

Integrating Non-monotonic Logical Reasoning and Inductive Learning with Deep Learning for Explainable Visual Question Answering

Heather Riley¹ and Mohan Sridharan^{2,*}

¹*Electrical and Computer Engineering, The University of Auckland, Auckland 1023, New Zealand*

²*Intelligent Robotics Lab, School of Computer Science, University of Birmingham, Birmingham B15 2TT, United Kingdom*

Correspondence*:

Mohan Sridharan

m.sridharan@bham.ac.uk

2 ABSTRACT

3 State of the art algorithms for many pattern recognition problems rely on data-driven deep network
4 models. Training these models requires a large labeled dataset and considerable computational
5 resources. Also, it is difficult to understand the working of these learned models, limiting their
6 use in some critical applications. Towards addressing these limitations, our architecture draws
7 inspiration from research in cognitive systems, and integrates the principles of commonsense
8 logical reasoning, inductive learning, and deep learning. As a motivating example of a task that
9 requires explainable reasoning and learning, we consider Visual Question Answering in which,
10 given an image of a scene, the objective is to answer explanatory questions about objects in the
11 scene, their relationships, or the outcome of executing actions on these objects. In this context, our
12 architecture uses deep networks for extracting features from images and for generating answers
13 to queries. Between these deep networks, it embeds components for non-monotonic logical
14 reasoning with incomplete commonsense domain knowledge, and for decision tree induction.
15 It also incrementally learns and reasons with previously unknown constraints governing the
16 domain's states. We evaluated the architecture in the context of datasets of simulated and
17 real-world images, and a simulated robot computing, executing, and providing explanatory
18 descriptions of plans and experiences during plan execution. Experimental results indicate that in
19 comparison with an "end to end" architecture of deep networks, our architecture provides better
20 accuracy on classification problems when the training dataset is small, comparable accuracy with
21 larger datasets, and more accurate answers to explanatory questions. Furthermore, incremental
22 acquisition of previously unknown constraints improves the ability to answer explanatory questions,
23 and extending non-monotonic logical reasoning to support planning and diagnostics improves
24 the reliability and efficiency of computing and executing plans on a simulated robot.

25 **Keywords:** Nonmonotonic logical reasoning, inductive learning, deep learning, visual question answering, commonsense reasoning,
26 human-robot collaboration

1 INTRODUCTION

27 Deep neural network architectures and the associated algorithms represent the state of the art for many
28 perception and control problems in which their performance often rivals that of human experts. These
29 architectures and algorithms are increasingly being used for a variety of tasks such as object recognition,
30 gesture recognition, object manipulation, and obstacle avoidance, in domains such as healthcare,
31 surveillance, and navigation. Common limitations of deep networks are that they are computationally
32 expensive to train, and require a large number of labeled training samples to learn an accurate mapping
33 between input(s) and output(s) in complex domains. It is not always possible to satisfy these requirements,
34 especially in dynamic domains where previously unseen situations often change the mapping between
35 inputs and outputs over time. Also, it is challenging to understand or provide an explanatory description
36 of the observed behavior of a learned deep network model. Furthermore, it is difficult to use domain
37 knowledge to improve the computational efficiency of learning these models or the reliability of the
38 decisions made by these models. Consider a self-driving car on a busy road. Any error made by the car,
39 e.g., in recognizing or responding to traffic signs, can result in serious accidents and make humans more
40 reluctant to use such cars. In general, it is likely that humans interacting with a system designed for complex
41 domains, with autonomy in some components, will want to know why and how the system arrived at
42 particular conclusions; this “explainability” will help designers improve the underlying algorithms and
43 their performance. Understanding the operation of these systems will also help human users build trust in
44 the decisions made by these systems. Despite considerable research in recent years, providing explanatory
45 descriptions of decision making and learning continues to be an open problem in AI.

46 We consider Visual Question Answering (VQA) as a motivating example of a complex task that inherently
47 requires explanatory descriptions of reasoning and learning. Given a scene and a natural language question
48 about an image of the scene, the objective of VQA is to provide an accurate answer to the question.
49 These questions can be about the presence or absence of particular objects in the image, the relationships
50 between these objects, or the potential outcome of executing particular actions on objects in the scene.
51 For instance, a system recognizing and responding to traffic signs on a self-driving car may be posed
52 questions such as “what is the traffic sign in the image?”, or “what is the meaning of this traffic sign?”,
53 and a system controlling a robot arm constructing stable arrangements of objects on a tabletop may be
54 asked “why is this structure unstable?” or “what would make the structure stable?”. We assume that
55 any such questions are provided as (or transcribed into) text, and that answers to questions are also
56 generated as text (that may be converted to speech) using existing software. Deep networks represent
57 the state of the art for VQA, but are characterized by the known limitations described above. We seek
58 to address these limitations by drawing inspiration from research in cognitive systems, which indicates
59 that reliable, efficient, and explainable reasoning and learning can be achieved in complex problems by
60 jointly reasoning with commonsense domain knowledge and learning from experience. Specifically, the
61 architecture described in this paper tightly couples knowledge representation, reasoning, and learning, and
62 exploits the complementary strengths of deep learning, inductive learning, and non-monotonic logical
63 reasoning with incomplete commonsense domain knowledge. We describe the following characteristics of
64 the architecture:

- 65 • For any input image of a scene of interest, Convolutional Neural Networks (CNNs) extract concise
66 visual features characterizing the image.
- 67 • Non-monotonic logical reasoning with the extracted features and incomplete commonsense domain
68 knowledge is used to classify the input image, and to provide answers to explanatory questions about
69 the classification and the scene.

- 70 • Feature vectors that the non-monotonic logical reasoning is unable to classify are used to train a
71 decision tree classifier that is also used to answer questions about the classification during testing.
- 72 • Feature vectors not classified by non-monotonic logical reasoning, along with the output of the decision
73 tree classifier, train a Recurrent Neural Network (RNN) that is used to answer explanatory questions
74 about the scene during testing.
- 75 • Feature vectors not classified by non-monotonic logical reasoning are also used to inductively learn,
76 and subsequently reason with, constraints governing domain states; and
- 77 • Reasoning with commonsense knowledge is expanded (when needed) to support planning, diagnostics,
78 and the ability to answer related explanatory questions.

79 This architecture builds on our prior work on combining commonsense inference with deep learning (Riley
80 and Sridharan, 2018a; Mota and Sridharan, 2019) by introducing the ability to learn and reason with
81 constraints governing domain states, and extending explainable inference with commonsense knowledge to
82 also support planning and diagnostics to achieve any given goal.

83 Although we use VQA as a motivating example, it is not the main focus of our work. State of the
84 art algorithms for VQA focus on generalizing to images from different domains, and are evaluated on
85 benchmark datasets of several thousand images drawn from different domains (Shrestha et al., 2019). Our
86 focus, on the other hand, is on transparent reasoning and learning in any given domain in which a large,
87 labeled dataset is not readily available. Towards this objective, our approach explores the interplay between
88 non-monotonic logical reasoning, incremental inductive learning, and deep learning. We thus neither
89 compare our architecture and algorithms with state of the art algorithms for VQA, nor use large benchmark
90 VQA datasets for evaluation. Instead, we evaluate our architecture's capabilities in the context of: (i)
91 estimating the stability of configurations of simulated blocks on a tabletop; (ii) recognizing different traffic
92 signs in a benchmark dataset of images; and (iii) a simulated robot delivering messages to the intended
93 recipients at different locations. The characteristics of these tasks and domains match our objective. In
94 both domains, we focus on answering explanatory questions about images of scenes and the underlying
95 classification problems (e.g., recognizing traffic signs). In addition, we demonstrate how our architecture
96 can be adapted to enable a robot assisting humans to compute and execute plans, and to answer questions
97 about these plans. Experimental results show that in comparison with an architecture based only on deep
98 networks, our architecture provides: (i) better accuracy on classification problems when the training dataset
99 is small, and comparable accuracy on larger datasets; and (ii) significantly more accurate answers to
100 explanatory questions about the scene. We also show that the incremental acquisition of state constraints
101 improves the ability to answer explanatory questions, and to compute minimal and correct plans.

102 We begin with a discussion of related work in Section 2. The architecture and its components are
103 described in Section 3, with the experimental results discussed in Section 4. Section 5 then describes the
104 conclusions and directions for further research.

2 RELATED WORK

105 State of the art approaches for VQA are based on deep learning algorithms (Jiang et al., 2015; Malinowski
106 et al., 2017; Masuda et al., 2016; Pandhre and Sodhani, 2017; Shrestha et al., 2019; Zhang et al., 2017).
107 These algorithms use labeled data to train neural network architectures with different arrangements of
108 layers and connections between them, capturing the mapping between the inputs (e.g., images, text
109 descriptions) and the desired outputs (e.g., class labels, text descriptions). Although deep networks have
110 demonstrated the ability to model complex non-linear mappings between inputs and outputs for different

111 pattern recognition tasks, they are computationally expensive and require large, labeled training datasets.
112 They also make it difficult to understand and explain the internal representations, identify changes that
113 will improve performance, or to transfer knowledge acquired in one domain to other related domains. In
114 addition, it is challenging to accurately measure performance or identify dataset bias, e.g., deep networks
115 can answer questions about images using question-answer training samples without even reasoning about
116 the images (Jabri et al., 2016; Teney and van den Hengel, 2016; Zhang et al., 2017). There is on-going
117 research on each of these issues, e.g., to explain the operation of deep networks, reduce training data
118 requirements and bias, reason with domain knowledge, and incrementally learn the domain knowledge. We
119 review some of these approaches below, primarily in the context of VQA.

120 Researchers have developed methods to understand the internal reasoning of deep networks and other
121 machine learning algorithms. Selvaraju et al. (2017) use the gradient in the last convolutional layer of a
122 CNN to compute the relative contribution (importance weight) of each neuron to the classification decision
123 made. However, the weights of neurons do not provide an intuitive explanation of the CNN’s operation
124 or its internal representation. Researchers have also developed general approaches for understanding the
125 predictions of any given machine learning algorithm. For instance, Koh and Liang (2017) use second-
126 order approximations of influence diagrams to trace any model’s prediction through a learning algorithm
127 back to the training data in order to identify training samples most responsible for any given prediction.
128 Ribeiro et al. (2016) developed a framework that analyzes any learned classifier model by constructing a
129 interpretable simpler model that captures the essence of the learned model. This framework formulates the
130 task of explaining the learned model, based on representative instances and explanations, as a submodular
131 optimization problem. In the context of VQA, Norcliffe-Brown et al. (2018) provide interpretability by
132 introducing prior knowledge of scene structure as a graph that is learned from observations based on the
133 question under consideration. Object bounding boxes are graph nodes while edges are learned using an
134 attention model conditioned on the question. Mascharka et al. (2018) augment a deep network architecture
135 with an image-space attention mechanism based on a set of composable visual reasoning primitives that
136 help examine the intermediate outputs of each module. Li et al. (2018) introduce a captioning model to
137 generate an image’s description, reason with the caption and the question to construct an answer, and
138 use the caption to explain the answer. However, these algorithms do not support the use of commonsense
139 reasoning to (i) provide meaningful explanatory descriptions of learning and reasoning; (ii) guide learning
140 to be more efficient; or (iii) provide reliable decisions when large training datasets are not available.

141 The training data requirements of a deep network can be reduced by directing attention to data relevant
142 to the tasks at hand. In the context of VQA, Yang et al. (2016) use a Long Short-Term Memory (LSTM)
143 network to map the question to an encoded vector, extract a feature map from the input image using a CNN,
144 and use a neural network to compute weights for feature vectors based on their relevance to the question. A
145 stacked attention network is trained to map the weighted feature vectors and question vector to the answer,
146 prioritizing feature vectors with greater weights. Schwartz et al. (2017) use learned higher-order correlations
147 between various data modalities to direct attention to elements in the data modalities that are relevant to
148 the task at hand. Lu et al. (2016) use information from the question to identify relevant image regions and
149 uses information from the image to identify relevant words in the question. A co-attentional model jointly
150 and hierarchically reasons about the image and the question at three levels, embedding words in a vector
151 space, using one-dimensional CNNs to model information at the phrase level, and using RNNs to encode
152 the entire question. A generalization of this work, a Bilinear Attention Network, considers interactions
153 between all region proposals in the image with all words in the (textual) question (Kim et al., 2018). A
154 Deep Attention Neural Tensor Network for VQA, on the other hand, uses tensor-based representations
155 to discover joint correlations between images, questions, and answers (Bai et al., 2018). The attention

156 module is based on a discriminative reasoning process, and regression with KL-divergence losses improves
157 scalability of training and convergence. Recent work by Anderson et al. (2018) combines top-down and
158 bottom-up attention mechanisms, with the top-down mechanism providing an attention distribution over
159 object proposals provided by the bottom-up mechanism.

160 In addition to reducing the training data requirements, researchers have focused on reducing the number
161 of annotated samples needed for training, and on minimizing the bias in deep network models. In the
162 context of VQA, Lin et al. (2014) iteratively revise a model trained on an initial training set by expanding
163 the training set with image-question pairs involving concepts it is uncertain about, with an “oracle” (human
164 annotator) providing the answers. This approach reduces annotation time, but the database includes just as
165 many images and questions as before. Goyal et al. (2017) provide a balanced dataset with each question
166 associated with a pair of images that require different answers, and provide a counterexample based
167 explanation for each image-question pair. Agrawal et al. (2018), on the other hand, separate the recognition
168 of visual concepts in an image from the identification of an answer to any given question, and include
169 inductive biases to prevent the learned model from relying predominantly on priors in the training data.

170 In computer vision, robotics and other applications, learning from data can often be made more efficient
171 by reasoning with prior knowledge about the domain. In the context of VQA, Wang et al. (2017) reason
172 with knowledge about scene objects to answer common questions about these objects, significantly
173 expanding the range of natural language questions that can be answered without making the training data
174 requirements impractical. However, this approach does not reduce the amount of data required to train the
175 deep network. Furbach et al. (2010) directly use a knowledge base to answer questions and do not consider
176 the corresponding images as inputs. Wagner et al. (2018), on the other hand, use physics engines and prior
177 knowledge of domain objects to realistically simulate and explore different situations. These simulations
178 guide the training of deep network models that anticipate action outcomes and answer questions about
179 all situations. Based on the observation that VQA often requires reasoning over multiple steps, Wu et al.
180 (2018) construct a chain of reasoning for multi-step and dynamic reasoning with relations and objects. This
181 approach iteratively forms new relations between objects using relational reasoning operations, and forms
182 new compound objects using object refining operations, to improve VQA performance. Given the different
183 components of a VQA system, Teney and van den Hengel (2018) present a meta learning approach to
184 separate question answering from the information required for the task, reasoning at test time over example
185 questions and answers to answer any given question. Two meta learning methods adapt a VQA model
186 without the need for retraining, and demonstrate the ability to provide novel answers and support vision
187 and language learning. Rajani and Mooney (2018) developed an ensemble learning approach, Stacking
188 With Auxiliary Features, which combines the results of multiple models using features of the problem as
189 context. The approach considers four categories of auxiliary features, three of which are inferred from
190 image-question pairs while the fourth uses model-specific explanations.

191 Research in cognitive systems indicates that reliable, efficient, and explainable reasoning and learning can
192 be achieved by reasoning with domain knowledge and learning from experience. Early work by Gil (1994)
193 enabled an agent to reason with first-order logic representations and incrementally refined action operators.
194 In such methods, it is difficult to perform non-monotonic reasoning, or to merge new, unreliable information
195 with existing beliefs. Non-monotonic logic formalisms have been developed to address these limitations,
196 e.g., Answer Set Prolog (ASP) has been used in cognitive robotics (Erdem and Patoglu, 2012) and other
197 applications (Erdem et al., 2016). ASP has been combined with inductive learning to monotonically learn
198 causal laws (Otero, 2003), and methods have been developed to learn and revise domain knowledge
199 represented as ASP programs (Balduccini, 2007; Law et al., 2018). Cognitive architectures have also

200 been developed to extract information from perceptual inputs to revise domain knowledge represented in
201 first-order logic (Laird, 2012), and to combine logic and probabilistic representations to support reasoning
202 and learning in robotics (Zhang et al., 2015). However, approaches based on classical first-order logic
203 are not expressive enough, e.g., modeling uncertainty by attaching probabilities to logic statements is not
204 always meaningful. Logic programming methods, on the other hand, do not support one or more of the
205 desired capabilities such as efficient and incremental learning of knowledge, reasoning efficiently with
206 probabilistic components, or generalization as described in this paper. These challenges can be addressed
207 using interactive task learning, a general knowledge acquisition framework that uses labeled examples or
208 reinforcement signals obtained from observations, demonstrations, or human instructions (Chai et al., 2018;
209 Laird et al., 2017). Sridharan and Meadows (2018) developed such a framework to combine non-monotonic
210 logical reasoning with relational reinforcement learning and inductive learning to learn action models to be
211 used for reasoning or learning in dynamic domains. In the context of VQA, there has been interesting work
212 on reasoning with learned symbolic structure. For instance, Yi et al. (2018) present a neural-symbolic VQA
213 system that uses deep networks to infer structural object-based scene representation from images, and to
214 generate a hierarchical (symbolic) program of functional modules from the question. An executor then
215 runs the program on the representation to answer the question. Such approaches still do not (i) integrate
216 reasoning and learning such that they inform and guide each other; or (ii) use the rich domain-specific
217 commonsense knowledge that is available in any application domain.

218 In summary, deep networks represent the state of the art for VQA and many other pattern recognition
219 tasks. Recent surveys on VQA methods indicate that despite considerable research, it is still difficult to
220 use these networks to support efficient learning, intuitive explanations, or generalization to simulated and
221 real-world images (Pandhre and Sodhani, 2017; Shrestha et al., 2019). Our architecture draws on principles
222 of cognitive systems to address these limitations. It tightly couples deep networks with components for
223 non-monotonic logical reasoning with commonsense domain knowledge, and for learning incrementally
224 from samples over which the learned model makes errors. This work builds on our proof of concept
225 architecture that integrated deep learning with commonsense inference for VQA (Riley and Sridharan,
226 2018a). It also builds on work in our research group on using commonsense inference and learned state
227 constraints to guide deep networks that estimate object stability and occlusion in images (Mota and
228 Sridharan, 2019). In comparison with our prior work, we introduce a new component for incrementally
229 learning constraints governing domain states, expand reasoning with commonsense knowledge to support
230 planning and diagnostics, explore the interplay between the architecture’s components, and discuss detailed
231 experimental results.

3 ARCHITECTURE

232 Figure 1 is an overview of our architecture that provides answers to explanatory questions about images of
233 scenes and an underlying classification problem. The architecture seeks to improve accuracy and reduce
234 training effort, i.e., reduce training time and the number of training samples, by embedding non-monotonic
235 logical reasoning and inductive learning in a deep network architecture. We will later demonstrate how
236 the architecture can be adapted to address planning problems on a simulated robot—see Section 3.5. The
237 architecture may be viewed as having four key components that are tightly coupled with each other.

- 238 1. A component comprising CNN-based feature extractors, which are trained and used to map any given
239 image of a scene under consideration to a vector of image features.

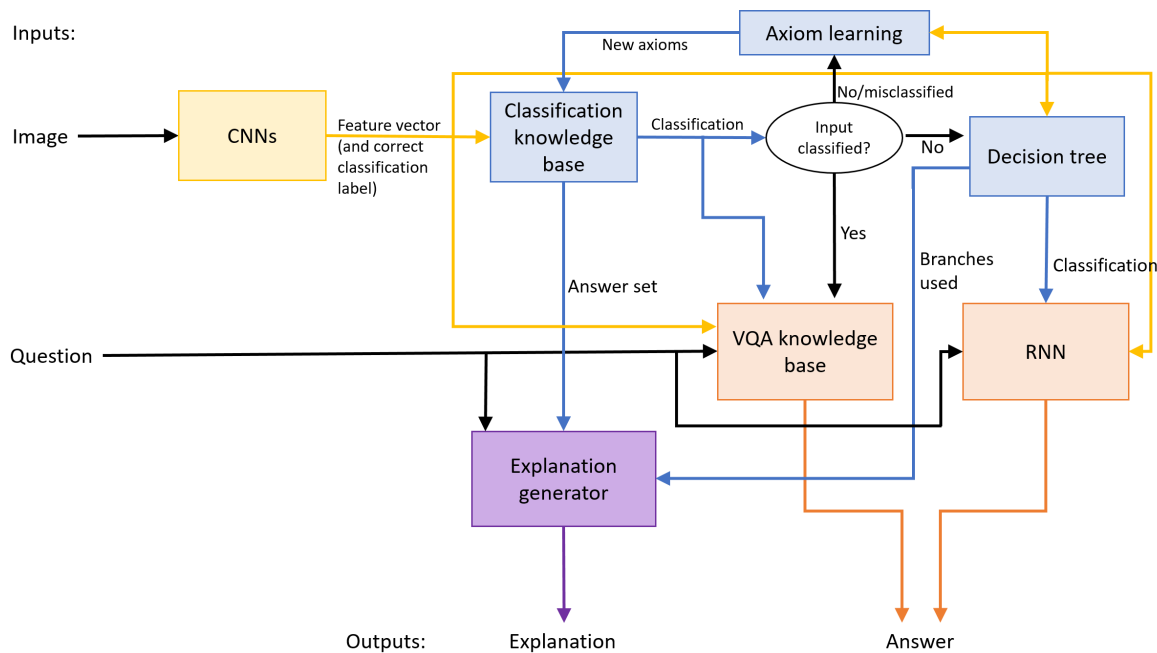


Figure 1. Overview of the architecture that combines the principles of deep learning, non-monotonic logical reasoning, and decision-tree induction.

- 240 2. A component that uses one of two methods to classify the feature vector. The first method uses
 241 non-monotonic reasoning with incomplete domain knowledge and the features to assign a class label
 242 and explain this decision. If the first method cannot classify the image, the second method trains and
 243 uses a decision tree to map the feature vector to a class label and explain the classification.
- 244 3. A component that answers explanatory questions. If non-monotonic logical reasoning is used for
 245 classification, it is also used to provide answers to these questions. If a decision tree is instead used for
 246 classification, an RNN is trained to map the decision tree’s output, the image features, and the question,
 247 to the corresponding answer.
- 248 4. A component that uses the learned decision tree and the existing knowledge base to incrementally
 249 construct and validate constraints on the state of the domain. These constraints revise the existing
 250 knowledge that is used for subsequent reasoning.

251 This architecture exploits the complementary strengths of deep learning, non-monotonic logical reasoning,
 252 and incremental inductive learning with decision trees. Reasoning with commonsense knowledge guides
 253 learning, e.g., the RNN is trained on (and processes) input data that cannot be processed using existing
 254 knowledge. The CNNs and RNN can be replaced by other methods for extracting image features and
 255 answering explanatory questions (respectively). Also, although the CNNs and RNN are trained in an initial
 256 phase in this paper, these models can be revised over time if needed. We hypothesize that embedding non-
 257 monotonic logical reasoning with commonsense knowledge and the incremental updates of the decision
 258 tree, between the CNNs and the RNN, makes the decisions more transparent, and makes learning more time
 259 and sample efficient. Furthermore, the overall architecture and methodology can be adapted to different
 260 domains. In this paper, we will use the following two domains to illustrate and evaluate the architecture’s
 261 components and the methodology.

- 262 1. **Structure Stability (SS):** this domain has different structures, i.e., different arrangements of simulated
 263 blocks of different colors and sizes, on a tabletop—see Figure 2 for some examples. We generated

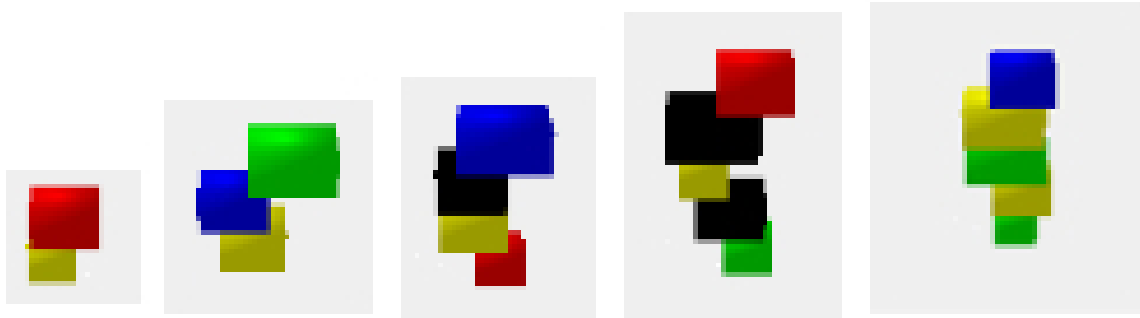


Figure 2. Illustrative images of structures of blocks of different colors and sizes; these images were obtained from a physics-based simulator for the SS domain.



Figure 3. Illustrative images of traffic signs from the BelgiumTS dataset (Timofte et al., 2013).

264 2500 such images using a physics-based simulator. The relevant features of the domain include the
 265 number of blocks, whether the structure is on a lean, whether the structure has a narrow base, and
 266 whether any block is placed such that it is not well balanced on top of the block below. The objective
 267 in this domain is to classify structures as being stable or unstable, and to answer explanatory questions
 268 such as “why is this structure unstable?” and “what should be done to make this structure stable?”.

269 2. **Traffic Sign (TS):** this domain focuses on recognizing traffic signs from images—see Figure 3 for
 270 some examples. We used the BelgiumTS benchmark dataset (Timofte et al., 2013) with ≈ 7000
 271 real-world images (total) of 62 different traffic signs. This domain’s features include the primary
 272 symbol of the traffic sign, the secondary symbol, the shape of the sign, the main color in the middle,
 273 the border color, the sign’s background image, and the presence or absence of a cross (e.g., some
 274 signs have a red or black cross across them to indicate the end of a zone, with the absence of the cross
 275 indicating the zone’s beginning). The objective is to classify the traffic signs and answer explanatory
 276 questions such as “what is the sign’s message?” and “how should the driver respond to this sign?”.

277 In addition to these two domains, Section 3.5 will introduce the **Robot Assistant (RA)** domain, a simulated
 278 domain to demonstrate the use of our architecture for computing and executing plans to achieve assigned
 279 goals. In the *RA* domain, a simulated robot reasons with existing knowledge to deliver messages to target
 280 people in target locations, and to answer explanatory questions about the plans and observed scenes.

281 The focus of our work is on understanding and using the interplay between deep learning, commonsense
 282 reasoning, and incremental learning, in the context of *reliable and efficient scene understanding in any*
 283 *given dynamic domain*. The benchmark VQA datasets and the algorithms, on the other hand, focus on
 284 generalizing across images from different scenarios in different domains, making it difficult to support the
 285 reasoning and learning capabilities of our architecture. We thus do not use these datasets or algorithms in
 286 our evaluation.

287 3.1 Feature Extraction using CNNs

288 The first component of the architecture trains CNNs to map input images to concise features representing
289 the objects of interest in the images. For the SS domain and TS domain, semi-automated annotation was
290 used to label the relevant features in images for training and testing. The selection of these features for
291 each domain was based on domain expertise. In the SS domain, the features of interest are:

- 292 • Number of blocks in structure (number $\in [1, 5]$);
- 293 • Whether the structure is on a lean (true, false);
- 294 • Width of the base block (wide, narrow); and
- 295 • Whether any block is displaced, i.e., not well balanced on top of the block below (true, false).

296 In the TS domain, the features of interest are:

- 297 • Primary symbol in the middle of the traffic sign; 39 primary symbols such as *bumpy_road*,
298 *slippery_road*, *stop*, *left_turn*, and *speed_limit*;
- 299 • Secondary symbol in the traffic sign; 10 secondary symbols such as *disabled*, *car* and *fence*;
- 300 • Shape of the sign; *circle*, *triangle*, *square*, *hexagon*, *rectangle*, *wide rectangle*, *diamond*, or
301 *inverted triangle*;
- 302 • Main color in the middle of the sign; *red*, *white*, or *blue*;
- 303 • Border color at the edge of the sign; *red*, *white*, or *blue*;
- 304 • Background image, e.g., some symbols are placed over a square or a triangle; and
- 305 • Presence of a red or black cross across a sign to indicate a zone's end or invalidity; the sign without the
306 cross indicates the zone's beginning or validity, e.g., a parking sign with a cross implies no parking.

307 To reduce the training data requirements and simplify the training of CNNs, we (i) train a separate CNN for
308 each feature to be extracted from an image; and (ii) start with a basic model for each CNN and incrementally
309 make it more complex as needed. The number of CNNs is thus equal to the number of features to be
310 extracted from each image for any given domain, and the CNN trained for each feature may be different
311 even within a particular domain. The basic CNN model we begin with has an input layer, a convolutional
312 layer, a pooling layer, a dense layer, a dropout layer, and a logit layer, as seen on the left of Figure 4.
313 Additional convolutional and pooling layers are added until the feature extraction accuracy converges or
314 exceeds a threshold (e.g., $\geq 90\%$). Our architecture also includes the option of fine-tuning previously
315 trained CNN models instead of starting from scratch. The right side of Figure 4 shows a CNN model
316 learned in our example domains, which has three convolutional layers and pooling layers. We trained and
317 validated these CNNs in an initial phase, and used them for evaluation. Our code for constructing these
318 CNNs for features (in our example domains) is in our repository (Riley and Sridharan, 2018b).

319 3.2 Classification using Non-monotonic Logical Reasoning or Decision Trees

320 The feature vector extracted from an image is used for decision making. In the SS domain and TS domain,
321 decisions take the form of assigning a class label to each feature vector¹. The second component of our
322 architecture performs this task using one of two methods: (i) non-monotonic logical inference using ASP;
323 or (ii) a classifier based on a learned decision tree. We describe these two methods below.

¹ In the RA domain discussed in Section 3.5, decision making also includes planning and diagnostics.

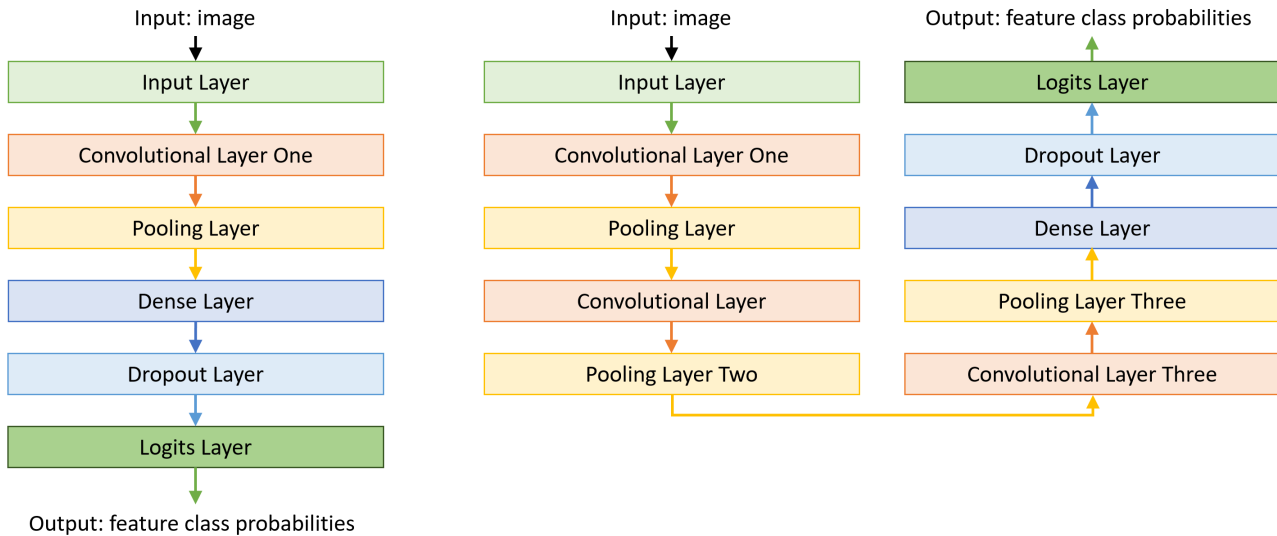


Figure 4. Basic CNN model used for extracting each feature in our architecture. CNNs for individual features may end up with different numbers of convolutional and pooling layers.

ASP-based Inference with Commonsense Knowledge: The first step in reasoning with incomplete commonsense domain knowledge is the representation of this knowledge. In our architecture, an *action language* is used to describe the dynamics of any domain under consideration. Action languages are formal models of parts of natural language used for describing transition diagrams of dynamic systems. Our architecture uses action language \mathcal{AL}_d (Gelfond and Incezan, 2013), with a *sorted signature* Σ that can be viewed as the *vocabulary* used to describe the domain’s transition diagram. The signature Σ comprises *basic sorts*, which are similar to *types* in a programming language, *statics*, i.e., domain attributes whose values do not change over time, *fluents*, i.e., domain attributes whose values can change over time, and *actions*. The domain’s fluents can be *basic*, i.e., those that obey the laws of inertia and are changed directly by actions, or *defined*, i.e., those that do not obey the laws of inertia and are defined by other attributes. A domain attribute or its negation is a *literal*; if all its variables are ground, it is a *ground literal*. \mathcal{AL}_d allows three types of statements: *causal law*, *state constraint* and *executability condition*.

$$\begin{aligned}
 a \text{ causes } l_b \text{ if } p_0, \dots, p_m & \quad (\text{Causal law}) \\
 l \text{ if } p_0, \dots, p_m & \quad (\text{State constraint}) \\
 \text{impossible } a_0, \dots, a_k \text{ if } p_0, \dots, p_m & \quad (\text{Executability condition})
 \end{aligned}$$

324 where a is an action, l is a literal, l_b is a basic literal, and p_0, \dots, p_m are domain literals.

325 The domain representation (i.e., the knowledge base) comprises a *system description* \mathcal{D} , which is a
 326 set of statements of \mathcal{AL}_d , and a *history* \mathcal{H} . \mathcal{D} comprises a sorted signature Σ and axioms describing
 327 the domain dynamics. For instance, in the SS domain, Σ includes basic sorts such as *structure*, *color*,
 328 *size*, and *attribute*; the basic sorts of the TS domain include *main_color*, *other_color*, *main_symbol*,
 329 *other_symbol*, *shape*, *cross* etc. The sort *step* is also in Σ to support temporal reasoning over time steps.
 330 The statics and fluents in the SS domain include:

This is a provisional file, not the final typeset article

$$\begin{aligned} & num_blocks(structure, num), block_color(block, color), block_size(block, size) \\ & block_displaced(structure), stable(structure) \end{aligned} \quad (1)$$

331 which correspond to the image features extracted in the domain, and are described in terms of their
332 arguments' sorts. In a similar manner, statics and fluents of the TS domain include:

$$\begin{aligned} & primary_symbol(sign, main_symbol), primary_color(sign, main_color) \\ & secondary_symbol(sign, other_symbol), secondary_color(sign, other_color) \\ & sign_shape(shape), background_image(image) \end{aligned} \quad (2)$$

333 In both domains, signature Σ includes a predicate $holds(fluent, step)$, which implies that a particular
334 fluent holds true at a particular time step. As stated above, Σ for a dynamic domain typically includes
335 actions that cause state transitions, but this capability is not needed to answer explanatory questions about
336 specific scenes and the underlying classification problem in our (SS, TS) domains. For ease of explanation,
337 we thus temporarily disregard the modeling of actions, and their preconditions and effects. We will revisit
338 actions in Section 3.5 when we consider planning tasks in the RA domain.

Given a signature Σ for a domain, a *state* of the domain is a collection of ground literals, i.e., statics, fluents, actions and relations with values assigned to their arguments—for more details, please see (Gelfond and Kahl, 2014; Sridharan et al., 2019). The axioms of \mathcal{D} are defined in terms of the signature and govern domain dynamics; this typically includes a distributed representation of the constraints related to domain actions, i.e., causal laws and executability conditions that define the preconditions and effects of actions, and constraints related to states, i.e., state constraints. In the SS domain and TS domain, axioms govern the belief about domain states; we will discuss axioms related to actions in Section 3.5 when we discuss the RA domain. Specifically, the axioms of the SS domain include state constraints such as:

$$\neg stable(S) \text{ if } block_displaced(S) \quad (3a)$$

$$stable(S) \text{ if } num_blocks(S, 2), \neg structure_type(S, lean) \quad (3b)$$

339 where Statement 3(a) says that any structure with a block that is displaced significantly is unstable, and
340 Statement 3(b) says that any pair of blocks without a significant lean is stable.

Axioms of the TS domain include statements such as:

$$\begin{aligned} sign_type(TS, no_parking) \text{ if } primary_color(TS, blue), primary_symbol(TS, blank), \\ cross(TS), shape(TS, circle) \end{aligned} \quad (4a)$$

$$\begin{aligned} sign_type(TS, stop) \text{ if } primary_color(TS, red), primary_symbol(TS, stoptext), \\ shape(TS, octagon) \end{aligned} \quad (4b)$$

341 where Statement 4(a) implies that a blue, blank, circular traffic sign with a cross across it is a no parking
342 sign. Statement 4(b) implies that a red, octagon-shaped traffic sign with the text “stop” is a stop sign.

343 The history \mathcal{H} of a dynamic domain is usually a record of fluents observed to be true or false at a particular
 344 time step, i.e., $obs(fluent, boolean, step)$, and the successful execution of an action at a particular time
 345 step, i.e., $hpd(action, step)$; for more details, see (Gelfond and Kahl, 2014). The domain knowledge in
 346 many domains often includes default statements that are true in all but a few exceptional circumstances.
 347 For example, we may know in the SS domain that “structures with two blocks of the same size are usually
 348 stable”. To encode such knowledge, we use our recent work that expanded the notion of history to represent
 349 and reason with defaults describing the values of fluents in the initial state (Sridharan et al., 2019).

350 Key tasks of an agent equipped with a system description \mathcal{D} and history \mathcal{H} include reasoning with this
 351 knowledge for inference, planning and diagnostics. In our architecture, these tasks are accomplished
 352 by translating the domain representation to a program $\Pi(\mathcal{D}, \mathcal{H})$ in CR-Prolog, a variant of ASP that
 353 incorporates consistency restoring (CR) rules (Balduccini and Gelfond, 2003). In this paper, we use the
 354 terms “ASP” and “CR-Prolog” interchangeably. ASP is a declarative programming paradigm designed
 355 to represent and reason with incomplete commonsense domain knowledge. It is based on stable model
 356 semantics, and supports *default negation* and *epistemic disjunction*. For instance, unlike “ $\neg a$ ”, which
 357 implies that *a is believed to be false*, “*not a*” only implies *a is not believed to be true*. Also, unlike “ $p \vee \neg p$ ”
 358 in propositional logic, “*p or not p*” is not tautological. Each literal can thus be true, false or unknown, and
 359 the agent reasoning with domain knowledge does not believe anything that it is not forced to believe.
 360 ASP can represent recursive definitions, defaults, causal relations, special forms of self-reference, and
 361 language constructs that occur frequently in non-mathematical domains, and are difficult to express in
 362 classical logic formalisms (Baral, 2003; Gelfond and Kahl, 2014). Unlike classical first-order logic, ASP
 363 supports non-monotonic logical reasoning, i.e., it can revise previously held conclusions or equivalently
 364 reduce the set of inferred consequences, based on new evidence—this ability helps the agent recover from
 365 any errors made by reasoning with incomplete knowledge. ASP and other paradigms that reason with
 366 domain knowledge are often criticized for requiring considerable (if not complete) prior knowledge and
 367 manual supervision, and for being unwieldy in large, complex domains. However, modern ASP solvers
 368 support efficient reasoning in large knowledge bases with incomplete knowledge, and are used by an
 369 international research community for cognitive robotics (Erdem and Patoglu, 2012; Zhang et al., 2015)
 370 and other applications (Erdem et al., 2016). For instance, recent work has demonstrated that ASP-based
 371 non-monotonic logical reasoning can be combined with: (i) probabilistic reasoning for reliable and efficient
 372 planning and diagnostics (Sridharan et al., 2019); and (ii) relational reinforcement learning and active
 373 learning methods for interactively learning or revising commonsense domain knowledge based on input
 374 from sensors and humans (Sridharan and Meadows, 2018).

In our architecture, the automatic translation from statements in \mathcal{AL}_d to the program Π is based on a custom-designed script². The resultant program Π includes the signature and axioms of \mathcal{D} , inertia axioms, reality checks, closed world assumptions for defined fluents and actions, and observations, actions, and defaults from \mathcal{H} . For instance, Statements 3(a-b) are translated to:

$$\neg stable(S) \leftarrow block_displaced(S) \quad (5a)$$

$$stable(S) \leftarrow num_blocks(S, 2), \neg structure_type(S, lean) \quad (5b)$$

375 In addition, features extracted from an input image (to be processed) are encoded as the initial state of
 376 the domain in Π . Each *answer set* of $\Pi(\mathcal{D}, \mathcal{H})$ then represents the set of beliefs of an agent associated

² An independent group of researchers has developed a general-purposed software to automatically translate any description in \mathcal{AL}_d to the corresponding CR-Prolog program; we expect this software to be made publicly available soon.

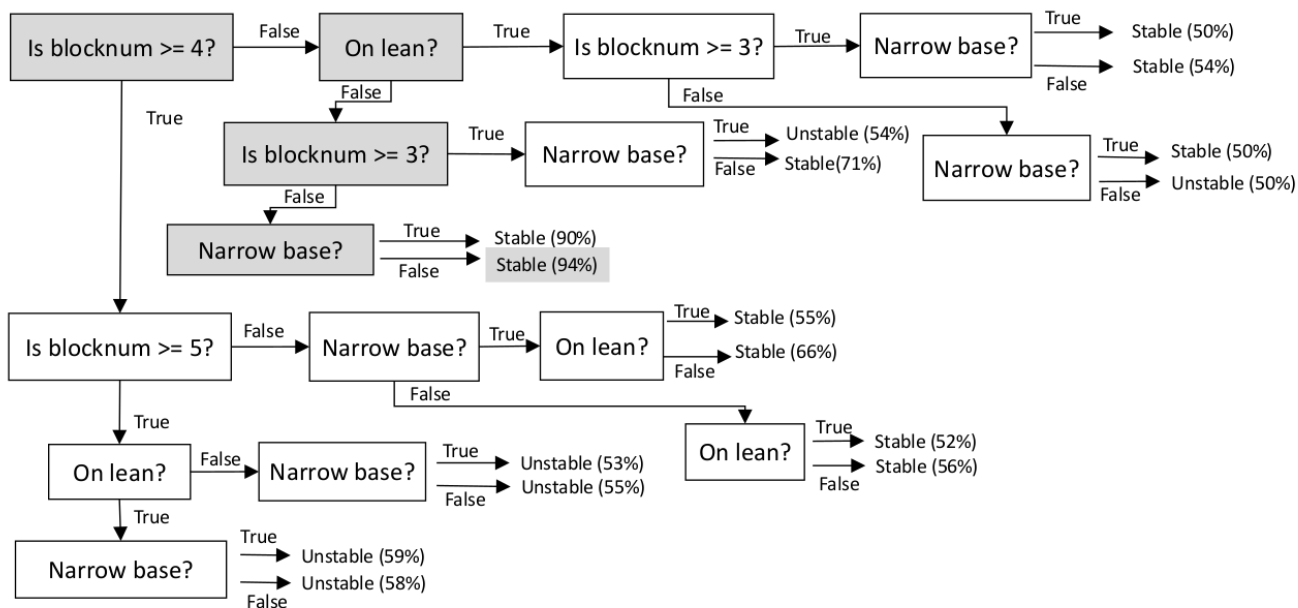


Figure 5. Example of part of a decision tree constructed from labeled samples and used for classification in the SS domain. The nodes used to classify a particular example are highlighted. Each leaf shows a class label and indicates the proportion of the labeled examples (at the leaf) that correspond to this label.

377 with this program. Algorithms for computing entailment, and for planning and diagnostics, reduce these
 378 tasks to computing answer sets of CR-Prolog programs. We compute answer sets of CR-Prolog programs
 379 using the SPARC system (Balai et al., 2013). The CR-Prolog programs for our example domains are in our
 380 open-source software repository (Riley and Sridharan, 2018b). For the classification task in our example
 381 domains, the relevant literals in the answer set provide the class label and an explanatory description of
 382 the assigned label (see Section 3.3); we will consider the planning task in Section 3.5. The accuracy of
 383 the inferences drawn from the encoded knowledge depends on the accuracy and extent of the knowledge
 384 encoded, but encoding comprehensive domain knowledge is difficult. The decision of what and how much
 385 knowledge to encode is made by the designer.

386 **Decision Tree Classifier:** If ASP-based inference cannot classify the feature vector extracted from an
 387 image, the feature vector is mapped to a class label using a decision tree classifier learned from labeled
 388 training examples. In a decision tree classifier, each node is associated with a question about the value of a
 389 particular feature, with the child nodes representing the different answers to the question, i.e., the possible
 390 values of the feature. Each node is also associated with samples that satisfy the corresponding values of the
 391 features along the path from the root node to this node. We use a standard implementation of a decision
 392 tree classifier (Duda et al., 2000). This implementation uses the Gini measure to compute information
 393 gain (equivalently, the reduction in entropy) that would be achieved by splitting an existing node based on
 394 each feature that has not already been used to create a split in the tree. Among the features that provide a
 395 significant information gain, the feature that provides the maximum information gain is selected to split the
 396 node. If none of the features would result in any significant information gain, this node becomes a leaf
 397 node with a class label that matches a majority of the samples at the node.

398 The decision tree’s search space is quite specific since it only considers samples that could not be
 399 classified by ASP-based reasoning. The decision tree does not need to generalize as much as it would have
 400 to if it had to process every training (or test) sample in the dataset. Also, although overfitting is much less

401 likely, we still use pruning to minimize the effects of overfitting. Figure 5 shows part of a learned decision
 402 tree classifier; specific nodes used to classify a particular example are highlighted to indicate that 94% of
 403 the observed examples of structures that have fewer than three blocks, do not have a significant lean, and
 404 do not have a narrow base, correspond to stable structures. These “active” nodes along any path in the
 405 decision tree that is used to classify an example can be used to explain the classification outcome in terms
 406 of the values of particular features that were used to arrive at the class label assigned to a specific image
 407 under consideration.

408 3.3 Answering Explanatory Questions

409 The third component of the architecture provides two methods for answering explanatory questions. The
 410 available inputs are the (i) question; (ii) vector of features extracted from the image under consideration;
 411 and (iii) classification output. The human designer also provides pre-determined templates for questions
 412 and their answers. In our case, we use a controlled vocabulary, templates based on language models and
 413 parts of speech for sentences, and existing software for natural language processing. Any given question
 414 is transcribed using the controlled vocabulary, parsed (e.g., to obtain parts of speech), and matched with
 415 the templates to obtain a relational representation. Recall that questions in the SS domain are of the form:
 416 “is this structure stable/unstable?” and “what is making this structure stable/unstable?”. These questions
 417 can be translated into relational statements such as $stable(S)$ or $\neg stable(S)$ and used as a question, or
 418 as the desired consequence, during inference or in a search process. In a similar manner, questions in
 419 the TS domain such as: “what sign is this?” and “what is the sign’s message?” can be translated into
 420 $sign_type(S, sign)$ and used for subsequent processing.

421 The first method for answering explanatory questions is based on the understanding that if the feature
 422 vector extracted from the image is processed successfully using ASP-based reasoning, it is also possible
 423 to reason with the existing knowledge to answer explanatory questions about the scene. To support such
 424 question answering, we need to revise the signature Σ in the system description \mathcal{D} of the domain. For
 425 instance, we add sorts such as $query_type$, $answer_type$, and $query$ to encode different types of queries
 426 and answers. We also introduce suitable relations to represent questions, answers to these questions, and
 427 more abstract attributes, e.g., of structures of blocks, traffic signs etc.

In addition to the signature, we also augment the axioms in \mathcal{D} to support reasoning with more abstract attributes, and to help construct answers to questions. For instance, we can include an axiom such as:

$$\begin{aligned} many_blocks(S) \leftarrow unstable(S), \neg base(S, narrow), \\ \neg struc_type(S, lean), \neg block_displaced(S) \end{aligned} \quad (6)$$

428 which implies that if a structure (of blocks) is not on a narrow base, does not have a significant lean, and
 429 does not have blocks significantly displaced, any instability in the structure implies (and is potentially
 430 because) there are too many blocks in the structure. Once the ASP program $\Pi(\mathcal{D}, \mathcal{H})$ has been revised as
 431 described above, we can compute answer set(s) of this program to obtain the beliefs of the agent associated
 432 with this program. For any given question, the answer set(s) are parsed based on the known controlled
 433 vocabulary and templates (for questions and answers) to extract relevant literals—these literals are included
 434 in the corresponding templates to construct answers to explanatory questions. These answers can also be
 435 converted to speech using existing software.

436 The second method for answering explanatory questions is invoked if the decision tree is used to process
 437 (i.e., classify in the context of the SS domain and TS domain) the vector of image features. The inability

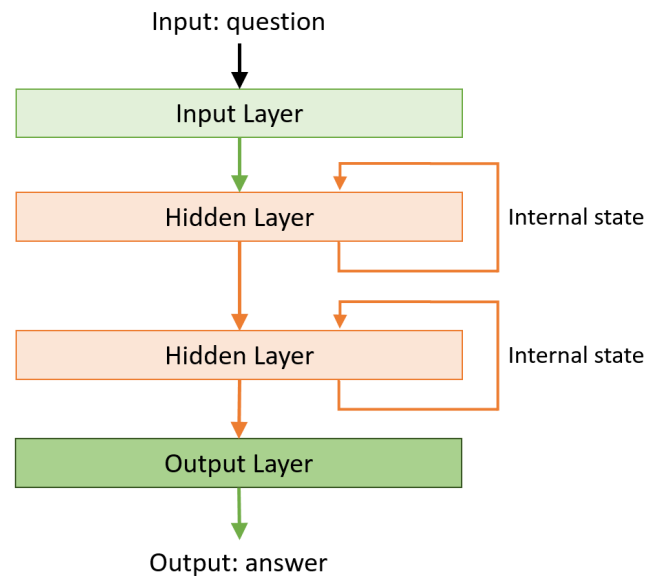


Figure 6. Example of the basic RNN used to construct explanations. The RNN learned for the example domains has 26 – 30 hidden layers.

438 to classify the feature vector through ASP-based reasoning is taken to imply that the encoded domain
 439 knowledge is insufficient to answer explanatory questions about the scene. In this case, an LSTM network-
 440 based RNN is trained and used to answer the explanatory questions. The inputs are the feature vector,
 441 classification output, and a vector representing the transcribed and parsed query. The output (provided
 442 during training) is in the form of answers in the predetermined templates. Similar to the approach used
 443 in Section 3.2, the RNN is built incrementally during training. We begin with one or two hidden layer(s),
 444 as shown in Figure 6, and add layers as long as it results in a significant increase in the accuracy. We
 445 also include the option of adding a stack of LSTMs if adding individual layers does not improve accuracy
 446 significantly. In our example domains, the RNN constructed to answer explanatory questions had as many
 447 as 26 – 30 hidden layers and used a softmax function to provide one of about 50 potential answer types.
 448 An example of the code used to train the RNN is available in our repository (Riley and Sridharan, 2018b).

449 3.4 Learning State Constraints

450 The components of the architecture described so far support reasoning with commonsense knowledge,
 451 learned decision trees, and deep networks, to answer explanatory questions about the scene and an
 452 underlying classification problem. In many practical domains, the available knowledge is incomplete, the
 453 number of labeled examples is small, or the encoded knowledge changes over time. The decisions made
 454 by the architecture can thus be incorrect or sub-optimal, e.g., a traffic sign can be misclassified or an
 455 ambiguous answer may be provided to an explanatory question. The fourth component of our architecture
 456 seeks to address this problem by supporting incremental learning of domain knowledge. Our approach is
 457 inspired by the inductive learning methods mentioned in Section 2, e.g., Sridharan and Meadows (2018) use
 458 relational reinforcement learning and decision tree induction to learn domain axioms. The work described
 459 in this paper uses decision tree induction to learn constraints governing domain states. The methodology
 460 used in this component, in the context of VQA, is as follows:

- 461 1. Identify training examples that are not classified, or are classified incorrectly, using the existing
 462 knowledge. Recall that this step is accomplished by the component described in Section 3.2, which

463 processes each training example using the existing knowledge encoded in the CR-Prolog program, in
 464 an attempt to assign a class label to the example.

465 2. Train a decision tree using the examples identified in Step-1 above. Recall that this step is also
 466 accomplished by the component described in Section 3.2.

3. Identify paths in the decision tree (from root to leaf) such that (i) there are a sufficient number of examples at the leaf, e.g., 10% of the training examples; and (ii) all the examples at the leaf have the same class label. Since the nodes correspond to checks on the values of domain features, the paths will correspond to combinations of partial state descriptions and class labels that have good support among the labeled training examples. Each such path is translated into a candidate axiom. For instance, the following are two axioms identified by this approach in the SS domain:

$$\neg stable(S) \leftarrow num_blocks(S, 3), base(S, wide), \quad (7a)$$

$$struc_type(S, lean)$$

$$\neg stable(S) \leftarrow num_blocks(S, 3), base(S, narrow), \quad (7b)$$

$$struc_type(S, lean)$$

4. Generalize candidate axioms if possible. For instance, if one candidate axiom is a over-specification of another existing axiom, the over-specified version is removed. In the context of the axioms in Statement 7(a-b), the second literal represents redundant information, i.e., if a structure with three blocks has a significant lean, it is unstable irrespective of whether the base of the structure is narrow or wide. Generalizing over these two axioms results in the following candidate axiom:

$$\neg stable(S) \leftarrow num_blocks(S, 3), struc_type(S, lean) \quad (8)$$

467 which only includes the literals that encode the essential information.

468 5. Validate candidate axioms one at a time. To do so, the candidate axiom is added to the CR-Prolog
 469 program encoding the domain knowledge. A sufficient number of training examples (e.g., 10% of the
 470 dataset, as before) relevant to this axiom, i.e., the domain features encoded by the examples should
 471 satisfy the body of the axiom, are drawn randomly from the training dataset. If processing these
 472 selected examples with the updated CR-Prolog program results in misclassifications, the candidate
 473 axiom is removed from further consideration.

474 6. Apply sanity checks to the validated axioms. The validated axioms and existing axioms are checked to
 475 remove over-specifications and retain the most generic version of any axiom. Axioms that pass these
 476 sanity checks are added to the CR-Prolog program and used for subsequent reasoning.

477 Section 4.3 examines the effect of such learned constraints on classification and VQA performance.

478 3.5 Planning with Domain Knowledge

479 The description of the architecture’s components has so far focused on classification and VQA, and
 480 reasoning has been limited to inference with knowledge. However, the architecture is also applicable to
 481 planning (and diagnostics) problems. Consider the **RA domain** in which a simulated robot has to navigate
 482 and deliver messages to particular people in particular places, and to answer explanatory questions, i.e., the
 483 domain includes aspects of planning and VQA. Figure 7 depicts this domain and a simulated scenario in it;
 484 semantic labels of the offices and rooms are shown in the upper half.



Figure 7. Block diagram and a simulated scenario in the RA domain in which the robot has to deliver messages to people in target locations.

485 A robot planning and executing actions in the real world has to account for the uncertainty in sensing and
 486 actuation. In other work, we addressed this issue by coupling ASP-based coarse-resolution planning with
 487 probabilistic fine-resolution planning and execution (Sridharan et al., 2019). In this paper, we temporarily
 488 abstract away such probabilistic models of uncertainty to thoroughly explore the interplay between
 489 reasoning and learning, including the effect of added noise in sensing and actuation (in simulation).

To support planning, the signature Σ of system description \mathcal{D} has basic sorts such as: *place*, *robot*, *person*, *object*, *entity*, *status*, and *step*, which are arranged hierarchically, e.g., *robot* and *person* are subsorts of *agent*, and *agent* and *object* are subsorts of *entity*. Σ also includes ground instances of sorts, e.g., *office*, *workshop*, *kitchen*, and *library* are instances of *place*, and *Sarah*, *Bob*, *John*, and *Sally* are instances of *person*. As before, domain attributes and actions are described in terms of the sorts of their arguments. The fluents include $loc(agent, place)$, which describes the location of the robot and people in the domain, and $message_status(message_id, person, status)$, which denotes whether a particular message has been delivered (or remains undelivered) to a particular person. Static attributes include relations such as $next_to(place, place)$ and $work_place(person, place)$ to encode the arrangement of places and the work location of people (respectively) in the domain. Actions of the domain include:

$$\begin{aligned}
 &move(robot, place) \\
 &deliver(robot, message_id, person)
 \end{aligned} \tag{9}$$

490 which move the robot to a particular place, and cause a robot to deliver a particular message to a particular
 491 person (respectively). For ease of explanation, we assume that the locations of people are defined fluents
 492 whose values are determined by external sensors, and that the locations of objects are static attributes; as a
 493 result, we do not consider actions that change the value of these attributes. The signature Σ also includes
 494 (as before) the relation $holds(fluent, step)$ to imply that a particular fluent is true at a particular time step.

Axioms of the RA domain capture the domain’s dynamics. These axioms include causal laws, state constraints and executability conditions encoded as statements in \mathcal{AL}_d such as:

$$\text{move}(\text{rob}_1, L) \text{ causes } \text{loc}(\text{rob}_1, L) \quad (10a)$$

$$\text{deliver}(\text{rob}_1, ID, P) \text{ causes } \text{message_status}(ID, P, \text{delivered}) \quad (10b)$$

$$\text{loc}(P, L) \text{ if } \text{work_place}(P, L), \text{ not } \neg \text{loc}(P, L) \quad (10c)$$

$$\neg \text{loc}(T, L_2) \text{ if } \text{loc}(T, L_1), L_1 \neq L_2 \quad (10d)$$

$$\text{impossible } \text{deliver}(\text{rob}_1, ID, P) \text{ if } \text{loc}(\text{rob}_1, L_1), \text{loc}(P, L_2), \\ L_1 \neq L_2 \quad (10e)$$

$$\text{impossible } \text{move}(\text{rob}_1, L) \text{ if } \text{loc}(\text{rob}_1, L) \quad (10f)$$

495 where Statement 10(a) states that executing a move action causes the robot’s location to be the target place;
 496 Statement 10(b) states that executing a deliver action causes the message to be delivered to the desired
 497 person; Statement 10(c) is a constraint stating that unless told otherwise the robot expects (by default) a
 498 person to be in her/his place of work; Statement 10(d) is a constraint stating that any thing can be in one
 499 place at at time; Statement 10(e) implies that a robot cannot deliver a message to an intended recipient if
 500 the robot and the person are not in the same place; and Statement 10(f) states that a robot cannot move to a
 501 location if it is already there.

502 As described in Section 3.2, the domain history is a record of observations (of fluents), the execution
 503 of actions, and the values of fluents in the initial state. Also, planning (similar to inference) is reduced to
 504 computing answer set(s) of the program $\Pi(\mathcal{D}, \mathcal{H})$ after including some helper axioms for computing a
 505 minimal sequence of actions; for examples, please see (Gelfond and Kahl, 2014; Sridharan et al., 2019). If
 506 the robot’s knowledge of the domain is incomplete or incorrect, the computed plans may be suboptimal or
 507 incorrect. The approach described in Section 3.4 can then be used to learn the missing constraints; we will
 508 explore the interplay between learning and planning in Section 4.4.

4 EXPERIMENTAL SETUP AND RESULTS

509 In this section, we describe the results of experimentally evaluating the following hypotheses about the
 510 capabilities of our architecture:

- 511 • **H1:** for the underlying classification problem, our architecture outperforms an architecture based on
 512 just deep networks for small training datasets, and provides comparable performance as the size of the
 513 dataset increases;
- 514 • **H2:** in the context of answering explanatory questions, our architecture provides significantly better
 515 performance in comparison with an architecture based on deep networks;
- 516 • **H3:** our architecture supports reliable and incremental learning of state constraints, which improves
 517 the ability to answer explanatory questions; and
- 518 • **H4:** our architecture can be adapted to planning tasks, with the incremental learning capability
 519 improving the ability to compute minimal plans.

520 These hypotheses were evaluated in the context of the domains (SS, TS and RA) introduced in Section 3.
 521 Specifically, hypotheses $H1$, $H2$ and $H3$ are evaluated in the SS domain and TS domain in the context
 522 of VQA. As stated in Section 1, VQA is used in this paper only as an instance of a complex task that

523 requires explainable reasoning and learning. We are primarily interested in exploring the interplay between
524 reasoning with commonsense domain knowledge, incremental learning, and deep learning, in any given
525 domain in which large labeled datasets are not readily available. State of the art VQA algorithms, on the
526 other hand, focus instead on generalizing across different domains, using benchmark datasets of several
527 thousand images. Given the difference in objectives between over work and the existing work on VQA,
528 we thus do not compare with state of the art algorithms, and do not use the benchmark VQA datasets.
529 Furthermore, we evaluated hypothesis $H4$ in the RA domain in which the robot's goal was to deliver
530 messages to appropriate people and answer explanatory questions about this process.

531 We begin by describing some execution traces in Section 4.1 to illustrate the working of our architecture.
532 This is followed by Sections 4.2- 4.4, which describe the results of experimentally evaluating the
533 classification, VQA, axiom learning, and planning capabilities, i.e., hypotheses $H1-H4$. We use accuracy
534 (precision) as the primary performance measure. Classification accuracy was measured by comparing the
535 assigned labels with the ground truth values, and question answering accuracy was evaluated heuristically
536 by computing whether the answer mentions all image attributes relevant to the question posed. This
537 relevance was established by a human expert, the lead author of this paper. Unless stated otherwise, we
538 used two-thirds of the available data to train the deep networks and other computational models, using the
539 remaining one-third of the data for testing. For each image, we randomly chose from the set of suitable
540 questions for training the computational models. We repeated this process multiple times and report the
541 average of the results obtained in these trials. For planning, accuracy was measured as the ability to
542 compute minimal and correct plans for the assigned goals. Finally, Section 4.5 discusses the reduction in
543 computational effort achieved by our architecture in comparison with the baselines.

544 4.1 Execution Traces

545 The following execution traces illustrate our architecture's ability to reason with commonsense knowledge
546 and learned models to provide intuitive answers for explanatory questions.

547 **Execution Example 1.** [*Question Answering, SS domain*] Consider a scenario in the SS Domain in
548 which the input (test) image is the one on the extreme right in Figure 2.

- 549 • First **classification-related question** posed: "*is this structure unstable?*"
550 The architecture's **answer**: "*no*".
- 551 • The **explanatory question** posed: "*what is making this structure stable?*"
552 The architecture's **answer**: "*the structure has five blocks and a narrow base, it is standing straight,
553 and there is no significant lean*".
- 554 • This answer was based on the following features extracted by CNNs from the image: (i) five blocks;
555 (ii) narrow base; (iii) standing straight; and (iv) no significant lean, i.e, all blocks in place.
- 556 • The extracted features were converted to literals. ASP-based inference provided an answer about the
557 stability of the arrangement of objects in the scenario. Relevant literals in the corresponding answer set
558 were then inserted into a suitable template to provide the answers described above.
- 559 • Since the example was processed successfully using ASP-based inference, it was not processed using
560 the decision tree (for classification) or the RNN (for answering the explanatory question).

561 **Execution Example 2.** [*Question Answering, TS domain*] Consider a scenario in the TS Domain with
562 the input (test) image is the one on the extreme right in Figure 3.

- 563 • The **classification question** posed was: “*what is the sign’s message?*”
564 The architecture’s **answer**: “*uneven surfaces ahead*”.
- 565 • When asked to explain this answer (“*Please explain this answer*”), the architecture identified that the
566 CNNs extracted the following features of the sign in the image: (i) it is triangle-shaped; (ii) main color
567 is white and other (i.e., border) color is red; (iii) it has no background image; (iv) it has a bumpy-road
568 symbol and no secondary symbol; and (v) it has no cross.
- 569 • These features were converted to literals and used in ASP-based inference based on existing knowledge
570 in the TS domain. ASP-based inference is unable to provide an answer, i.e., unable to classify the sign.
- 571 • The extracted features were processed using the trained decision tree, which only used the colors in the
572 sign to assign the class label. The main (or border) color is normally insufficient to accurately classify
573 signs. However, recall that the decision tree is trained to classify signs that cannot be classified by
574 reasoning with existing knowledge.
- 575 • The decision tree output, image feature vector, and input question, were processed by the previously
576 trained RNN to provide the answer type and the particular answer described above.

577 These (and other such) execution traces illustrate the working of our architecture, especially that:

- 578 • The architecture takes advantage of (and perform non-monotonic logical inference with) the existing
579 commonsense domain knowledge to reliably and efficiently address the decision-making problem
580 (classification in the examples above) when possible. In such cases, it is also able to answer explanatory
581 questions about the classification decision and the underlying scene.
- 582 • When the desired decision cannot be made using non-monotonic logical inference with domain
583 knowledge, the architecture smoothly transitions to training and using a decision-tree to make and
584 explain the classification decision. In such cases, the architecture also learns and uses an RNN to
585 answer explanatory questions about the scene.

586 4.2 Experimental Results: Classification + VQA

587 To quantitatively evaluate hypotheses $H1$ and $H2$, we ran experimental trials in which we varied the size
588 of the training dataset. In these trials, the baseline performance was provided by a CNN-RNN architecture,
589 with the CNNs processing images to extract and classify features, and the RNN providing answers to
590 explanatory questions. The number of questions considered depends on the complexity of the domain, e.g.,
591 we included eight different types of questions in the SS domain and 248 different types of questions in
592 the TS domain. We repeated the trials 50 times (choosing the training set randomly each time) and the
593 corresponding average results are summarized in Figures 8 and 9 for the SS domain, and in Figures 10
594 and 11 for the TS domain. We make some observations based on these figures:

- 595 1. The classification performance of our architecture depends on the domain. In the relatively simpler SS
596 domain, the baseline deep network architecture is at least as accurate as our architecture, even with a
597 small training set—see Figure 8. This is because small differences in the position and arrangement of
598 blocks (which could almost be considered as noise) influence the decision about stability. For instance,
599 two arrangements of blocks that are almost identical end up receiving different ground truth labels for
600 stability, and it is not possible to draft rules based on abstract image features to distinguish between
601 these cases. The baseline deep network architecture, which generalizes from data, is observed to be
602 more sensitive to these small changes than our architecture. Exploring the reason for this performance
603 is an interesting direction for further research.

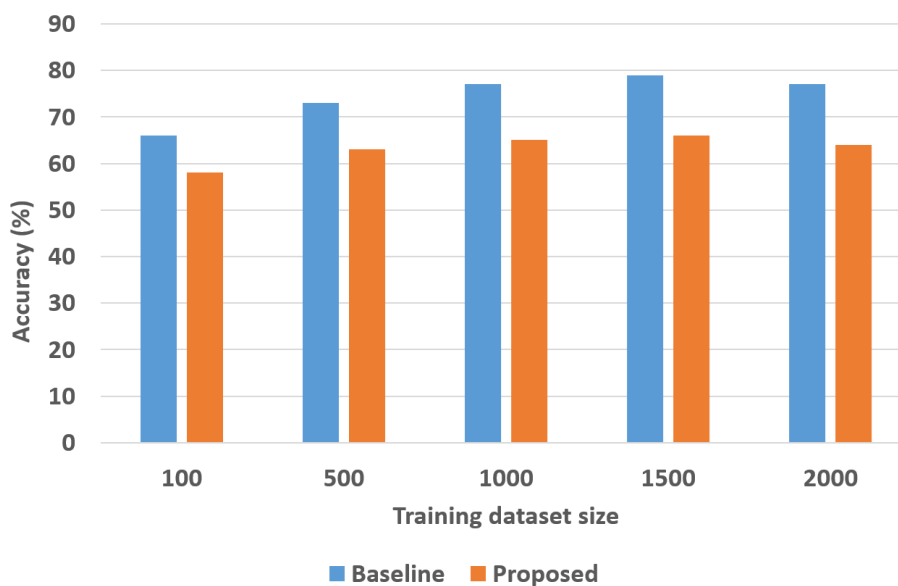


Figure 8. Classification accuracy as a function of the number of training samples in the SS domain.

- 604 2. In the more complex TS domain, our architecture provides better classification accuracy than the
 605 baseline architecture based on just deep networks, especially when the size of the training set is
 606 small—see Figure 10. The classification accuracy increases with the size of the training set³, but our
 607 architecture is always at least as accurate as the baseline architecture.
- 608 3. Our architecture is much more capable of answering explanatory questions about the classification
 609 decisions than the baseline architecture. When the answer provided by our architecture does not
 610 match the ground truth, we are able to examine why that decision was made. We were thus able to
 611 understand and explain the lower classification accuracy of our architecture in the SS domain. The
 612 baseline architecture does not provide this capability.
- 613 4. Unlike classification, the VQA performance of our architecture is much better than that of the baseline
 614 architecture in both domains. Also performance does not improve just by increasing the size of training
 615 set, even in simpler domains, e.g., see Figure 9. This is because VQA performance also depends on
 616 the complexity of the explanatory questions. For more complex domains, the improvement in VQA
 617 accuracy provided by our architecture is much more pronounced, e.g., see Figure 11.

618 We explored the statistical significance of the observed performance by running paired t-tests. We observed
 619 that the VQA performance of the proposed architecture was significantly better than that of the baseline
 620 architecture; this is more pronounced in the TS domain that is more complex than the SS domain. Also,
 621 although the baseline architecture provides better classification performance in the SS domain, the difference
 622 is not always statistically significant.

623 To further explore the observed results, we obtained a “confidence value” from the logits layer of each
 624 CNN used to extract a feature from the input image. For each CNN, the confidence value is the largest
 625 probability assigned to any of the possible values of the corresponding feature, i.e., it is the probability
 626 assigned to the most likely value of the feature. These confidence values are considered to be a measure of
 627 the network’s confidence in the corresponding features being a good representation of the image. We trained
 628 a version of our architecture in which if the confidence value for any feature was low, the image features

³ We limit ourselves to training sets that are not too large in order to match the focus of our paper.

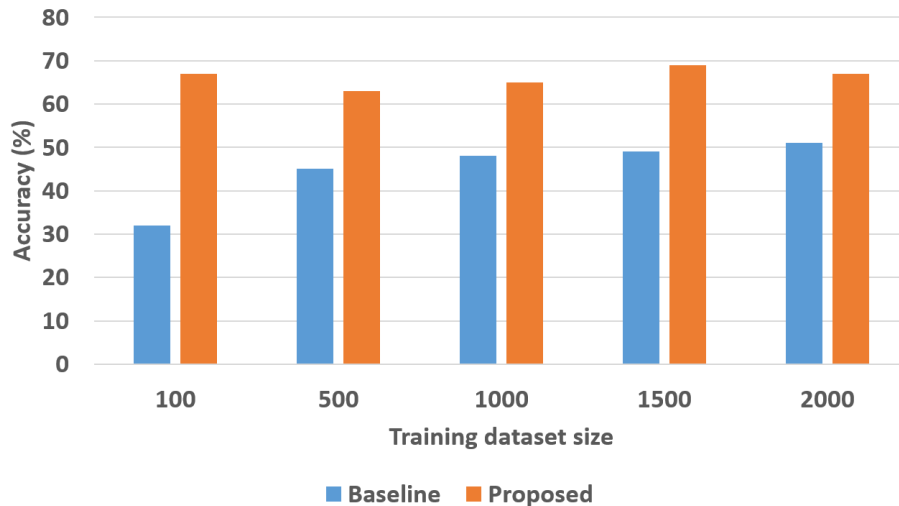


Figure 9. VQA accuracy as a function of the number of training samples in the SS domain.

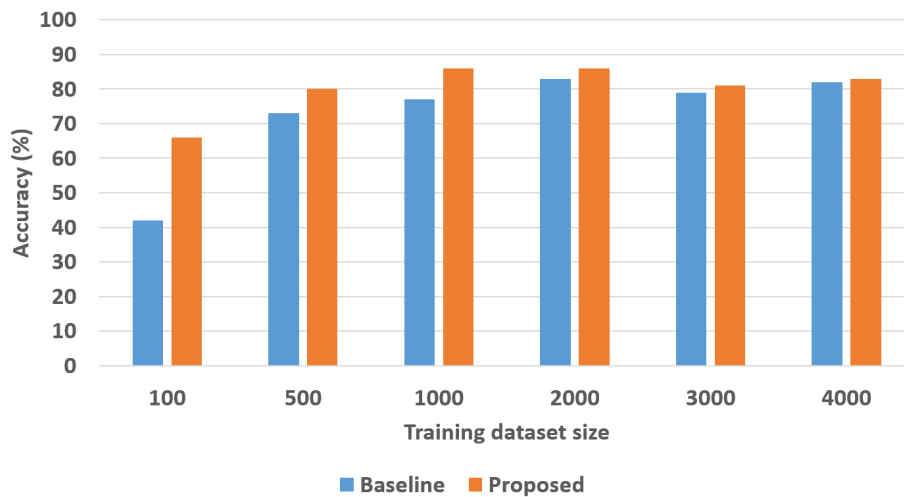


Figure 10. Classification accuracy as a function of number of training samples in TS domain.

629 were only used to revise the decision tree (during training), or were processed using the decision tree
 630 (during testing). In other words, features that do not strongly capture the essence of the image are not used
 631 for non-monotonic logical reasoning; the deep network architectures provide much better generalization
 632 to noise. We hypothesized that this approach would improve the accuracy of classification and question
 633 answering, but it did not make any significant difference in our experimental trials. We believe this is
 634 because the extracted features were mostly good representations of the objects of interest in the images.
 635 We thus did not use such networks (that compute the confidence value) in any other experiments.

636 4.3 Experimental Results: Learn Axiom + VQA

637 Next, we experimentally evaluated the ability to learn axioms, and the effect of such learning on the
 638 classification and VQA performance. For the SS domain, we designed a version of the knowledge base
 639 with eight axioms related to stability or instability of the structures. Out of these, four were chosen
 640 (randomly) to be removed and we examined the ability to learn these axioms, and the corresponding
 641 accuracy of classification and VQA, as a function of the number of labeled training examples (ranging
 642 from 100 to 2000). We repeated these experiments 30 times and the results (averaged over the 30 trials)

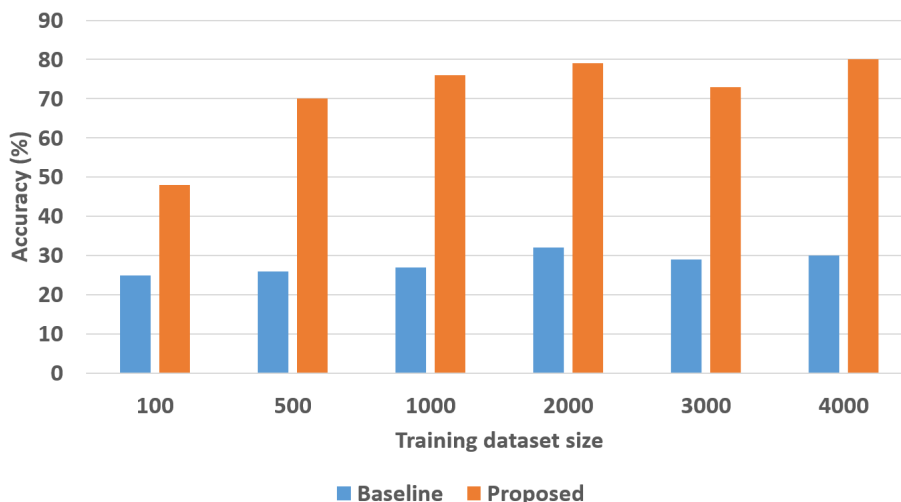


Figure 11. VQA accuracy as a function of number of training samples in TS domain.

643 are summarized in Figures 12-13. In the TS domain, the methodology for experimental evaluation was
 644 the same. However, since the domain was more complex, there were many more axioms in the domain
 645 description (for classification and VQA); we also had access to more labeled training examples. In each
 646 experimental trial, a quarter of the available axioms were thus selected and commented out, and the
 647 accuracy of classification and VQA were evaluated with the number of labeled training examples varying
 648 from 100 to 4000. The results averaged over 30 such trials are summarized in Figures 14-15.



Figure 12. Comparison of classification accuracy in the SS domain with and without axiom learning. In both cases, some axioms were missing from the knowledge base.

649 In these figures, “Original KB” (depicted in blue) represents the baseline with some axioms missing
 650 from the system description, e.g., four in the SS domain and one quarter of the axioms in the TS domain.
 651 The results obtained by using the available labeled examples to learn the axioms that are then used for
 652 classification and answering explanatory questions about the scene, are shown as “Learned KB” in orange.
 653 We observe that our approach supports incremental learning of the domain axioms, and that using the
 654 learned axioms improves the classification accuracy and the accuracy of answering explanatory questions,
 655 in comparison with the baseline. This improvement was found to be statistically significant using paired
 656 tests at 95% level of significance. These results support hypothesis *H3*.

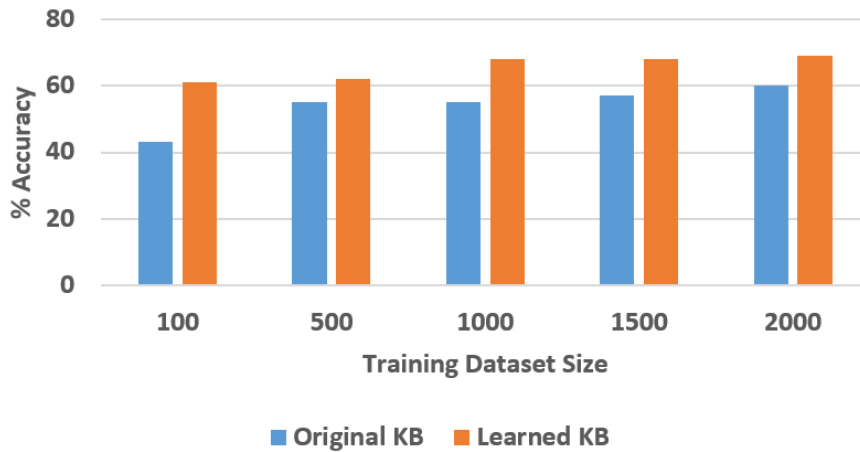


Figure 13. Comparison of VQA accuracy in the SS domain with and without axiom learning. In both cases, some axioms were missing from the knowledge base.

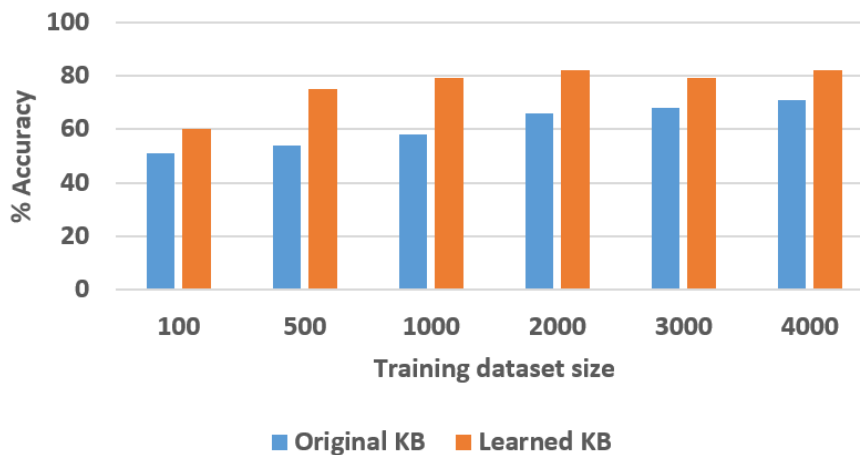


Figure 14. Comparison of classification accuracy in the TS domain with and without axiom learning. In both cases, some axioms were missing from the knowledge base.

657 4.4 Experimental Results: Learn Axiom + Plan

658 Next, we experimentally evaluated the ability to learn axioms and the effect of the learned axioms on
 659 planning, in the RA domain. The simulated robot was equipped with domain knowledge for planning,
 660 classification, and question answering. It uses this knowledge to navigate through an office building, locate
 661 the intended recipient of a message, deliver the message, detect and reason about objects in its surroundings,
 662 and answer questions about the rooms it has visited. We considered 24 different types of questions in
 663 this domain. As stated in Section 3.5, we limit uncertainty in sensing and actuation on robots to noise
 664 added in simulation. Average results from 100 trials indicates a VQA accuracy of $\approx 85\%$ after training
 665 the architecture’s components with just 500 labeled images. The domain knowledge includes learned
 666 axioms—the corresponding experimental results and the planning performance are discussed later in this
 667 section. We begin with an execution trace in this domain.

668 **Execution Example 3.** [*Question Answering, RA Domain*] Consider the scenario in the RA domain
 669 (Figure 7) in which the robot’s goal was to deliver a message from John to Sally, and return to John to
 670 answer questions.

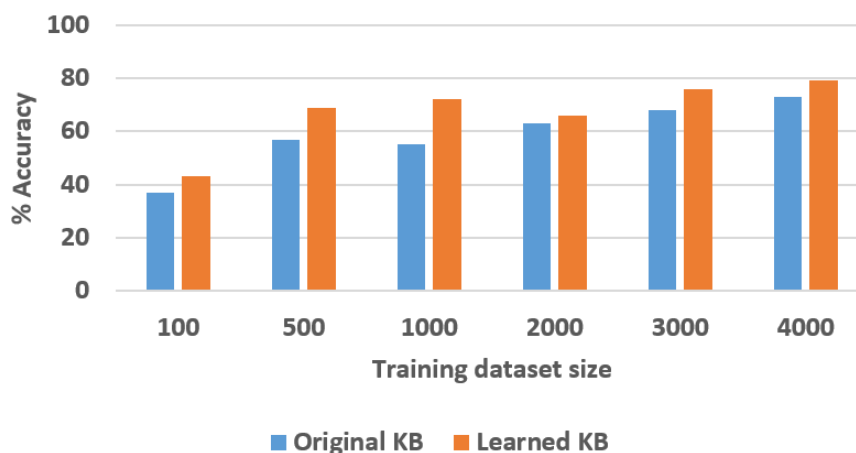


Figure 15. Comparison of VQA accuracy in the TS domain with and without axiom learning. In both cases, some axioms were missing from the knowledge base.

- 671 • The robot was initially in John’s office. It computed a plan that comprises moving to Sally’s office
 672 through the library and the kitchen, delivering the message to Sally, and returning to John’s office
 673 through the same route to answer questions.
- 674 • During plan execution, the robot periodically takes images of the scenes in the domain, which are used
 675 for planning, classification and question answering.
- 676 • After returning to John’s office, the robot and the human had an exchange about the plan constructed
 677 and executed, and the observations received. The exchange includes instances such as:
- 678 **John’s question:** “is Sally’s location cluttered?”
- 679 **Robot’s answer:** “Yes”.
- 680 When asked, the robot provides an **explanation** for this decision: “Sally is in her office. Objects
 681 detected are Sally’s chair, desk, and computer, and a cup, a large box, and a sofa. The room is cluttered
 682 because the cup, large box, and sofa are not usually in that room”.

The RA domain was also used to evaluate the effects of axiom learning. There were four employees in offices in the simulated scenario, as shown in Figure 7, and the robot was asked to find particular individuals and deliver particular messages to them. Employees are initially expected to be in their assigned workplace (i.e., their office), and spend most of their time in these offices, unless this default knowledge has been negated by other knowledge or observations. This information is encoded as follows:

$$\text{holds}(\text{loc}(P, L), 0) \leftarrow \text{not default_negated}(P, L), \text{work_place}(P, L)$$

where $\text{work_place}(P, L)$ specifies the default location of each person, and $\text{default_negated}(P, L)$ is used to encode that a particular person may not be in their default location. These exceptions to the defaults can be encoded as follows:

$$\text{default_negated}(P, L) \leftarrow \text{obs}(\text{loc}(P, L1), \text{true}, I), L \neq L1 \quad (11a)$$

$$\text{default_negated}(P, L) \leftarrow \text{obs}(\text{loc}(P, L), \text{false}, I) \quad (11b)$$

- 683 Statement 11(a) implies that the default assumption should be ignored if the person in question is observed
 684 to be in a location other than their workplace, and Statement 11(b) implies that a default assumption

685 should be ignored if the corresponding person is not observed in their workplace. Including such default
 686 knowledge (and exceptions) in the reasoning process allows the robot to compute better plans and execute
 687 the plans more efficiently, e.g., when trying to deliver a message to a particular person. However, this
 688 knowledge may not be known in advance, the existing knowledge may be inaccurate or change with time
 689 (e.g., humans can move between the different places), or the observations may be incorrect. Our axiom
 690 learning approach was used in this domain to acquire previously unknown information about the default
 691 location of people and exceptions to these defaults. In all the trials, the simulated robot was able to learn
 692 the previously unknown axioms.

Axiom learning	Plans (per trial)	Actions (per trial)	Execution time (per trial)	Planning time (per trial)	Planning time (per plan)
Before	4	2.3	1.6	6.0	1.6
After	1	1	1	1	1

Table 1. Planning performance in a scenario in the RA domain (see Figure 7) before and after axiom learning. Results averaged over 100 paired trials indicate that reasoning with previously unknown axioms results in fewer plans with fewer actions in each trial, and significantly reduces the time taken to compute and execute the plans.

693 We then conducted 100 paired trials to explore the effects of the learned axioms on planning, with the
 694 corresponding results summarized in Table 1. In each trial, we randomly chose a particular goal and initial
 695 conditions, and measured planning performance before and after the previously unknown axioms had been
 696 learned and used for reasoning. Since the initial conditions are chosen randomly, the object locations, the
 697 initial location of the robot, and the goal, may vary significantly between trials. Under these circumstances,
 698 it is not meaningful to average the results obtained in the individual trials for performance measures such
 699 as planning time and execution time. Instead, the results obtained without including the learned axioms
 700 were computed as a ratio of the results obtained after including the learned axioms; the numbers reported
 701 in Table 1 are the average of these computed ratios. Before axiom learning, the robot often explored an
 702 incorrect location (for a person) based on other considerations (e.g., distance to the room) and ended
 703 up having to replan. After the previously unknown axioms were included in the reasoning process, the
 704 robot went straight to the message recipient’s most likely location, which also happened to be the actual
 705 location of the recipient in many trials. As a result, we observe a (statistically) significant improvement
 706 in planning performance after the learned axioms are used for reasoning. Note that in the absence of the
 707 learned axioms, the robot computes four times as many plans taking six times as much time in any given
 708 trial (on average) as when the learned axioms are included in reasoning. Even the time taken to compute
 709 each plan (with potentially multiple such plans computed in each trial) is significantly higher in the absence
 710 of the learned axioms. This is because the learned axioms enable the robot to eliminate irrelevant paths in
 711 the transition diagram from further consideration. In a similar manner, reasoning with intentional actions
 712 enables the robot to significantly reduce the plan execution time by terminating or revising existing plans
 713 when appropriate, especially in the context of unexpected successes and failures. These results provide
 714 evidence in support of hypothesis $H4$.

715 Finally, we conducted some initial proof of concept studies exploring the use of our architecture on
 716 physical robots. We considered a robot collaborating with a human to jointly describe structures of blocks
 717 on a tabletop (similar to the SS domain described in this paper). We also considered a mobile robot finding
 718 and moving objects to desired locations in an indoor domain (similar to the RA domain). These initial
 719 experiments provided some promising outcomes. The robot was able to provide answers to explanatory

720 questions, compute and execute plans to achieve goals, and learn previously unknown constraints. In the
721 future, we will conduct a detailed experimental analysis on robots in different domains.

722 4.5 Computational Effort

723 In addition to the improvement in accuracy of classification and VQA, we also explored the reduction in
724 computational effort provided by our architecture in comparison with the baselines. Measuring this time
725 quantitatively is challenging because it depends on various factors such as the task being performed (e.g.,
726 classification, VQA), the knowledge encoded in the knowledge base, the size and order of samples in the
727 training set, and the parameters of the deep networks. However, we were able to gain the following insights.

728 The computation time includes the training time and the testing (i.e., execution) time, and we first
729 considered the training time. Depending on the task being performed (e.g., classification, VQA, and/or
730 planning), this time includes the time taken to encode and draw inferences from the knowledge base,
731 process queries and construct answers, and train the deep network models. Encoding the incomplete
732 domain knowledge is a one-time exercise for any given domain. The time taken to reason with this
733 knowledge, and the time taken to process queries and construct answers, are negligible in comparison
734 with the time taken to learn the deep network models. Also, the use of CNNs to extract features from
735 images is common to both our architecture and the baselines, and these networks (for the most part) do
736 not need to be retrained multiple times for any given domain. The key difference between our architecture
737 and the baselines is observed in the context of answering explanatory questions about the scenes and the
738 underlying classification problem. Recall that with our architecture, only examples that cannot be processed
739 by ASP-based reasoning are processed by decision-trees and the RNNs for VQA. In our experimental
740 trials, $\approx 10 - 20\%$ of a training set is used (on average) to train the RNNs with our architecture, whereas
741 the entire training set is used for training the RNNs with the baseline architectures. This difference often
742 translates to an order of magnitude difference in the training time, e.g., a few minutes for each training set
743 (in a particular domain) with our architecture compared with hours or days with the baseline architectures.
744 Note that accuracy of our architecture is still much better than that of the baselines, e.g., see Figure 9 and
745 Figure 11, i.e., any given accuracy is achieved using a much smaller number of training samples.

746 The execution time of our architecture is comparable with that of the baselines and is often less. Once
747 the deep network models have been learned, using them for the different tasks does not take much time,
748 e.g., a few seconds to process the input and provide a decision and/or the answer to a query. However,
749 similar to the situation during training, only test samples that cannot be processed by ASP-based reasoning
750 are processed by the decision trees and RNNs with our architecture. Also, since the deep networks in our
751 architecture only need to disambiguate between a small(er) number of training examples, they often have a
752 much simpler structure than the deep networks in the baseline architectures.

753 Note that in addition to classification and VQA, our architecture also supports explainable reasoning
754 for planning and incremental learning of previously unknown constraints. Providing similar capabilities
755 using just deep network architectures will (at the very least) require a large number of training examples of
756 planning under different conditions; it is often not possible to provide such training examples in dynamic
757 domains. We thus conclude that our architecture significantly reduces the computational effort while
758 supporting a range of capabilities in comparison with the baseline architectures comprising deep networks.

5 DISCUSSION AND FUTURE WORK

759 Visual question answering (VQA) combines challenges in computer vision, natural language processing,
760 and explainability in reasoning and learning. Explanatory descriptions of decisions help identify errors,
761 and to design better algorithms and frameworks. In addition, it helps improve trust in the use of reasoning
762 and learning systems in critical application domains. State of the art algorithms for VQA are based on deep
763 networks and the corresponding learning algorithms. Given their focus on generalizing across different
764 domains, these approaches are computationally expensive, require large training datasets, and make it
765 difficult to provide explanatory descriptions of decisions. We instead focus on enabling reliable and efficient
766 operation in any given domain in which a large number of labeled training examples may not be available.
767 Inspired by research in cognitive systems, our architecture tightly couples representation, reasoning, and
768 interactive learning, and exploits the complementary strengths of deep learning, non-monotonic logical
769 reasoning with commonsense knowledge, and decision tree induction. Experimental results on datasets of
770 real world and simulated images indicate that our architecture provides the following benefits in comparison
771 with a baseline architecture for VQA based on deep networks:

- 772 1. Better accuracy, improved sample efficiency, and reduced computational effort on classification
773 problems when the training dataset is small, and comparable accuracy with larger datasets while still
774 using only a subset of these samples for training;
- 775 2. Ability to provide answers to explanatory questions about the scenes and the underlying decision
776 making problems (e.g., classification, planning);
- 777 3. Incremental learning of previously unknown domain constraints, whose use in reasoning improves the
778 ability to answer explanatory questions; and
- 779 4. Ability to adapt the complementary strengths of non-monotonic logical reasoning with commonsense
780 domain knowledge, inductive learning, and deep learning, to address decision-making (e.g., planning)
781 problems on a robot.

782 Our architecture opens up multiple directions of future work, which will address the limitations of existing
783 work and significantly extend the architecture’s capabilities. We discuss some of these extensions below:

- 784 1. The results reported in this paper are based on image datasets (simulated, real-world) chosen or
785 constructed to mimic domains in which a large, labeled dataset is not readily available. One direction of
786 future work is to explore the use of our architecture in other domains that provide datasets of increasing
787 complexity, i.e., with a greater number of features and more complex explanatory questions. This
788 exploration may require us to consider larger datasets, and to examine the trade-off between the size of
789 the training dataset, the computational effort involved in processing such a dataset with many labeled
790 examples, and the effort involved in encoding and reasoning with the relevant domain knowledge.
- 791 2. In our architecture, we have so far used variants of existing network structures as the deep network
792 components (i.e., CNN, RNN). In the future, we will explore different deep network structures in
793 our architecture, using the explanatory answers to further understand the internal representation of
794 these network structures. Towards this objective, it would be particularly instructive to construct and
795 explore deep networks and logic-based domain representations that provide similar behavior on a set
796 of tasks, or provide different behavior when operating on the same dataset. As stated in the discussion
797 in Section 4.2, such an exploration may help us better understand (and improve) the design and use of
798 deep network models for different applications.

799 3. This paper used VQA as a motivating problem to address key challenges in using deep networks in
800 dynamic domains with limited labeled training examples. We also described the use of our architecture
801 (with tightly-coupled reasoning and learning components) for planning on a simulated robot. In the
802 future, we will combine this architecture with other architectures we have developed for knowledge
803 representation, reasoning, and interactive learning in robotics (Sridharan and Meadows, 2018; Sridharan
804 et al., 2019). The long-term goal will be to support explainable reasoning and learning on a physical
805 robot collaborating with humans in complex domains.

ACKNOWLEDGMENTS

806 The authors thank Ales Leonardis for feedback on the architecture described in this paper. An earlier
807 version of this architecture appeared in a conference paper (Riley and Sridharan, 2018a). In this paper,
808 we have introduced a new component for incrementally learning constraints governing domain states,
809 expanded reasoning to support planning and diagnostics, and discussed more detailed experimental results.

CONFLICT OF INTEREST STATEMENT

810 The authors declare that the research was conducted in the absence of any commercial or financial
811 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

812 HR and MS designed the algorithms and architecture. HR implemented the algorithms and architecture
813 with feedback from MS. MS and HR designed the experimental setup. HR conducted the evaluation, and
814 HR and MS analyzed the results. MS and HR wrote the paper.

FUNDING

815 This work was supported in part by the US Office of Naval Research Science of Autonomy award N00014-
816 17-1-2434, and the Asian Office of Aerospace Research and Development award FA2386-16-1-4071.

DATA AVAILABILITY STATEMENT

817 The datasets generated or analyzed for this study, and the software implementation of the architecture
818 and algorithms, can be found in the following online repository: [https://github.com/hril230/
819 masters_code](https://github.com/hril230/masters_code)

REFERENCES

- 820 Agrawal, A., Batra, D., Parikh, D., and Kembhavi, A. (2018). Don't Just Assume; Look and Answer:
821 Overcoming Priors for Visual Question Answering. In *International Conference on Computer Vision
822 and Pattern Recognition* (Salt Lake City, USA)
- 823 Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., et al. (2018). Bottom-up and
824 Top-down Attention for Image Captioning and Visual Question Answering. In *International Conference
825 on Computer Vision and Pattern Recognition* (Salt Lake City, USA)
- 826 Bai, Y., Fu, J., Zhao, T., and Mei, T. (2018). Deep Attention Neural Tensor Network for Visual Question
827 Answering. In *European Conference on Computer Vision (ECCV)* (Munich, Germany)

- 828 Balai, E., Gelfond, M., and Zhang, Y. (2013). Towards Answer Set Programming with Sorts. In
829 *International Conference on Logic Programming and Nonmonotonic Reasoning* (Corunna, Spain)
- 830 Balduccini, M. (2007). Learning Action Descriptions with A-Prolog: Action Language C. In *AAAI Spring*
831 *Symposium on Logical Formalizations of Commonsense Reasoning*
- 832 Balduccini, M. and Gelfond, M. (2003). Logic Programs with Consistency-Restoring Rules. In *AAAI*
833 *Spring Symposium on Logical Formalization of Commonsense Reasoning*. 9–18
- 834 Baral, C. (2003). *Knowledge Representation, Reasoning and Declarative Problem Solving* (Cambridge
835 University Press)
- 836 Chai, J. Y., Gao, Q., She, L., Yang, S., Saba-Sadiya, S., and Xu, G. (2018). Language to Action:
837 Towards Interactive Task Learning with Physical Agents. In *International Joint Conference on Artificial*
838 *Intelligence (IJCAI)* (Stockholm, Sweden)
- 839 Duda, R. O., Hart, P. E., and Stork, D. G. (2000). *Pattern Classification (Second Edition)* (Wiley-Blackwell)
- 840 Erdem, E., Gelfond, M., and Leone, N. (2016). Applications of Answer Set Programming. *AI Magazine*
841 37, 53–68
- 842 Erdem, E. and Patoglu, V. (2012). Applications of Action Languages to Cognitive Robotics. In *Correct*
843 *Reasoning* (Springer-Verlag). 229–246
- 844 Furbach, U., Glöckner, I., Helbig, H., and Pelzer, B. (2010). Logic-Based Question Answering. *KI -*
845 *Künstliche Intelligenz* 24, 51–55
- 846 Gelfond, M. and Incelesan, D. (2013). Some Properties of System Descriptions of AL_d . *Journal of Applied*
847 *Non-Classical Logics, Special Issue on Equilibrium Logic and Answer Set Programming* 23, 105–120
- 848 Gelfond, M. and Kahl, Y. (2014). *Knowledge Representation, Reasoning and the Design of Intelligent*
849 *Agents* (Cambridge University Press)
- 850 Gil, Y. (1994). Learning by Experimentation: Incremental Refinement of Incomplete Planning Domains.
851 In *International Conference on Machine Learning* (New Brunswick, USA), 87–95
- 852 Goyal, Y., Khot, T., Stay, D. S., Batra, D., and Parikh, D. (2017). Making the V in VQA Matter: Elevating
853 the Role of Image Understanding in Visual Question Answering. In *International Conference on*
854 *Computer Vision and Pattern Recognition* (Honolulu, USA), 6325–6334
- 855 Jabri, A., Joulin, A., and van der Maaten, L. (2016). Revisiting Visual Question Answering Baselines. In
856 *European Conference on Computer Vision* (Amsterdam)
- 857 Jiang, A., Wang, F., Porikli, F., and Li, Y. (2015). *Compositional Memory for Visual Question Answering*.
858 Tech. rep., <https://arxiv.org/abs/1511.05676>
- 859 Kim, J.-H., Jun, J., and Zhang, B.-T. (2018). Bilinear Attention Networks. In *Neural Information*
860 *Processing Systems* (Montreal, Canada)
- 861 Koh, P. W. and Liang, P. (2017). Understanding Black-box Predictions via Influence Functions. In
862 *International Conference on Machine Learning (ICML)* (Sydney, Australia), 1885–1894
- 863 Laird, J. E. (2012). *The Soar Cognitive Architecture* (The MIT Press)
- 864 Laird, J. E., Gluck, K., Anderson, J., Forbus, K. D., Jenkins, O. C., Lebiere, C., et al. (2017). Interactive
865 Task Learning. *IEEE Intelligent Systems* 32, 6–21
- 866 Law, M., Russo, A., and Broda, K. (2018). The Complexity and Generality of Learning Answer Set
867 Programs. *Artificial Intelligence* 259, 110–146
- 868 Li, Q., Fu, J., Yu, D., Mei, T., and Luo, J. (2018). *Tell-and-Answer: Towards Explainable Visual Question*
869 *Answering using Attributes and Captions*. Tech. rep., <https://arxiv.org/abs/1801.09041>
- 870 Lin, T., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., et al. (2014). Microsoft COCO:
871 Common Objects in Context. In *European Conference on Computer Vision* (Zurich, Switzerland)

- 872 Lu, J., Yang, J., Batra, D., and Parikh, D. (2016). Hierarchical Question-Image Co-Attention for Visual
873 Question Answering. In *Advances in Neural Information Processing Systems* (Barcelona, Spain)
- 874 Malinowski, M., Rohrbach, M., and Fritz, M. (2017). Ask Your Neurons: A Deep Learning Approach to
875 Visual Question Answering. *International Journal of Computer Vision* 125, 110–135
- 876 Mascharka, D., Tran, P., Soklaski, R., and Majumdar, A. (2018). Transparency by Design: Closing the
877 Gap between Performance and Interpretability in Visual Reasoning. In *International Conference on*
878 *Computer Vision and Pattern Recognition* (Salt Lake City, USA)
- 879 Masuda, I., de la Puente, S. P., and Nieto, X. G. (2016). Open-Ended Visual Question-Answering. In
880 *International Conference on Computer Vision and Pattern Recognition* (Las Vegas, USA)
- 881 Mota, T. and Sridharan, M. (2019). Commonsense Reasoning and Knowledge Acquisition to Guide Deep
882 Learning on Robots. In *Robotics Science and Systems* (Freiburg, Germany)
- 883 Norcliffe-Brown, W., Vafeais, E., and Parisot, S. (2018). Learning Conditioned Graph Structures for
884 Interpretable Visual Question Answering. In *Neural Information Processing Systems* (Montreal, Canada)
- 885 Otero, R. P. (2003). Induction of the Effects of Actions by Monotonic Methods. In *International Conference*
886 *on Inductive Logic Programming* (Szeged, Hungary), 299–310
- 887 Pandhre, S. and Sodhani, S. (2017). *Survey of Recent Advances in Visual Question Answering*. Tech. rep.,
888 <https://arxiv.org/abs/1709.08203>
- 889 Rajani, N. F. and Mooney, R. J. (2018). Stacking With Auxiliary Features for Visual Question Answering.
890 In *16th Annual Conference of the North American Chapter of the Association for Computational*
891 *Linguistics: Human Language Technologies* (New Orleans, USA), 2217–2226
- 892 Ribeiro, M., Singh, S., and Guestrin, C. (2016). Why Should I Trust You? Explaining the Predictions of
893 Any Classifier. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*
894 *(KDD)* (San Francisco, USA), 1135–1144
- 895 Riley, H. and Sridharan, M. (2018a). Non-monotonic Logical Reasoning and Deep Learning for Explainable
896 Visual Question Answering. In *International Conference on Human-Agent Interaction* (Southampton,
897 UK)
- 898 [Dataset] Riley, H. and Sridharan, M. (2018b). Software for Architecture combining Non-monotonic
899 Logical Reasoning, Inductive Learning and Deep Learning for VQA. Accessed January 2019 from
900 https://github.com/hril230/masters_code
- 901 Schwartz, I., Schwing, A., and Hazan, T. (2017). High-Order Attention Models for Visual Question
902 Answering. In *Advances in Neural Information Processing Systems (NIPS)* (Long Beach, USA),
903 3664–3674
- 904 Selvaraju, R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). Grad-CAM: Visual
905 Explanations from Deep Networks via Gradient-based Localization. In *International Conference on*
906 *Computer Vision* (Venice, Italy), 618–626
- 907 Shrestha, R., Kafle, K., and Kanan, C. (2019). Answer Them All! Toward Universal Visual Question
908 Answering Models. In *International Conference on Computer Vision and Pattern Recognition* (Long
909 Beach, USA)
- 910 Sridharan, M., Gelfond, M., Zhang, S., and Wyatt, J. (2019). REBA: A Refinement-Based Architecture for
911 Knowledge Representation and Reasoning in Robotics. *Journal of Artificial Intelligence Research* 65,
912 87–180
- 913 Sridharan, M. and Meadows, B. (2018). Knowledge Representation and Interactive Learning of Domain
914 Knowledge for Human-Robot Collaboration. *Advances in Cognitive Systems* 7, 69–88
- 915 Teney, D. and van den Hengel, A. (2016). *Zero-Shot Visual Question Answering*. Tech. rep., <https://arxiv.org/abs/1611.05546>
- 916

- 917 Teney, D. and van den Hengel, A. (2018). Visual Question Answering as a Meta Learning Task. In
918 *European Conference on Computer Vision (ECCV)* (Munich, Germany)
- 919 Timofte, R., Mathias, M., Benenson, R., and Gool, L. V. (2013). Traffic Sign Recognition - How far are we
920 from the Solution? In *International Joint Conference on Neural Networks (IJCNN)* (Dallas, USA), 1–8
- 921 Wagner, M., Basevi, H., Shetty, R., Li, W., Malinowski, M., Fritz, M., et al. (2018). Answering
922 Visual *What-If* Questions: From Actions to Predicted Scene Descriptions. In *Visual Learning and*
923 *Embodied Agents in Simulation Environments (VLEASE) Workshop at ECCV* (Munich, Germany).
924 <https://arxiv.org/abs/1809.03707>
- 925 Wang, P., Wu, Q., Shen, C., van den Hengel, A., and Dick, A. R. (2017). Explicit Knowledge-based
926 Reasoning for Visual Question Answering. In *International Joint Conference on Artificial Intelligence*
927 (Melbourne, Australia)
- 928 Wu, C., Liu, J., Wang, X., and Dong, X. (2018). Chain of Reasoning for Visual Question Answering. In
929 *Advances in Neural Information Processing Systems (NeurIPS)* (Montreal, Canada), 273–283
- 930 Yang, Z., He, X., Gao, J., Deng, L., and Smola, A. J. (2016). Stacked Attention Networks for Image
931 Question Answering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Las
932 Vegas, USA), 21–29
- 933 Yi, K., Wu, J., Gan, C., Torralba, A., Kohli, P., and Tenenbaum, J. B. (2018). Neural-Symbolic VQA:
934 Disentangling Reasoning from Vision and Language Understanding. In *Neural Information Processing*
935 *Systems* (Montreal, Canada)
- 936 Zhang, S., Sridharan, M., and Wyatt, J. (2015). Mixed Logical Inference and Probabilistic Planning for
937 Robots in Unreliable Worlds. *IEEE Transactions on Robotics* 31, 699–713
- 938 Zhang, T., Dai, D., Tuytelaars, T., Moens, M.-F., and Gool, L. V. (2017). *Speech-Based Visual Question*
939 *Answering*. Tech. rep., <https://arxiv.org/abs/1705.00464>