LEVERAGING METAPATHS FOR LEARNING FROM KNOWLEDGE GRAPHS IN THE CONTEXT OF VISION BASED CLASSIFICATION OF OBJECT STATES

Anonymous authors

Paper under double-blind review

ABSTRACT

Zero-Shot Object State Classification (ZS-OSC) aims to recognize unseen object states without any visual training examples. Existing methods typically rely on Knowledge Graphs (KGs) to provide semantic information about states, but they often treat KGs as homogeneous, overlooking the rich relational knowledge encoded in their structure. We propose a novel approach to ZS-OSC¹ that leverages meta-paths to capture complex relationships between object states in a KG. Our method learns to project semantic information from the KG into the visual space via meta-path learning, generating discriminative visual embeddings for unseen state classes. To the best of our knowledge, this is the first work to utilize metapaths for ZS-OSC. We conduct extensive experiments on four benchmark datasets, demonstrating the superior performance of our approach compared to SoTA zeroshot learning methods and a graph-based baseline. Our ablation study further provides insights into the impact of key design choices on the effectiveness of our method.

025 026 027

028 029

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

Knowledge graphs (KGs) have become increasingly important in addressing various computer vision (CV) tasks, such as object classification (Marino et al., 2017), zero-shot recognition (Wang 031 et al., 2018b) and Visual Question Answering (Krishna et al., 2017). This surge in KG utilization 032 can be attributed to their ability to provide rich semantic information and contextual knowledge that 033 can enhance the understanding of visual data. However, current approaches often under-utilize the 034 full potential of KGs. Several factors contribute to this sub-optimal utilization. First, KGs used in Computer Vision (CV) are often large and contain irrelevant or erroneous information, which can introduce noise and hinder performance. Second, many methods fail to consider the diversity of edge types in KGs, treating all edges as homogeneous and overlooking valuable relational knowl-037 edge. This simplistic approach limits the ability to effectively exploit the rich semantic information 038 embedded within KGs.

040 A number of approaches such as filtering mechanisms (Wang et al., 2014; Domingos & Richardson, 041 2007), ad-hoc KG construction (Dong et al., 2014; Gouidis et al., 2024) and random walking (Per-042 ozzi et al., 2014; Grover & Leskovec, 2016) were proposed as strategies against these shortcomings. An alternative approach focusing on the overcoming of these limitations concerned the concept of 043 meta-paths (Sun et al., 2011; Dong et al., 2017). The utilization of meta-paths essentially involves 044 the learning of the relative importance of the different paths between nodes within the KG via the assignment of weights to them. Meta-path-based approaches offer several advantages when applied 046 to heterogeneous information networks (HINs). First, meta-paths enable the modeling of complex 047 semantic relationships by explicitly defining sequences of node and edge types (Figure 1), making 048 them particularly suitable for capturing the diverse interactions in multi-typed networks (Sun et al., 049 2011; Shi et al., 2014). This provides a significant advantage over traditional graph learning methods (Sun et al., 2011), which treat all nodes and edges as homogeneous, thus failing to leverage 051 the rich information embedded in HINs. Furthermore, meta-paths allow for task-specific traversal 052 and filtering, enabling more accurate representations for applications such as link prediction, recom-

¹The code can be found at https://anonymous.4open.science/r/Metapaths-7811/

054

056

059

066

067

068 069



Figure 1: Unlike standard graph learning methods that prioritize local connections, meta-paths can capture distant relationships. In this toy example of a household objects graph, a meta-path can detect the stronger connection in specific contexts between "bottle" and "glass" (linked indirectly) than between "bottle" and "table", "bottle" and "kitchen" and "bottle" and "water" (directly linked).

mendation systems, and node classification. The ability to define domain-specific meta-paths also allows for more precise similarity measures, improving performance in tasks like clustering and search. Importantly, meta-paths enhance the interpretability of models by making the relationships between entities more transparent, providing domain experts with insights into how predictions are generated (Xiong et al., 2017).

076 Motivated by the great potential that the meta-paths learning seems to hold, this work attempts 077 to explore the utilization of this approach in the context of the Object State Classification (OSC) task (Isola et al., 2015; Gouidis et al., 2022; Souček et al., 2022; Saini et al., 2023), which is a CV 079 task attracting growing research attention over the last few years. OSC concerns the recognition of object states appearing in images and videos and is closely related to the more popular problems 081 of Object Recognition and Action Recognition. OSC is an important problem whose solution is of significant impact. The recognition of object states and state changes is crucial for determining an 083 object's condition and the interaction that was performed upon or could be performed in the future on it (Jamone et al., 2016). Moreover, the capacity for efficient OSC is of primary importance in 084 AI systems that support tasks such as learning object affordances (Chuang et al., 2018), recognizing 085 interactions (Wang et al., 2016b; Isola et al., 2015; Liu et al., 2017; Mancini et al., 2022), reasoning 086 to achieve an object state change (Farhadi et al., 2009), recognizing the completion or failure of goals 087 and recovery from possible mistakes during procedural activities (Schoonbeek et al., 2024) and many 088 others. Meanwhile, large-scale video datasets (Grauman et al., 2022; Saini et al., 2023) concerning 089 human-object interactions provide rich annotation data which are related to object state changes enabling the definition of new problems and the establishment of benchmarks and challenges related 091 to object state detection and classification (Grauman et al., 2022). 092

This work addresses the challenging task of Zero-Shot Object State Classification (ZS-OSC), where 093 the goal is to classify object states without any visual training examples. The key challenge lies 094 in leveraging auxiliary non-visual information to enable the classification of these unseen states. 095 Existing state-of-the-art methods typically utilize KGs as sources of structured semantic informa-096 tion (Gouidis et al., 2023), but they often treat KGs as homogeneous, potentially overlooking valu-097 able relational information. We propose a novel approach that leverages meta-paths for more effec-098 tive learning of zero-shot representations. Our method learns to project semantic information from 099 a KG into the visual space via meta-path learning, generating visual embeddings for unseen state classes. To the best of our knowledge, this is the first method to utilize meta-paths in the context of 100 Zero-Shot Visual Classification, and therefore ZS-OSC, through embedding generation. 101

Beyond their demonstrated utility in ZS-OSC, we posit that the generation of meta-path-based embeddings as a primary research objective holds significant promise. This direction offers several compelling advantages. First, meta-path-based embeddings can capture both local and global semantic relationships in heterogeneous networks, providing universal entity representations applicable across diverse tasks and domains without requiring task-specific fine-tuning. This fosters the development of multi-purpose, generalizable embeddings. Second, generating embeddings with a focus on representation quality can facilitate knowledge transfer across domains. Third, while cur-

108 rent embedding techniques, such as (Dong et al., 2017), are often task-specific, a framework that 109 prioritizes the generation of meta-path embeddings as an independent objective would offer a gen-110 eralized tool for analyzing heterogeneous information networks. This would mitigate the need for 111 task-dependent tuning and enable researchers to investigate embeddings without being constrained 112 by a particular task. Finally, since meta-paths effectively capture complex, higher-order relationships within networks, prioritizing the efficient generation of meta-path-based embeddings can lead 113 to more compact and informative representations, optimizing storage and computational efficiency, 114 particularly in large-scale heterogeneous networks. 115

- ¹¹⁶ This work makes the following key contributions:
 - We introduce a novel method for generating embeddings by leveraging meta-paths within Graph Neural Networks (GNNs). This approach enables GNNs to effectively harness the rich information encoded in KGs. To the best of our knowledge, this is the first work to explore meta-path utilization for embedding generation in this context.
 - We conduct a comprehensive ablation study to analyze the impact of various design choices and parameters on the performance of our proposed method. This analysis provides valuable insights into the interplay between meta-path learning and embedding generation.
 - We perform an extensive experimental evaluation on four diverse datasets, comparing our method against established KG-based baselines and SoTA Large Pre-trained Models. The results demonstrate that our approach achieves superior performance by a significant margin.
- 127 128 129

117

118

119

120

121

122

123

124

125

126

2 RELATED WORK

130 131

132 Meta-path Learning: Meta-paths, a fundamental concept in heterogeneous information networks 133 (HINs) (Shi et al., 2016), have been widely studied in applications such as similarity search, recom-134 mendation systems, and link prediction. Meta-paths were introduced to capture complex relation-135 ships in HINs, enabling the study of object proximities and connectivity patterns. The work by Sun 136 et al. (2011) proposed PathSim that used meta-paths to measure similarity between objects based on 137 shared relationships, with success in applications like similarity search and clustering. Extending this, Shi et al. (2014) proposed HeteSim, which computes relevance between objects via meta-paths, 138 incorporating directionality and node types for enhanced flexibility. Meta-paths have also enhanced 139 recommendation systems in HINs. Yu et al. (2013) developed a collaborative filtering algorithm 140 incorporating meta-path-based similarities between users and items, improving recommendation 141 accuracy. Similarly, Wang et al. (2016a) used meta-path-based features in matrix factorization for 142 item recommendation, leveraging HIN structure to model user-item interactions more effectively. 143

Learning optimal meta-path weights for specific tasks has been a key research area. Dong et al. 144 (2017) introduced MetaPath2Vec, which learns node embeddings through meta-path-based random 145 walks, showing improved performance in classification and clustering tasks. Fu et al. (2020) ex-146 tended this with a scalable meta-path-guided graph neural network, learning meta-path importance 147 for tasks in large-scale HINs. In link prediction, meta-paths have been used to predict missing 148 or future links in networks. Liu et al. (2018) proposed a meta-path-based link prediction method, 149 capturing complex node interactions and outperforming traditional algorithms in networks with mul-150 tiple types of nodes. Recent advancements in graph neural networks (GNNs) have further improved 151 meta-path-based link prediction. Zhang et al. (2019) introduced a heterogeneous graph neural net-152 work (HGNN) that integrates meta-path-based features into node aggregation, achieving better link 153 prediction.

154 Sun & Han (2013) introduced the concept of meta-paths for mining HINs, capturing semantic re-155 lationships across data types. Automatic discovery methods for meta-paths were later explored 156 by Meng et al. (2015) to address challenges in manually retrieving meta-paths. Ferrini et al. (2024) 157 proposed a novel approach to enhance GNN accuracy through effective meta-path identification, 158 while Noori et al. (2023) explored meta-paths for flexible similarity search in biological knowledge 159 graphs. Additionally, Yun et al. (2022) discussed learning new graph structures using meta-paths, demonstrating their role in enhancing GNN performance. Recent studies have also focused on chal-160 lenges such as automatic meta-path discovery (Huang et al., 2020), integration with deep learning 161 models (Wang et al., 2019), and the application of meta-paths in dynamic HINs (Trivedi et al., 2019).

The main novelty of our work distinguishing it from the aforementioned works concerns the utilization of meta-paths in the context of embeddings generation. Although many actual works focusing on meta-paths utilize embeddings, they serve as a means for another goal, i.e. downstream task, such as link prediction or entity alignment. In contrast, in our work, the generation of embeddings is a final objective.

167 Object State and Attribute Classification: Visual attributes are commonly defined as visual con-168 cepts that are both machine-detectable and human-understandable (Duan et al., 2012). The prevail-169 ing approach to learning attributes mirrors that of object classes, where a convolutional neural net-170 work is trained with discriminative classifiers using annotated image datasets (Singh & Lee, 2016). 171 However, existing labeled attribute datasets suffer from limitations such as smaller scale compared 172 to object datasets, a restricted number of generic attributes, and limited category coverage (Lampert et al., 2009; Isola et al., 2015; Patterson & Hays, 2016; Yu & Grauman, 2017; Mancini et al., 173 2022). Furthermore, research specifically focusing on state classification remains limited (Gouidis 174 et al., 2022), with most existing approaches relying on similar assumptions and techniques as those 175 employed for attribute classification. This highlights a need for dedicated methods tailored to the 176 unique challenges of state recognition. 177

Zero-shot Classification: Zero-shot learning (ZSL) has attracted significant attention due to its 178 179 ability to address the practical challenge of classifying objects without any training examples (Xian et al., 2018a). This is particularly crucial in real-world scenarios where obtaining labeled data for 180 every possible class is often infeasible. Several approaches have been proposed for zero-shot ob-181 ject classification, including semantic embedding-based methods (Wang et al., 2018a; Xian et al., 182 2018b), attribute-based methods (Lampert et al., 2014), generative models (Xian et al., 2018b; 183 Changpinyo et al., 2016), and learning compatibility functions between image and class embed-184 dings (Akata et al., 2015). Semantic embedding methods utilize compact semantic spaces or at-185 tribute sets to bridge the gap between seen and unseen object classes. Attribute-based methods leverage descriptive attributes to infer the class of unseen objects. Generative models synthesize 187 images of unseen classes based on similarities to seen classes. Additionally, recent work has ex-188 plored the use of knowledge graphs to capture semantic relationships between objects and facilitate 189 ZSL (Kampffmeyer et al., 2019; Nayak & Bach, 2022). While ZSL has been extensively studied for object recognition, its application to Zero State Classification (ZSC) remains relatively unex-190 plored. With the exception of Gouidis et al. (2023), which focuses exclusively on state classification 191 without relying on prior knowledge about object classes, existing ZSC methods primarily address 192 the Compositional Zero-Shot Learning (CZSL) variant of the problem. CZSL aims to generalize 193 to unseen combinations of object and state primitives by learning their compositionality from the 194 training set (Misra et al., 2017; Nagarajan & Grauman, 2018; Yang et al., 2020). 195

195

3 Methodology

197 198

199 This work introduces a novel approach for generating embeddings in heterogeneous graphs by syn-200 ergistically combining KGs structures with meta-path-based GNNs. These generated embeddings 201 are then employed for the task of ZS-OSC. Our methodology harnesses the strengths of both KGs embeddings and the rich relational information encapsulated by meta-paths. The core idea is to 202 leverage a designated set of KG nodes as guides for meta-path learning. Specifically, we utilize the 203 visual embeddings of these nodes as ground truth and train a Graph Transformer Network (GTN) 204 to assign weights to different KG edge types. This is achieved by generating embeddings for the 205 guide nodes that align with their visual embeddings. The GTN architecture is inspired by Yun et al. 206 (2022), while the training procedure draws inspiration from Kampffmeyer et al. (2019); Gouidis 207 et al. (2023). Figure 2 provides a general overview of our method. 208

209 210

3.1 PRELIMINARIES

Before proceeding with describing our method, it is necessary to present related backgrounds related to our work. Additionally, Table S1 presented in the supplementary section lists commonly used notations in this paper for quick reference.

Heterogeneous Graph: A heterogeneous graph is denoted as $G = (V, E, \phi, \psi)$ and is associated with a node type mapping function $\phi : V \to T_v$ and an edge type mapping function $\psi : E \to T_e$,



of *P*. A neighbor that is connected through two different metapath instances is treated as two distinct nodes in $N_P(v)$. If *P* is symmetric, $N_P(v)$ includes the node *v* itself.

273 Metapath-based Graph: For a given metapath P, the metapath-based graph G_P is the graph con-274 structed using all metapath-based neighbor pairs from the original graph G. If P is symmetric, G_P 275 is homogeneous.

Heterogeneous Graph Embedding: Given a heterogeneous graph G = (V, E), with node attribute matrices $X_{A_i} \in \mathbb{R}^{|V_{A_i}| \times d_{A_i}}$ for each node type $A_i \in \mathcal{A}$, the goal of heterogeneous graph embedding is to learn *d*-dimensional node representations $h_v \in \mathbb{R}^d$ for each node $v \in V$, where $d \ll |V|$, capturing rich structural and semantic information from G.

280 281

3.2 Meta-path-based Embedding Generation

Meta-paths represent semantic connections between different node types in a heterogeneous graph.
GTNs automatically learn useful meta-paths by selecting and combining adjacency matrices of different edge types. This process allows the model to generate new graph structures that are useful and is used typically for downstream tasks such as node classification. In our case, the meta-paths learning is not mediated by a downstream task. Instead, meta-paths are learned via the generation of embeddings for the graph's nodes.

Following the notation introduced previously, we can use adjacency matrices, a different one for each edge type, to compute the different meta-paths. Specifically, the meta-path's adjacency matrix $A_P \in \{0, 1\}^{\mathbb{V} \times \mathbb{V}}$ of the graph is computed as:

291 292 293

295

289

290

$$A_P = \prod_{i=1}^{\ell} A_{t_i} \Longrightarrow A_P = A_{t_1} A_{t_2} \dots A_{t_{\ell}}$$

where A_{t_i} is the adjacency matrix corresponding to edge type t_i and l corresponds to the length of the meta-paths, with $a_{ij} = 1$ denoting a meta-path with length equal than l between nodes i and jand $a_{ij} = 0$ an absence of meta-path, respectively. In order to have also meta-paths with lengths less than l we add the identity matrix $I \in \mathbb{R}^{V \times V}$ to the adjacency matrices.

With the utilization of a soft selection mechanism, weights are assigned to the different types of edges. This technique enables the learning of the optimal combination of edge types for each metapath. Specifically, this is achieved with a 1×1 convolution with softmax activation over the adjacency matrices of different edge types:

304 305

306

307

308

where $\alpha_t^{(k)}$ is the learnt weight for edge type t at the k-th transformer layer. The resulting metapath adjacency matrices are normalized using degree matrices D and applied to perform multi-hop convolutions. More than one softmax convolutions could be also applied, in which case the objective is that each channel learns a different representation and the overall combination of the different representations results in more a robust outcome.

 $WA^{(k)} = \sum_{i=1}^{|T_e|} \alpha_t^{(k)} A_t,$

With the learnt softmax weights, the meta-path's adjacency matrix $WA_P \in \mathbb{R}^{\mathbb{V} \times \mathbb{V}}$ of the graph is computed as:

316 317 318

$$WA_P = \prod_{i=1}^{\ell} A_{t_i} \Longrightarrow \sum_{t=0}^{|T_e|} \alpha_t^{(k)} A_t = \sum_{t=1}^{|T_e|} \alpha^{(0)} A_{t_0} \sum_{t=1}^{|T_e|} \alpha^{(1)} A_{t_1} \sum_{t=1}^{|T_e|} \alpha^{(2)} A_{t_2} \cdots \sum_{t=1}^{|T_e|} \alpha^{(l)} A_{t_l}.$$

319 320 321

After having computed the weighted adjacency matrix WA_P we can compute the aggregation metapaths-based features matrix $\mathbb{X}_P \in \mathbb{R}^{|V_{A_i}| \times d_{A_i}}$ for the graphs nodes by multiplying the nodes features matrix matrices $\mathbb{X}_{A_i} \in \mathbb{R}^{|V_{A_i}| \times d_{A_i}}$ with WA_P :

$$\mathbb{X}_P = \mathbb{X}_A \times WA_P$$

In order to generate embeddings for the graphs nodes we use a multi-layer Graph Neural Network which learns to project the nodes features into the target space:

$$\mathbb{E}_P = f_\theta(\mathbb{X}_P)$$

The weights corresponding to the softmax layers and the stacked layers of the GNN are updated based on the minimization of a L2 distance function $\mathcal{L}_{\mathcal{E}}$ between the generated nodes embeddings and the ground truth nodes embeddings for the set of nodes which serve as ground truths:

$$\mathcal{L}_{\mathcal{E}} = \frac{1}{2N} \sum_{n \in \mathcal{N}} \sum_{d \in \mathcal{D}} (\mathbb{E}_P - \tilde{\mathbb{E}}_P)^2$$

where \mathcal{D} is the dimension of the embeddings and \mathcal{N} the number of the ground truth concepts, respectively.

3.3 ZERO-SHOT OBJECT STATE CLASSIFICATION

To achieve zero-shot object-state classification, it is crucial to learn representations for the unseen 342 target state classes. To this end, we construct a KG with state classes as nodes and utilize the 343 GTN trained in the previous step for meta-path learning to generate embeddings for these target 344 classes. These embeddings are then integrated into a visual classifier, specifically a pre-trained 345 CNN, to serve as visual representations of the visually known state classes. Following established 346 practices in transfer learning and zero-shot learning, we replace the final layer of the pre-trained 347 CNN with these generated embeddings. To ensure compatibility, the dimensionality (\mathcal{D}) of the 348 generated embeddings is set to match the dimension of the CN's last layer. This adaptation enables 349 the CNN to effectively classify the target state classes in a zero-shot manner.

350 351

324 325 326

327

328 329 330

331

332

333

338

339 340

341

4 EXPERIMENTAL EVALUATION

352 353 354

4.1 IMPLEMENTATION AND EVALUATION ISSUES

355 **Implementation Details:** The KG provided as input to the Graph Transformer Network (GTN) for 356 meta-path learning was the WordNet hierarchy of the 1000 classes from the ImageNet1000 dataset 357 (Russakovsky et al., 2015). These 1000 classes served as the ground truth set for training the GTN, 358 which was trained for 200 epochs with five different learning rates (see Section 4.2 for details). The 359 same KG was also used to generate the target object state classes. In experiments with multiple softmax channels, the generated embeddings were averaged. The convolutional neural network 360 (CNN) used for zero-shot object-state classification (ZS-OSC) was ResNet-101 (He et al., 2016) 361 pre-trained on the ImageNet1000 dataset. 362

Datasets: We utilized four publicly available datasets containing object state annotations: OSDD
 Gouidis et al. (2022), CGQA Mancini et al. (2022), MIT Isola et al. (2015), and VAW Pham et al.
 (2021). While OSDD is specifically designed for state detection, the other three are attribute datasets
 that include object states among their classes. We extracted subsets concerning exclusively object
 state classes. For OSDD and VAW, bounding boxes from the original images were extracted to
 create images suitable for the OSC task. Table S2 in the supplementary section provides details on
 the four datasets.

Metrics: Our evaluation follows the zero-shot method from (Purushwalkam et al., 2019). Following
 those guidelines we calculate accuracy per class and then average these instead of reporting overall
 accuracy. This ensures each class is equally important, regardless of its sample size.

373 374

- 4.2 ABLATION STUDY
- 375

This ablation study investigates the optimal configuration of key parameters related to meta-pathbased embedding generation, including: (a) the maximum length of meta-paths, (b) the number of softmax channels used for meta-path selection, and (c) the training learning rate.

070					
378	Meta-naths Max Length	OSDD	CGOA-States	MIT-States	VAW-States
379	The putties that Bengin	00000	eeq.i suites		
380	1	27.1	44.9	47.1	29.3
381	2	29.5	46.7	47.3	29.8
382	3	31.3	47.6	48.7	32.1
383	4	30.2	46.5	47.9	31.0

Table 1: Ablation results for maximum length of meta-paths. The number of softmax channels in the GTN was 3. The networks were trained with a learning rate equal to 1e - 2.

Number of Channels	OSDD	CGQA-States	MIT-States	VAW-States
1	27.2	44.3	46.2	28.5
2	28.9	46.1	46.4	29.7
3	31.3	47.6	48.7	32.1
4	30.5	46.9	47.0	30.7

Table 2: Ablation results for number of softmax channels. The maximum length of the meta-paths was 3. The networks were trained with a learning rate equal to 1e - 2.

Length of meta-path: Table 1 presents the results of varying the maximum length of meta-paths (1, 2, 3, and 4 hops). The best performance across all datasets was achieved with a maximum length of 3 hops. Performance generally improved with increasing meta-path length, likely because longer paths capture more global graph information. However, performance slightly decreased with a length of 4 hops, suggesting that excessively long paths might introduce noise, potentially by incorporating less relevant or spurious relationships. This finding highlights the importance of carefully selecting the appropriate meta-path length to balance information gain and noise reduction.

Number of Softmax Channels: Table 2 shows the impact of varying the number of softmax channels used for meta-path selection. The best performance was achieved with 3 channels, suggesting that multiple channels allow the GTN to learn diverse meta-path representations, thereby improving embedding quality. Increasing the number of channels allows the model to capture different aspects of the relationships encoded in the meta-paths, leading to richer and more informative embeddings. However, using too many channels, e.g., 4 in this case, might not provide further benefits and could potentially increase model complexity without a corresponding improvement in performance.

Training Learning Rate: Table 3 presents the results of using different learning rates (LR) during GTN training. The best performance was obtained with a learning rate of 1e-2. Smaller learning rates led to a significant performance drop, indicating the importance of this parameter for effective model training. This suggests that a learning rate of 1e-2 strikes a good balance between convergence speed and stability, allowing the model to effectively learn from the data without overshooting or getting stuck in local optima. The observed performance drop with smaller learning rates could be attributed to slower convergence and potential difficulties in escaping local optima.

4.3 COMPARISON TO COMPETING METHODS

This experiment had a two-fold objective. First, we compared our method against SoTA Large Pre-trained Models (LPMs) capable of ZS-OSC. We used six different prompts² related to the target states and report the mean average performance across all prompts, following standard convention. Second, we compared our approach to a graph-based method (Gouidis et al., 2023) specifically designed for ZS-OSC, which allows us to assess the impact of meta-paths on performance. For comparison with the baseline method (Gouidis et al., 2023), which relies on random walks, we used five different random seeds for initialization and report the mean performance over all seeds. The same KG was used as input to both the GCNs of the baseline method and our method. It is important to note that other ZSL methods, such as CZSL, are not applicable in this context, because they require information about object classes, which is not available in this zero-shot setting.

²The prompts are presented in the supplementary material.

0.0					
132	LR	OSDD	CGOA-States	MIT-States	VAW-States
33		ODDD	COQ11 States	MIII States	Min States
434	5e - 2	29.2	46.2	47.5	30.1
435	1e - 2	31.3	47.6	48.7	32.1
136	5e - 3	28.9	45.1	47.0	28.5
430	1e - 3	27.2	44.5	46.1	26.2
437	5e - 4	26.1	43.8	44 5	24.6
438	00 4	20.1	45.0		24.0

Table 3: Ablation results for learning rate. The maximum length of the meta-paths was 3, and the number of softmax channels was 3.

Method	OSDD	CGQA-States	MIT-States	VAW-States
Baseline (Gouidis et al., 2023)	27.3	45.1	43.3	25.6
CLIP-RN101 (Radford et al., 2021)	22.5	46.9	39.3	28.0
CLIP-VITBP16 (Radford et al., 2021)	28.8	44.9	46.4	30.1
CLIP-VITLP14 (Radford et al., 2021)	28.4	43.4	48.6	27.9
ALIGN (Jia et al., 2021)	29.5	40.0	44.2	28.4
BLIP (Li et al., 2022)	13.3	26.0	27.2	16.1
Ours	31.3	47.6	48.7	32.1
Our improvement over the SoTA	1.8	0.7	0.1	2.0
Our improvement over the Baseline	4.0	2.4	5.4	6.5

Table 4: Results against competing methods. Bold/Blue text indicates best/2nd best performance.

Table 4 presents the results. Our method achieved the best performance across all datasets, outper-forming both the baseline method and all LPMs. Notably, our method surpassed the baseline by a large margin (4.0%, 2.5%, 5.4%, and 6.5% for OSDD, CGQA, MIT, and VAW, respectively). This outcome highlights the effectiveness of meta-path learning compared to random walks for represen-tation learning in graph structures. The superior performance of our method can be attributed to the ability of meta-paths to capture and exploit specific semantic relationships within the KG, leading to more informative and discriminative embeddings for ZS-OSC. In contrast, random walk techniques might not effectively capture these relationships, resulting in less effective representations.

Furthermore, the fact that our method outperforms SoTA LPMs (1.8%, 0.7%, 0.1%, and 2.0% for OSDD, CGQA, MIT, and VAW, respectively) demonstrates the potential of incorporating KG in-formation and meta-path learning into ZS-OSC. LPMs, while powerful, might not fully capture the nuanced semantic relationships between objects and their states that are encoded in KGs. By leveraging meta-paths to extract and utilize this information, our method achieves a significant per-formance improvement.

CONCLUSION

This paper introduced a novel method for generating embeddings in heterogeneous graphs by lever-aging meta-paths within a graph neural network framework. Our approach utilizes KGs structures and visual embeddings to guide the learning of meta-paths, enabling the generation of informative and discriminative embeddings. We demonstrated the effectiveness of our method in the context of ZS-OSC, achieving superior performance compared to state-of-the-art LPMs and a graph-based baseline. Our ablation study provided insights into the impact of key parameters on the performance of our method.

Future work will focus on several promising directions. First, we aim to explore the applicability of our method to other zero shot CV tasks such as semantic segmentation, image captioning and visual question answering to further evaluate its generalizability. Second, we plan to investigate the integration of different types of KGs such as multimodal KGs and KGs containg different type of nodes and explore the impact of KG quality on embedding generation. Finally, we intend to extend our approach to incorporate more complex meta-path structures and explore alternative graph neural network architectures for enhanced performance.

486 REFERENCES

526

527

528

- Zeynep Akata, Scott Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. Evaluation of output embeddings for fine-grained image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2927–2936, 2015. 4
- 491 Soravit Changpinyo, Wei-Lun Chao, Boqing Gong, and Fei Sha. Synthesized classifiers for zero 492 shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
 493 pp. 5327–5336, 2016. 4
- Ching-Yao Chuang, Jiaman Li, Antonio Torralba, and Sanja Fidler. Learning to act properly: Predicting and explaining affordances from images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 975–983, 2018. 2
- Pedro Domingos and Matthew Richardson. Markov logic: A unifying framework for statistical
 relational learning. 2007. 1
- Xin Dong, Evgeniy Gabrilovich, Geremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmann, Shaohua Sun, and Wei Zhang. Knowledge vault: A web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 601–610, 2014. 1
- Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. Metapath2vec: Scalable representation
 learning for heterogeneous networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 135–144. ACM, 2017. 1, 3
- Kun Duan, Devi Parikh, David Crandall, and Kristen Grauman. Discovering localized attributes for
 fine-grained recognition. In 2012 IEEE conference on computer vision and pattern recognition,
 pp. 3474–3481. IEEE, 2012. 4
- Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes.
 2009 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2009, pp. 1778– 1785, 2009. doi: 10.1109/CVPRW.2009.5206772. 2
- Francesco Ferrini, Antonio Longa, Andrea Passerini, and Manfred Jaeger. Meta-path learning for multi-relational graph neural networks. In *Learning on Graphs Conference*, pp. 2–1. PMLR, 2024. 3
- Xiaoyang Fu, Jiawei Zhang, Zitao Meng, and S Yu Philip. Scalable graph neural networks for heterogeneous networks. *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1001–1010, 2020. 3
- F. Gouidis, T. Patkos, A. Argyros, and D. Plexousakis. Detecting object states vs detecting objects:
 A new dataset and a quantitative experimental study. In *Proceedings of the 17th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications* (VISAPP), volume 5, pp. 590–600, 2022. 2, 4, 7, 15
 - Filipos Gouidis, Theodore Patkos, Antonis Argyros, and Dimitris Plexousakis. Leveraging knowledge graphs for zero-shot object-agnostic state classification. *arXiv preprint arXiv:2307.12179*, 2023. 2, 4, 8, 9
- Filippos Gouidis, Konstantinos Papoutsakis, Theodore Patkos, Antonis Argyros, and Dimitris Plexousakis. Exploring the impact of knowledge graphs on zero-shot visual object state classification. *Proceedings Copyright*, 738:749, 2024. 1
- Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan,
 Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray,
 Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano
 Fragomeni, Qichen Fu, Abrham Gebreselasie, Cristina González, James Hillis, Xuhua Huang,
 Yifei Huang, Wenqi Jia, Weslie Khoo, Jáchym Kolář, Satwik Kottur, Anurag Kumar, Federico
 Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan

556

558

559

576

540 Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz, Merey Ramazanova, Leda Sari, 541 Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, 542 Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbeláez, David Crandall, Dima Damen, 543 Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawa-544 har, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torre-545 sani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the world in 3,000 hours of egocentric 546 video. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition 547 (CVPR), pp. 18995–19012, June 2022. 2 548

- Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings* of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 855–864, 2016. 1
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016. 7
 - Shaohua Huang, Xiao He, Ruifeng Jiang, Li Song, Xinyu Xie, and Hongyuan Ma. Meta-path guided heterogeneous graph neural network for intent recommendation. In *Proceedings of the 29th ACM International Conference on Information Knowledge Management*, pp. 2249–2256, 2020. **3**
- Phillip Isola, Joseph J. Lim, and Edward H. Adelson. Discovering states and transformations in image collections. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 07-12-June:1383–1391, 2015. ISSN 10636919. doi: 10.1109/CVPR. 2015.7298744. 2, 4, 7, 15
- Lorenzo Jamone, Emre Ugur, Angelo Cangelosi, Luciano Fadiga, Alexandre Bernardino, Justus
 Piater, and José Santos-Victor. Affordances in psychology, neuroscience, and robotics: A survey.
 IEEE Transactions on Cognitive and Developmental Systems, 10(1):4–25, 2016. 2
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan
 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
 with noisy text supervision. In *International conference on machine learning*, pp. 4904–4916.
 PMLR, 2021. 9
- Michael Kampffmeyer, Yinbo Chen, Xiaodan Liang, Hao Wang, Yujia Zhang, and Eric P. Xing. Rethinking knowledge graph propagation for zero-shot learning. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2019-June:11479–11488, 2019. ISSN 10636919. doi: 10.1109/CVPR.2019.01175. 4
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie
 Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017. 1
- ⁵⁸¹ Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 951–958. IEEE, 2009. 4
- Christoph H. Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for
 zero-shot visual object categorizationa. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(3):453–465, 2014. ISSN 01628828. doi: 10.1109/TPAMI.2013.140. 4
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International Conference* on Machine Learning, pp. 12888–12900. PMLR, 2022. 9
- Xiaolei Liu, Xiaodan Liang, Zhiqiang Li, Jie Tang, and Li Lin. Link prediction in heterogeneous
 networks using meta-paths. *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 437–446, 2018. 3

594 595 596	Yang Liu, Ping Wei, and Song-Chun Zhu. Jointly recognizing object fluents and tasks in egocentric videos. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , pp. 2924–2932, 2017. 2
598 599 600 601	Massimiliano Mancini, Muhammad Ferjad Naeem, Yongqin Xian, and Zeynep Akata. Learning Graph Embeddings for Open World Compositional Zero-Shot Learning. <i>IEEE Transactions on</i> <i>Pattern Analysis and Machine Intelligence</i> , 8828(c):1–15, 2022. ISSN 19393539. doi: 10.1109/ TPAMI.2022.3163667. 2, 4, 7, 15
602 603 604 605	Kenneth Marino, Ruslan Salakhutdinov, and Abhinav Gupta. The more you know: Using knowledge graphs for image classification. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pp. 2673–2681, 2017. 1
606 607 608	Changping Meng, Reynold Cheng, Silviu Maniu, Pierre Senellart, and Wangda Zhang. Discovering meta-paths in large heterogeneous information networks. In <i>Proceedings of the 24th international conference on world wide web</i> , pp. 754–764, 2015. 3
609 610 611	Ishan Misra, Abhinav Gupta, and Martial Hebert. From red wine to red tomato: Composition with context. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pp. 1792–1801, 2017. 4
613 614 615	Tushar Nagarajan and Kristen Grauman. Attributes as operators: factorizing unseen attribute-object compositions. In <i>Proceedings of the European Conference on Computer Vision (ECCV)</i> , pp. 169–185, 2018. 4
616 617 618	N. V. Nayak and S. H. Bach. Zero-shot learning with common sense knowledge graphs. <i>Transactions on Machine Learning Research (TMLR)</i> , 2022. 4
619 620 621	Ayush Noori, Michelle M Li, Amelia LM Tan, and Marinka Zitnik. Metapaths: similarity search in heterogeneous knowledge graphs via meta-paths. <i>Bioinformatics</i> , 39(5):btad297, 2023. 3
622 623 624	 Genevieve Patterson and James Hays. Coco attributes: Attributes for people, animals, and objects. In Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VI 14, pp. 85–100. Springer, 2016. 4
625 626 627	Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social repre- sentations. In <i>Proceedings of the 20th ACM SIGKDD international conference on Knowledge</i> <i>discovery and data mining</i> , pp. 701–710, 2014. 1
628 629 630 631	Khoi Pham, Kushal Kafle, Zhe Lin, Zhihong Ding, Scott Cohen, Quan Tran, and Abhinav Shrivas- tava. Learning to predict visual attributes in the wild. In <i>Proceedings of the IEEE/CVF CVPR</i> , pp. 13018–13028, June 2021. 7, 15
632 633 634 635	Senthil Purushwalkam, Maximilian Nickel, Abhinav Gupta, and Marc'Aurelio Ranzato. Task-driven modular networks for zero-shot compositional learning. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 3593–3602, 2019. 7
636 637 638 639	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021. 9
640 641 642 643	Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. <i>International journal of computer vision</i> , 115:211–252, 2015. 7
644 645 646	N. Saini, H. Wang, A. Swaminathan, V. Jayasundara, B. He, K. Gupta, and A. Shrivastava. Chop amp; learn: Recognizing and generating object-state compositions. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 20190–20201, Los Alamitos, CA, USA,

 647
 oct 2023. IEEE Computer Society. doi: 10.1109/ICCV51070.2023.01852. URL https: //doi.ieeecomputersociety.org/10.1109/ICCV51070.2023.01852. 2
 668

677

681

683

684

685

686

687 688

689

690

691

692

- 648 Tim J Schoonbeek, Tim Houben, Hans Onvlee, Fons van der Sommen, et al. Industreal: A dataset 649 for procedure step recognition handling execution errors in egocentric videos in an industrial-like 650 setting. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 651 pp. 4365–4374, 2024. 2 652 Chuan Shi, Bin Hu, Xiaowei Zhao, and Philip S Yu. Hetesim: A general framework for relevance 653 measure in heterogeneous networks. IEEE Transactions on Knowledge and Data Engineering, 654 26(10):2479-2492, 2014. 1, 3 655 656 Chuan Shi, Yitong Li, Jiawei Zhang, Yizhou Sun, and S Yu Philip. A survey of heterogeneous 657 information network analysis. IEEE Transactions on Knowledge and Data Engineering, 29(1): 658 17–37, 2016. 3 659 Krishna Kumar Singh and Yong Jae Lee. End-to-end localization and ranking for relative attributes. 660 In Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, Oc-661 tober 11-14, 2016, Proceedings, Part VI 14, pp. 753-769. Springer, 2016. 4 662 663 Tomáš Souček, Jean-Baptiste Alayrac, Antoine Miech, Ivan Laptev, and Josef Sivic. Multi-Task 664 Learning of Object State Changes from Uncurated Videos. 2022. URL http://arxiv.org/ 665 abs/2211.13500.2 666 667 Yizhou Sun and Jiawei Han. Meta-path-based search and mining in heterogeneous information networks. Tsinghua Science and Technology, 18(4):329–338, 2013. 3
- 669 Yizhou Sun, Jiawei Han, Xifeng Yan, Philip S Yu, and Tianyi Wu. Pathsim: Meta path-based 670 top-k similarity search in heterogeneous information networks. In Proceedings of the 17th ACM 671 SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 992–1000. 672 ACM, 2011. 1, 3 673
- 674 Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Dyngem: Deep embedding method 675 for dynamic heterogeneous information networks. In Proceedings of the 2019 World Wide Web 676 Conference, pp. 1132–1143, 2019. 3
- Daixin Wang, Peng Cui, and Wenwu Zhu. Structural deep network embedding. In Proceedings of 678 the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 679 pp. 1225–1234. ACM, 2016a. 3 680
- Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, and Philip S Yu. Heterogeneous 682 graph attention network. In The World Wide Web Conference, pp. 2022–2032, 2019. 3
 - Xiaolong Wang, Ali Farhadi, and Abhinav Gupta. Actions ~ Transformations. In Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, volume 2016-Decem, pp. 2658–2667. IEEE, jun 2016b. ISBN 9781467388504. doi: 10.1109/CVPR.2016.291. URL http://ieeexplore.ieee.org/document/7780660/.2
 - Xiaolong Wang, Yufei Ye, and Abhinav Gupta. Zero-Shot Recognition via Semantic Embeddings and Knowledge Graphs. Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 6857–6866, 2018a. ISSN 10636919. doi: 10.1109/CVPR. 2018.00717. 4
- Xiaolong Wang, Yufei Ye, and Abhinav Gupta. Zero-shot recognition via semantic embeddings 693 and knowledge graphs. In Proceedings of the IEEE conference on computer vision and pattern 694 recognition, pp. 6857-6866, 2018b. 1 695
- 696 Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by trans-697 lating on hyperplanes. In Proceedings of the AAAI conference on artificial intelligence, vol-698 ume 28, 2014. 1
- Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning-a 700 comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis* 701 and machine intelligence, 41(9):2251-2265, 2018a. 4

702 703 704	Yongqin Xian, Tobias Lorenz, Bernt Schiele, and Zeynep Akata. Feature generating networks for zero-shot learning. In <i>Proceedings of the IEEE conference on computer vision and pattern recog-</i> <i>nition</i> , pp. 5542–5551, 2018b. 4
705 706 707 708	Wenhan Xiong, Thien Hoang, and William Yang Wang. Deeppath: A reinforcement learning method for knowledge graph reasoning. In <i>Proceedings of the 2017 ACM on Conference on Information and Knowledge Management</i> , pp. 565–574, 2017. 2
709 710 711	Muli Yang, Cheng Deng, Junchi Yan, Xianglong Liu, and Dacheng Tao. Learning unseen concepts via hierarchical decomposition and composition. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 10248–10256, 2020. 4
712 713 714 715	Aron Yu and Kristen Grauman. Semantic jitter: Dense supervision for visual comparisons via synthetic images. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , pp. 5570–5579, 2017. 4
716 717 718	Xin Yu, Zhaochun Ren, Yizhou Sun, Quanquan Gu, Zhiqiang Luo, and Dawei Chen. Collaborative filtering with entity similarity regularization in heterogeneous information networks. In <i>IJCAI International Joint Conference on Artificial Intelligence</i> , pp. 21–29, 2013. 3
719 720 721	Seongjun Yun, Minbyul Jeong, Sungdong Yoo, Seunghun Lee, S Yi Sean, Raehyun Kim, Jaewoo Kang, and Hyunwoo J Kim. Graph transformer networks: Learning meta-path graphs to improve gnns. <i>Neural Networks</i> , 153:104–119, 2022. 3, 4
723 724 725	Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, and Nitesh V Chawla. Heteroge- neous graph neural network. In <i>Proceedings of the 25th ACM SIGKDD International Conference</i> <i>on Knowledge Discovery & Data Mining</i> , pp. 793–803. ACM, 2019. 3
726	
729	
720	
729	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
/50	
750	
752	
757	
755	

SUPPLEMENTARY MATERIAL

This section presents tables that due to space limitation were omitted from the main body of the paper. Specifically, Table S1 shows the notation symbols that are used, Table S2 presents the details about the datasets used in the experimental evaluation and Table S3 shows the prompts that were using for the LPMs respectively.

763	Symbol	Description
764	Symbol	Description
765	G = (V, E)	Heterogeneous graph
705	V	Set of nodes in G
766	E	Set of edges in G
767	$\phi: V \to T_v$	Node type mapping function
768	$\psi: E \to T_e$	Edge type mapping function
769	T_v	Set of node types
770	T_e	Set of edge types
771	$ T_v $	Number of node types
772	$ T_e $	Number of edge types
773	Р	Meta-path
774	p	Meta-path instance
775	v_i	Node in G
776	e_l	Edge in G
	$\tau_e(e_l)$	Edge type of edge e_l
///	$N_P(v)$	Set of metapath-based neighbors of node v under meta-path P
778	G_P	Metapath-based graph derived from meta-path P
779	X_A	Node feature matrix
780	$\mathcal{A}^{''}$	Adjacency matrix
781	$\mathcal{A}_{\mathcal{P}}$	Meta-path adjacency matrix
782	$\mathbb{E}_{P}^{'}$	Meta-path based generated embedding of node v
783	$\tilde{\mathbb{E}}_P$	Ground truth embedding of node v
784	\mathcal{D}	Dimensionality of node embeddings
785	\mathcal{N}	Number of ground truth concepts
786	$\mathcal{L}_{\mathcal{E}}$	L2 distance loss function

Table S1: Notation Table

Dataset	Train	Val	Test	States	Objects	VOSC	TOSC	S\O
OSDD (Gouidis et al., 2022)	6,977	1,124	5,275	9	14	35	126	2.36
CGQA-states (Mancini et al., 2022)	244	46	806	5	17	41	75	1.71
MIT-states (Isola et al., 2015)	170	34	274	5	14	20	70	1.57
VAW (Pham et al., 2021)	2,752	516	1,584	9	23	51	207	2.61

Table S2: Details about the four image datasets utilized in this work. Train/Val/Test: Number of Training/Validation/Testing Images. States: Number of State classes, Objects: Number of Object classes. VOSC/TOSC: Valid/Total Object-State combinations. S\O: Average number of states than an Object can be situated in.

	Prompt
1	An image of a {} object
2	The object in the image is {}
3	The state of the object in the image is {}
4	The object in the image is currently {}
5	An image of a object in a state of {}
6	The scene depicts a object that appears to be {}