

# Sparse Prototype Network for Explainable Pedestrian Behavior Prediction

Yan Feng<sup>\*1</sup>, Alexander Carballo<sup>2,3,4</sup>, and Kazuya Takeda<sup>1,2,4</sup>

**Abstract**—Predicting pedestrian behavior is challenging yet crucial for applications such as autonomous driving and smart city. Recent deep learning models have achieved remarkable performance in making accurate predictions, but they fail to provide explanations of their inner workings. One reason for this problem is the multi-modal inputs. To bridge this gap, we present Sparse Prototype Network (SPN), an explainable method designed to simultaneously predict a pedestrian’s future action, trajectory, and pose. SPN leverages an intermediate prototype bottleneck layer to provide sample-based explanations for its predictions. The prototypes are modality-independent, meaning that they can correspond to any modality from the input. Therefore, SPN can extend to arbitrary combinations of modalities. Regularized by mono-semanticity and clustering constraints, the prototypes learn consistent and human-understandable features and achieve state-of-the-art performance on action, trajectory and pose prediction on TITAN and PIE. Finally, we propose a metric named Top-K Mono-semanticity Scale to quantitatively evaluate the explainability. Qualitative results show the positive correlation between sparsity and explainability. Code available at <https://github.com/Equinoxxxx/SPN>.

## I. INTRODUCTION

Predicting pedestrian behavior in complex environments is a critical task for autonomous systems, with applications ranging from self-driving vehicles to intelligent surveillance. While traditional models have focused on predicting either the future trajectory, action, or pose of pedestrians, there has been a growing demand for models that can simultaneously predict pedestrian behaviors in different types, including trajectory, action class and pose.<sup>1</sup>

Recent efforts in applying deep learning methods to address the prediction of any one type of behavior are proven effective [1]–[4]. However, only a few works discussed joint prediction of multiple types of behaviors [5], [6]. Moreover, most of these works lack the mechanism of revealing their inner workings during inference, which causes additional testing costs when the model faces unseen scenarios, decreases the trustworthiness, and hinders developers and researchers from making further improvements. A major challenge that causes the problem is multi-modal inputs. Pedestrian behaviors can be inferred from various

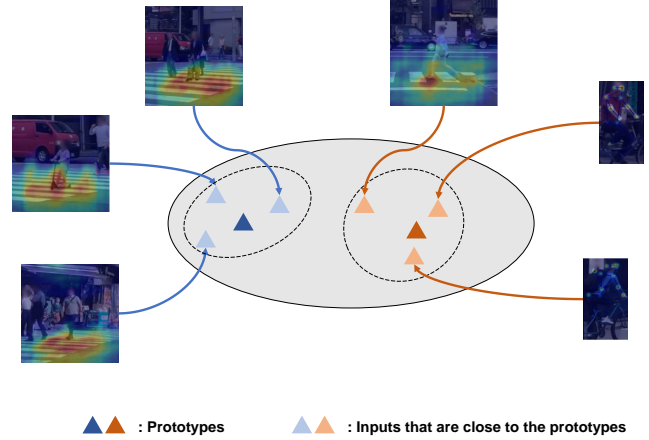


Fig. 1. Modality-independent prototypes. Multi-modal inputs are projected in a joint latent space, and the predictions are made based on their distance from the prototypes. Each prototype is represented by samples that are the closest to it.

clues, including the historical trajectories, poses, contextual elements, etc. Although many of these clues can originate from visual inputs, the lack of annotated data and the high dimensionality of raw inputs make it hard for neural networks to learn salient features directly from images and videos. Therefore, disentangling different types of data and regarding them as multi-modalities became a common practice [1], [2], [4], [5], [7]. However, most explaining techniques are modality-specific [8] or architecture-specific [9]–[12], which limits the scalability of the model. To bridge this gap, we present Sparse Prototype Network (SPN), a prototype-based framework designed to jointly predict multiple types of pedestrian behaviors, and provide explanations of its inferences based on the distance between the prototypes it learns and the input. The method is inspired by the idea that, multi-modalities derived from the same observation can be regarded as different fragments. By mapping the fragments into a joint latent space, a prototype vector can match any one of them, and is thus modality-independent, which enables the model to extend to arbitrary combinations of modalities. The modalities in the training data with the closest distance to a prototype are selected to represent that prototype. Figure ?? briefly illustrates the inner working of SPN.

In the wider range of literature on explainable AI (XAI), one major challenge in developing explainable models is poly-semanticity, i.e. a unit of explaining (e.g. a neuron or a vector) activates on multiple semantics in the input. Since SPN uses prototype vectors to explain its inference, it also

<sup>1</sup>Graduate School of Informatics, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, Aichi 464-8601, Japan.

<sup>2</sup>Institutes of Innovation for Future Society, Nagoya University, Furo-cho, Chikusa-ku, Nagoya, Aichi 464-8601, Japan.

<sup>3</sup>Faculty of Engineering and Graduate School of Engineering, Gifu University, 1-1 Yanagido, Gifu City, 501-1193, Japan

<sup>4</sup>Tier IV Inc., Nagoya University Open Innovation Center, 1-3, Mei-eki 1-chome, Nakamura-ward, Nagoya, 450-6610, Japan

<sup>1</sup>In this paper, the term “behavior” refers to pedestrian trajectory, action class and pose unless otherwise specified.

suffers from poly-semanticity. We argue that mono-semantic prototypes should activate sparsely on the whole dataset, since features that human understand are sparse [13], [14]. Motivated by this insight, we propose a sparsity loss term to encourage the prototypes to capture only sparse features, and thus to improve the mono-semanticity. Moreover, we propose a quantitative metric, named Top-K Mono-semanticity Scale, to evaluate the explainability of the prototypes without the influence of subjectivity brought by human annotators.

The contributions of this paper are as follows:

- We introduce SPN, a novel approach that simultaneously predicts pedestrian action, trajectory and pose while offering explainable predictions based on modality-independent prototypes. Therefore, SPN can be extended to arbitrary modalities.
- We propose a metric, the Top-K Mono-semanticity Scale, to quantitatively evaluate the explainability of the learned prototypes.
- We apply a complementary combination of regularization losses to the prototypes. The sparsity term encourages the prototypes to capture mono-semantic features, while the clustering term provides the prototypes with the commonality between different modalities.
- We validate SPN’s performance on two popular pedestrian behavior prediction datasets, TITAN and PIE, and show that it achieves state-of-the-art performance while maintaining a high level of explainability.
- We made the code available at <https://github.com/Equinoxxxxxx/SPN>.

## II. RELATED WORKS

### A. Pedestrian Behavior Prediction

Traditional approaches for predicting pedestrian actions and trajectories relied on handcrafted features and physical models, such as Kalman filters and social force models [15]. These methods, while inherently explainable, often struggle to capture the complexity of real-world pedestrian behavior, especially in crowded or dynamic environments.

More recently, the field of pedestrian behavior prediction has witnessed significant advances brought by deep learning [3]–[5], [16]–[18]. State-of-the-art works primarily employed multi-modal inputs such as appearance, motion, past trajectory, etc [1], [2], [5], [7], [18]–[20], and applied certain integration strategies, such as attention mechanisms, to make predictions. However, these models often lacked the ability to explain their predictions, limiting their applicability in safety-critical systems such as autonomous driving.

### B. Prototype-based Models

A large portion of efforts toward explainability in deep learning models have primarily focused on post-hoc explanation methods, where explanations are generated after the model has made its predictions [9]–[12]. Recently, large language models (LLMs) have been applied to generate, text-based explanations, thanks to their outstanding ability of reasoning [8], [21]. However, these post-hoc explanations are limited in their fidelity. Moreover, they are often specific to a

particular modality, such as visual data, and do not generalize well to multi-modal inputs.

On the other hand, prototype-based methods [?], [22] offer an alternative path toward inherently explainable models. These methods learn prototypes—representative features—from the data, which are then used to make predictions. While most of these methods were applied to few-shot learning problem [23]–[25], some were designed to make explainable predictions in simple tasks and small-scale datasets [22], [26]. The advantage of prototype-based approaches is that they provide sample-based explanations, i.e. using representative samples to explain the semantics of certain prototypes. However, these methods typically focus on single modalities and do not address the challenge of multi-modal explainability.

In order to bridge the above gaps for pedestrian behavior prediction, SPN leverages modality-independent prototypes, as well as a sparsity loss to promote mono-semanticity of the learned prototypes.

### C. Mono-semantic Explanation

Sample-based explanation has recently become one of the major explaining techniques of XAI methods, including prototype-based methods [22], [26]–[29] and mechanistic methods [8], [30], where the most related samples are used to represent the explaining units such as neurons and prototypes. One recent challenge faced by these methods is poly-semanticity [31], meaning that the units are represented by unrelated samples. On the contrary, mono-semantic units refer to those that are related to relatively consistent features. Although works from the field of language models [14], [32] used sample-wise metric to promote or inhibit the mono-semanticity of a single sample, there has not been a metric to evaluate the mono-semanticity of an explaining unit. Inspired by Mono-semantic Scale [32] where sparsity was used as a proxy of mono-semanticity, we propose Top-K Mono-semanticity Scale (Top-K MS) for short as a quantitative metric for explainability of prototypes.

## III. METHOD

In this section, we describe the architecture of the Sparse Prototype Network (SPN) and the explainability metric Top-K Mono-semanticity Scale (Top-K MS) used to evaluate the model. Our proposed method focuses on creating a predictive framework that not only achieves accurate pedestrian behavior prediction but also provides interpretable, human-understandable explanations for its decisions.

### A. Sparse Prototype Network

The SPCM employs multi-modal input and predicts three types of data to describe pedestrian behavior: action, trajectory, and pose. To achieve this, the model is composed of three main modules: input encoding, prototype layer, and prediction heads.

**Input encoding.** The input encoding module processes each modality independently, transforming the raw inputs into compact, high-dimensional feature vectors. Let  $\mathcal{X} =$

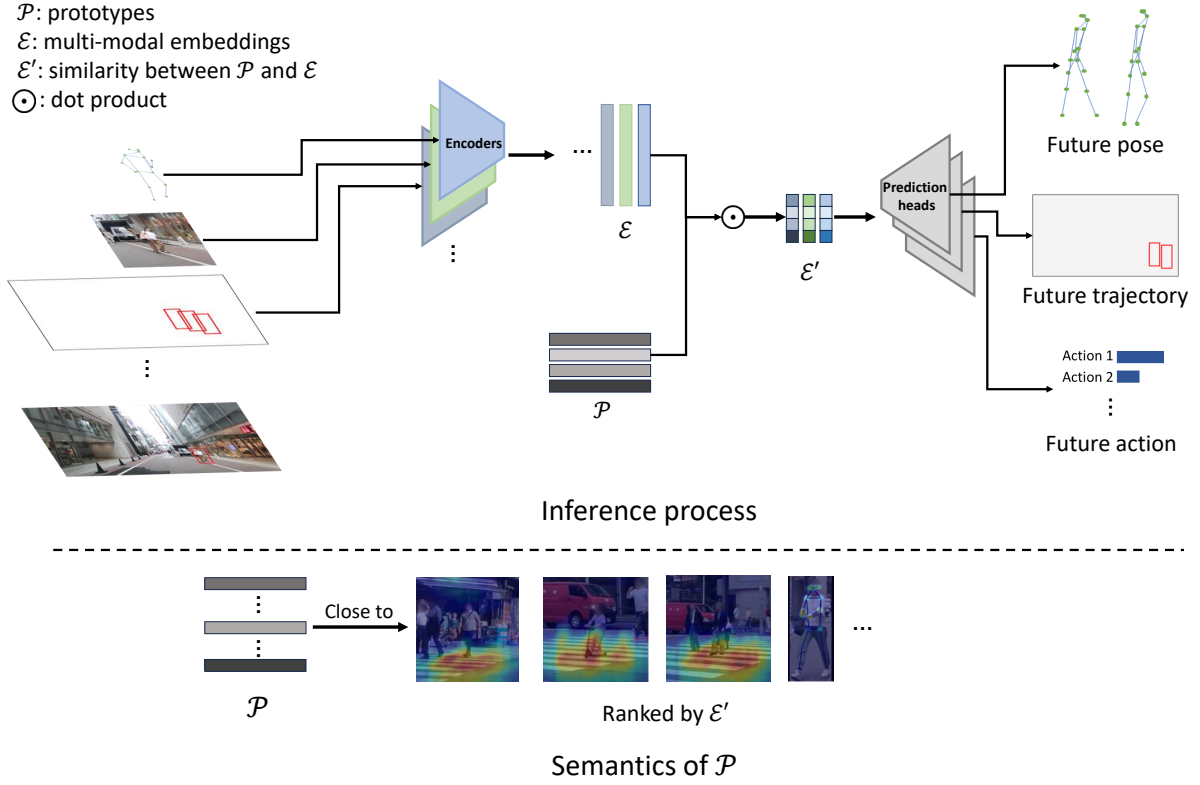


Fig. 2. Inference process of SPN and how the prototypes are explained.

$\{\mathbf{x}_m\}_{m=1}^M$  be the inputs from  $M$  modalities. The input encoding module consists of  $M$  encoders  $\mathcal{F} = \{f_m\}_{m=1}^M$ . Each encoder  $f_m$  corresponds to a modality  $\mathbf{x}_m$ , and is followed by an MLP layer to project the outputs into a joint latent space. The multi-modal embeddings are given by

$$\mathbf{e}_m = \text{ReLU}(W_e f_m(\mathbf{x}_m)) \quad (1)$$

where  $\mathcal{E} = \{\mathbf{e}_m\}_{m=1}^M$  represents the multi-modal embeddings, and  $W_e$  is the weights of the MLPs.

In our experiments, we exploit up to 5 different modalities. Their formats are as follows:

- 1) Local context: a single image cropped by enlarged bounding box around the pedestrian. Following the practice in [2], we add semantic segmentation map as additional channels to the image. The encoder is in the same architecture as in [2], which is a multi-layer 2D convolutional block.
- 2) Past pose: a sequence of the past skeletons of the pedestrian in the COCO format.
- 3) Past trajectory: a sequence of the past coordinates of bounding boxes of the target pedestrian in the image plane.
- 4) Ego motion: a sequence of acceleration of the ego vehicle.
- 5) Social relation: the relative location of other pedestrians and vehicles in the surroundings. Following the practice in [18], the format is given by

$$\mathcal{G}_k = [\log(\frac{|x_b - x_k|}{w_b}), \log(\frac{|y_b - y_k|}{h_b}), \log(\frac{w_k}{w_b}), \log(\frac{h_k}{h_b})] \quad (2)$$

where  $\mathcal{G}_k$  is the relative location of the  $k$ -th neighbor.

**Prototype layer.** In this module, the multi-modal embeddings  $\mathcal{E}$  is mapped to the similarity with a set of prototype vectors  $\mathcal{P} = \{\mathbf{p}_n\}_{n=1}^N$ , where each  $\mathbf{p}_n$  is in the same dimensions as  $\mathbf{e}_m$ . The matching process is given by

$$\mathcal{E}' = \text{ReLU}(\mathcal{P}^\top \mathcal{E}) \quad (3)$$

where  $\mathcal{E}' \in \mathbb{R}^{N \times M}$  is the matching results. The ReLU activation ensures that the similarity is non-negative and that the mapping is non-linear without losing interpretability [31]. Although there is previous work [22] applying purely linear functions to ensure strictly interpretable architecture, non-linear activation in a single layer provides more possibilities of further improvements and investigations, such as more prototypes than dimensions [31], [33].

**Task-specific prediction heads.** Given the matching  $\mathcal{E}'$  between prototypes and multi-modal embeddings, the predictions are made by a task-specific prediction head for each task. We frame trajectory prediction and pose prediction as generation task, and action prediction as classification task. For action prediction, we simply use a linear function to generate the action logits

$$\hat{\mathbf{y}}_{act} = \text{softmax}(W_{act} \mathcal{E}') \quad (4)$$

For trajectory and pose prediction, we apply a model  $g$  that supports conditioned generation and input  $\mathcal{E}'$  as the condition

that contains all past observations. The prediction is given by

$$\hat{\mathbf{y}}_{traj/pose} = g(\epsilon; \mathcal{E}') \quad (5)$$

where  $\epsilon$  is the Gaussian noise required for generation. More details about the implementation is introduced in Section IV-A.

### B. Loss Functions

The model is optimized using a combination of two groups of losses:

- **Task-specific Losses:** For each prediction task (action, trajectory, and pose), we compute the respective loss. For action prediction, we use cross-entropy loss; for trajectory prediction and pose prediction, mean squared error (MSE).
- **Prototype Regularization Losses:** Since the multi-modal inputs are fed to the prototype layer independently, no inter-modality interaction is maintained. Therefore, the underlying commonality between different modalities might be lost in the process. As compensation for this loss, a clustering term  $\mathcal{L}_{cluster}$  is applied to encourage the prototypes to capture the commonality between different modalities. Specifically, we follow the practice in [34], drawing the embeddings of the same sample together while pushing embeddings from different samples away from each other. The clustering loss is given by

$$\mathcal{L}_{cluster} = -\frac{1}{BM^2} \sum_{i=1}^B \sum_{m=1}^M \sum_{n=1}^M \log \frac{\exp(\mathbf{e}_{m,i} \cdot \mathbf{e}_{n,i} / \tau)}{\sum_{j=1}^B \exp(\mathbf{e}_{m,i} \cdot \mathbf{e}_{n,j} / \tau)} \quad (6)$$

where  $\mathbf{e}_{m,i}$  is the embedding for modality  $m$  of the  $i$ -th sample, and  $\tau$  is a temperature factor to normalize the dot product.

To avoid the multi-modal representations from collapsing to the same point, we also apply L1 loss  $\mathcal{L}_{l1}$  as a counterpart to improve the sparsity degree of  $\mathcal{E}'$ , and meanwhile as a proxy of mono-semanticity. The total loss is formulated as:

$$\mathcal{L} = \lambda_{cluster} \mathcal{L}_{cluster} + \lambda_{l1} \mathcal{L}_{l1} + \mathcal{L}_{task} \quad (7)$$

### C. Top-K Mono-semanticity Scale

Inspired by the Mono-semantic Scale [32] that uses sparsity as a proxy of mono-semanticity, we introduce the Top-K Mono-semanticity Scale to quantitatively evaluate the explainability of our model. This metric measures how well the top-K most activated are far from the average values. For each prototype, we rank the samples based on their activation strength and compute the mean relative variance of the top-K activations. Specifically, given the matching results between the prototypes and a batch of inputs of size  $B$ , the Top-K Mono-semanticity Scale (Top-K MS) is defined as

$$\psi_n = \frac{1}{K} \sum_{k=1}^K \frac{\mathcal{E}'_{n,k} - \bar{\mathcal{E}}'_n}{S^2} \quad (8)$$

where

$$\bar{\mathcal{E}}'_n = \frac{1}{BM} \sum_{i=1}^{BM} \mathcal{E}'_{n,i}, S^2 = \frac{1}{BM-1} \sum_{i=1}^{BM} (\mathcal{E}'_{n,i} - \bar{\mathcal{E}}'_n)^2 \quad (9)$$

and  $\mathcal{E}'_{n,k}$  ranked in the descending order. A high  $\psi_n$  value means that the  $n$ -th prototype has high activations on a limited number of samples, and the feature is thus sparse. Figure 3<sup>2</sup> illustrates prototypes with different degrees of Top-K MS and the samples with top activations. In the first row where the Top-K MS is high, the top samples are all in the same modality, and the highlighted regions focus on the crosswalks at a relatively far distance; in the second row, the prototype seems to focus on such a layout of agents where the target pedestrian is at the left edge of the view, and a surrounding vehicle or pedestrian on the opposite side is attended to; in the third row, the prototype focuses on a seemingly region of shadow; and in the last row, where the Top-K MS is the lowest, the prototype seems to capture both the head joints where the pedestrian is standing and a near-zero ego acceleration sequence. All in all, as the Top-5 MS values decrease, the semantics of the representative samples of a prototype tend to become less consistent and hard to understand.

## IV. EXPERIMENTS

### A. Implementation Details

The detailed settings of the multi-modal encoders we used in the experiments are as follows:

1) **Local context.** We use the backbone in [2], a multi-layer 2D convolutional block.

2) **Past pose, past trajectory, ego motion and social relation.** For these modalities, we use a one-layer transformer encoder [35] with 8 heads and 64 dimensions. For modalities with more than one sequence dimension, we flatten these dimensions together with the temporal dimension before feeding the input to the backbone.

For the prediction heads of trajectory prediction and pose prediction, we use DePOSit [36], a diffusion-based generation model originally designed for pose regression. We input  $\mathcal{E}'$  as a condition vector into DePOSit for pose and trajectory prediction.

We use Adam optimizer during training with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ . We set  $\lambda_{cluster} = 0.001$  and  $\lambda_{l1} = 0.01$  out of a random search. We use 50 prototypes and 512 dimensions for each prototype as a default setting. The model is trained for 50 epochs with a batch size of 64.

<sup>2</sup>See <https://github.com/Equinoxxxxxx/SPN> for more visualization results.

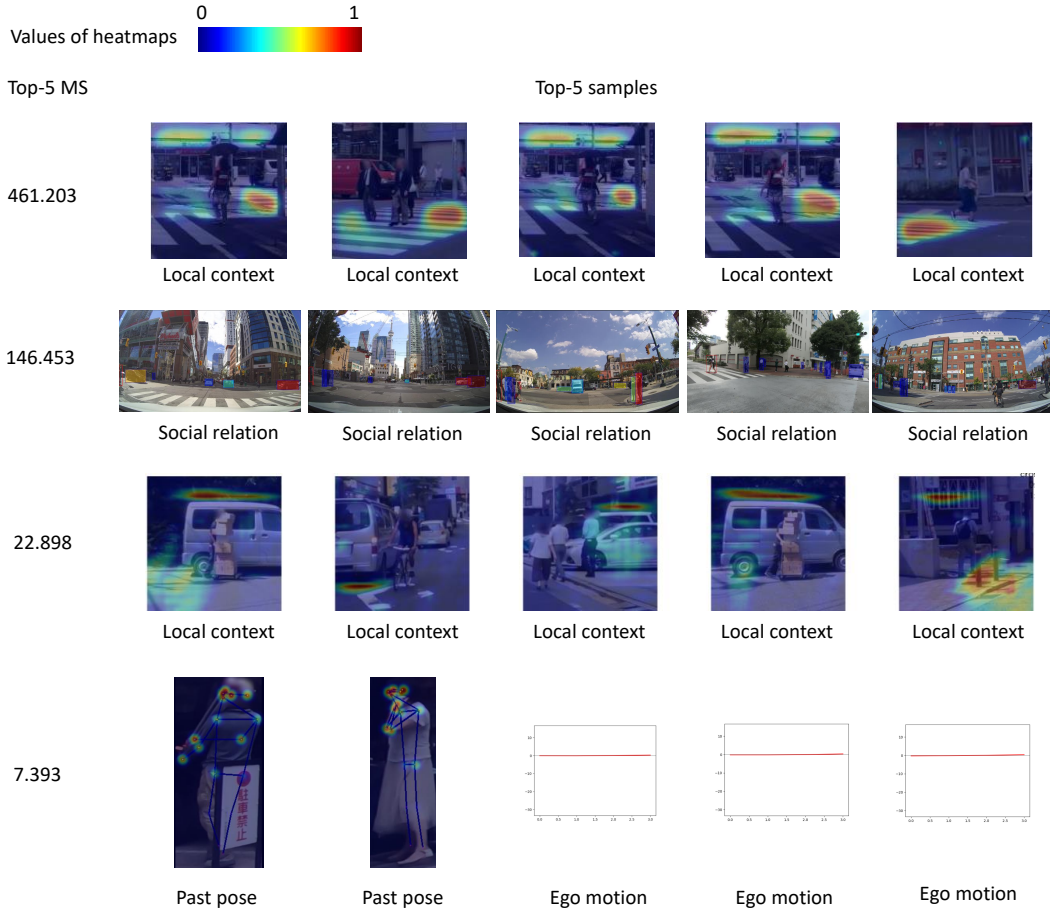


Fig. 3. Examples of prototypes with different Top-K MS. As the Top-K MS decreases, the semantics of the prototype become hard to recognize.

## B. Datasets

We evaluate SPN on PIE [5] and TITAN [37]. TITAN contains 10 hours of 60 FPS driving videos in Tokyo and multiple independent action sets at 10 Hz frequency: communicative actions, transportive actions, complex contextual actions, simple contextual actions and atomic actions. PIE is a commonly used dataset for pedestrian action and trajectory prediction. It contains 6 hours of 30 FPS driving videos in Toronto with annotations of binary crossing action labels at 30 Hz frequency. Since PIE only contains action labels of pedestrian crossing the road, we use crossing prediction to evaluate the action prediction performance.

## C. Evaluation Metrics

We evaluate SPN using standard prediction metrics. For action prediction, we use accuracy and F1-score; for trajectory and pose prediction, we use mean squared error (MSE). Since TITAN and PIE have different image sizes, we normalize the trajectory and pose coordinates by their corresponding image size, and calculate the MSE on the normalized data.

## D. Behavior Prediction Performance

We compare SPN on all three tasks with state-of-the-art methods, including SGNet [38], next [18], PCPA [1],

TABLE I  
BEHAVIOR PREDICTION RESULTS ON TITAN&PIE

Models	TITAN&PIE			
	Acc	F1	Traj MSE	Pose MSE
SGNet	-	-	1.2	-
SGNet-CVAE	-	-	1.0	-
DePOSit	-	-	-	0.16
Next	0.72	0.57	0.11	-
PCPA	<b>0.78</b>	0.45	-	-
PedGraph+	0.77	0.51	-	-
SPN(ours)	<b>0.78</b>	<b>0.59</b>	<b>0.09</b>	<b>0.07</b>

Pedestrian Graph+ [2] and DePOSit [36]. Note that we concatenate TITAN and PIE together to evaluate the ability of the models with maximum data availability. And since SGNet, DePOSit and SPN are stochastic when predicting trajectory and pose, we choose the best result out of 5 runs as the final result.

The results in Table I show that SPN outperforms other models on most metrics. Note that PCPA and Pedestrian Graph+ have high accuracy but low F1, indicating that they are biased to the dominant class, since neither TITAN or PIE are balanced on the crossing action. For trajectory prediction and pose prediction, SPN exceeds other models, especially

TABLE II  
ABLATION EXPERIMENTS ON REGULARIZATION TERMS

Models	TITAN&PIE				
	Acc	F1	Traj MSE	Pose MSE	Mean Top-5 MS
$\lambda_{cluster} = 0$ $\lambda_{l1} = 0$	<b>0.79</b>	0.55	0.095	0.067	68.91
$\lambda_{cluster} = 0.001$ $\lambda_{l1} = 0$	0.49	0.47	0.110	0.069	34.83
$\lambda_{cluster} = 0$ $\lambda_{l1} = 0.01$	0.32	0.32	0.221	0.210	78.55
$\lambda_{cluster} = 0.001$ $\lambda_{l1} = 0.01$	<b>0.79</b>	<b>0.59</b>	<b>0.093</b>	<b>0.066</b>	<b>306.03</b>

autoregressive models such as SGNet and DePOSit, to a considerable extent, indicating that the prototypes effectively encoded the rich information from multi-modal inputs and improved the prediction based on this information, despite the compactness of the prototype-based representations.

In Table II we list the results of ablation experiments on the regularization terms by removing one or both of them. We also add the Top-5 MS value as a reference of the explainability. It can be seen from the second and third row that, either the clustering term or the L1 term alone would cause training collapse, meaning the model falls to a suboptimal. Only when both terms are present, the performance continue to improve. Beside the improvement on the prediction performance, the difference brought by the combination of two loss terms can also be seen in Figure 4, where prototypes with high Top-5 MS from the two cases are shown. While the upper part shows prototypes with relatively consistent and clear semantics, the lower part, where no regularization terms are applied, shows the tendency of focusing on patterns that are hard to understand, as well as lower Top-5 MS values. The comparison shows the effect of the regularization terms in not only improving prediction but also encouraging the prototypes to focus on more meaningful semantics.

#### E. Evaluation of Partial Prototypes

To further test the functionality of the prototypes with high sparsity, we conduct an experiment to evaluate the performance of SPN with only a small portion of prototypes functioning. Specifically, we select 5 out of 50 prototypes from a trained SPN, and set the rest prototypes to 0. The 5 prototypes are selected by two different criteria:

- Top-5 MS. Prototypes with the highest Top-5 MS values are selected.
- L1 norm of the corresponding columns in  $W_{act}$ . Since  $W_{act}$  transforms the prototype matching results  $\mathcal{E}'$  into action logits, each column of  $W_{act}$  corresponds to a certain prototype, where each element represents the relation between the prototype and one specific action class. By selecting columns of  $W_{act}$  with the high L1 norm values, we are actually selecting prototypes with the highest importance to action prediction.

The results of the partial prototypes selected by the above two criteria are listed in Table III. While the performances of trajectory and pose prediction are the same,

TABLE III  
EXPERIMENTS ON PARTIAL PROTOTYPES

Selection metric	TITAN&PIE			
	Acc	F1	Traj MSE	Pose MSE
Top-5 MS	<b>0.75</b>	<b>0.58</b>	<b>0.092</b>	<b>0.066</b>
Linear weights	0.63	0.49	<b>0.092</b>	0.067

it is interesting that prototypes selected by Top-5 MS have better performances. This could indicate that prototypes with high sparsity actually contains more effective action-related information than the prototypes with high importance in  $W_{act}$ , which are supposed to be the optimal sets for action prediction. Such a phenomenon can be evidence that sparse features could contain more meaningful semantics, even though not paid much attention by the model itself.

#### V. LIMITATIONS

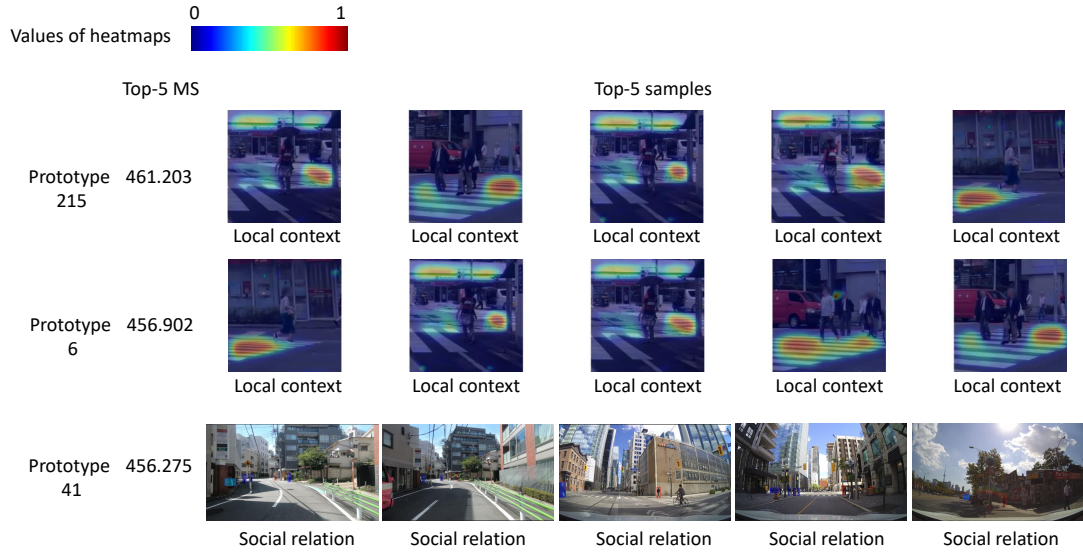
Despite the improvements in prediction performance and in the mono-semanticity of explanations, the above results also reveal limitations of SPN. First of all, mono-semanticity is only one aspect of general explainability, and applying Top-K MS and sparsity loss cannot totally diminish the subjectivity from the evaluation of explainability. Secondly, despite the fact that the classification task can be made completely interpretable by using linear functions as the prediction head, the transformation from prototypes to the generation results, i.e. trajectory and pose, still remains a black box. Although the prototypes can provide more transparent conditions than previous methods, more efforts are required to improve the explainability of the generation and regression process. Last but not least, although SPN shows progress with a relatively small number of prototypes, the potential still remains to be discovered to have more prototypes than dimensions such that the prototypes can learn further disentangled and fine-grained features from multi-modalities.

#### VI. CONCLUSIONS

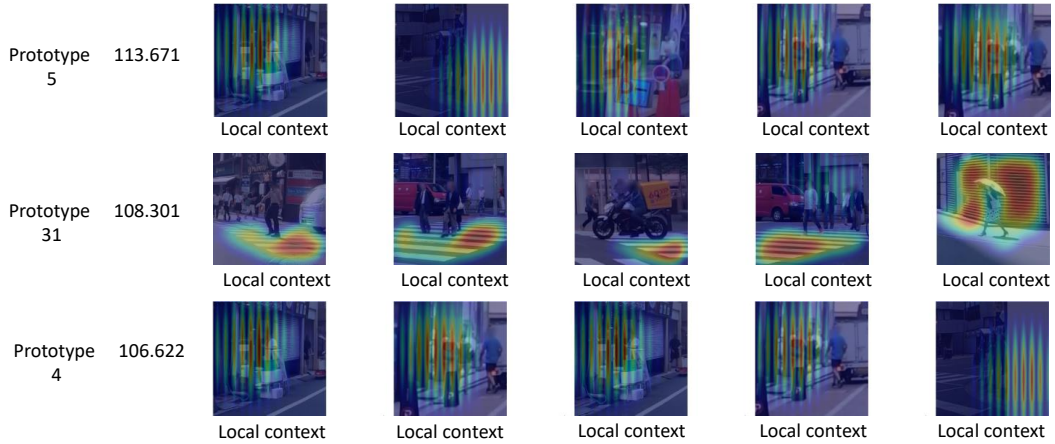
In this paper, we introduced the Sparse Prototype Network (SPN), an explainable architecture designed to address the challenge of explainable pedestrian behavior prediction. SPN surpasses existing models in three tasks, i.e. action prediction, trajectory prediction and pose prediction. By mapping multi-modal inputs to modality-independent prototypes, SPN can trace its own decisions to human-understandable concepts. To compensate for the lost commonality between modalities, we also apply a complementary combination of regularization terms, preventing the prototypes from collapsing and also encouraging the model to learn mono-semantic features.

Furthermore, we introduced the Top-K Mono-semanticity Scale, a new metric to quantitatively evaluate the explainability of not only prototype-based models, but also other explainable methods relying on sample-based explanations. This metric enables a systematic assessment of how well each prototype aligns with a single understandable concept.





(a) Prototypes when  $\lambda_{l1} = 0.01, \lambda_{cluster} = 0.001$



(b) Prototypes when  $\lambda_{l1} = 0, \lambda_{cluster} = 0$

Fig. 4. Qualitative comparison between SPN with the regularization terms and SPN without them.

In multiple experiments, both the quantitative and qualitative results demonstrate the SPN can make accurate predictions and meanwhile provide explanations of its decision-making.

Despite the strengths of SPN, there are several limitations that require further investigation, including a more comprehensive evaluation of general explainability, a more transparent generation module, and larger scales of prototypes. On top of that, our future work will focus on:

1) Extending SPN to a larger scale of data. So far, SPN only shows promising results with the TITAN and PIE datasets, due to the limitation of data with rich action labels. Future research should focus on extending SPN to larger and more diverse datasets, as well as exploring unsupervised learning to utilize data without manual labels.

2) Exploiting knowledge from well-trained models, such as LLMs. Due to the flexibility of natural language, there have been relatively mature evaluation metrics for the ex-

plainability of LLMs. Although it is difficult to directly apply these metrics to multi-modal tasks, it is still feasible and worth investigating to develop explainable multi-modal methods on top of a transparent and trustworthy language model,

## REFERENCES

- [1] I. Kotseruba, A. Rasouli, and J. K. Tsotsos, “Benchmark for evaluating pedestrian action prediction,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 1258–1268.
- [2] P. R. G. Cadena, Y. Qian, C. Wang, and M. Yang, “Pedestrian graph+: A fast pedestrian crossing prediction model based on graph convolutional networks,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 11, pp. 21 050–21 061, 2022.
- [3] A. Rasouli, I. Kotseruba, and J. K. Tsotsos, “Are they going to cross? a benchmark dataset and baseline for pedestrian crosswalk behavior,” in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 206–213.
- [4] J. Lorenzo, I. P. Alonso, R. Izquierdo, A. L. Ballardini, Á. H. Saz, D. F. Llorca, and M. Á. Sotelo, “Capformer: Pedestrian crossing action prediction using transformer,” *Sensors*, vol. 21, no. 17, p. 5694, 2021.

- [5] A. Rasouli, I. Kotseruba, T. Kunic, and J. K. Tsotsos, "Pie: A large-scale dataset and models for pedestrian intention estimation and trajectory prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6262–6271.
- [6] H. Girase, H. Gang, S. Malla, J. Li, A. Kanehara, K. Mangalam, and C. Choi, "Loki: Long term and key intentions for trajectory prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9803–9812.
- [7] B. Liu, E. Adeli, Z. Cao, K.-H. Lee, A. Shenoi, A. Gaidon, and J. C. Niebles, "Spatiotemporal relationship reasoning for pedestrian intent prediction," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3485–3492, 2020.
- [8] S. Bills, N. Cammarata, D. Mossing, H. Tillman, L. Gao, G. Goh, I. Sutskever, J. Leike, J. Wu, and W. Saunders, "Language models can explain neurons in language models," <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>, 2023.
- [9] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2921–2929.
- [10] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [11] H. Wang, Z. Wang, M. Du, F. Yang, Z. Zhang, S. Ding, P. Mardziel, and X. Hu, "Score-cam: Score-weighted visual explanations for convolutional neural networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 24–25.
- [12] S. Abnar and W. Zuidema, "Quantifying attention flow in transformers," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 4190–4197.
- [13] J. Wang, S. Di, L. Chen, and C. W. W. Ng, "Learning from emergence: A study on proactively inhibiting the monosemantic neurons of artificial neural networks," in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2024, pp. 3092–3103.
- [14] H. Yan, Y. Xiang, G. Chen, Y. Wang, L. Gui, and Y. He, "Encourage or inhibit monosemanticity? revisit monosemanticity from a feature decorrelation perspective," *arXiv preprint arXiv:2406.17969*, 2024.
- [15] D. Helbing and P. Molnar, "Social force model for pedestrian dynamics," *Physical review E*, vol. 51, no. 5, p. 4282, 1995.
- [16] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social lstm: Human trajectory prediction in crowded spaces," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 961–971.
- [17] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social gan: Socially acceptable trajectories with generative adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2255–2264.
- [18] J. Liang, L. Jiang, J. C. Niebles, A. G. Hauptmann, and L. Fei-Fei, "Peeking into the future: Predicting future person activities and locations in videos," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5725–5734.
- [19] Y. Ling, Z. Ma, Q. Zhang, B. Xie, and X. Weng, "Pedast-gcn: Fast pedestrian crossing intention prediction using spatial-temporal attention graph convolution networks," *IEEE Transactions on Intelligent Transportation Systems*, 2024.
- [20] F. Marchetti, T. Mordan, F. Becattini, L. Seidenari, A. Del Bimbo, and A. Alahi, "Crossfeat: Semantic cross-modal attention for pedestrian behavior forecasting," *IEEE Transactions on Intelligent Vehicles*, 2024.
- [21] S. Huang, S. Mamidanna, S. Jangam, Y. Zhou, and L. H. Gilpin, "Can large language models explain themselves? a study of llm-generated self-explanations," *arXiv preprint arXiv:2310.11207*, 2023.
- [22] C. Chen, O. Li, D. Tao, A. Barnett, C. Rudin, and J. K. Su, "This looks like that: deep learning for interpretable image recognition," *Advances in neural information processing systems*, vol. 32, 2019.
- [23] W. Xu, Y. Xian, J. Wang, B. Schiele, and Z. Akata, "Attribute prototype network for zero-shot learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 969–21 980, 2020.
- [24] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," *Advances in neural information processing systems*, vol. 30, 2017.
- [25] M. Lee, S. Cho, S. Lee, C. Park, and S. Lee, "Unsupervised video object segmentation via prototype memory network," in *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2023, pp. 5924–5934.
- [26] M. Nauta, J. Schlötterer, M. Van Keulen, and C. Seifert, "Pip-net: Patch-based intuitive prototypes for interpretable image classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 2744–2753.
- [27] P. W. Koh, T. Nguyen, Y. S. Tang, S. Mussmann, E. Pierson, B. Kim, and P. Liang, "Concept bottleneck models," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5338–5348.
- [28] Z. Chen, Y. Bei, and C. Rudin, "Concept whitening for interpretable image recognition," *Nature Machine Intelligence*, vol. 2, no. 12, pp. 772–782, 2020.
- [29] A. Sarkar, D. Vijaykeerthy, A. Sarkar, and V. N. Balasubramanian, "A framework for learning ante-hoc explainable models via concepts," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 286–10 295.
- [30] Z. He, X. Ge, Q. Tang, T. Sun, Q. Cheng, and X. Qiu, "Dictionary learning improves patch-free circuit discovery in mechanistic interpretability: A case study on othello-gpt," *arXiv preprint arXiv:2402.12201*, 2024.
- [31] N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds, R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg, and C. Olah, "Toy models of superposition," *Transformer Circuits Thread*, 2022. [Online]. Available: [https://transformer-circuits.pub/2022/toy\\_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html)
- [32] J. Wang, S. Di, L. Chen, and C. W. W. Ng, "Learning from emergence: A study on proactively inhibiting the monosemantic neurons of artificial neural networks," in *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2024, pp. 3092–3103.
- [33] T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison, A. Askell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell, N. Joseph, Z. Hatfield-Dodds, A. Tamkin, K. Nguyen, B. McLean, J. E. Burke, T. Hume, S. Carter, T. Henighan, and C. Olah, "Towards monosemanticity: Decomposing language models with dictionary learning," *Transformer Circuits Thread*, 2023, <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- [34] Y. Feng, A. Carballo, Y. Niu, and K. Takeda, "Contrasting disentangled partial observations for pedestrian action prediction," in *2024 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2024, pp. 2828–2833.
- [35] A. Vaswani, "Attention is all you need," *Advances in Neural Information Processing Systems*, 2017.
- [36] S. Saadatnejad, A. Rasekh, M. Mofayez, Y. Medghalchi, S. Rajabzadeh, T. Mordan, and A. Alahi, "A generic diffusion-based approach for 3d human pose prediction in the wild," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 8246–8253.
- [37] S. Malla, B. Dariush, and C. Choi, "Titan: Future forecast using action priors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 186–11 196.
- [38] C. Wang, Y. Wang, M. Xu, and D. J. Crandall, "Stepwise goal-driven networks for trajectory prediction," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 2716–2723, 2022.