

Multimodal Model for Diagram Question Answering

George Washington University
Capstone Fall 22 - DATS 6501_80
Alexis Kaldany, Joshua Ting
[GitHub Repo](#)
[Final Presentation](#)

Table of Contents

Multimodal Model for Diagram Question Answering	1
Table of Contents	2
1. Introduction	3
2. Background	4
2.1 Prior Work	4
2.2 Our Goals and Planned Contributions	5
2.3 Data Description	5
2.4 Selected Framework and Environment Setup	6
2.5 Metrics Selection	6
3. Exploratory Data Analysis	9
4. Data Preprocessing	11
4.1 Data Acquisition	11
4.2 Splitting Data	12
4.3 Types of Input Scenarios	12
4.4 Embeddings	13
4.5 Image Annotations	14
4.6 String Annotations	15
5. Modeling	16
5.1 Model Selection	16
5.2 Model Training	16
6. Results	19
7. Conclusions	22
7.1 Lessons Learned	22
8. Appendix	24
8.1 References	24

1. Introduction

Our Capstone project can be classified as a *Visual Question and Answering* task, a multimodal problem, where a text question and image are used as inputs for the model. The model then generates a text output, the answer to the input question. The images are elementary school level diagrams from textbooks, with appropriate level of complexity questions and answers.

In this paper, we explore several methodologies to perform the task of Diagram Question and Answering. Specifically, we attempt to accurately predict the correct multiple choice answer given a diagram, multiple choice question, and annotations as inputs. Additionally in this paper, we outline the methodologies for:

- Background Research
- Data Overview and Acquisition
- Exploratory Data Analysis
- Data Processing
- Model Selection, Training, and Metrics Evaluation

We decided to work with this data as we had never worked on multi-modal models before. We both had also taken the Data Science of Networks class with Professor Johnson, and correspondingly believed we had the knowledge base to integrate Natural Language Processing, Computer Vision, and Graph-Networks into a model.

Additionally, upon seeing the scores of the paper achieved, and the relative age of the paper (2016 was a long time ago!), we thought we could achieve similar results with newer deep learning models, specifically Transformers based models. We also wanted to reduce the complexity of the model by eliminating the Diagram Parse Graph (described in following sections) and using a single model, start to finish.

We found the complexity of the data and the richness of the annotations took nearly a month to fully integrate into scripts which made the data fully available to experiment with. Then we spent a month working on various iterations of annotation usage. We never received positive results at any point, so we experimented with different visual embedding systems, different annotation uses, different training parameters, and so on, to no avail.

2. Background

2.1 Prior Work

Kembhavi et al. (2016)'s paper, “*A Diagram Is Worth A Dozen Images*”³, extends Visual Question and Answering by first performing it on diagrams rather than on natural images. In their own words,

“From a computer vision perspective, these [diagram] illustrations are inherently different from natural images and offer a unique and interesting set of problems...visual illustrations may depict complex phenomena and higher-order relations between objects.”³

The paper introduces “Diagram Parse Graphs” (DPG), Figure 1, which encodes diagram constituents and their relationships in a graph. They construct these graphs using annotations of the images, which include relationships between various image objects. The graphs were then combined with the question and answers as inputs into their custom LSTM-based architecture, DQA-NET in Figure 2, which was then trained to predict the correct multiple choice question.

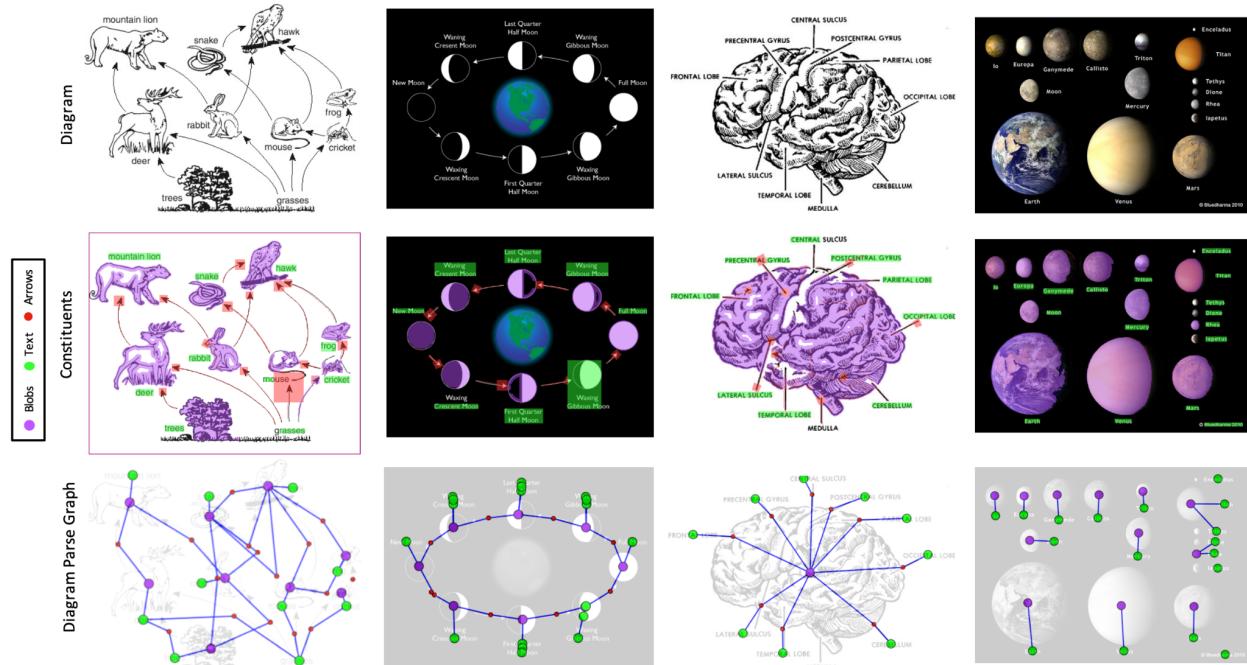


Figure 1: Diagram Parse Graphs³.

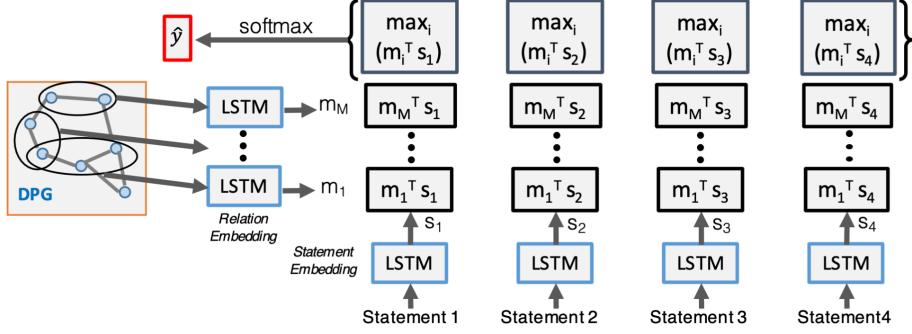


Fig. 3. An overview of the DQA-NET solution to diagram question answering. The network encodes the DPG into a set of facts, learns to attend on the most relevant fact, given a question and then answers the question.

Figure 2: DQA-NET Solution³.

The results of the above methods were shown by the paper to be 38.47% accuracy for the DQA-NET and 32.90% accuracy using the VQA method.

Method	JIG Score
GREEDY SEARCH	28.96
A* SEARCH	41.02
DSDP-NET	51.45

Method	Training set	Accuracy
VQA	VQA	29.06
VQA	AI2D	32.90
DQA-NET	AI2D	38.47

Table 2. (left) Syntactic parsing results, (right) Question answering results

Figure 3: DQA-Net results³.

2.2 Our Goals and Planned Contributions

Since the original paper, Transformers and Auto Encoder-Decoder architectures have been used in highly complex domains such as natural language understanding and image understanding to achieve impressive results in various tasks under each domain respectively.

Our objective is to explore using contemporary pretrained transformer models on this particular dataset and task. We aim to utilize techniques in both NLP and computer vision to generate text and visual embeddings as inputs for our chosen transformer model. The creation of DPGs required extensive feature engineering and training of a separate model to just create the DPGs. We hope to achieve similar accuracy scores to the original paper while reducing the complexity needed to create DPGs and create a more streamlined end-to-end trainable model.

2.3 Data Description

The dataset is composed of 5000 scientific diagrams of roughly elementary school level complexity. Roughly three questions were created per image, with 4 possible answers per question totaling 15,000 diagrams and question-answer pairs. The figure below shows an example of a diagram and question-answer pair.

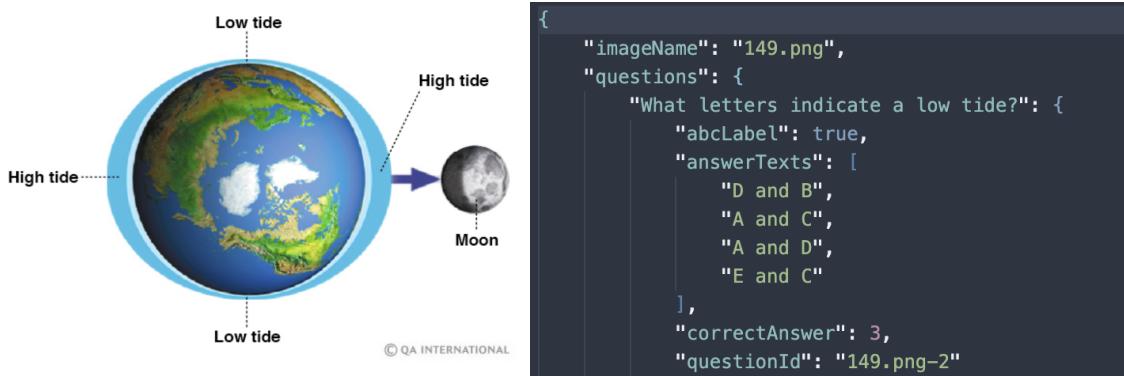


Figure 4: Sample diagram and question-answer example³.

Certain instances of question-answer pairs require annotations to answer the associated questions. Each image has an associated json file containing extensive annotations for each image. These annotations map to areas on the image with bounding boxes.

```
"relationships": {
  "T0+A0+B1": {
    "category": "intraObjectLinkage",
    "connector": "A0",
    "destination": "B1",
    "hasDirectionality": false,
    "id": "T0+A0+B1",
    "origin": "T0"
  },
}
```

Figure 5: Sample annotations example³.

2.4 Selected Framework and Environment Setup

We used PyTorch and Transformers from HuggingFace for our deep learning framework. Our programs are intended to be run on an Ubuntu server containing a GPU (CUDA). A requirements.txt file exists with our codebase containing a list of required dependencies for the virtual environment.

Our team setup 2 separate AWS EC2 instances to process our data and train our models, with each instance having access to GPU. We used GitHub for version control and documentation of our work.

2.5 Metrics Selection

Since this is a multiclass classification problem, we identified our metrics to judge our performance, with the main metric being accuracy. We also chose to include other metrics that

can give us useful insights into the performance of the model such as F1 Score, Recall, Precision, Specificity, and a Confusion Matrix.

Metric
*Accuracy
F1 Score
Recall
Precision
Specificity

Figure 6: Metrics considerations. * main metric to compare against the original paper's results.

2.6 VisualBERT

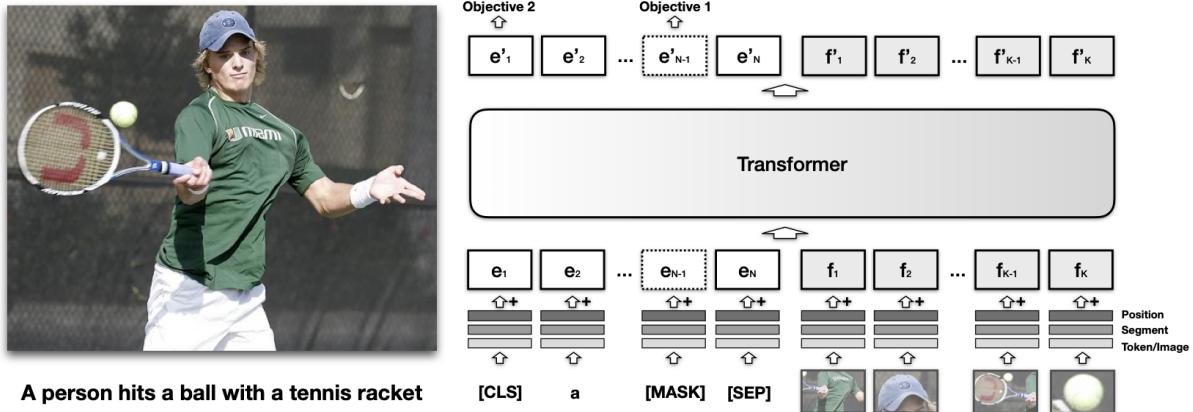


Figure 7: VisualBERT architecture.

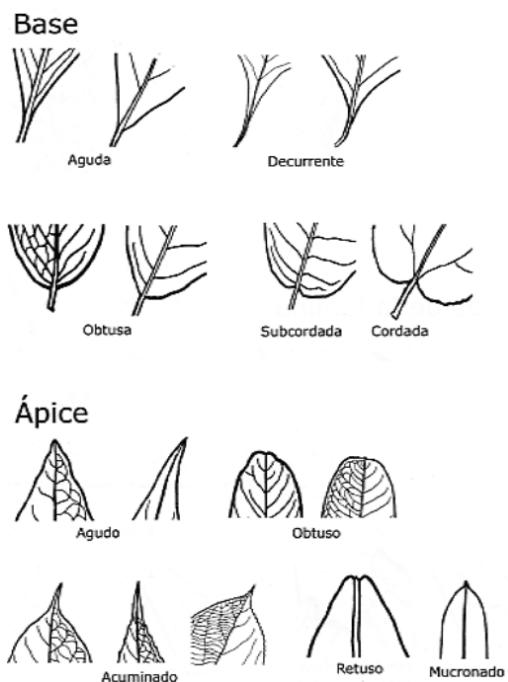
VisualBERT consists of layers of Transformer images that align elements of the input text and regions of the images using self-attention. VisualBERT can connect language and image regions without direct supervision and is somewhat sensitive to syntactic relationships, and can recognize associations between verbs and image regions. Those interactions between words and objects are of great use for Diagram Question and Answering.

3. Exploratory Data Analysis

The AI2D dataset² is hosted for free and open source access, a Creative Commons license, on the AWS Marketplace within a S3 bucket. The dataset consists of 4,817 illustrative diagrams in various dimensions. Within a json file, each diagram is accompanied by at least one question prompt with each question having 4 answer choices, and 1 correct answer. Lastly, each diagram includes human labeled annotations in a json file format. Each illustrative diagram having multiple Q/A results in over 10k samples to allow for training.

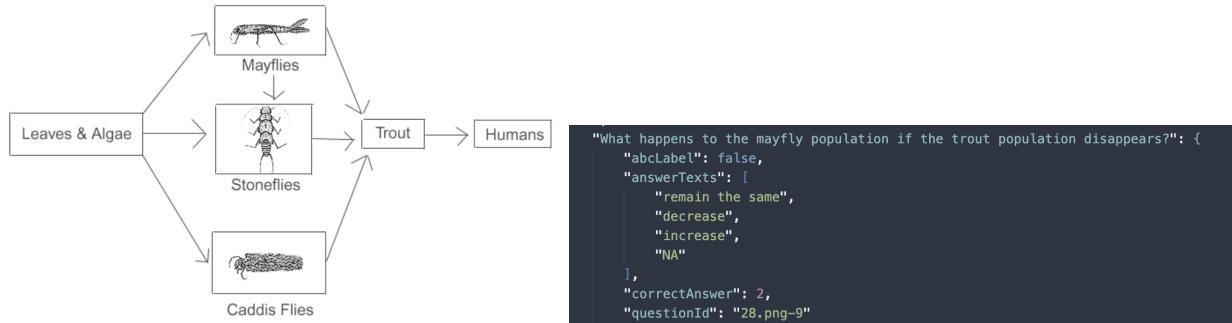
On the surface, this task seems reasonably simple in that the output space is just multiclass classification with 4 classes, but diving in more, the complexity of the problem begins to show itself. For each sample, the output classes differ and are not related between samples. For example, a B label in one sample has entirely different underlying meaning to a B in a different sample. Additionally, the input space is quite complex with both images, annotations, and Q/A pairs.

During EDA, we explored the variety of images and Q/A pairs. The main thing we noticed was that the complexity of each illustrative diagram and Q/A pair differs quite considerably. Some questions merely asked to identify objects within an image while others asked abstract questions that required critical understanding of not just the objects within an image, but also the higher level relationships between objects, then to use logical reasoning to deduce the correct answer.



```
"questions": {
    "Which of these is not apice?": {
        "abcLabel": false,
        "answerTexts": [
            "obtuso",
            "cordada",
            "none of the above",
            "agudo"
        ],
        "correctAnswer": 1,
        "questionId": "4713.png-0"
    }
}
```

Figure 8: Although not trivial, especially for computer vision, this shows an example of a less complex sample from the data.

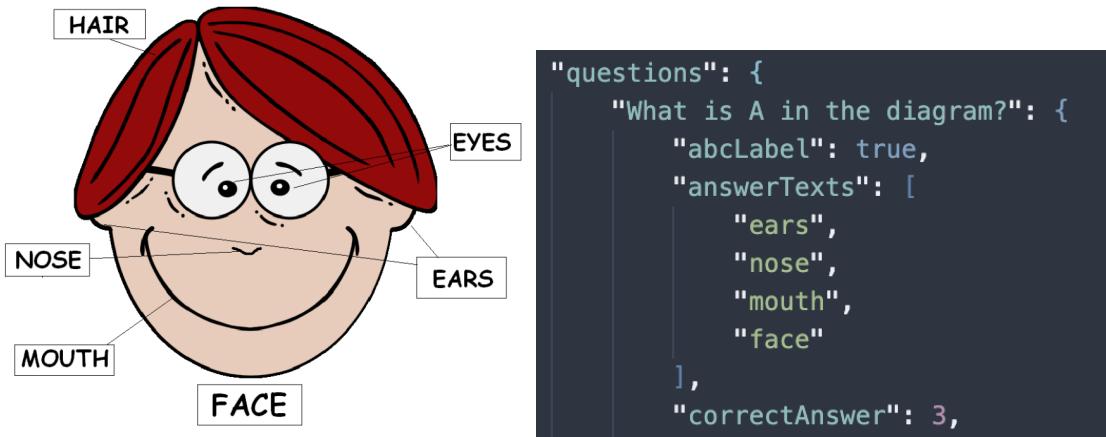


```

"What happens to the mayfly population if the trout population disappears?": {
  "abcLabel": false,
  "answerTexts": [
    "remain the same",
    "decrease",
    "increase",
    "NA"
  ],
  "correctAnswer": 2,
  "questionId": "28.png-9"
}
  
```

Figure 9: Example data sample of a complex Q/A that requires understanding of relationships and logical reasoning.

Additionally, some samples asked questions pertaining to labels within the annotations that are not present on the original image. This meant that the annotations were absolutely crucial for certain samples but slightly less so for others.



```

"questions": {
  "What is A in the diagram?": {
    "abcLabel": true,
    "answerTexts": [
      "ears",
      "nose",
      "mouth",
      "face"
    ],
    "correctAnswer": 3,
    "questionId": "28.png-10"
  }
}
  
```

Figure 10: Example data sample where questions refer to objects within annotations and not in the original image.

4. Data Preprocessing

During preprocessing, we have built several key modules to support our model:

1. **Data Downloader:** Downloads the ai2d dataset to our data folder, unzips the file, and deletes the zip file.
2. **Data Preprocessor:** First there is a function which generates a json containing the local paths to images, annotations, and questions. From that json a different function generates a dataframe with each row being a separate question (images can contain multiple sets of question/answer pairs).
3. **Train/Val/Test Split:** A module to randomly shuffle our dataset then split to training, validation, and blind test datasets at 70%, 20%, and 10% portions respectively.
4. **Image Annotator:** The image annotator module processes the annotations in a json format to draw and label onto the original diagram images using the included coordinates. This module is useful to include both the original diagram image and any annotations within the same input space.
5. **Visual Embedder:** Currently we use Resnet18 to generate visual embeddings of the dimensions needed for *VisualBERT for Multiple Choice*⁵. Generating these embeddings was one of the first challenges that was encountered, and this method of generating embeddings may not be ideal. Other models are being investigated.
6. **Tokenizer:** A pre-trained BERT tokenizer is used to generate our text embeddings. Each sample of our data includes 4 instances of differing text representing each of our different class choices. Each instance is tokenized to generate a separate embedding.
7. **Data Loader:** A custom class to support our unique multimodal data. The main steps of the Data Loader is to yield a batch of data to our model while performing each of the preprocessing modules above to output text and visual embeddings.
8. **Annotation to String:** A series of functions which generate strings from the annotations. This was created after we discovered how poorly the model was working. The intent was to supply the model with additional information by concatenating these string “annotations” to the question.

4.1 Data Acquisition

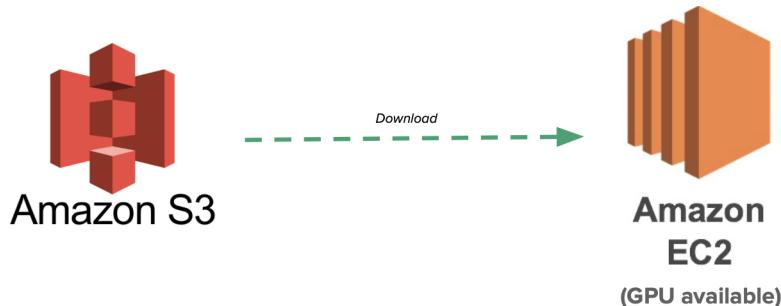


Figure 11: Data acquisition schematic.

Our data is hosted by AWS under the “Open Data Sponsorship Program”. It can easily be downloaded for its S3 bucket. The Allen Institute for AI paid for the labeling of the data and continues to sponsor the dataset.

4.2 Splitting Data

Dataset Split	Choice A	Choice B	Choice C	Choice D	TOTAL
Training	25.3% (2,831)	25.5% (2,845)	25.6% (2,856)	23.6% (2,639)	72.1% (11,171)
Validation	23.8% (296)	26.5% (329)	27.7% (345)	21.9% (272)	8.0% (1,242)
Testing	24.5% (758)	26.0% (802)	26.9% (832)	22.5% (696)	19.9% (3,088)
TOTAL	25.1% (3,885)	25.6% (3,976)	26.0% (4,033)	23.3% (3,607)	15,501

Figure 12: Data acquisition schematic.

We used the same indices as the original paper, so the training, validation, and testing data we used is exactly the same as that used in the paper. There is a slight class imbalance with more correct A values, and less D values being correct. We discovered that our model was choosing A at a much higher frequency than B, B much more than C, C much more than D. We think this may be due to the relative distribution of correct answers inside the training data.

4.3 Types of Input Scenarios

As part of our iterative experimentation we designed multiple input scenarios. All three scenarios differed in the way we attempt to incorporate the annotations into our model. The first approach is to generate text embeddings from the question-answer pairs and a visual embedding from the diagram images using the methods described in sections 4.4 to 4.6. This scenario refrains altogether from incorporating the annotations. The second scenario is to add annotations into the inputs by drawing the annotations straight onto the diagram images before turning them into the aforementioned visual embeddings. The last scenario is to turn the annotations into strings that we concatenate to the end of the question-answer pairs before turning them into text embeddings.

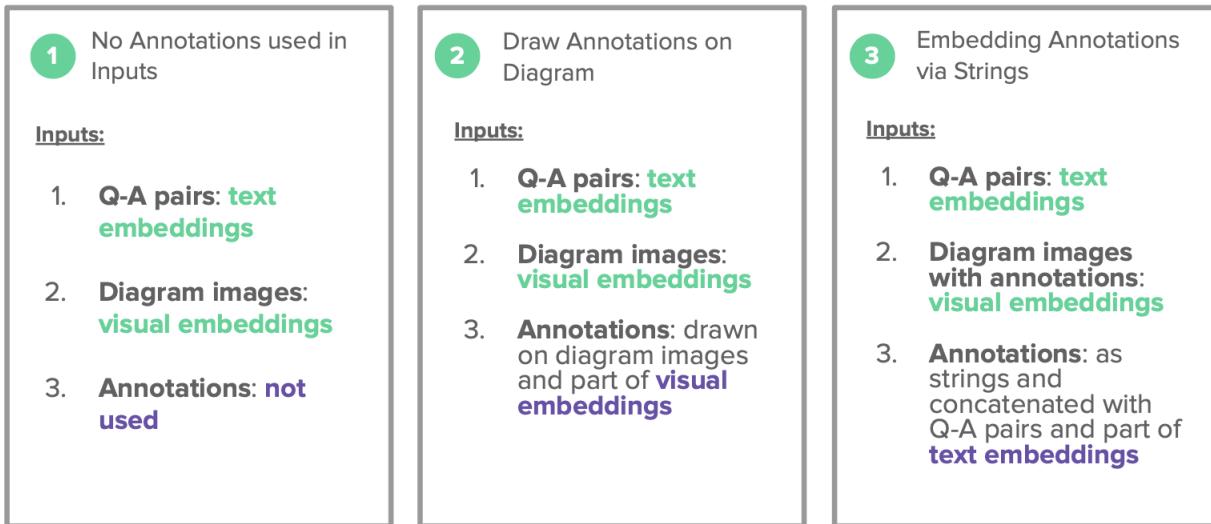


Figure 13: Three different model input scenarios.

Our team felt that processing the annotations were key to the effectiveness of our model to learn and produce results so we decided to try a few different methods in order to see the effects on the overall process.

4.4 Embeddings

For text embeddings we used the BERT Tokenizers. This was necessary as VisualBERT requires BERT tokens as input. The text input is composed of the question, and for scenario three with the annotation strings, also the annotations in the form of the string.

For image embeddings we used ResNET 18, taking the average pooling layer as our visual embeddings. We resize the images to 224x224 and feed them into the ResNET model right before they are entered into the main VisualBERT model.

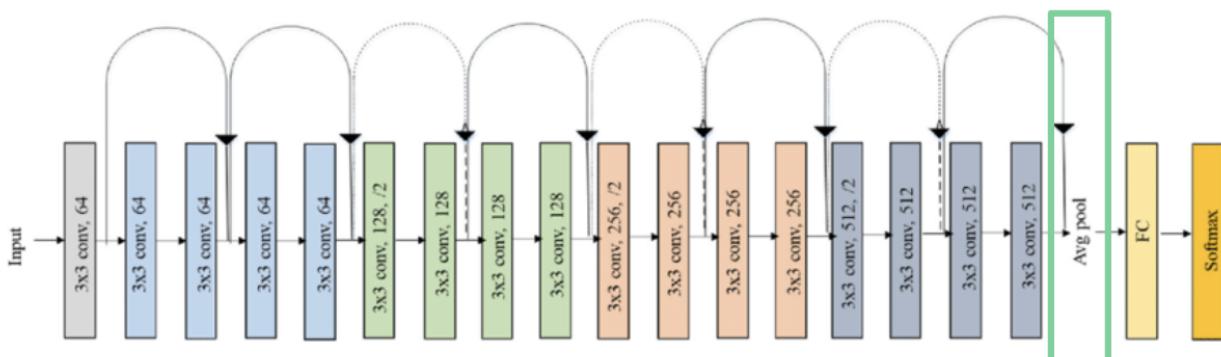


Figure 14: ResNET18 Model using last feature map before fully connected layers as embeddings.

4.5 Image Annotations

Our team worked on applying the annotations straight onto the diagram images before turning them into visual embeddings. A sample of the annotated diagrams can be found below.

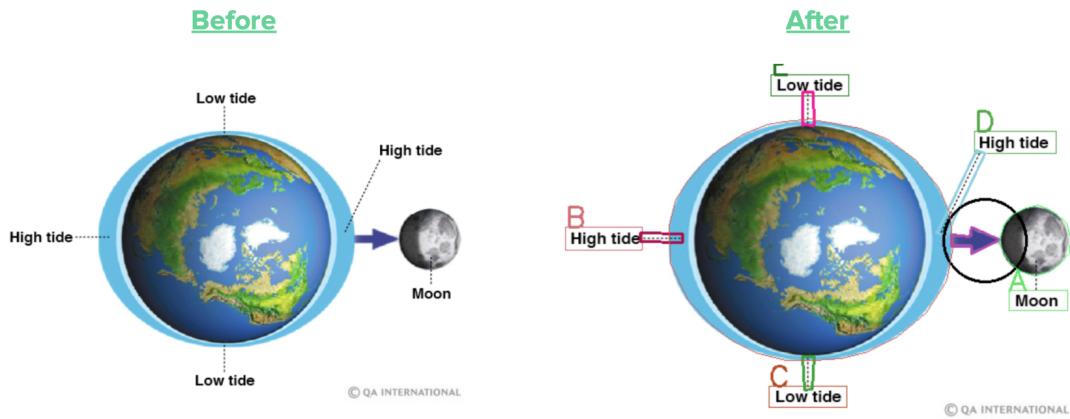
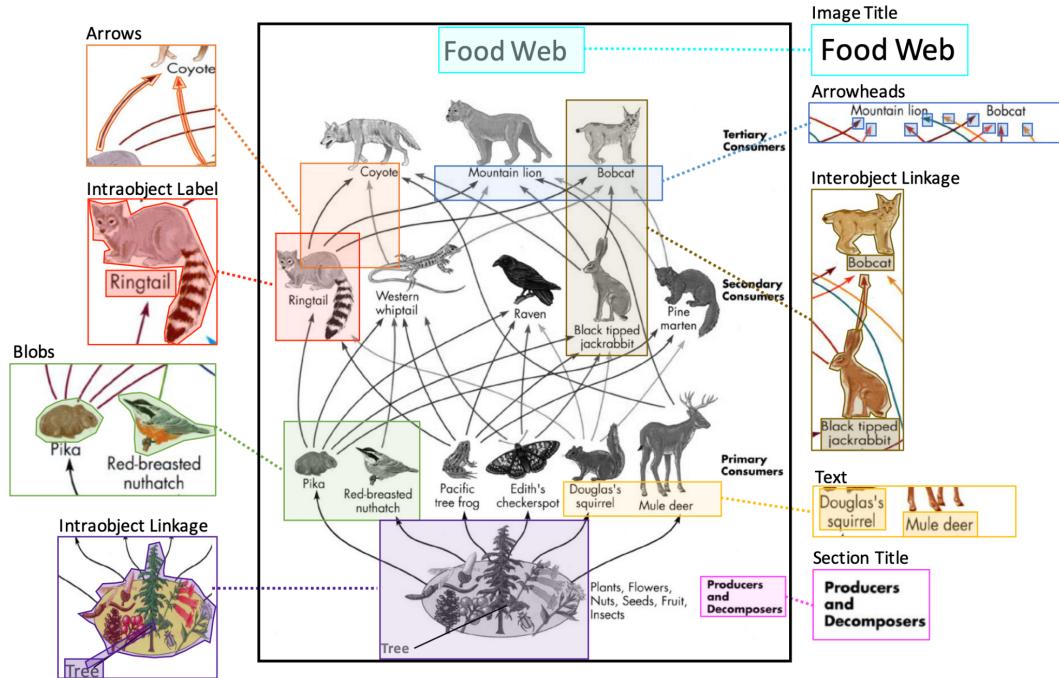


Figure 15: Before and after annotations are applied onto the diagram image

In this input situation we draw the annotations directly on the image. Each of the below objects had a bounding box within the image which could be drawn as shown from the original paper below.



Multiple Choice Question: From the above food web diagram, what will lead to an increase in the population of deer? a) increase in lion b) decrease in plants c) **decrease in lion** d) increase in pika

Figure 16: Sample data example with diagram, annotations, and graph relationships³

Each representation of the annotations has a specific meaning and can be summarized in the list below. There are objects that represent regions, labels, linkages, and captions, and other rich information that is crucial for our model.

-
- Intra-Object Label (R₁)**: A text box naming the entire object.
- Intra-Object Region Label (R₂)**: A text box referring to a region within an object.
- Intra-Object Linkage (R₃)**: A text box referring to a region within an object via an arrow.
- Inter-Object Linkage (R₄)**: Two objects related to one another via an arrow.
- Arrow Head Assignment (R₅)**: An arrow head associated to an arrow tail.
- Arrow Descriptor (R₆)**: A text box describing a process that an arrow refers to.
- Image Title (R₇)**: The title of the entire image.
- Image Section Title (R₈)**: Text box that serves as a title for a section of the image.
- Image Caption (R₉)**: A text box that adds information about the entire image, but does not serve as the image title.
- Image Misc (R₁₀)**: Decorative elements in the diagram.
-

Table 1. Different types of relationships in our diagram parse graphs.

Figure 17: Types of annotations³

4.6 String Annotations

Rather than drawing the annotations directly onto the image, the “*A Diagram Is Worth A Dozen Images*”³ team transformed the annotations into strings, which were subsequently tokenized and embedded in their model. Considering the structure of *VisualBERT for Multiple Choice*⁵, namely that we can only enter one string alongside the visual embeddings of the image, we intend to combine the question string with the annotation string and tokenize them together. This has potential to yield better results than the configurations we have already tried.

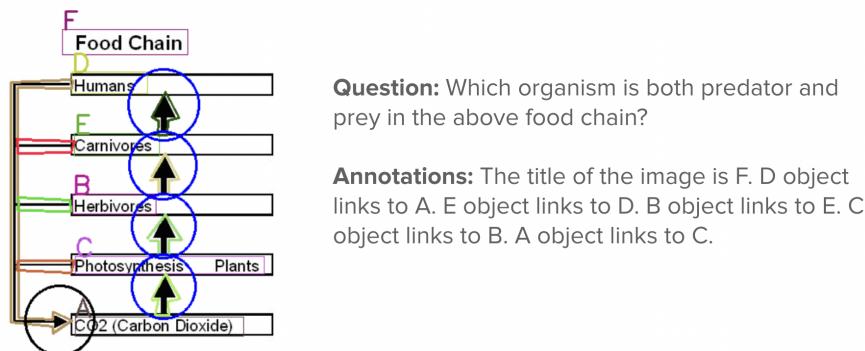


Figure 18: Sample annotations as a string.

5. Modeling

5.1 Model Selection

During model selection, we explored several model options readily available on Hugging Face before finally selecting Li et. al's *VisualBERT for Multiple Choice*⁵. This model choice seemed the most promising as it allowed input embeddings for both text and images while also having a prebuilt Multiple Choice task head. A couple of other models we considered included *VisualBERT* with a Q&A head, Kim et al's *ViLT*⁶, and Huang et al's *LayoutMV3*⁷.

The table below shows the different considerations we explored regarding each of the model's pre-training objectives and available key heads before selecting our final model. All three models have a masked language model pre-training objective while they have different methods to try and learn the joint representation of text and image. VisualBERT uses sentence-image prediction on captions, ViLT uses both image text matching and word patch alignment, while LayoutMV3 uses masked image modeling and word patch alignment.

Name	Type	Pre-Training Objectives	Key Heads
VisualBERT ⁵	Vision-Language Transformer Model	1. Bidirectional masked language model 2. Sentence-image prediction on caption data	1. Multiple Choice 2. Question Answering 3. Visual Reasoning
ViLT ⁶	V-L Transformer Model without Convolution for Visual Embedding	1. Masked language model with whole word masking 2. Image text matching 3. Word patch alignment: predict masked image patches of a text word	1. Question Answering
LayoutMV3 ⁷	Document Transformer Model	1. Bidirectional masked language model 2. Masked image model 3. Word patch alignment	1. Question Answering

Figure 19: Model selection considerations.

5.2 Model Training

As mentioned in *Section 4.3*, our team had three different input configurations for our model. We finetuned the VisualBERT model on each of the input configurations for 16 epochs each.

We have selected the following architecture choices and hyperparameter selections:

1. **Optimizer:** AdamW()¹⁶
2. **Criterion:** torch.nn.CrossEntropy()¹⁷

3. Learning Rate: 5e-5

Each epoch of training took roughly an hour each. Therefore, we built a custom module to allow for model checkpointing and the ability to resume training from a saved checkpoint in order to allow for iterative experimentation during training.

As seen below by the training curves for each of the input configuration scenarios, the validation accuracy for the first configuration improved as training went on but not for the other two configurations. We surmise that the addition of annotations in the format that we tried complicated the input space and prevented the model from learning.

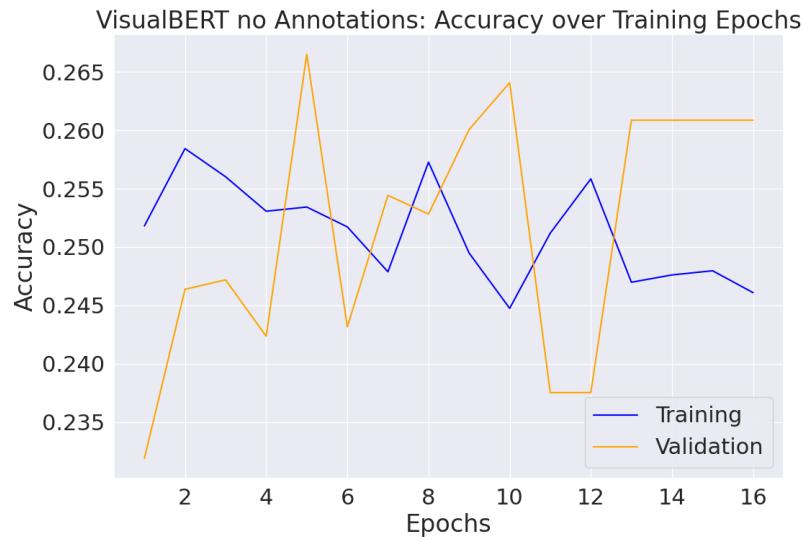


Figure 20: Training curves - Input configuration 1.

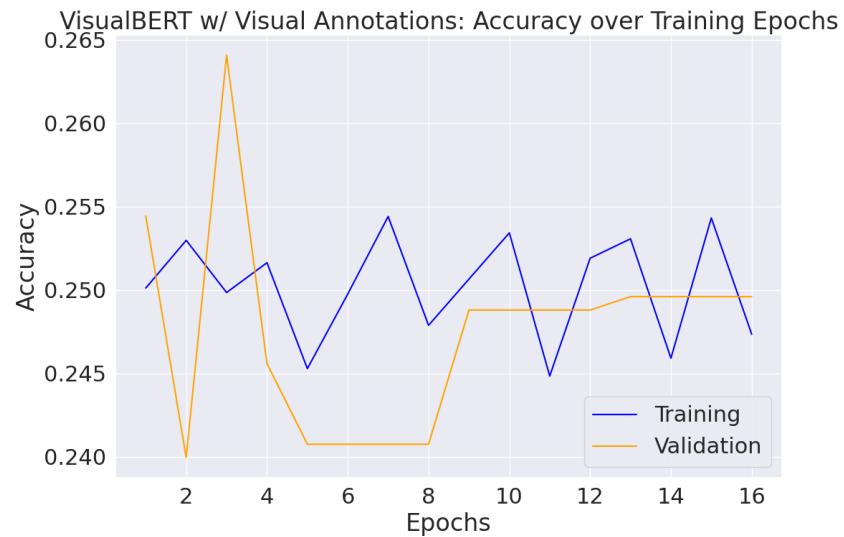


Figure 21: Training curves - Input configuration 2.

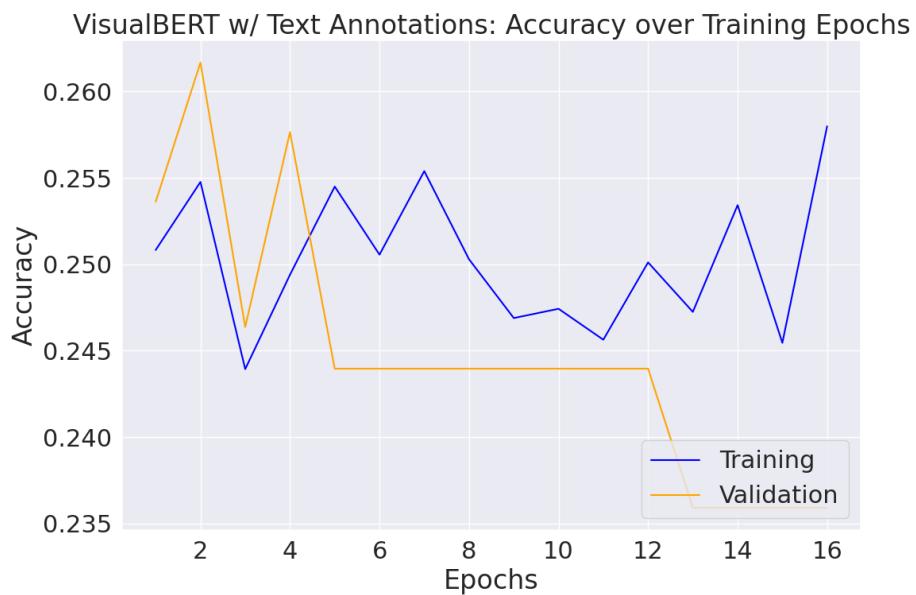


Figure 22: Training curves - Input configuration 3.

6. Results

Rather than draw After training completion, our team tested each of the three models on the same testing dataset that the original paper used. Our models' testing accuracies can be found in the figure below. The y-axis represents the accuracy while each bar represents a model. The yellow bars are our VisualBERT trained models representing each of the three different input configurations while the purple represents the benchmark models from the original paper.

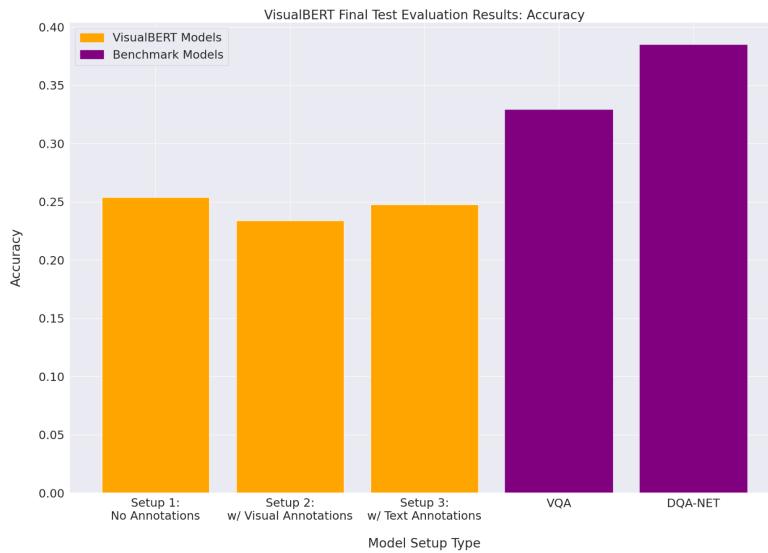


Figure 23: Final testing accuracies compared against benchmark models.

As shown, all three input scenarios had results nearing the random output of 25% and under the benchmark models of VQA and DQA-NET. It suffices to say that the novel experimentation setup that we tried did not result in good results. The model was not able to learn the data as we had hoped.

Diving further into the class specific results, we have the confusion matrix for each of the input configuration scenarios in the figures below. As displayed, the model seems to want to classify the answer A most often. Our team dived into why that may have been the case but could not find a specific reason beyond the data input was too complicated for the model to successfully learn.

Confusion Matrix: Setup 1

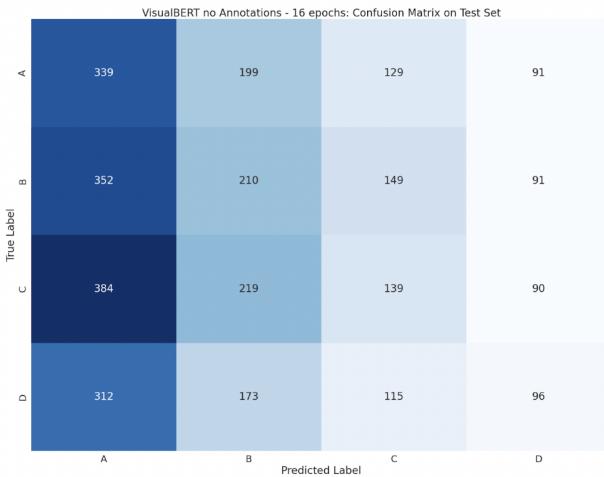


Figure 24: Confusion matrix for input setup 1.

Confusion Matrix: Setup 2

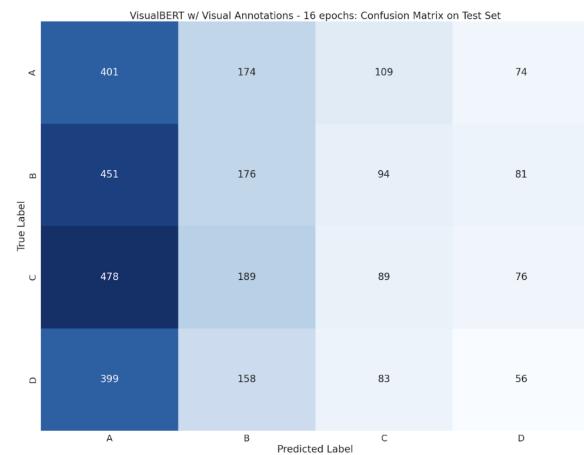


Figure 25: Confusion matrix for input setup 2.

Confusion Matrix: Setup 3

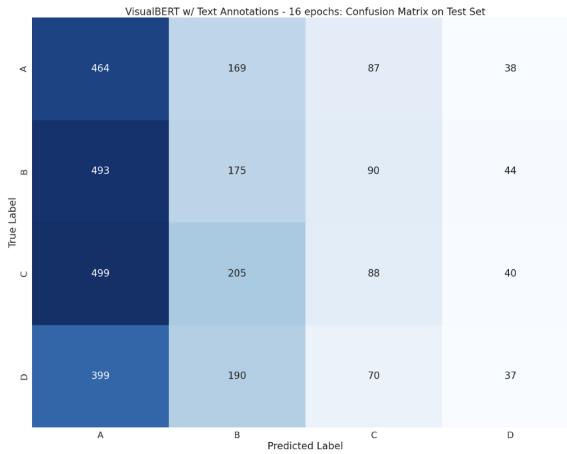


Figure 26: Confusion matrix for input setup 3.

The complexity of the dataset can be broken down into two areas. The first being that diagrams are not like regular natural images where arrows, relationships between objects, and text on the images are crucial to answering the respective questions. The second area being that the annotations are difficult to process into an input format that was easily usable and bridge the signal gaps between the image and text. Our group tried several methods to encode the annotations but these methods were not particularly successful. There were lots of signals between the different data modes but it was difficult to ensure the model can appropriately learn a good joint representation of all those inputs.

7. Conclusions

7.1 Lessons Learned

Visual-Diagram reasoning is at the frontier of both computer vision and language tasks. Within this space, our team tried something novel based on the research performed. As is the case for trying something novel, things don't always work out as planned. We hope that the work described in this report will inform others as to what to try next in solving this particular problem task. Our team learned and proved that large pretrained transformer models may not always perform better than more data-customized architectures, at least for diagram question-answering.

The crux of this task is working with large and complex input spaces where the really important signals can be difficult to prioritize, such as structuring the input for the model in a configuration where recognizing objects within a diagram and understanding that an arrow between 2 objects signifies a type of relationship is achievable.

As far as we know nobody has built a model that had an accuracy rating above 40%. Considering the baseline is 25%, there is clearly a lot of work to be done on these multimodal-knowledge graph tasks. Working on problems at the limit of current technologies is a new experience and challenging.

7.2 Future Improvements

If we had more time there are a few avenues we could pursue.

1. Move away from exploiting transfer learning and build a custom model which has modern characteristics (transformer layers) but has a more flexible embedding than the “off-the-shelf” models we used (VisualBert).
2. Devote far more time to build a really sophisticated, stand-alone visual embedding model which can supply embeddings to a second model which combines the visual embeddings with the textual embeddings.
3. Try using the HuggingFace “Processors” which generate visual and text embeddings from the same model, enabling us to use a pretrained model that has been designed for Visual Question and Answering and other multimodal tasks.
4. Filter the questions and/or images and focus on specific subsets. Both have more variation than we had previously thought, and it may be the case that we need to build different models for different subsets of questions and images because of this variation.
5. Try chaining different models together to solve the problem, using pre-trained models at all steps.. For example, create a model which generates a summary of an image, (Visual Summarization), then perform a simple question and answering task on the summary.

Additionally, if we have enough time, we hope to explore further image processing techniques to help present only the most useful signals for our model. The current visual embedding technique is rudimentary. The difficulty of enhancing the visual embeddings is that the answers

don't necessarily have a bounding box on this image where the answer can be identified, so improved identification of objects in the image may not necessarily improve overall accuracy.

8. Appendix

8.1 References

1. [Github Repo](#)
2. [AI2 Diagram Dataset \(AI2D\)](#)
AI2 Diagram Dataset (AI2D) was accessed on 9/5/2022 from <https://registry.opendata.aws/allenai-diagrams>.
3. [A Diagram is Worth a Thousand Words](#)
@article{Kembhavi2016ADI,
title={A Diagram is Worth a Dozen Images},
author={Aniruddha Kembhavi and Michael Salvato and Eric Kolve and Minjoon Seo and
Hannaneh Hajishirzi and Ali Farhadi},
journal={ArXiv},
year={2016},
volume={abs/1603.07396}}
4. Paper code:
5. [VISUALBERT: A SIMPLE AND PERFORMANT BASELINE FOR VISION AND LANGUAGE](#), Li et al., 2019
6. [VILT: Vision-and-Language Transformer Without Convolution or Region Supervision](#). Kim et al., 2021
7. [LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking](#). Huang et al., 2022
8. [VisualBERT for Multiple Choice Hugging Face](#)
9. [VisualBERT for Question Answering Hugging Face](#)
10. [VILT for Question Answering Hugging Face](#)
11. [LayoutMV3 Hugging Face](#)
12. [VisualBERT Demo](#)
13. [BERT Multiple Choice Sample](#)
14. [Fine Tuning on Multiple Choice Task](#)
15. [Hugging Face](#)
16. [PyTorch AdamW Optimizer](#)
17. [PyTorch Cross Entropy](#)