

A Neuro-Symbolic Approach with Reinforcement Learning for Explainable Anomaly Detection in Pedestrian Video Sequence

Jaeil Park and Sung-Bae Cho

Department of Computer Science, Yonsei University, Seoul 03722, South Korea
{wodlf603, sbcho}@yonsei.ac.kr

Abstract

Video anomaly detection in pedestrian streets requires to explain the anomaly because of its danger, such as a car moving on a pedestrian road, and to interact with supervisors with question answering. To explain the anomaly, the methods based on neural networks such as SHAP have been investigated, but they have a limitation that only takes account of the properties of the abnormal objects and is not interactive with the supervisor. This paper proposes a video anomaly detection method supporting question answering with a reinforcement learning-based neuro-symbolic approach. After converting a question into executable programs, it is operated on a scene graph with the video anomaly detection result to provide an answer for the question. After that, it executes reinforcement learning through a comparison between the result of the model and the ground-truth feedback from the supervisor. A question-answering experiment on UCSD dataset confirms that the proposed method answers the questions about anomalies, confirming 99% accuracy and demonstrating the causal inference through case analysis.

1 Introduction

Due to the extensive usage of surveillance cameras and the limitations of manpower, there is a growing demand for an automated video surveillance system [Fleck and Straßer, 2010]. One of the primary challenges that autonomous video surveillance systems face with is the automatic detection of anomalies, defined as unusual, uncommon, or irregular events occurring in complex and crowded environments [Cong et al., 2011; Xu et al., 2017].

In addition to performing detection using black box models, the anomaly detection model should provide explanations regarding the causes, outcomes, and necessary precautions for identifying visual scenarios that encompass real complex situations in a logical manner [Amarasinghe et al., 2018]. Previous research on explanation has been predominantly focused on the properties of abnormal objects, thereby neglecting the comprehensive associations between objects that are linked to the risk of such abnormalities [Szymanowicz et al.,

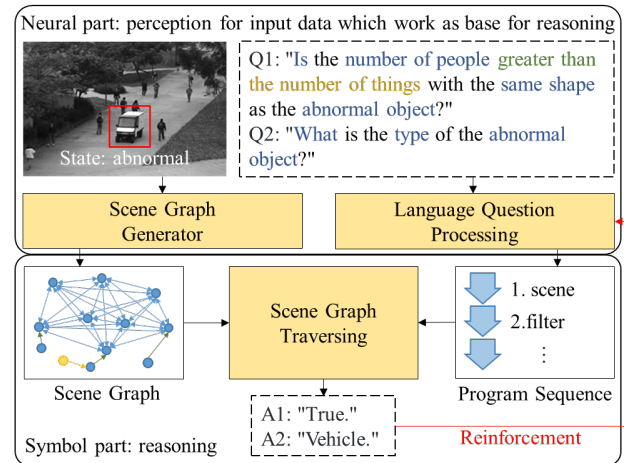


Figure 1. Solving visual question answering tasks in a pedestrian video environment with a combination of reasoning on scene graph traversing and cognition using reinforcement learning.

2022; Szymanowicz et al., 2023; Wu et al., 2021]. This limited scope inhibits the ability to interact with diverse explanations and formulate strategies for risk mitigation. Consequently, the implementation of anomaly detection through question-answering grounded in visual data presents a critical challenge in assessing the explainable video anomaly detection systems.

In the context of question-answering, integrating neural networks for recognition and reasoning between recognized symbols has been effectively employed across various tasks, such as phishing URL detection [Park et al., 2021]. The effectiveness of this strategy is evidenced by the enhanced accuracy in synthetic photographic scenarios, characterized by their simplified data collection and object relationships [Yi et al., 2018]. However, while the inference-based question-answering model has demonstrated validity within a restricted domain, its inference in complex real-world datasets is regarded as a significant problem in cognitive neural networks and question-answering [Amizadeh et al., 2020]. For instance, a broader definition of the perception and question scope is required in pedestrian video environments, where the diversity of object types, relationships, and potential situations is increased. Furthermore, effective integration with previously established recognition methods should be considered.

66 In this paper, the range of possible questions and answers
67 within pedestrian video surveillance environments is repre-
68 sented by predefined programs, each composed of question
69 queries expressed in a domain-specific language. A program
70 must balance usefulness and mapping performance. Before
71 establishing the program, the properties, and relationships be-
72 tween the objects in the pedestrian image, forming the foun-
73 dation of the program, are outlined for all images. This paper
74 defines object properties and relations as shown in Table 1.
75 Questions are translated into a sequence of the program com-
76 mands $\{p_1, \dots, p_s\}$ for simpler execution. The programs are

77 classified into six categories, as illustrated in Table 2, and the
78 format for each input and output is explicitly defined.
79 Moreover, we present a neuro-symbolic approach integrat-
80 ing reinforcement learning for scene graph construction and
81 constant curvature manifold (CCM)-based anomaly detection
82 to resolve the problem (see Figure 1). During image anomaly
83 detection, the image is transposed to a latent space founded
84 on a non-Euclidean framework, enabling the detection of
85 anomalies within video content. Images are transformed into
86 a graph structure, and an explicit inference process for the

Name	Description
Object property	
Shape	Person, bicycle, car, skateboard, wheelchair, cart, truck, others
Size	Small, large
Position	x-y location
Velocity	Computed by comparing x-y coordinates within frames
Object relation	
Relative position	Left, right, in front, behind, over, under
Relative size	Larger than, smaller than
Relative velocity	Faster than, slower than
Numbers	More than, less than
Equal	Same shape, same size, same position, same abnormal

Table 1: Definition of object property and relation on a pedestrian video environment.

Name	Description		Name	Description	
Function	Input	Output	Function	Input	Output
Basic program					
scene	-	Object list	count	Object list	Integer
unique	Object list	Object	exist	Object list	Boolean
relate	Object list	Object list	get_frame	Integer	scene
Filter program					
filter_size	List, size	List	filter_object	List, abnormal	List
filter_shape	List, shape	List	filter_scene	List, Int	List
filter_position	Position	List	filter_frame	List, Int	List
filter_velocity	List, integer	List			
Query program					
query_size	Object	Size	query_velocity	Object	Integer
query_shape	Object	Shape	query_type	Object	List
query_position	Object	Position			
Logic program					
AND	List, List	Object list	OR	List, List	Object list
Sameness program					
same_size	Object	Object list	same_position	Object	Object list
same_shape	Object	Object list	same_velocity	Object	Object list
Compare function					
equal_integer	Int, Int	Boolean	equal_shape	Shape, Shape	Boolean
equal_size	Size, Size	Boolean	less_then	Int, Int	Boolean
equal_color	Col, Col	Boolean	greater_then	Int, Int	Boolean

Table 2: Definition of domain-specific language set for visual question and answering.

Approach	Visual perception	QA processing	Environment
End-to-end neural network	Convolutional neural network	Long short-term memory [Antol <i>et al.</i> , 2015]	Synthetic visual scenes (CLEVR)
		Modular network with encoder-decoder [Hu <i>et al.</i> , 2017]	General objects (MS-COCO)
Neuro-symbolic	Mask R-CNN (Object tables)	Domain-specific language [Yi <i>et al.</i> , 2018]	Synthetic visual scenes (CLEVR)
	Mask R-CNN (Feature vectors)	Quasi-symbolic program execution [Mao <i>et al.</i> , 2019]	
	Faster R-CNN	Differentiable first-order logic [Amizadeh <i>et al.</i> , 2020]	General objects with scene graph (GQA)

Table 3: Previous research for combining deep learning and inference algorithms for visual question and answering.

87 detected anomaly generates an output using an algorithm de-
88 signed to traverse the transformed graph. A neuro-symbolic
89 system then takes on the task of scene-graph reasoning, inte-
90 grating the anomaly detection results from pedestrian video
91 data with programs translated from question sets. Scene
92 graph reasoning, predicated on object properties and relation-
93 ships, generates complex inferences suitable for question-an-
94 swering in realistic settings. Furthermore, our proposed
95 model is designed to learn with reinforcement feedback, us-
96 ing both predicted and original answers. This allows model
97 for tuning toward more accurate answer prediction. We have
98 verified the model's ability to handle complex queries.

99 We illustrate the superior performance of our method com-
100 pared to extant visual question-answering techniques through
101 graph reasoning. Our proposed method is evaluated on five
102 distinct types of questions, finding that it outperforms the
103 convnetional methods in terms of efficacy. Drawing from our
104 experimental results, we posit that the integration of a neuro-
105 symbolic system for scene-graph reasoning with a deep learn-
106 ing-based question-answering mechanism furnishes a level of
107 inference that is highly suited to real-world environments.

108 2 Related Works

109 **Video Anomaly Detection (VAD).** Numerous studies
110 have explored anomaly detection, typically supervised or un-
111 supervised methods. Despite the challenges in data collection,
112 supervised anomaly detection models have been studied due
113 to their superior performance. Shin and Cho, for instance, de-
114 veloped a data augmentation method using a generative ad-
115 versarial network (GAN) [Shin and Cho, 2018]. Conversely,
116 unsupervised anomaly detectors overcome some limitations
117 inherent in supervised models. Zhao *et al.* proposed a model
118 that identifies unusual events in videos via dynamic sparse
119 coding [Zhao *et al.*, 2011], while Liu *et al.* devised a future
120 frame prediction model for anomaly detection [Liu *et al.*,
121 2018]. In the latter model, predicted frames are compared
122 with actual future frames, with large differences indicating an
123 anomaly. Further advancements in the field include end-to-
124 end architecture for one-class classification [Sabokrou *et al.*,
125 2018], and a modified GAN method that learns an encoder
126 simultaneously during training to develop an anomaly detec-
127 tion method [Zenati *et al.*, 2018]. They constructed an adver-
128 sorially learned one-class classifier (ALOCC) composed of
129 an encoder, decoder, and discriminator. However, defining

130 data representation as a simple distribution can result in un-
131 seen data easily following that distribution, potentially caus-
132 ing novel data to be incorrectly classified as normal. To ad-
133 dress this issue, we propose a one-class anomaly detection
134 model based on a constant curvature manifold, a type of non-
135 Euclidean space.

136 **Scene Graph Generation.** Numerous generative methods
137 such as conditional random field (CRF), CNN, RNN, LSTM,
138 and graph neural networks have been developed for scene
139 graphs. CRF-based models like SG-CRF effectively model
140 statistical correlation in visual relationships [Cong *et al.*,
141 2018]. With the advent of neural models for scene graph gen-
142 eration, CNN- and RNN-based models have been explored.
143 BAR-CNN, a CNN-based model, incorporates an attention
144 mechanism but may still suffer from limited receptive neuron
145 regions [Kolesnikov *et al.*, 2019]. The RNN-based Zoom-Net
146 model successfully recognizes complex visual relationships
147 through deep message propagation and interaction between
148 local object features and global predicate features without a
149 linguistic dictionary [Yin *et al.*, 2018]. Despite the success of
150 these models, GCN has proven to be highly effective in graph
151 reasoning tasks, leading to numerous researchers exploring
152 scene graph generation methods based on the graph [Goller
153 and Kuchler, 1996; Gori *et al.*, 2005]. Graph R-CNN, for ex-
154 ample, trims the original scene graph to generate sparse can-
155 didate graph structures [Yang *et al.*, 2018]. In this paper, we
156 adopt Graph-RCNN, considering its efficiency and effective-
157 ness in generating scene graphs within complex scenarios.

158 **Question-Answering.** Table 3 outlines the methods com-
159 bining deep learning and inference algorithms for visual
160 question-answering, categorized by approach, method, and
161 environment. Initial attempts to implement image recognition
162 and processing, as well as mapping with neural networks, de-
163 fined visual question-answering tasks within a synthetic en-
164 vironment [Antol *et al.*, 2015]. Several methods using modu-
165 lar neural networks demonstrated the necessity of distin-
166 guishing between recognition and natural language pro-
167 cessing tasks [Hu *et al.*, 2017]. In the neuro-symbolic ap-
168 proach, which combines inference algorithms with deep
169 learning, symbol grounding and inference methods of objects
170 were examined [Yi *et al.*, 2018]. The research aiming to de-
171 velop domain-specific languages and symbolic processes for
172 query and relationship representation [Amizadeh *et al.*, 2020]

173 demonstrated high-performance question-answering com-
 174 pared to human respondents in synthetic environments [Mao
 175 *et al.*, 2019]. Based on previous studies, this paper redefines
 176 objects and query range to extend the neuro-symbolic ap-
 177 proach to more complex pedestrian video surveillance envi-
 178 ronments and enhances the practicality by incorporating an
 179 anomaly detection module with neural networks.

180 3 Methodology

181 Figure 2 illustrates the proposed method in this paper. An au-
 182 toencoder employing a constant curvature manifold detects
 183 anomalies, and a scene graph is formulated by integrating
 184 anomaly detection outcomes with object detection results
 185 from pedestrian video data. The input questions and associ-
 186 ated programs are mapped onto a supervised long short-term
 187 memory (LSTM) encoder-decoder framework. The set of
 188 programs, extracted from the input questions, executes a fil-
 189 tration process with scene graph traversal. This methodology
 190 produces the outcomes by applying a specific program to a
 191 group of nodes within a scene graph. After that, the model is
 192 trained to generate suitable responses via reinforcement
 193 learning.

194 3.1 Anomaly Detection with Autoencoder

195 Variational autoencoders (VAEs) or generative adversarial
 196 networks (GANs) may not be well-suited for learning com-
 197 plex data representations. We aim to address this issue using
 198 a constant curvature manifold in the latent space. As a result,
 199 even when a novel anomaly appears, it is readily simulated
 200 by the normal variance. This phenomenon can be easily ob-
 201 served in videos with minor changes, where the background
 202 remains fixed while only the object changes. The representa-
 203 tion that our model learns is based on a constant curvature
 204 manifold, which belongs to a class of non-Euclidean spaces.
 205 The d -dimensional CCM \mathcal{T} is a Riemannian manifold
 206 characterized by a constant curvature $\kappa \in \mathbb{R}$. It can be defined
 207 as follows:

$$208 \quad \mathcal{T} = \{\mathbf{x} \in \mathbb{R}^{d+1} \mid \langle \mathbf{x}, \mathbf{x} \rangle = \kappa^{-1}\} \quad (1)$$

209 where $\langle \cdot, \cdot \rangle$ denotes a scalar product. In the CCM, it is de-
 210 fined from the pseudo-Euclidean scalar product:

$$211 \quad \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \begin{pmatrix} I_{d \times d} & \mathbf{0} \\ \mathbf{0} & -1 \end{pmatrix} \mathbf{y} \quad (2)$$

212 where $I_{d \times d}$ is the identity matrix with size of d and T means
 213 transpose operator.

214 The three components are trained to define the data repre-
 215 sentation as the constant curvature manifold. An encoder g is
 216 trained to project the input data into latent space while the
 217 features of data are maintained. A discriminator D learns to
 218 distinguish features $g(x')$ of normal data from other ex-
 219 tracted features. The encoder is forced to project x and x' to
 220 the same point which follows a CCM. Compared to the previ-
 221 ous works [Cruz-Esquivel and Guzman-Zavaleta, 2022;
 222 Wang *et al.*, 2022; Chang *et al.*, 2020], our discriminator has
 223 compressed features as input, resulting in the small size of the
 224 model. In this process, the encoder is forced to project x and
 225 x' to the same point which follows a CCM as shown in Figure
 226 3, and the discriminator is trained to classify $g(x)$, $g(x')$, and
 227 z . Therefore, to explicitly verify whether the trained latent

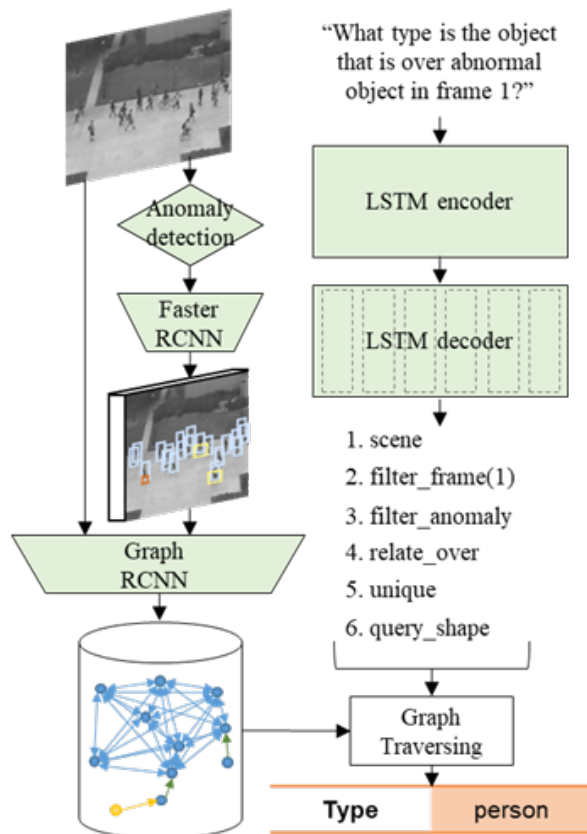


Figure 2. An illustration of the proposed method. Anomalies have been detected with autoencoder with CCM, which is added to the scene graph generated from graph-RCNN. Questions are translated into executable programs with LSTM, and neuro-symbolic integration is applied with scene graph traversal.

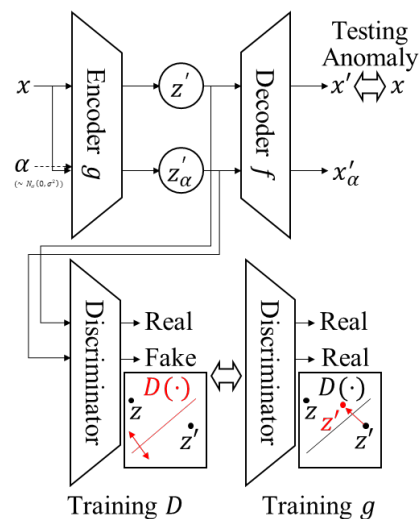


Figure 3. Structure of the anomaly detector with CCM.

228 space forms CCM, we add a membership function $\mu(\cdot)$ as
 229 follows:

Algorithm 1: Anomaly Detection Training Process

Data: hyperparameters

Result: discriminator D , autoencoder AE , encoder g , and decoder f

for $i = 1, \dots, M$ **do**

for $j = 1, \dots, N$ **do**

 Sample x and x' from X and $X + N_\sigma$

 Sample z from CCM

$\mathcal{L}_D \leftarrow \mathcal{L}_D - \frac{\partial \mathcal{L}_D}{\partial D}(x, x', z)$

$\mathcal{L}_{AE} \leftarrow \mathcal{L}_{AE} - \frac{\partial \mathcal{L}_{AE}}{\partial D}(x)$

$\mathcal{L}_g \leftarrow \mathcal{L}_g - \frac{\partial \mathcal{L}_g}{\partial g}(x, x')$

$\mathcal{L}_f \leftarrow \mathcal{L}_f - \frac{\partial \mathcal{L}_f}{\partial f}(x, x')$

$\mathcal{L}_D \leftarrow \mathcal{L}_D - \frac{\partial \mathcal{L}_D}{\partial D}(x, x', z)$

end

end

return f, g and D

$$\mu(z) = \exp\left(-\frac{(\langle z, z \rangle - \kappa^{-1})^2}{2\sigma^2}\right) \quad (3)$$

230 where σ is the hyperparameter to control the scale of CCM.
231 The final forms of the objective function \mathcal{L}_g and \mathcal{L}_D for the
232 encoder and the discriminator are as follows:
233

$$\mathcal{L}_g = \mathbb{E}_{x \sim X} [l(x, f(g(x)))] + \mathbb{E}_{x' \sim X + N_\sigma} [\log(1 - (D(g(x')) + \alpha\mu(x')))] \quad (4)$$

$$\mathcal{L}_D = \mathbb{E}_{x \sim X, z \sim \text{CCM}} \left[\log\left(1 - \frac{D(g(x)) + D(z)}{2} + \alpha\mu(x)\right) \right] + \mathbb{E}_{x' \sim X + N_\sigma} [\log(D(g(x')) + \alpha\mu(x'))] \quad (5)$$

236 l is a binary function to measure the difference between the
237 input data and the reconstructed data. α is a hyperparameter
238 for balance between the outputs of the discriminator (implicit
239 verification) and the membership function (explicit verifica-
240 tion), and α is a hyperparameter for balance.

$$\mathcal{L}_f = \mathbb{E}_{x \sim X, \alpha \sim N_\sigma} [l(x, f(g(x))), l(x, f(g(x + \alpha)))] \quad (6)$$

$$\mathcal{L}_{AE} = \mathbb{E}_{x \sim X} [l(x, f(g(x)))] \quad (7)$$

243 The final objective function for the proposed one-class
244 anomaly detection model is shown in equation (8). To bal-
245 ance each term, we use the hyperparameters β , γ , and δ . Al-
246 gorithm 1 shows the whole training process.

$$\mathcal{L} = \mathcal{L}_D + \beta\mathcal{L}_g + \gamma\mathcal{L}_f + \delta\mathcal{L}_{AE} \quad (8)$$

248 where M is the number of epochs and N is the number of
249 batches.

250 3.2 Scene Graph Generation using Graph R-CNN

251 The scene graph describes the properties and relationships of
252 objects. Given a set of object property categories $\mathcal{C} =$
253 $\{C_1, \dots, C_m\}$ and a set of object relationship categories R , a
254 scene graph is a tuple (O, E) where $O = \{o_1, \dots, o_n\}$ is a set
255 of objects with each o_i , an object that $o_i = \{c_{i1}, c_{i2}, \dots, c_{im}\}$
256 where $c_{ij} \in C_j$, and $E \subseteq O \times R \times O$ is a set of directed

257 edges of the form (o_i, r, o_j) where $o_i, o_j \in O$ and $r \in R$. For
258 this scene graph, object property and relation are defined as
259 the same as those of the questions (Table 1 and Table 2).

260 This paper details the transformation of images from pe-
261 destrian video sequences into scene graphs via a three-step
262 procedure based on the defined scene graph and object prop-
263 erty. Initially, a 3D convolution operation-based autoencoder,
264 factoring in a time axis, determines the normality of the cor-
265 responding image. Subsequently, objects within images are
266 identified where anomalies have been detected with Faster R-
267 CNN. Lastly, the detection outcome is produced through
268 Graph R-CNN in conjunction with the original image.

269 Faster R-CNN undertakes object detection for scene graph
270 parsing. The model proposed integrates anomaly detection
271 results from an autoencoder with constant curvature manifold
272 and an image to detect objects, with an accompanying repre-
273 sentation of their normality.

274 Graph R-CNN, a leading method among scene graph gen-
275 eration algorithms, successfully elucidates the relationships
276 between objects more effectively. It employs a relationship
277 proposition network (RePN) that efficiently manages second-
278 ary potential relationships between image objects and a graph
279 convolutional network (GCN). In this paper, the images of
280 pedestrian video frames are input into corresponding algo-
281 rithm models, trained with the VQA dataset, preserving valid
282 information correlating to the predefined object properties.
283 The resultant scene graph facilitates knowledge representa-
284 tion that can more distinctly express object relationships
285 while safeguarding information on the objects.

286 3.3 Neuro-Symbolic QA with Reinforcement

287 In the proposed model, a question is translated into a se-
288 quence of programs $\{p_1, \dots, p_s\}$ via an LSTM encoder-de-
289 coder structure. Scene graph traversing is performed using
290 the corresponding translated program and the resultant scene
291 graph from section 3.1. Each program operates on a set of
292 nodes in the scene graph. For instance, the "scene" program
293 returns all objects in the current scene. Programs other than
294 "relation" and "scene" do not require any relational infor-
295 mation. Each of these objects in the set is processed by an "if-
296 else" operation, and the resulting output is calculated.

297 Programs associated with relationships require infor-
298 mation about the relationships between objects. This model
299 employs a method of searching through the edges of the scene
300 graph. It verifies whether an edge, corresponding to a con-
301 nection for filtering, is connected to each node for a set of
302 nodes that are used as input when the edge is present. Then,
303 the program for connection proceeds by calculating a set of
304 nodes comprised of target nodes and outputs it. This process
305 is illustrated in Algorithm 2. Through this algorithm, logical
306 reasoning for each program stage becomes feasible.

307 In this paper, a two-stage procedure is implemented to train
308 LSTM, with the aim of elucidating the mapping between a
309 question and its corresponding program. Initially, a few
310 ground truth question-program pairs are extracted from the
311 training set to pretrain the model under direct supervision.
312 Subsequently, the model is paired with a deterministic pro-
313 gram executor. Reinforcement learning is then employed to

Algorithm 2: Scene Graph Traversing Algorithm

Data: scene graph $G = \{O, E\}$ and program sequence $P = \{p_1, p_2, \dots, p_n\}$
Result: traversing result – answer to question
for $p_i \in P$ **do**
 if p_i is “scene” **do**
 $S.push(\phi)$
 else
 if p_i is in “relation” **then**
 $O_{org} = S.pop()$
 $O_{new} = \phi$
 for $o_j \in O_{org}$ **do**
 for $e_{jk} \in E$ where $e_{jk} = \overline{o_j o_k}$ **do**
 if e_k is relation in p_i **then**
 $O_{new} = O_{new} + \{e_k\}$
 end
 end
 end
 $S.push(O_{new})$
 else
 $S = p_i(S)$
 end
 end
end
return $S.pop()$

314 fine-tune the LSTM, utilizing a larger dataset of question-an-
315 swer pairs. Notably, only the accuracy of the execution result
316 is used as the reward signal in this reinforcement learning
317 phase.

318 Employing reinforcement learning for question and answer
319 pairs contributes to generating more precise responses to in-
320 quiries. The decision to respond to the input image and query
321 serves as a reward signal r , wherein the value of $r - b$ is
322 propagated for model learning by establishing a baseline b to
323 inhibit decay. The value of b is initially set to zero and is sub-
324 sequently updated whenever a reward value manifests, shown
325 as equation (9), thereby modulating the learning of extant
326 models.

$$327 \quad \mathbf{b} \leftarrow (\mathbf{1} - \alpha_{decay})\mathbf{r} + \alpha_{decay}\mathbf{b} \quad (9)$$

328

329 4 Experiments

330 4.1 Real-World Pedestrian Video Dataset

331 In order to evaluate the efficacy of the proposed method, we
332 employ the UCSD pedestrian datasets, which are collected
333 from stationary CCTV footage. This data comprises pedestri-
334 ans and various moving objects captured moving in both di-
335 rections. As in Table 1, the object attribute table generated
336 from this data includes combinations of two to four-wheeled
337 vehicles (bicycles, cars, skateboards, wheelchairs, carts, and
338 trucks), along with various backgrounds (wood, roads, and
339 grass).

340 In this paper, question-program pairs are formulated based
341 on the objects within an image. The program is comprised of
342 a sequence of domain-specific languages, as specified in Ta-
343 ble 2, and each pair originates from a predefined template.
344 The queries have been categorized into five types, each typi-
345 fied by its distinct properties.

346 “Querying Attribute” refers to inquiries about an object’s
347 characteristics, including queries concerning the attributes of
348 anomaly objects. “Compare Attribute” involves the compar-
349 ison of attributes between two objects and contains queries
350 that can also determine anomaly attributes. “Exist” and
351 “count” are demarcated as queries about the existence and
352 quantity of specific objects, respectively. Lastly, “compare
353 Number” is classified as a query that contrasts the number of
354 objects across various sets.

355 4.2 Question and Answering Performance

356 Table 4 compares the performance of our method with the
357 conventional question-answering methods, segregated by
358 program type. In scenarios where image recognition using
359 convolutional neural networks is coupled with question pro-
360 cessing using LSTM, and mapped using simple supervised
361 learning, our method seldom misclassifies, exhibiting an ac-
362 curacy of 0.9971. This contrasts starkly with the considerably
363 lower accuracy of 0.6457 when the number of objects is pre-
364 cisely specified in the table of perceived object properties.
365 Furthermore, in complex environments such as pedestrian
366 video sequences, our method, in combination with inference
367 capabilities, outperforms the encoder-decoder approaches
368 based on modular neural networks, achieving an accuracy of
369 0.9991 against the latter’s 0.9232.

Method	Count	Exist	Compare number	Compare attribute	Query attribute	Overall accuracy
CNN-LSTM [Antol et al., 2015]	64.57%	87.44%	53.78%	77.47%	77.47%	72.15%
Mask R-CNN	85.23%	92.93%	83.45%	90.68%	92.68%	88.99%
Module network with Encoder-decoder [Hu et al., 2017]	86.77%	96.61%	86.48%	96.51%	95.27%	92.32%
Ours	99.71%	99.97%	99.96%	99.93%	99.98%	99.91%

Table 4: 10-fold cross-validation of accuracy with other methods by query type.

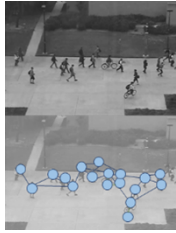
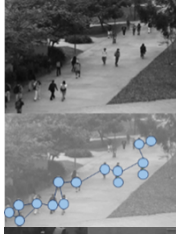
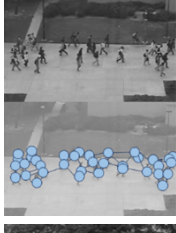

Scene	Query	Question	Program Representation	Answer (P&L)/ Answer (Graph)
	Count	What number of large normal persons are behind the small man?	scene filter size[small] unique relate[behind] filter size[large] filter anomaly[normal] filter shape[person] count	14/15
	Exist	Are there any things in front of the small normal person?	scene filter size[small] filter anomaly[normal] filter shape[person] unique relate[front] exist	False / True
	Compare number	Are there more humans on the left side of the scene than on the right?	scene filter position[left] filter shape[person] scene filter position[right] filter shape[person] greater than	True / False
	Count	What number of large normal persons are behind the small man?	scene filter size[small] unique relate[behind] filter size[large] filter anomaly[normal] filter shape[person] count	7/8

Table 5: Program representation and scene-graph for each case of correct response.

370 For every classification, our method exhibits the highest
371 accuracy. Notably, our method achieves a remarkable accu-
372 racy of 0.9996 in the "compare number" classification, the
373 most challenging category that records the lowest figure for
374 all other algorithms. This demonstrates the potential of the
375 neuro-symbolic approach in tackling problems that could
376 yield varied and complex values, such as numerical compar-
377 ison, and affirms the role of the scene graph in bolstering this
378 capability. In addition, we also report higher accuracy in
379 "query attribute" and "count" categories, which can lead to
380 complex results and require precise determination, respec-
381 tively.

382 Table 5 shows the instances where questions and answers
383 fail when employing data inclusive of object properties and
384 locations, but succeed when scene graph data is employed. In
385 cases where relational information is required for questions,
386 misidentification of relationships frequently occurs based on
387 data with object property and location. However, when the
388 proposed method is adopted for image information represen-
389 tation, the relational information can be more accurately han-
390 dled even with more complex programs.

391 5 Conclusions

392 In this paper, we propose a neuro-symbolic visual question-
393 answering method tailored for pedestrian anomaly video se-
394 quences, which closely resemble real-world environments.
395 This method is facilitated by defining object properties, rela-
396 tionships, and question coverage and incorporating a scene
397 graph generator as well as an anomaly detector. The proposed
398 method demonstrates considerable accuracy of 0.9978 across
399 five types of queries.

400 However, the proposed method's inference algorithm, de-
401 signed to map questions and answers, is implemented as a
402 basic filter algorithm operation. This approach needs valida-
403 tion in the general image field, where object relationships are
404 more complex than in pedestrian video sequences. Particu-
405 larly, as the emergence of various objects tends to complicate
406 the scene graph, thereby increasing computational demand, a
407 learning method that considers computational optimization
408 will be required in the future work.

409 Acknowledgements

410 This work was supported by the Yonsei Fellow Program
411 funded by Lee Youn Jae, and Institute of Information & Com-
412 munications Technology Planning & Evaluation (IITP) grant
413 funded by the Korean government (MSIT) (No. 2020-0-
414 01361, Artificial Intelligence Graduate School Program
415 (Yonsei University); No.2021-0-02068, Artificial Intelli-
416 gence Innovation Hub).

417 References

418 [Amarasinghe *et al.*, 2018] K. Amarasinghe, K. Keney, and
419 M. Manic. Toward explainable deep neural network based
420 anomaly detection. *11th Int. Conf. on Human System In-*
421 *teraction*, pp. 311–317. IEEE, 2018.

422 [Amizadeh *et al.*, 2020] S. Amizadeh, H. Palangi, A.
423 Polozov, Y. Huang, and K. Koishida. Neuro-symbolic
424 visual reasoning: Disentangling. *Int. Conf. on Machine*
425 *Learning*, pp. 279–290. PMLR, 2020.

426 [Antol *et al.*, 2015] S. Antol, A. Agrawal, J. Lu, M. Mitchell,
427 D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual ques-
428 tion answering. *IEEE Int. Conf. on Computer Vision*, pp.
429 2425–2433, 2015.

430 [Chang *et al.*, 2020] Y. Chang, Z. Tu, W. Xie, and J. Yuan.
431 Clustering driven deep autoencoder for video anomaly de-
432 tection. *16th European Conf. on Computer Vision, Part*
433 *XV 16*, pp. 329–345. Springer, 2020.

434 [Cong *et al.*, 2011] Y. Cong, J. Yuan, and J. Liu. Sparse re-
435 construction cost for abnormal event detection. *IEEE*
436 *Conf. on Computer Vision and Pattern Recognition*, pp.
437 3449–3456. IEEE, 2011.

438 [Cong *et al.*, 2018] W. Cong, W. Wang, and W.-C. Lee.
439 Scene graph generation via conditional random fields.
440 *arXiv preprint arXiv:1811.08075*, 2018.

441 [Cruz-Esquivel and Guzman-Zavaleta, 2022] E. Cruz-Es-
442 quivel and Z. J. Guzman-Zavaleta. An examination on au-
443 toencoder designs for anomaly detection in video surveil-
444 lance. *IEEE Access*, 10:6208–6217, 2022.

445 [Fleck and Straßer, 2010] S. Fleck and W. Straßer. Privacy
446 sensitive surveillance for assisted living—A smart camera
447 approach. *Handbook of Ambient Intelligence and Smart*
448 *Environments*, pp. 985–1014, 2010.

449 [Goller and Kuchler, 1996] C. Goller and A. Kuchler. Learn-
450 ing task-dependent distributed representations by back-
451 propagation through structure. *Int. Conf. on Neural Net-*
452 *works*, vol. 1, pp. 347–352. IEEE, 1996.

453 [Gori *et al.*, 2005] M. Gori, G. Monfardini, and F. Scarselli.
454 A new model for learning in graph domains. *IEEE Int.*
455 *Joint Conf. on Neural Networks*, vol. 2, pp. 729–734,
456 2005.

457 [Hu *et al.*, 2017] R. Hu, J. Andreas, M. Rohrbach, T. Darrell,
458 and K. Saenko. Learning to reason: End-to-end module
459 networks for visual question answering. *IEEE Int. Conf.*
460 *on Computer Vision*, pp. 804–813, 2017.

461 [Kolesnikov *et al.*, 2019] A. Kolesnikov, A. Kuznetsova, C.
462 Lampert, and V. Ferrari. Detecting visual relationships us-
463 ing box attention. *IEEE/CVF Int. Conf. on Computer Vi-*
464 *sion Workshops*, pp. ???–???, 2019.

465 [Li *et al.*, 2018] Y. Li, W. Ouyang, B. Zhou, J. Shi, C. Zhang,
466 and X. Wang. Factorizable net: An efficient subgraph-
467 based framework for scene graph generation. *European*
468 *Conf. on Computer Vision*, pp. 335–351, 2018.

469 [Liu *et al.*, 2018] W. Liu, W. Luo, D. Lian, and S. Gao. Future
470 frame prediction for anomaly detection—A new baseline.
471 *IEEE Conf. on Computer Vision and Pattern Recognition*,
472 pp. 6536–6545, 2018.

473 [Mao *et al.*, 2019] J. Mao, C. Gan, P. Kohli, J. B. Tenenbaum,
474 and J. Wu. The neurosymbolic concept learner: Interpret-
475 ing scenes, words, and sentences from natural supervi-
476 sion. *arXiv preprint arXiv:1904.12584*, 2019.

477 [Park *et al.*, 2021] K.-W. Park, S.-J. Bu, and S.-B. Cho. Evo-
478 lutionary optimization of neuro-symbolic integration for
479 phishing URL detection. *Int. Conf. on Hybrid Artificial*
480 *Intelligence Systems*, pp. 88–100. Springer, 2021.

481 [Sabokrou *et al.*, 2018] M. Sabokrou, M. Khalooei, M. Fathy,
482 and E. Adeli. Adversarially learned one-class classifier
483 for novelty detection. *IEEE Conf. on Computer Vision*
484 *and Pattern Recognition*, pp. 3379–3388, 2018.

485 [Shin and Cho, 2018] W. Shin and S.-B. Cho. CCTV image
486 sequence generation and modeling method for video
487 anomaly detection using generative adversarial network.
488 *Int. Conf. on Intelligent Data Engineering and Automated*
489 *Learning*, pp. 457–467. Springer, 2018.

490 [Szymanowicz *et al.*, 2021] S. Szymanowicz, J. Charles, and
491 R. Cipolla. X-man: Explaining multiple sources of anom-
492 alies in video. *IEEE/CVF Conf. on Computer Vision and*
493 *Pattern Recognition*, pp. 3224–3232, 2021.

494 [Szymanowicz *et al.*, 2022] S. Szymanowicz, J. Charles, and
495 R. Cipolla. Discrete neural representations for explainable
496 anomaly detection. *IEEE/CVF Winter Conf. on Applica-*
497 *tions of Computer Vision*, pp. 148–156, 2022.

498 [Tang *et al.*, 2019] K. Tang, H. Zhang, B. Wu, W. Luo, and
499 W. Liu. Learning to compose dynamic tree structures for
500 visual contexts. *IEEE/CVF Conf. on Computer Vision and*
501 *Pattern Recognition*, pp. 6619–6628, 2019.

502 [Wang *et al.*, 2022] L. Wang, H. Tan, F. Zhou, W. Zuo, and
503 P. Sun. Unsupervised anomaly video detection via a dou-
504 ble-flow ConvLSTM variational autoencoder. *IEEE Ac-*
505 *cess*, 10:44278–44289, 2022.

506 [Wu *et al.*, 2021] C. Wu, S. Shao, C. Tunc, P. Satam, and S.
507 Hariri. An explainable and efficient deep learning frame-
508 work for video anomaly detection. *Cluster Computing*,
509 pp. 1–23, 2021.

510 [Xu *et al.*, 2017] D. Xu, Y. Yan, E. Ricci, and N. Sebe. De-
511 tecting anomalous events in videos by learning deep rep-
512 resentations of appearance and motion. *Computer Vision*
513 *and Image Understanding*, 156:117–127, 2017.

- 514 [Yang *et al.*, 2018] J. Yang, J. Lu, S. Lee, D. Batra, and D.
515 Parikh. Graph R-CNN for scene graph generation. *Euro-*
516 *pean Conf. on Computer Vision*, pp. 670–685, 2018.
- 517 [Yi *et al.*, 2018] K. Yi, J. Wu, C. Gan, A. Torralba, P. Kohli,
518 and J. Tenenbaum. Neural-symbolic VQA: Disentangling
519 reasoning from vision and language understanding. *Ad-*
520 *vances in Neural Information Processing Systems*, 31,
521 2018.
- 522 [Yin *et al.*, 2018] G. Yin, L. Sheng, B. Liu, N. Yu, X. Wang,
523 J. Shao, and C. C. Loy. Zoom-net: Mining deep feature
524 interactions for visual relationship recognition. *European*
525 *Conf. on Computer Vision*, pp. 322–338, 2018.
- 526 [Zenati *et al.*, 2018] H. Zenati, C. S. Foo, B. Lecouat, G.
527 Manek, and V. R. Chandrasekhar. Efficient GAN-based
528 anomaly detection. *arXiv preprint arXiv:1802.06222*,
529 2018.
- 530 [Zhao *et al.*, 2011] B. Zhao, L. Fei-Fei, and E. P Xing. Online
531 detection of unusual events in videos via dynamic sparse
532 coding. *IEEE Conf. on Computer Vision and Pattern*
533 *Recognition*, pp. 3313–3320. IEEE, 2011.