
IconQA: A New Benchmark for Abstract Diagram Understanding and Visual Language Reasoning

Pan Lu¹, Liang Qiu¹, Jiaqi Chen², Tony Xia¹, Yizhou Zhao¹,
Wei Zhang³, Zhou Yu⁴, Xiaodan Liang², Song-Chun Zhu¹

¹Center for Vision, Cognition, Learning and Autonomy, UCLA

²Sun Yat-sen University, ³East China Normal University, ⁴Columbia University

Abstract

Current visual question answering (VQA) tasks mainly consider answering human-annotated questions for natural images. However, aside from natural images, abstract diagrams with semantic richness are still understudied in visual understanding and reasoning research. In this work, we introduce a new challenge of Icon Question Answering (IconQA) with the goal of answering a question in an icon image context. We release IconQA, a large-scale dataset that consists of 107,439 questions and three sub-tasks: *multi-image-choice*, *multi-text-choice*, and *filling-in-the-blank*. The IconQA dataset is inspired by real-world diagram word problems that highlight the importance of abstract diagram understanding and comprehensive cognitive reasoning. Thus, IconQA requires not only perception skills like object recognition and text understanding, but also diverse cognitive reasoning skills, such as geometric reasoning, commonsense reasoning, and arithmetic reasoning. To facilitate potential IconQA models to learn semantic representations for icon images, we further release an icon dataset Icon645 which contains 645,687 colored icons on 377 classes. We conduct extensive user studies and blind experiments and reproduce a wide range of advanced VQA methods to benchmark the IconQA task. Also, we develop a strong IconQA baseline Patch-TRM that applies a pyramid cross-modal Transformer with input diagram embeddings pre-trained on the icon dataset. IconQA and Icon645 are available at <https://iconqa.github.io>.

1 Introduction

We are witnessing an exciting development of visual question answering (VQA) research in recent years. The long-standing goal of the VQA task is to exploit systems that can answer natural questions that correspond to visual information. Several datasets have been released to evaluate the systems' visual and textual content understanding abilities [3, 57, 14, 21, 18, 52]. One of the underlying limitations of current VQA datasets is that they are focusing on answering visual questions for natural images. However, aside from natural pictures, abstract diagrams with visual and semantic richness account for a large proportion of the visual world. For instance, it is shown that emojis can express rich human sentiments [26, 10], and diagrams like icons can map the physical worlds into symbolic and aesthetic representations [31, 40, 24].

Some pioneering works attempt to propose datasets that are capable of answering questions for abstract diagrams. However, these datasets either address domain-specific charts, plots, and illustrations [26, 22], or are generated from limited templates [55, 48, 21]. These limitations impede their practical applications in real-world scenarios. For example, in elementary school, abstract diagrams in math world problems are involved with diverse objects and various reasoning skills [25].

To address these shortcomings, we introduce Icon Question Answering (IconQA), a new challenge for *abstract diagram* visual reasoning and question answering. The task, stemming from math word

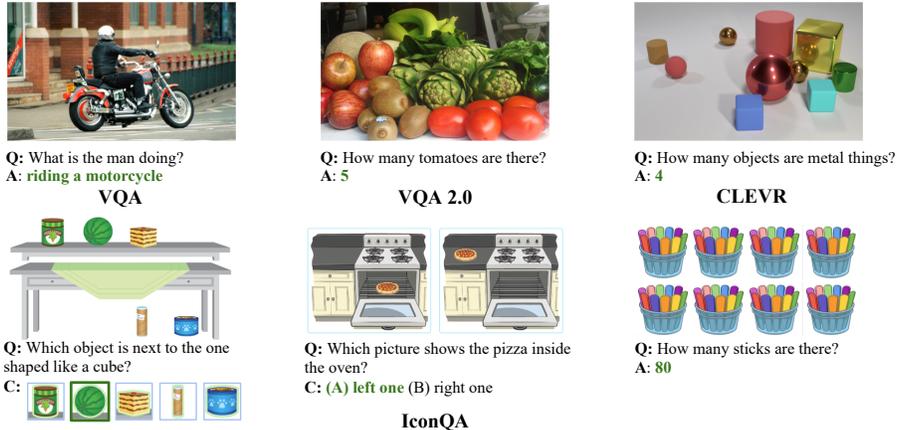


Figure 1: **Top**: Examples in three popular VQA datasets: VQA [3], VQA 2.0 [14], and CLEVR [21]. **Bottom**: Examples of three sub-tasks in our IconQA dataset. For answering these icon questions, it requires diagram recognition and text understanding, as well as diverse cognitive reasoning skills.

problems for children [41], exhibits a promising potential to develop education assistants. We name the proposed task as IconQA because the images depict icons, which simplify recognition and allow us to focus on reasoning skills for further research. We release IconQA, a large-scale dataset that contains 107,439 QA pairs and covers three different sub-tasks: *multiple-image-choice*, *multiple-text-choice* and *filling-in-the-blank*. A typical IconQA problem is provided with an icon image and a question, and the answer is in the form of either a short piece of text or a choice from multiple visual or textual choices. Correctly answering IconQA questions needs diverse human intelligence skills. As the examples in Figure 1 show, IconQA poses new challenges for abstract diagram understanding like recognizing objects and identifying attributes. Besides, it is critical to develop diverse cognitive reasoning skills, including counting objects, comparing attributes, performing arithmetic operations, making logical inferences, completing spatial reasoning, or leveraging external commonsense to answer IconQA questions. More examples from the dataset are shown in Appendix A.1.

We use the IconQA dataset to benchmark various VQA approaches in the IconQA task, including four attention-based multimodal pooling methods [2, 28, 54, 11] and four Transformer-based pre-trained methods [33, 6, 53, 29], as illustrated in Figure 6. Also, we conduct extensive user studies to evaluate the performance differences between the algorithms and human beings. Three blind studies show that the IconQA dataset is robust against biased shortcuts when answering icon questions. We further develop a strong baseline called pyramid patch cross-modal Transformer (Patch-TRM), which effectively learns implicit visual and linguistic relationships in IconQA. Patch-TRM parses the diagrams into patch sequences in a spatial pyramid structure and learns a joint embeddings within a multimodal Transformer. Along with the IconQA dataset, we collect an auxiliary icon dataset, Icon645, that features 645,687 colored icons on 377 object classes. The icon dataset is used to pre-train the diagram embedding module in Patch-TRM to enhance abstract diagram understanding.

Our contributions can be summarized as 1) we propose a new challenge, IconQA, that requires abstract diagram understanding of icon images and diverse visual reasoning skills; 2) we establish two large-scale datasets: IconQA, a question answering dataset in the icon domain, and Icon645, an icon dataset for model pre-training; 3) we benchmark the IconQA dataset extensively via experiments on eight existing methods and develop a strong multimodal Transformer-based baseline.

2 Related Works

VQA Datasets. There have been efforts to develop datasets for the visual question answering (VQA) task since the first large-scale benchmark was introduced in [3]. Early released datasets [14, 30, 47, 52] contain natural images and related questions, where understanding the visual and textual contents is essential for question answering. Some recent datasets introduce questions that involve more diverse visual scenes or require external knowledge to answer, which leads to more complex visual and semantic reasoning for question answering. For example, CLEVR [21] is a synthetic dataset that serves as a diagnostic test for a range of visual reasoning abilities over combinations of three object

shapes. However, these datasets are limited to the natural image domain and pay little attention to abstract diagrams, which also have informative semantics and wide applications.

Diagram QA Datasets. To address the need for vision-and-language reasoning for diagrams, several abstract diagram QA datasets have been developed. For example, abstract VQA [3, 55] considers the task of answering questions on abstract scenes. Similarly, NLVR [48], FigureQA [23], and DVQA [22] feature diagrams of scientific plots that are generated with several figure types and question templates. However, questions and diagrams in these datasets are generated from limited templates, leading to the existence of unintended linguistic shortcuts for question answering. Some more works have proposed datasets of middle school math or science problems in more practical and complex scenarios [46, 27, 43, 44, 37]. A central limitation of the subject QA datasets is that they require complex domain-specific knowledge, which makes disentangling visual reasoning and domain knowledge difficult. Herein, we address these limitations by introducing the IconQA dataset, where only elementary commonsense is required. Through IconQA, we aim to provide a new benchmark for abstract scene understanding and learning different visual reasoning skills in *real-world* scenarios.

VQA Methods. Early VQA approaches usually combine multi-modal inputs by applying attention mechanisms over image regions or question words [28, 39, 38, 12, 54, 11]. Inspired by the semantic nature of VQA images, a line of approaches adopt object proposals from pre-trained object detectors and learn their semantic relationships [28, 54, 11]. As Transformers achieve excellent performance on vision tasks, pioneering works have attempted to use pre-trained models to learn visual representations for natural images in the VQA task [36, 33, 6, 29] and achieve significant improvements. However, current VQA models are not capable of extracting meaningful visual representations from abstract diagrams, as they require image embeddings or object proposals learned from natural images. Instead, we develop a strong baseline that feeds spatial patch sequences into a Transformer encoder that is powered by the embedding module pre-trained on our Icon645 dataset.

3 The IconQA Dataset

The IconQA dataset provides diverse questions that require abstract diagram recognition, comprehensive visual reasoning skills, and basic commonsense knowledge. IconQA consists of 107,439 questions split across three different sub-tasks. To the best of our knowledge, IconQA is the largest VQA dataset that focuses on real-world problems with icon images while involving multiple human intelligence reasoning abilities (see Table 4).

3.1 Data Collection

We aim to collect icon-based question answering pairs that involve multiple reasoning skills, such as visual reasoning and commonsense reasoning. To construct the IconQA dataset, which stems from real-world math word problems, we search for open-source math textbooks with rich icon images and diverse topics. Of those, we choose *IXL Math Learning* which compiles popular textbooks aligned to California Common Core Content Standards¹. We ask well-trained crowd workers to collect problems that cover content from pre-K to third grade, as these problems usually contain abstract images and involve little to none complex domain knowledge. With the driven interest of visual reasoning over abstract images, we filter out the questions that do not accompany icon images or only have images in black and white. Redundant or repetitive data instances are also removed. Question choices are randomly shuffled to ensure a balanced answer distribution. See Appendix A for full details of the dataset collection and usage.

3.2 Data Analysis

Finally, we collect 107,439 IconQA data instances, where each data point contains a colored icon image, a natural language question, optional image or text choices, as well as a correct answer. The IconQA dataset consists of 107,439 questions and is divided into train, validation, and test splits with a ratio of 6:2:2, as shown in Table 1. The dataset consists of three sub-tasks: *multi-image-choice*, *multi-text-choice*, and *filling-in-the-blank*. The *multi-image-choice* sub-task is defined as choosing the correct image from a list of image candidates based on a given diagram and its corresponding question. Similarly, the *multi-text-choice* sub-task is defined as a multiple choice question with 2-5

¹<https://www.ixl.com/standards/california/math>

and on average, a question requires 1.63 skills. The detailed statistics are demonstrated in Table 2. In general, the *filling-in-the-blank* sub-task consists of questions that require the most number of skills, averaging 1.81 skills per question. 9.25% of the *filling-in-the-blank* questions require 3 skills. As the examples from IconQA shown in Figure 1, the first and second questions require the skills of *scene* understanding and *spatial* reasoning. The third example asks how many sticks exist in the diagram, requiring the basic ability of *counting* and basic *algebra* operations. As stated before, the IconQA dataset requires a wide range of skills for a model to perform well on IconQA.

Skill types	Description
Geometry	Identify shapes, symmetry, transformations
Counting	Count objects, shapes
Comparing	Compare object attributes
Spatial	Identify spatial positions and relations
Scene	Understand abstract scenes
Pattern	Identify next and different patterns
Time	Identify time of clocks, events
Fraction	Perform fraction operations
Estimation	Estimate lengths, large numbers
Algebra	Perform algebraic operations
Measurement	Measure widths, lengths, heights
Commonsense	Apply external knowledge
Probability	Perform probability and statistics operations

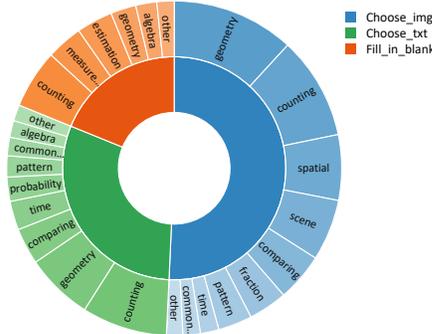


Table 3: Definition of reasoning skill types. Figure 5: Skill distribution in IconQA questions.

Comparisons to Other Datasets. We compare our IconQA dataset with two datasets on natural images and five datasets on abstract diagrams in Table 4. To summarize, IconQA is different from these datasets in various aspects. Unlike natural images (VQA [3], CLEVR [21]) or abstract diagrams like scenes, charts, plots, and illustrations (VQA-Abstract [3], DVQA [22], NLVR [48], AI2D [26], Geometry3K [37]), IconQA features icon images and covers the largest object set of 388 classes. As questions in IconQA stem from real-world math problems and they may describe complex problem scenarios, IconQA has the longest question length among all related datasets. Furthermore, IconQA requires both commonsense and arithmetic reasoning due to its origin from real-world problems. Lastly, IconQA contains more QA task types including answering questions with image choices.

Table 4: Statistics for the IconQA dataset and comparisons with existing datasets.

	#QA	#Image	AvgQ	MaxQ	Image Type	QSource	#Object	#Task	VisualAns	CommonSen	Arithmetic
VQA [3]	614,163	204,721	6.1	23	Natural	Annotated	-	2		✓	
CLEVR [21]	999,968	100,000	18.4	43	Natural	Generated	3	1			
VQA-Abstract [3]	150,000	50,000	6.0	21	Scene	Annotated	131	2			
DVQA [22]	2,325,316	300,000	10.3	23	Bar chart	Generated	-	1			✓
NLVR [48]	92,244	92,244	11.2	25	Scatter plot	Generated	3	1			
Geometry3K [37]	3,002	2,342	10.1	46	Diagram	Real-world	4	1			✓
AI2D [26]	4,563	4,903	9.8	64	Illustration	Real-world	-	1		✓	
IconQA (Ours)	107,439	96,817	8.4	73	Icon image	Real-world	388	3	✓	✓	✓

3.3 Impact and Ethics

Impact & Usage. IconQA is useful for not only follow-up research projects but also real-world applications (e.g. K-6 education applications like tutoring assistants). Moreover, visual recognition in the abstract domain is essential to general AI agents, but rarely investigated in the community, posing new challenges in abstract and symbolic visual reasoning – a natural ability of human.

Social Ethics. Unlike VQA datasets in the natural image domain, IconQA is completely built upon abstract icon images. Therefore, it is less likely to be used in surveillance systems that might infringe on people’s privacy. Moreover, due to the abstract nature of the dataset, IconQA does not contain any sensitive personal information such as gender and race, nor does it contain data that might exacerbate biases against under-represented communities. Therefore, after careful examinations of our dataset, we think the dataset is unlikely to be used to harm people directly.

4 The Icon645 Dataset

As discussed in Section 3.2, IconQA questions are accompanied by abstract diagrams that cover a wide range of icon objects. Using existing backbone networks to extract image representations for

Table 5: Collected icon examples in the Icon645 dataset.

Icons	Examples	Icons	Examples
Bed		Bucket	
Cake		Car	
Castle		Dog	
Giraffe		Kite	
Soda		Tree	

these icon images is inadequate, as most of these networks are pre-trained on natural images. To overcome the limitation, we develop a new large-scale icon dataset for pre-training existing vision backbone networks. We use the collected icon data to pre-train the current backbone networks, which can be applied to extract diagram representations in IconQA.

We retrieve the 388 icon classes mentioned in the question texts from Flaticon², the largest database of free vector icons. After removing 11 classes that can't be retrieved, we construct an icon dataset containing 377 classes, called Icon645. As summarized in Table 10 (Appendix), the Icon645 dataset includes 645,687 colored icons with a minimum size of 64 by 64 and a maximum size of 256 by 256. Examples in Table 5 show that our collected icons include a wide variety of colors, formats and styles. On top of pre-training encoders, the large-scale icon data could also contribute to future research on abstract aesthetics and symbolic visual understanding. In this work, we use the icon data to pre-train backbone networks on the icon classification task in order to extract semantic representations from abstract diagrams in IconQA. See Appendix B for the details of data collection and analysis.

5 Benchmarks

In this section, we first develop a patch cross-modal Transformer model (Patch-TRM) as a strong baseline for the IconQA task. To benchmark the IconQA dataset, we consider multi-modal pooling methods with attention mechanisms [2, 28, 11, 54], Transformer-based VQA approaches [36, 6, 53, 29], and three blind study methods as benchmark models, as summarized in Figure 6. Additionally, a user study is conducted to explore the performances of human beings in different age groups. In the sections below, we discuss the main principles of the core networks in the benchmarks we performed.

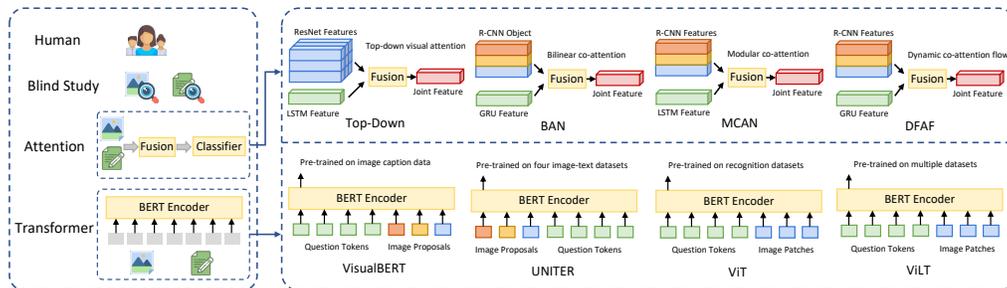


Figure 6: An overview of benchmark baselines on the IconQA task.

5.1 Our Baseline Model

Inspired by recent advances Transformer has achieved in vision-language tasks [33, 36], we develop a cross-modal Transformer model Patch-TRM for icon question answering. Taking the *multi-image choice* sub-task as an example, the overall architecture is shown in Figure 7. The diagram is first parsed into ordered patches in a hierarchical pyramid layout. These patches are then encoded by a pre-trained ResNet and passed through a vision Transformer. Question text is encoded by a language Transformer and fused with patch embeddings via the attention mechanism. The encoded image choices are concatenated with the joint diagram-question representation and then fed to a

²Flaticon: <https://www.flaticon.com/>

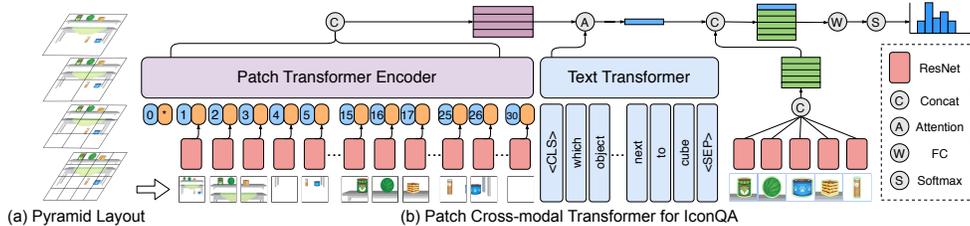


Figure 7: Our IconQA baseline Patch-TRM. Patch-TRM takes patches parsed from a hierarchical pyramid layout and embeds them through ResNet pre-trained on our Icon645 dataset. The joint diagram-question feature is learned via cross-modal Transformers followed by the attention module.

classifier for question answering. The other two sub-tasks utilize similar network architectures, except that in the *multi-text-choice* sub-task, we use an LSTM encoder [17] for choice embedding, while *filling-in-the-blank* does not need a choice encoder.

Current dominant VQA methods either rely heavily on the ResNet backbone network to extract image features or depend on the Transformer encoders to learn image embeddings. However, these networks are pre-trained on natural images and are likely to fail to extract meaningful representations or reasonable object proposals when processing the diagrams in IconQA. Instead, we pre-train the ResNet network on the icon classification task with the icon dataset we compiled (Section 4). Patch-TRM hierarchically parses the diagram into patches that retain complete objects to a large extent, and the parsed patches are embedded by the pre-trained ResNet network before being fed into the vision Transformer. The hierarchical parsing structure, along with the ResNet pre-trained on icon data facilitate our Patch-TRM to learn informative diagram representations for the IconQA task. More details of the pre-training task are discussed in Section 6.4.

5.2 Benchmark Methods

Attention models. We construct four attention models for benchmarking. The first model implements Top-Down attention [2] for VQA, which is a strong attention method that applies free-form based attention on image representations from a pre-trained ResNet-101 network. The remaining three models utilize the bottom-up attention mechanism with the help of object detection proposals from Faster-RCNN [42]. Specifically, BAN [28] proposes a method that utilizes bilinear attention distributions to learn joint vision-language information. DFAF [11] is an advanced model that applies self-attention and cross-modal attention and introduces the information flow to help the model focus on target question words and image regions. The last approach, MCAN [54], learns the self-attention on the questions and images and the question-guided-attention of images jointly.

Transformer models. Four Transformer-based models are also implemented as benchmarks. ViL-BERT [36] and UNITER [6] are two Transformer-based approaches that take image object proposals from Faster-RCNN [42] and question tokens as inputs. Specifically, ViLBERT learns the joint representation of the image content and the natural language content from image proposal regions and question tokens, while UNITER processes multimodal inputs simultaneously for joint visual and textual understanding. The last two benchmarks ViL [53] and ViLT [29] are more recently proposed Transformer models that take image patch tokens instead of object proposals as inputs when representing the image.

Blind study models. We develop three models to check for possible data biases in the IconQA dataset. A random baseline picks up one from the given choice candidates for the *multiple-choice* sub-tasks while predicts the answer by randomly selecting one from all possible answers in the train data for the *filling-in-the-blank* sub-task. Q-Only is set up similar to the Top-Down [2] model, but it only considers textual inputs. This baseline learns the question bias in the training set. I-Only also has a Top-Down architecture, but it only takes abstract diagrams as inputs, and tests the distribution biases in the images and answers in IconQA.

User study. To assess human performances in the IconQA task, we post the test set of IconQA on Amazon Mechanical Turk (AMT) and ask people to provide answers to the questions in the test set. We also ask the participants to provide us with their age group anonymously. We strongly encourage parents who have young children to let their children complete the questionnaires, as their answers

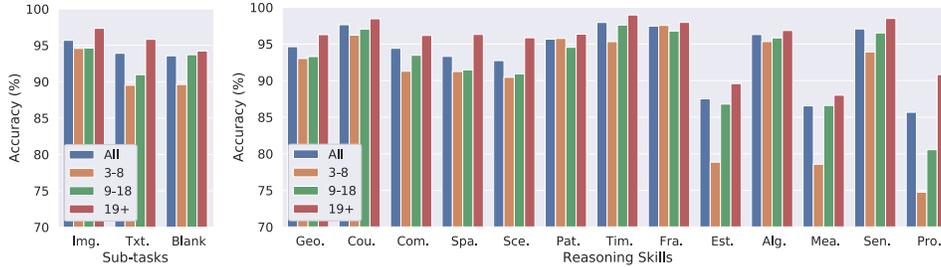


Figure 8: Performance of humans in different age groups for the IconQA task. **Left:** Accuracy over three sub-tasks; **Right:** Accuracy over thirteen reasoning skills.

Table 6: Results on the IconQA dataset.

Method	Sub-tasks (3)			Reasoning skills (13)												
	Img.	Txt.	Blank	Geo.	Cou.	Com.	Spa.	Sce.	Pat.	Tim.	Fra.	Est.	Alg.	Mea.	Sen.	Pro.
Human	95.69	93.91	93.56	94.63	97.63	94.41	93.31	92.73	95.66	97.94	97.45	87.51	96.29	86.55	97.06	85.67
Random	41.70	36.87	0.29	30.30	18.38	41.20	36.49	34.25	34.81	35.82	34.84	3.62	11.12	0.36	45.16	38.81
Q-Only	41.64	36.86	28.45	38.03	33.63	48.19	37.14	35.37	33.66	48.09	33.06	40.46	28.02	38.07	45.25	40.76
I-Only	41.56	36.02	46.65	38.71	37.64	45.26	37.52	35.47	36.29	47.37	32.48	62.29	31.73	64.02	45.25	37.51
Top-Down [2]	75.92	68.51	73.03	80.07	65.01	80.65	45.78	58.22	55.01	68.28	72.43	99.54	50.00	99.46	84.54	83.75
BAN [28]	76.33	70.82	75.54	79.99	67.56	82.12	53.20	66.92	55.67	66.50	73.77	97.06	47.46	96.50	82.12	82.45
ViLBERT [33]	76.66	70.47	77.08	80.05	71.05	75.60	49.46	58.52	62.78	66.72	74.09	99.22	50.62	99.07	81.78	70.94
MCAN [54]	77.36	71.25	74.52	79.86	68.94	82.73	49.70	62.49	54.79	68.00	76.20	99.08	47.32	98.99	83.25	84.87
DFAF [11]	77.72	72.17	78.28	81.80	70.68	81.69	51.42	67.01	56.60	67.72	77.60	99.02	50.27	98.83	84.11	85.70
UNITER [6]	78.71	72.39	78.53	81.31	71.01	83.67	48.34	61.25	60.81	69.77	78.37	99.41	49.18	99.38	86.10	87.84
ViT [53]	79.15	72.34	78.92	82.60	70.84	82.12	54.64	68.80	58.46	68.66	77.41	98.95	51.10	98.76	84.72	86.07
ViLT [29]	79.67	72.69	79.27	82.61	71.13	84.95	53.38	66.72	59.22	69.99	75.81	99.02	50.55	98.91	86.10	87.65
Patch-TRM (Ours)	82.66	75.19	83.62	81.87	77.81	87.00	55.62	62.39	68.75	77.98	82.13	98.24	56.73	97.98	92.49	95.73

give us insights to how the designed audience of these questions perform. Further details about the user study are included in Appendix D.

6 Experiments

6.1 Training Details

Following prior work [3], all the baselines are trained on the IconQA training set, tuned on the validation set, and finally evaluated on the test set. Similar to [3], we choose accuracy as the evaluation metric. For the two *multi-choice* sub-tasks, the answer is regarded as correct only if it matches the ground truth. On the other hand, as the collected answers for *filling-in-blank* are short numbers, correct answers are expanded to include both the digital number and its corresponding words. More details of the benchmark setups and implementations can be found in Appendix E.1.

Our benchmarks and baselines are implemented using PyTorch. All experiments are run on one Nvidia RTX 3090 GPU. We use the Adamax optimizer with optimal learning rates of 7×10^{-4} , 8×10^{-4} , and 2×10^{-3} on the three sub-tasks respectively. We apply a binary cross-entropy loss to train the multi-class classifier with a batch size of 64 and a maximum epoch of 50. The early stopping strategy is used when the validation accuracy stops improving for five consecutive epochs. It takes about 50, 30, and 10 minutes to train our baseline Patch-TRM on three sub-tasks respectively.

6.2 Experimental Results

Table 6 demonstrates the results of the benchmark methods and our baseline on the IconQA test set. The first three columns of the results represent the three sub-tasks: *multi-image-choice*, *multi-text-choice*, and *filling-in-the-blank* respectively. The remaining 13 columns illustrate the results of these approaches over problems that require different reasoning skills, as defined in Table 3.

Human performance. Out of the 54,896 collected answers, 9,620 are made by young children from age 3 to 8, 19,040 are made by adolescents from age 9 to 18, and 26,236 are made by adults.

Q: Which object is next to the one shaped like a cube?



Figure 9: Text-to-image attention visualization.

Method	Img.	Txt.	Blank
Patch-TRM w/o pre	82.01	72.72	81.67
Patch-TRM w/o att	80.63	68.00	80.29
Patch-TRM w/o pos	81.27	64.98	80.68
Patch-TRM V-CLS	80.15	63.90	70.27
Pyramid 1+4+9+16	82.45	68.76	82.19
Pyramid 1+4+9	80.61	67.42	81.36
Full model	82.96	75.21	83.10

Table 7: Ablation study in IconQA.

The human performance over the three sub-tasks and thirteen skills is illustrated in Figure 8. As expected, young children do not answer the questions as well as adolescents or adults, suggesting that most participants answer their ages correctly. Moreover, the result shows that humans perform more consistently on all sub-tasks compared to machine algorithms. Interestingly, humans are outperformed by models quite significantly in questions that require numerical reasoning skills like *probability*, *measurement*, and *estimation*.

Analysis by Task Types. Humans outperform all benchmarks consistently over there sub-tasks and most reasoning skills. There is still a large gap to fill for future research of abstract diagram understanding and visual reasoning on the icon domain. The results achieved in blind studies of Q-only and I-only are close to random, showing that the IconQA dataset is robust and reliable in distribution. Our proposed Patch-TRM baseline outperforms current state-of-the-art VQA models in all three sub-tasks. These improvements mainly come from two insights: pre-training ResNet on icon images and taking a hierarchical approach with attention mechanism.

Analysis by Reasoning Types. Similarly, the Patch-TRM baseline obtains better results than the benchmarks over most reasoning skill types. Interestingly, in some skills such as *estimation*, *measurement*, and *probability*, Patch-TRM performs better than average human beings. This implies that neural networks have a promising potential to develop the basic ability of mathematical reasoning.

Quantitative Analysis. We visualize one example with the cross-modal attention map generated by our baseline Patch-TRM in Figure 9. The visualized attention shows that our baseline is capable of attending to the corresponding patch regions with higher weights given the input question.

6.3 Ablation Study

To study the functions of individual components in our model, we conduct an ablation analysis. Table 7 presents the results of different simplifications of our full model, where each implementation is trained on the IconQA train set and tested on the validation set. Instead of ResNet101 pre-trained on the icon classification task, *Patch-TRM w/o pre* utilizes ResNet101 pre-trained on natural image data for patch feature extraction. The decreasing performance of 0.95-2.49% indicates that pre-training backbones on tasks within similar domains is critical to downstream tasks. The attention mechanism helps to combine the image and question representations and improves the model performance by up to 7% compared to using simple concatenation (denoted as *Patch-TRM w/o att*). Positional embeddings of the ordered diagram patches benefit the vision Transformer by enabling it to learn spatial relationships among the patches, compared to the baseline without position embeddings (*Patch-TRM w/o pos*). *Patch-TRM V-CLS* uses the output embedding of [CLS] token as the diagram feature instead, which leads to a drastic performance decline. We have also experimented with coarse-grained patch cropping (e.g., *Pyramid 1+4+9+16* denotes 30 patches, *Pyramid 1+4+9* denotes 14 patches), which results in a performance degradation of 0.51% to 7.79%.

6.4 Icon Classification for Pre-training

The Icon645 dataset is collected to pre-train the backbone network for patch feature extraction. The dataset has a long-tailed distribution, and thus we address the class-imbalanced issue following previous studies on specific loss functions such as CB loss [8], Focal loss [34], and LDAM loss [5]. The metric of Top-5 accuracy is used to evaluate different model setups and the evaluation results are summarized in Table 8.

Table 8: Results for icon classification.

Method	Total	Head	Medium	Tail
ResNet32 [16] + CB [8]	27.91	19.66	36.51	33.53
ResNet32 [16] + Focal Loss [34]	32.80	51.59	36.51	8.94
ResNet32 [16] + LDAM [5]	42.65	55.68	46.42	24.94
ResNet101 [16] + LDAM [5]	62.93	70.29	70.50	47.51

Following [35], to demonstrate performances on different data parts, we divide the dataset into three balanced clusters: Head, Medium, and Tail, corresponding to 132, 122, and 123 classes respectively. All classes in Head have at least 1,000 instances, all classes in Medium have 300 - 999 instances, and all classes in Tail have fewer than 300 instances. As the results show, the backbone network ResNet101 with a re-balanced LDAM loss function achieves the best result for icon classification on Icon645. Consequently, we adopt this pre-trained ResNet101 network to extract patch features in our baseline Patch-TRM for IconQA.

7 Conclusion

In this work, we introduce IconQA, an open-source dataset of icon question answering in real-world scenarios for assessing the abilities of abstract diagram understanding and visual language reasoning. IconQA features 107,439 questions, three sub-tasks, and thirteen types of cognitive reasoning skills. We benchmark the IconQA task extensively with a user study, three blind studies, as well as multiple existing attention-based and Transformer-based approaches. We further develop a strong baseline, Patch-TRM, which parses the diagram in a pyramid layout and applies cross-modal Transformers with attention mechanism to learn the meaningful joint diagram-question feature. Additionally, we introduce Icon645, a large-scale icon dataset that is useful to pre-train the diagram encoding network used in Patch-TRM for the IconQA task.

By releasing a new dataset of icon question answering for abstract diagram understanding and visual language reasoning, we envision that IconQA will facilitate a wide range of research in computer vision and natural language processing, as well as smart education applications like tutoring systems, to invent the future of AI for science education.

References

- [1] Vedika Agarwal, Rakshith Shetty, and Mario Fritz. Towards causal vqa: Revealing and reducing spurious correlations by invariant and covariant semantic editing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9690–9698, 2020.
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision (CVPR)*, pages 2425–2433, 2015.
- [4] Jack Bandy and Nicholas Vincent. Addressing "documentation debt" in machine learning research: A retrospective datasheet for bookcorpus. *arXiv preprint arXiv:2105.05241*, 2021.
- [5] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arachiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1567–1578, 2019.
- [6] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [7] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [8] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9268–9277, 2019.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT (NAACL)*, 2018.
- [10] Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1615–1625, 2017.

- [11] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li. Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6639–6648, 2019.
- [12] Peng Gao, Hongsheng Li, Shuang Li, Pan Lu, Yikang Li, Steven CH Hoi, and Xiaogang Wang. Question-guided hybrid convolution for visual question answering. In *The European Conference on Computer Vision (ECCV)*, pages 469–485, 2018.
- [13] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*, 2018.
- [14] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [15] Kotaro Hara, Abigail Adams, Kristy Milland, Saiph Savage, Chris Callison-Burch, and Jeffrey P Bigham. A data-driven analysis of workers’ earnings on amazon mechanical turk. In *Proceedings of the 2018 CHI conference on human factors in computing systems*, pages 1–14, 2018.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, pages 770–778, 2016.
- [17] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [18] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6700–6709, 2019.
- [19] Huaizu Jiang, Ishan Misra, Marcus Rohrbach, Erik Learned-Miller, and Xinlei Chen. In defense of grid features for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10267–10276, 2020.
- [20] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. *arXiv preprint arXiv:1909.10351*, 2019.
- [21] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2901–2910, 2017.
- [22] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5648–5656, 2018.
- [23] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. *arXiv preprint arXiv:1710.07300*, 2017.
- [24] Takuro Karamatsu, Gibran Benitez-Garcia, Keiji Yanai, and Seiichi Uchida. Iconify: Converting photographs into icons. In *Proceedings of the 2020 Joint Workshop on Multimedia Artworks Analysis and Attractiveness Computing in Multimedia*, pages 7–12, 2020.
- [25] Sevilay Karamustafaoglu. Improving the science process skills ability of science student teachers using i diagrams. *International Journal of Physics & Chemistry Education*, 3(1):26–38, 2011.
- [26] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Min Joon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [27] Aniruddha Kembhavi, Minjoon Seo, Dustin Schwenk, Jonghyun Choi, Ali Farhadi, and Hannaneh Hajishirzi. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4999–5007, 2017.
- [28] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1571–1581, 2018.

- [29] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pages 5583–5594. PMLR, 18–24 Jul 2021.
- [30] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, pages 32–73, 2017.
- [31] Manuel Lagunas, Elena Garces, and Diego Gutierrez. Learning icons appearance similarity. *Multimedia Tools and Applications*, 78(8):10733–10751, 2019.
- [32] Li-Jia Li, Hao Su, Li Fei-Fei, and Eric P Xing. Object bank: A high-level image representation for scene classification & semantic feature sparsification. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1378–1386, 2010.
- [33] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [34] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision (ICCV)*, pages 2980–2988, 2017.
- [35] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2537–2546, 2019.
- [36] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 13–23, 2019.
- [37] Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-Chun Zhu. Intergps: Interpretable geometry problem solving with formal language and symbolic reasoning. In *The 59th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2021.
- [38] Pan Lu, Lei Ji, Wei Zhang, Nan Duan, Ming Zhou, and Jianyong Wang. R-vqa: learning visual relation facts with semantic attention for visual question answering. In *The ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, pages 1880–1889, 2018.
- [39] Pan Lu, Hongsheng Li, Wei Zhang, Jianyong Wang, and Xiaogang Wang. Co-attending free-form regions and detections with multi-modal multiplicative feature embedding for visual question answering. In *The AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [40] Spandan Madan, Zoya Bylinskii, Matthew Tancik, Adrià Recasens, Kimberli Zhong, Sami Alsheikh, Hanspeter Pfister, Aude Oliva, and Fredo Durand. Synthetically trained icon proposals for parsing and summarizing infographics. *arXiv preprint arXiv:1807.10441*, 2018.
- [41] Maria Martiniello. Language and the performance of english-language learners in math word problems. *Harvard Educational Review*, 78(2):333–368, 2008.
- [42] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems (NeurIPS)*, pages 91–99, 2015.
- [43] Mrinmaya Sachan, Kumar Dubey, and Eric Xing. From textbooks to knowledge: A case study in harvesting axiomatic knowledge from textbooks to solve geometry problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 773–784, 2017.
- [44] Mrinmaya Sachan, Kumar Avinava Dubey, Tom M Mitchell, Dan Roth, and Eric P Xing. Learning pipelines with limited data and domain knowledge: A study in parsing physics problems. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 140–151, 2018.
- [45] Mike Schuster and Kaisuke Nakajima. Japanese and korean voice search. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5149–5152. IEEE, 2012.
- [46] Minjoon Seo, Hannaneh Hajishirzi, Ali Farhadi, Oren Etzioni, and Clint Malcolm. Solving geometry problems: Combining text and diagram interpretation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1466–1476, 2015.

- [47] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8317–8326, 2019.
- [48] Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A corpus of natural language for visual reasoning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 217–223, 2017.
- [49] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962*, 2019.
- [50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems (NeurIPS)*, pages 5998–6008, 2017.
- [51] Shuo Wang, Yizhou Wang, and Song-Chun Zhu. Learning hierarchical space tiling for scene modeling, parsing and attribute tagging. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(12):2478–2491, 2015.
- [52] Xinyu Wang, Yuliang Liu, Chunhua Shen, Chun Chet Ng, Canjie Luo, Lianwen Jin, Chee Seng Chan, Anton van den Hengel, and Liangwei Wang. On the general value of evidence, and bilingual scene-text visual question answering. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [53] Ildoo Kim Wonjae Kim, Bokyung Son. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. under review.
- [54] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6281–6290, 2019.
- [55] Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Yin and Yang: Balancing and answering binary visual questions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [56] Jun Zhu, Tianfu Wu, Song-Chun Zhu, Xiaokang Yang, and Wenjun Zhang. A reconfigurable tangram model for scene representation and categorization. *IEEE Transactions on Image Processing*, 25(1):150–166, 2015.
- [57] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] Please see Appendix A.13.
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] Please see Section 3.3.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments (e.g. for benchmarks)...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] Please see our project page at <https://iconqa.github.io>.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Please see Section 6.1 for training details. For details on benchmark model settings, see Appendix E.1.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] Please see Section 6.1.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes] See <https://github.com/lupantech/IconQA#license>.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] All datasets are available on the IconQA website <https://iconqa.github.io>, or the github repository <https://github.com/lupantech/IconQA>.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] As we discuss in Section 3.3, our datasets do not contain identifiable or offensive content.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [Yes] See Appendix D.
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [Yes] See Appendix D.3.