-means, they do not require a predetermined number of clusters, which is crucial for SR extraction.

These approaches define SR extraction more as spatial clustering rather than time series segmentation. Therefore, the constraint that the clustered points must be consecutive is not always enforced. Many extensions have been proposed to utilize the temporal information as well [16], [17].

C. Others

The authors in [18] and [19] classify single location points as stay points when a GPS connection loss is detected. The algorithm proposed by [20] extracts SRs by searching for local minima of speed and zero crossings in acceleration within the trajectory.

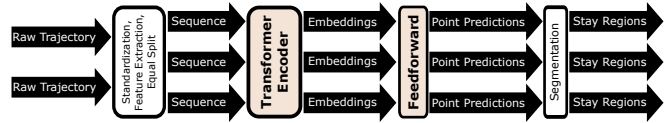


Fig. 2: Overall architecture of DeepStay. Light brown colored boxes indicate trainable models.

IV. METHODOLOGY

A. Architectural Overview

Figure 2 shows the overall architecture of our model DeepStay and the intermediate results of the processing pipelines. First, the raw trajectories are standardized and split into sequences of equal size. Furthermore, additional features are extracted to improve the performance of the subsequent transformer encoder. This encoder receives a sequence of constant length and outputs an embedding vector for each point comprising latent features about the point within its sequence. The following feedforward layer acts as a decoder and predicts a probability for each vector to be part of a stay. In the next step, all consecutive points with a predicted probability above a certain threshold are grouped as SRs.

1) *Preprocessing of raw GNSS trajectories*: All datasets in this work contain GNSS coordinates, such as GPS. In the first step, we project all coordinates into a 2D Cartesian system (x, y) using an appropriate UTM zone [21].

Since trajectories may have varying sample rates, we use the time difference Δt between each point and its predecessor as an additional feature. Our preliminary experiments indicate that this approach leads to better results than using linear interpolation as proposed by [22]. We also add the current velocity v as the ratio of the Euclidean distance and Δt between two consecutive points as another feature.

All trajectories are chunked into sequences of equal length $n = 256$. This allows the transformer encoder to be trained with multiple sequences in a single batch of size $B = 64$.

Furthermore, we standardize the features Δt and v separately based on their distribution in the training set to obtain a mean of 0 and a standard deviation of 1. The standardization of the location features x and y is done jointly. Each sequence is subtracted from its mean $(\bar{x}_{\text{seq}}, \bar{y}_{\text{seq}})$ and divided by the common standard deviation of the entire training set $\sigma_{x,y_{\text{train}}}$ to prevent the model from memorizing specific regions. To further reduce overfitting, we rotate every sequence uniformly at random with respect to its origin $(0, 0)$ before feeding it to the model. The final features of the i -th data point in sequence seq are shown in 4.

$$\text{seq}_i = [x_i \quad y_i \quad \Delta t_i \quad v_i], \quad \text{seq} \in \mathbb{R}^{n \times 4} \quad (4)$$

2) *Transformer Encoder*: We choose the encoder of the transformer model [1] to learn latent embeddings emb_i for each sequence point seq_i . This allows us to predict the class probabilities pointwise instead of segmentwise. Thus, by design, segmentation and classification are performed jointly. We stick with the original setting of the base encoder

including the projection and positional encoding as described in [1] to get the final embeddings $emb \in \mathbb{R}^{n \times d_{\text{model}}}$.

3) *Decoder*: A feedforward layer with sigmoid activation decodes the embeddings and predicts the probability for each point emb_i to be part of a stay:

$$\hat{c}_i = \sigma(emb_i W_d^T + b_d), \quad \hat{c}_i \in [0, 1]^n \quad (5)$$

Now the segmentation can be done by simply grouping consecutive points where $\hat{c}_i < 0.5$ for non-SRs and $\hat{c}_i > 0.5$ for SRs, respectively.

4) *Supervision*: In the case of available SR labels, we can compute the pointwise binary cross entropy (BCE) between the prediction \hat{c} and the ground truth c :

$$\text{BCE}(\hat{c}_i, c_i) = -c_i \log \hat{c}_i - (1 - c_i) \log(1 - \hat{c}_i) \quad (6)$$

The distribution of the binary labels can be highly imbalanced. To prevent the model leaning towards one of the classes, we apply class weighting based on the mean \bar{c}_{train} within the training set:

$$\text{BCE}_w(\hat{c}_i, c_i, \bar{c}_{\text{train}}) = \left(\frac{c_i}{\bar{c}_{\text{train}}} + \frac{1 - c_i}{1 - \bar{c}_{\text{train}}} \right) \text{BCE}(\hat{c}_i, c_i) \quad (7)$$

Now the total loss $\mathcal{L}_{\text{super}}$ is the average weighted BCE over all points in all N_{train} training sequences:

$$\mathcal{L}_{\text{super}} = \frac{1}{N_{\text{train}} \cdot n} \sum_{j=1}^{N_{\text{train}}} \sum_{i=1}^n \text{BCE}_w(\hat{c}_i^{(j)}, c_i^{(j)}, \bar{c}_{\text{train}}) \quad (8)$$

B. Weakly Supervised SR Extraction

Since the vast amount of publicly available location trajectories does not contain SR labels, we apply *programmatically* weak supervision [23] by generating weak labels based on other data sources. These labels are often inaccurate. However, since we can generate them on a large scale, and since the error generally does not correlate with the input, we expect our model to still learn useful latent representations.

For that, we define a function f_{weak} that returns the estimated probability $c_{i_{\text{weak}}}$ that the location point g_i is part of a stay, and a confidence score $w_{i_{\text{weak}}}$ for that prediction:

$$f_{\text{weak}}(g_i) = (c_{i_{\text{weak}}}, w_{i_{\text{weak}}}) \quad (9)$$

Here, $c_{i_{\text{weak}}}$ replaces the ground truth value c_i , while $w_{i_{\text{weak}}}$ is used to weight the influence of the weak label on the total loss. Thus, the model learns more from weak labels, where the labeling function is more certain. The total loss is then:

$$\mathcal{L}_{\text{weak}} = \sum_{j=1}^{N_{\text{train}}} \sum_{i=1}^n \frac{w_{i_{\text{weak}}^{(j)}}}{N_{\text{train}} \cdot n} \text{BCE}_w(\hat{c}_i^{(j)}, c_{i_{\text{weak}}}^{(j)}, \bar{c}_{\text{train}_{\text{weak}}}) \quad (10)$$

Furthermore, the mean label $\bar{c}_{\text{train}_{\text{weak}}}$ is also weighted by the confidence score:

$$\bar{c}_{\text{train}_{\text{weak}}} = \frac{\sum_{j=1}^{N_{\text{train}}} \sum_{i=1}^n c_{i_{\text{weak}}}^{(j)} \cdot w_{i_{\text{weak}}}^{(j)}}{\sum_{j=1}^{N_{\text{train}}} \sum_{i=1}^n w_{i_{\text{weak}}}^{(j)}}$$

f_{weak} works as an ensemble of separate labeling functions that predict whether a location point is part of a stay or

not. All labeling functions implement simple heuristics and may conflict with each other. Here, they predict a pair $(c_{\text{weak}}, w_{\text{weak}}(g))$ with a constant value for c_{weak} and a confidence weight w_{weak} depending on the input data g . In total, four different functions are defined:

- f_{build} predicts a stay with high confidence if a location lies within a building.
- f_{am} predicts a stay with high confidence if a location lies within small amenities.
- f_{street} predicts a non-stay with high confidence if a location is close to a street.
- $f_{\text{transport}}$ predicts a non-stay based on available transportation mode labels.

The data source for the first three functions is OSM. Similar to [19], the coordinates of the location g_i are used to query additional information from the map service.

1) *Stay Labeling Functions*: f_{build} checks, if g_i lies within a building $b \in \mathcal{B}_{\text{OSM}}$. If so, it returns a confidence weight of 1, since points that fall inside a building have a high chance of being part of a stay:

$$w_{\text{build}}(g_i) = \begin{cases} 1, & \text{if } \exists b \in \mathcal{B}_{\text{OSM}} \mid b \cap g_i \neq \{\} \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

Similarly, f_{am} returns $w_{\text{am}} > 0$, if g_i lies within an amenity $a \in \mathcal{A}_{\text{OSM}}$. This is an OSM category for facilities like hospitals or airports that can encapsulate multiple buildings. For larger amenities, since it is less certain that people will stay in a single location, we model the confidence weight as a function of their geographic area:

$$w_{\text{am}}(g_i) = \begin{cases} \max_{a \in \mathcal{A}_{\text{OSM}} \mid a \cap g_i \neq \{\}} \exp\left(-\frac{\text{area}(a)}{|\mathcal{A}_{\text{OSM}}| \sum_j \text{area}(\mathcal{A}_{\text{OSM}_j})}\right) & \text{if } \exists a \mid a \cap g_i \neq \{\} \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

Here, the fraction is the ratio of the area of an encapsulating amenity to the average area of all amenities in the dataset. By using a negative exponent, the weight starts at 1 for an area of 0 and decreases as the area increases. We use the maximum value, when multiple amenities encapsulate g_i .

2) *Non-Stay Labeling Functions*: f_{street} checks if g_i is near a street, since those points have a high probability of being part of a non-stay. Thus, if a street $s \in \mathcal{S}_{\text{OSM}}$ intersects with a centered bounding box $bb_i(l_s)$ around g_i , f_{street} returns a confidence weight of 1. The box has a shape of $d(l_s) \times d(l_s)$ with l_s denoting the importance level of the street s (e.g. highway). Formally, this can be defined as:

$$w_{\text{street}}(g_i) = \begin{cases} 1, & \text{if } \exists s \in \mathcal{S}_{\text{OSM}} \mid s \cap bb_i(l_s) \neq \{\} \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

$f_{\text{transport}}$ is designed for the GeoLife (GL) dataset, which will be introduced in Chapter V in more detail. Its trajectories include additional transportation mode labels, which are: walking, running, biking, motorcycle, car, taxi, bus, train, subway, boat, and airplane. While the dataset lacks a separate stay mode, we can use all motorized modes (i.e., all modes

except walking, running, and biking) as a heuristic for non-stays. The confidence weight is formalized as:

$$w_{\text{transport}}(g_i) = \begin{cases} 1, & \text{if } \text{label}(g_i) \in \text{modes}_{\text{motorized}} \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

3) *Combining the Labeling Functions:* We combine the results of all heuristics \mathcal{H} by averaging the predicted probabilities and adding up all confidence weights as follows:

$$f_{\text{weak}}(g_i) = (c_{i_{\text{weak}}}, w_{i_{\text{weak}}}) = \left(\frac{\sum_{j \in \mathcal{H}} c_j \cdot w_j(g_i)}{\sum_{j \in \mathcal{H}} w_j(g_i)}, \sum_{j \in \mathcal{H}} w_j(g_i) \right), \quad (15)$$

where $\mathcal{H} = \{\text{build}, \text{am}, \text{street}, \text{transport}\}$.

Thus, each labeling function f_j has a linear influence on the total confidence weight $w_{i_{\text{weak}}}$, independent of the output of other labeling functions. On the one hand, this combination is similar to an ensemble with model averaging, whereas, on the other hand, this resembles also a Mixture of Experts, where the weights depend on the input g_i .

C. Self-Supervised Encoder

Many points are not captured by any heuristic and receive a total confidence weight of 0. Self-supervised learning (SSL) could still leverage those data in a (weakly) semi-supervised manner and further strengthen the model’s robustness to inaccurate training data [24]. Since [25], [26] show good results by using forecasting as a pretext task for time series data, we adopted their approach.

1) *Forecasting Task:* We choose the velocity as one forecast target, which is less dependent on the sample rate compared to the location. Additionally, the bearing angle is forecasted as a second target because it is not directly included in the input and requires the model to encode more informative embeddings. More specifically, we predict the sine and cosine values of the angle to capture the periodicity.

Given the encoder output emb , we concatenate its sequence mean \overline{emb} and last embedding vector emb_n as an aggregated embedding vector $emb_{\text{agg}} \in \mathbb{R}^{2d_{\text{model}}}$ for the whole sequence. This vector is then passed to two separate feedforward layers. No activation function is used for the velocity, while for the sine and cosine prediction we apply tanh to bind the output between -1 and 1:

$$\hat{v}_{n+1} = emb_{\text{agg}} W_{\text{vel}}^T + b_{\text{vel}} \quad (16)$$

$$\begin{bmatrix} \hat{\sin}_{\alpha_{n+1}} \\ \hat{\cos}_{\alpha_{n+1}} \end{bmatrix} = \tanh(emb_{\text{agg}} W_{\text{ang}}^T + b_{\text{ang}}) \quad (17)$$

2) *Multitask Loss:* The loss for each pretext task is defined by the Mean Squared Error (MSE) between the prediction and the ground truth:

$$\text{MSE}(\hat{y}, y) = \frac{1}{N_{\text{train}}} \sum_{j=1}^{N_{\text{train}}} \|\hat{y}_{n+1}^{(j)} - y_{n+1}^{(j)}\| \quad (18)$$

TABLE I: Summary of weak labels derived from GL dataset.

weak label	heuristic	total sum of confidence weights
stays ($c_{\text{weak}} = 1$)	building	1.0 M
	amenity	0.2 M
non-stays ($c_{\text{weak}} = 0$)	street	5.0 M
	transport	2.7 M

$$\mathcal{L}_{\text{vel}} = \text{MSE}(\hat{v}, v) \quad (19)$$

$$\mathcal{L}_{\text{ang}} = \text{MSE}(\hat{\sin}_{\alpha}, \sin \alpha) + \text{MSE}(\hat{\cos}_{\alpha}, \cos \alpha) \quad (20)$$

We follow [25] and approach SSL as multitask learning with the sum of the downstream loss $\mathcal{L}_{\text{weak}}$ and the pretext losses:

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{weak}} + \lambda_{\text{vel}} \mathcal{L}_{\text{vel}} + \lambda_{\text{ang}} \mathcal{L}_{\text{ang}}, \quad (21)$$

where λ denotes tunable hyperparameters. If ground truth labels are available, $\mathcal{L}_{\text{weak}}$ is replaced with $\mathcal{L}_{\text{super}}$.

V. DATA

For this study, we select two datasets: GeoLife (GL) by [6], [12], [27] and ExtraSensory (ES) by [28]. GL contains two orders of magnitude more location points than ES but lacks proper SR labels. ES is chosen because of its activity labels, from which we can infer ground truth SR labels.

Similar to [29] and [30], we remove outliers based on unrealistic velocity values and split a user’s trajectory if the time difference between two consecutive points exceeds 20 minutes or if an unrealistic location jump is detected.

A. GeoLife

Instead of ground truth SR labels, the GL dataset contain time-segmented transportation mode labels from 69 of all 182 participants. These labels are used to derive weak labels and for our experiment on TMD.

To reduce network traffic and memory, we gather OSM data for points that fall within the 15%/85% percentile of longitude and latitude, which covers about 62% of the total dataset. An overview of the total sum of confidence weights used for weak supervision can be found in Table I. The remaining unlabeled data is still used for SSL instead. The sample rate of GL is non-constant and varies between 1 and 6 seconds. For the UTM projection, we choose zone 50N.

B. ExtraSensory

We use the ES dataset to fine-tune and evaluate DeepStay. Besides GNSS points, this dataset contains other sensor data, which we ignore. It was collected for the task of activity recognition. Participants should self-report their current activities such as “biking” or “watching TV”. Some activity modes clearly indicate stays and non-stays. Thus, we define a function that maps these modes to SR labels.

In the second step, we remove suspicious stays, where the velocity is higher than the average velocity of non-stays. The final number of points and derived labels are listed in Table II. The sample rate of ES is nearly constant at $\frac{1}{\text{min}}$. To achieve reasonable results with an encoder pre-trained

TABLE II: Summary of the cleaned ES dataset.

	total number
GNSS points	306 k
stays ($c_i = 1$)	223 k
non-stays ($c_i = 0$)	28 k

on GL, we linearly interpolate the location trajectory at a rate of 0.5 Hz. However, for the final test results, only the predictions for the real, non-interpolated labels are evaluated. The prediction value is taken from the nearest interpolation point. For the map projection, we use the UTM zone 11N.

VI. EXPERIMENTS

In the first experiment, we train and test DeepStay on SR labels. The second experiment shows the ability of our architecture to be used for the more general task of TMD.

A. Experiment 1: Stay Region Extraction

For this experiment, DeepStay is pre-trained on weak labels from the GL dataset and then fine-tuned and tested together with traditional baselines on the ES dataset, where it achieves the best overall results among all methods.

1) *Baselines*: We implement the following algorithms as baselines and test them on the ES dataset:

- **Kang et al. [4]**: Threshold-based clustering. It collects consecutive points until a distance threshold to the points' centroid is exceeded. Then the time criterion 2 is checked and if the minimum duration is reached, the collected points form a SR. Although the authors only proposed a POI extraction algorithm, it also implicitly incorporates SR extraction, which can be outsourced.
- **D-Star [17]**: Density-based clustering. It is based on DBSCAN, but instead of solely clustering the location points spatially, it considers only neighboring points along the trajectory and tries to exclude outliers. D-Star seems to be state-of-the-art.
- **CB-SMoT [16]**: Density-based clustering. While the algorithm is similar to D-Star, the resulting SRs contain only consecutive points, which is more in line with our definition. It can incorporate prior known POIs. However, for a fair comparison, we exclude this data.

We optimize the hyperparameters of Kang et al. and CB-SMoT using a 3×3 grid search based on the values reported in the original publications. D-Star has 4 parameters to adjust, hence we perform a random search with 10 different constellations. Each parameter search is incorporated in a 5-fold cross-validation based on the F_1 score. We split the ES data in the same way as for DeepStay.

2) *Training, Validation, and Test*: The training and testing pipeline for DeepStay can be summarized in three steps:

- 1) **Hyperparameter optimization**: Training on about 80 % of the GL dataset with weak labels and optimization of hyperparameters on the remaining 20 % in respect to the loss $\mathcal{L}_{\text{weak}}$. These hyperparameters are the number of training epochs, the weight decay, the learning rate, and the SSL weights λ_{vel} and λ_{ang} .

TABLE III: Final results on the ES dataset.

Method	F_1	Acc.	Precision	Recall
DeepStay (ours)	0.788	0.954	0.822	0.757
D-Star [17]	0.753	0.951	0.877	0.660
CB-SMoT [16]	0.548	0.909	0.619	0.491
Kang et al. [4]	0.453	0.796	0.325	0.748
constant $\hat{c}_i = 0$	0.203	0.113	0.113	1.000
constant $\hat{c}_i = 1$	0.000	0.887	-	0.000

- 2) **Pre-training**: Creating a pre-trained DeepStay model by reinitiating the training on the full GL dataset and using the best-known hyperparameters.
- 3) **Fine-tuning and test**: Fine-tuning the decoder of the pre-trained model on the ES dataset and freezing all other model weights including the encoder layers. We apply 5-fold cross-validation, i.e. each iteration about 80 % of the data is used for training, and validation and 20 % for testing. Of this 80 %, 10 % is used for a second hyperparameter optimization.

We follow [31] and split the data by the participants of the respective study, to avoid leakage between training, validation and test set. During both the pre-training and the fine-tuning, we apply an Adam optimizer [32] and SSL.

3) *Metrics*: A common metric in time series segmentation is the pointwise accuracy, i.e. the ratio of correctly classified points to the total number of labels.

In addition, we measure the pointwise calculated recall and precision. The definition of the positive class is crucial for both metrics. Since the final test dataset, i.e. ES, is highly imbalanced and contains many more stays than non-stays (see Table II), it is more important to detect a non-stay than a stay. This also resembles everyday life, where people mostly stay in one place and only move from time to time. Therefore, we choose non-stays as the positive class. The derived F_1 score is used as the main metric to evaluate all SR extraction algorithms.

4) *Results*: The final results are shown in Table III. All reported values are calculated over all 5 ES test data splits. In addition to the three baselines, two simplistic baselines predict a constant value (either always non-stay $\hat{c}_i = 0$ or always stay $\hat{c}_i = 1$).

DeepStay achieves higher overall scores than all implemented baselines, while the results for D-Star are comparable in terms of accuracy. Kang et al. use an approach with hard thresholds, which seems to be disadvantageous compared to a density-based approach. Even though CB-SMoT achieves relatively high accuracy, its F_1 score is significantly worse than the similar D-Star algorithm. This may be due to the missing outlier detection in CB-SMoT.

5) *Ablation Study*: We compare the contribution of different training components in Table IV, where we analyze the effect of training DeepStay first without any SSL and second without any pre-training, i.e. solely trained on the ES dataset. For the latter, the original sample rate of $\frac{1}{\text{min}}$ was used instead of interpolation. In addition to the previous metrics, we also measure the area under the PR curve (PR-AUC). It can be seen that the effect of SSL is relatively

TABLE IV: Ablation study of DeepStay tested on ES.

Method	F_1	Acc.	PR-AUC	Precision	Recall
DeepStay (full)	0.788	0.954	0.821	0.822	0.757
w/o SSL	0.780	0.953	0.809	0.837	0.729
w/o pre-training	0.557	0.850	0.787	0.418	0.838

small. However, the pre-training has a significant impact on the performance, showing that the model correctly handles the noise coming from the weak labels and learns reasonable latent representations of the SRs.

B. Experiment 2: Transportation Mode Detection

To further demonstrate the broader applicability of DeepStay and to contribute our findings to a broader field of research, we apply the same encoder for TMD.

There has been some work on transformer-based TMD for data other than trajectories, such as accelerometer, gyroscope, and magnetometer data [33]. However, these sensors are sampled at a much higher rate (> 20 Hz) and thus the input sequences cover only a few seconds. In this case, the transportation mode is mostly constant, so the segmentation part is dropped from the TMD task and only the classification part remains.

For TMD *from location trajectories*, sequences typically cover several minutes and therefore mode changes are likely to occur. Nevertheless, most of the related work presupposes a correct segmentation and simply classifies each of the segments as one of the available modes [22]. This is problematic for real-world applications, where a correct segmentation is never given in advance. Here, the advantage of using the transformer encoder is the joint segmentation and classification of transportation modes by simply predicting the pointwise class probabilities and grouping consecutive predicted points of the same modes together. The baseline model SECA [29] may be the state-of-the-art approach of those models that segment *and* classify from raw GNSS trajectories. Although their published code lacks the segmentation part, we compare their self-reported results on the GL dataset with our own results and show that our approach significantly outperforms SECA.

1) *Model Adaptations*: The only adaptations made to DeepStay are in the decoder and the supervision. The decoder’s weights are expanded and a softmax activation predicts the pointwise probability $\hat{c}_{i,m}$ for each of the $M = 5$ transportation modes:

$$\hat{c}_{i,m} = \text{softmax}(emb_i W_{d'}^T + b_{d'})_m, \quad \hat{c} \in [0, 1]^{n \times M} \quad (22)$$

Now BCE_w in 8 is replaced by the weighted cross entropy (CE_w) in 23 between prediction and ground truth $c_{i,m}$, where $\bar{c}_{\text{train},m}$ denotes the percentage of labels of the m -th class within the training set. For segmentation, we can simply group consecutive points with the same most probable class.

$$\text{CE}_w(\hat{c}_i, c_i, \bar{c}_{\text{train}}) = - \sum_{m=1}^M c_{i,m} \frac{\log(\hat{c}_{i,m})}{M \cdot \bar{c}_{\text{train},m}} \quad (23)$$

TABLE V: Total number of GNSS points after preprocessing.

	unlabeled (Training)	labeled (Training)	labeled (Test)
This work	16.29 M	3.74 M	4.76 M
SECA [29]	15.43 M	4.76 M	4.76 M

TABLE VI: Final results for TMD tested on the GL dataset.

Method	F_1	Acc.
DeepStay (ours)	0.830	0.831
SECA with ground truth segments	0.764	0.768
SECA with predicted segments	0.717	0.721

2) *Baseline and Comparable Datasets*: The SECA model [29] is used as the only baseline. The authors perform a change point search by using the PELT method [34] to first segment the trajectory. Second, they use a convolutional neural network (CNN) to predict the mode of each segment and integrate an autoencoder for semi-supervision.

We compare the size of the dataset after our own preprocessing with that of SECA in Table V. It shows that we train DeepStay with significantly fewer labels compared to SECA. However, in total more unlabeled data is available. Overall, the test sets are quite similar, which allows us to compare the final results of DeepStay and SECA.

3) *Training, Validation, and Test*: Both SECA and DeepStay are trained semi-supervised. While SECA uses an autoencoder, DeepStay applies SSL (see Section IV-C).

Unlike in Experiment 1, we randomly assign each sequence seq to one of the training or test sets, regardless of the including participants, to match the setup of SECA. We also apply 5-fold cross-validation. In addition, 20% of the training data is used to adjust the same hyperparameters as in Experiment 1. We optimize our model using Adam [32].

4) *Results*: We report the weighted F_1 score and the accuracy. This F_1 score is the average of the per-class F_1 scores weighted by the number of labels per class. SECA is performing segmentwise classification and DeepStay pointwise classification, thus the following results are not fully comparable.

Nevertheless, the final results in Table VI clearly demonstrate the significant performance improvement of DeepStay. A major reason may be the pointwise predictions, which do not require a prior segmentation, but intrinsically segment the data for the classification task. However, even when SECA is given ground truth segments, DeepStay still achieves better results. One reason may be that, unlike SECA, the input sequence for our model is not limited to a single transportation mode, i.e., it can also learn the transition between modes. E.g., it is intuitively more likely to see a transition from bus to train than from bus to car. Furthermore, the autoencoder in SECA only tries to reconstruct the trajectory, while SSL can provide proxy labels for DeepStay, which may be more informative. In addition, the transformer model with its attention mechanism seems to be superior in comparison to the CNN layers for this task.

VII. CONCLUSION AND FUTURE WORK

In this work, we show, how to derive programmatically weak labels for SR extraction and how to successfully train a transformer encoder with these data. We demonstrate the effectiveness of this model on ground truth data for SR extraction and TMD, where it outperforms state-of-the-art methods. This work should be seen as a starting point for new data-driven approaches to SR extraction and provides useful training and test data. Ideas for future work are:

1) *More data augmentation*: Instead of always training on the same sequences, all trajectories could be shifted by a number of points in each epoch. This results in slightly different sequences and SSL targets and reduces overfitting.

2) *Modeling dependencies*: We treat all labeling functions independently, although there are clear dependencies. E.g., w_{build} correlates strongly with w_{am} , because buildings are often part of amenities. Other work suggests that the performance benefits significantly from incorporating these dependencies [35].

3) *Pre-training on multiple datasets*: In this study, we stick with the GL dataset for pre-training. However, there are many public unlabeled GNSS trajectory datasets. All of them could be weakly labeled with our approach and carefully combined to have an even larger training set.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [2] Y. Zheng, Y. Chen, Q. Li, X. Xie, and W.-Y. Ma, "Understanding transportation modes based on gps data for web applications," *ACM Transactions on the Web (TWEB)*, vol. 4, no. 1, pp. 1–36, 2010.
- [3] J. Suzuki, Y. Suhara, H. Toda, and K. Nishida, "Personalized visited-poi assignment to individual raw gps trajectories," *ACM Transactions on Spatial Algorithms and Systems (TSAS)*, vol. 5, no. 3, pp. 1–28, 2019.
- [4] J. H. Kang, W. Welbourne, B. Stewart, and G. Borriello, "Extracting places from traces of locations," *ACM SIGMOBILE Mobile Computing and Communications Review*, vol. 9, no. 3, pp. 58–68, 2005.
- [5] Q. Li, Y. Zheng, X. Xie, Y. Chen, W. Liu, and W.-Y. Ma, "Mining user similarity based on location history," in *Proceedings of the 16th ACM SIGSPATIAL international conference on Advances in geographic information systems*, 2008, pp. 1–10.
- [6] Y. Zheng, L. Zhang, X. Xie, and W.-Y. Ma, "Mining interesting locations and travel sequences from gps trajectories," in *Proceedings of the 18th international conference on World wide web*, 2009, pp. 791–800.
- [7] T. Andrade, B. Cancela, and J. Gama, "Discovering locations and habits from human mobility data," *Annals of Telecommunications*, vol. 75, no. 9, pp. 505–521, 2020.
- [8] M. Umair, W. S. Kim, B. C. Choi, and S. Y. Jung, "Discovering personal places from location traces," in *16th International Conference on Advanced Communication Technology*. IEEE, 2014, pp. 709–713.
- [9] T. Bhattacharya, L. Kulik, and J. Bailey, "Extracting significant places from mobile user gps trajectories: a bearing change based approach," in *Proceedings of the 20th International Conference on Advances in Geographic Information Systems*, 2012, pp. 398–401.
- [10] M. Pavan, S. Mizzaro, I. Scagnetto, and A. Beggiato, "Finding important locations: A feature-based approach," in *2015 16th IEEE International Conference on Mobile Data Management*, vol. 1. IEEE, 2015, pp. 110–115.
- [11] A. Thomason, N. Griffiths, and V. Sanchez, "Identifying locations from geospatial trajectories," *Journal of Computer and System Sciences*, vol. 82, no. 4, pp. 566–581, 2016.
- [12] Y. Zheng, Q. Li, Y. Chen, X. Xie, and W.-Y. Ma, "Understanding mobility based on gps data," in *Proceedings of the 10th international conference on Ubiquitous computing*, 2008, pp. 312–321.
- [13] M. Pavan, S. Mizzaro, and I. Scagnetto, "Mining movement data to extract personal points of interest: A feature based approach," in *Information Filtering and Retrieval*. Springer, 2017, pp. 35–61.
- [14] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *kdd*, vol. 96, no. 34, 1996, pp. 226–231.
- [15] M. Ankerst, M. M. Breunig, H.-P. Kriegel, and J. Sander, "Optics: Ordering points to identify the clustering structure," *ACM Sigmod record*, vol. 28, no. 2, pp. 49–60, 1999.
- [16] A. T. Palma, V. Bogorny, B. Kuijpers, and L. O. Alvarez, "A clustering-based approach for discovering interesting places in trajectories," in *Proceedings of the 2008 ACM symposium on Applied computing*, 2008, pp. 863–868.
- [17] K. Nishida, H. Toda, and Y. Koike, "Extracting arbitrary-shaped stay regions from geospatial trajectories with outliers and missing points," in *Proceedings of the 8th ACM SIGSPATIAL International Workshop on Computational Transportation Science*, 2015, pp. 1–6.
- [18] D. Ashbrook and T. Starner, "Using gps to learn significant locations and predict movement across multiple users," *Personal and Ubiquitous computing*, vol. 7, no. 5, pp. 275–286, 2003.
- [19] X. Cao, G. Cong, and C. S. Jensen, "Mining significant semantic locations from gps data," *Proceedings of the VLDB Endowment*, vol. 3, no. 1-2, pp. 1009–1020, 2010.
- [20] G. Stylianou, "Stay-point identification as curve extrema," *arXiv preprint arXiv:1701.06276*, 2017.
- [21] R. B. Langley, "The utm grid system," *GPS world*, vol. 9, no. 2, pp. 46–50, 1998.
- [22] H. Moreau, A. Vassilev, and L. Chen, "The devil is in the details: An efficient convolutional neural network for transport mode detection," *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- [23] A. J. Ratner, C. M. De Sa, S. Wu, D. Selsam, and C. Ré, "Data programming: Creating large training sets, quickly," *Advances in neural information processing systems*, vol. 29, 2016.
- [24] D. Hendrycks, M. Mazeika, S. Kadavath, and D. Song, "Using self-supervised learning can improve model robustness and uncertainty," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [25] S. Jawed, J. Grabocka, and L. Schmidt-Thieme, "Self-supervised learning for semi-supervised time series classification," in *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, 2020, pp. 499–511.
- [26] S. Tipirneni and C. K. Reddy, "Self-supervised transformer for sparse and irregularly sampled multivariate clinical time-series," *ACM Transactions on Knowledge Discovery from Data (TKDD)*, vol. 16, no. 6, pp. 1–17, 2022.
- [27] Y. Zheng, X. Xie, W.-Y. Ma *et al.*, "Geolife: A collaborative social networking service among user, location and trajectory," *IEEE Data Eng. Bull.*, vol. 33, no. 2, pp. 32–39, 2010.
- [28] Y. Vaizman, K. Ellis, and G. Lanckriet, "Recognizing detailed human context in the wild from smartphones and smartwatches," *IEEE pervasive computing*, vol. 16, no. 4, pp. 62–74, 2017.
- [29] S. Dabiri, C.-T. Lu, K. Heaslip, and C. K. Reddy, "Semi-supervised deep learning approach for transportation mode identification using gps trajectory data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 5, pp. 1010–1023, 2019.
- [30] Y. Zheng, L. Liu, L. Wang, and X. Xie, "Learning transportation mode from raw gps data for geographic applications on the web," in *Proceedings of the 17th international conference on World Wide Web*, 2008, pp. 247–256.
- [31] M. Etemad, "Transportation modes classification using feature engineering," *arXiv preprint arXiv:1807.10876*, 2018.
- [32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [33] Y. Tian, D. Hettiarachchi, and S. Kamijo, "Transportation mode detection combining cnn and vision transformer with sensors recalibration using smartphone built-in sensors," *Sensors*, vol. 22, no. 17, p. 6453, 2022.
- [34] R. Killick, P. Fearnhead, and I. A. Eckley, "Optimal detection of changepoints with a linear computational cost," *Journal of the American Statistical Association*, vol. 107, no. 500, pp. 1590–1598, 2012.
- [35] P. Varma, B. D. He, P. Bajaj, N. Khandwala, I. Banerjee, D. Rubin, and C. Ré, "Inferring generative model structure with static analysis," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.