

Improved RAMEN: Towards Domain Generalization for Visual Question Answering

Bhanuka Manesha Samarasekara Vitharana Gamage
School of Information Technology
Monash University
Bandar Sunway, Malaysia
bsam0002@student.monash.edu

Lim Chern Hong
School of Information Technology
Monash University
Bandar Sunway, Malaysia
lim.chernhong@monash.edu

Abstract—Currently nearing human-level performance, Visual Question Answering (VQA) is an emerging area in artificial intelligence. Established as a multi-disciplinary field in machine learning, both computer vision and natural language processing communities are working together to achieve state-of-the-art (SOTA) performance. However, there is a gap between the SOTA results and real world applications. This is due to the lack of model generalisation. The RAMEN model [1] aimed to achieve domain generalization by obtaining the highest score across two main types of VQA datasets. This study provides two major improvements to the early/late fusion module and aggregation module of the RAMEN architecture, with the objective of further strengthening domain generalization. Vector operations based fusion strategies are introduced for the fusion module and the transformer architecture is introduced for the aggregation module. Improvements of up to five VQA datasets from the experiments conducted are evident. Following the results, this study analyses the effects of both the improvements on the domain generalization problem. The code is available on GitHub through the following link <https://github.com/bhanukaManesha/ramen>.

Index Terms—visual question answering, computer vision, natural language processing, attention, generalisation, RAMEN, early fusion, late fusion, transformer

I. INTRODUCTION

Visual Question Answering (VQA) is a multi-disciplinary problem in machine learning that exists at the intersection of the computer vision, natural language processing and knowledge representation fields [2]. Recently, the task of VQA has been classified as an AI-complete task due to the complexity of it. This problem requires the semantic understanding of each of the three fields as well as the relationship between each one of them [3]. One of the main issues in VQA is that the state-of-the-art (SOTA) results on the datasets do not translate on to real-world applications. This has directed the VQA field towards generalization.

The datasets in the field of VQA can be separated into two main categories [1]. The first type focuses on answering questions by understanding the objects on natural real world images and the other focuses on using synthetic images to test reasoning questions. The problem with this categorization is that the algorithms tend to focus on one or the other and not generalize on both. This known as the domain generalization problem, because the VQA models generalize on both types of dataset either through training from scratch or fine tuning to

the domains and not overfitting on one type. [1] addressed this issue by introducing a framework for domain generalization. This framework allows to train models of both domains with similar visual and textual features to evaluate their generalization ability.

They also introduced the RAMEN model architecture which was able to outperform all the other models compared in the study in terms of domain generalization. However, this model uses a simple architecture with a potential for improvement and exploration. Therefore, this study proposes improvements to the architecture of the RAMEN model while analyzing the effect of these changes to the overall problem of domain generalization.

The main contributions of this study includes the following:

- Improvements to the domain generalization performance of the RAMEN model architecture by proposing modifications to the fusion and aggregation modules.
- A broad comparison of the vector based fusion operations for early and late fusion pertaining to domain generalization.
- Implementation and analysis of a transformer based aggregation module to map the relationships between bi-modal embeddings of the regional proposals in the RAMEN model.

The rest of the paper is organized as follows: Section II provides more context to the domain generalization problem in VQA, the RAMEN model and the transformer architecture in VQA. The proposed improvements to the RAMEN model is detailed in Section III, which is followed by the experiment strategy used in Section IV. A comprehensive analysis of the results is conducted in Section V which is then summarized in Section VI.

II. RELATED WORK

This section first summarizes the main VQA datasets used in this study, followed by the RAMEN model. Next, the background of the transformer architecture in VQA is explored.

A. VQA Datasets

The dataset is the most important part of the VQA pipeline as it determines what the model learns. If the dataset contains inherent biases, the model will learn these and the performance

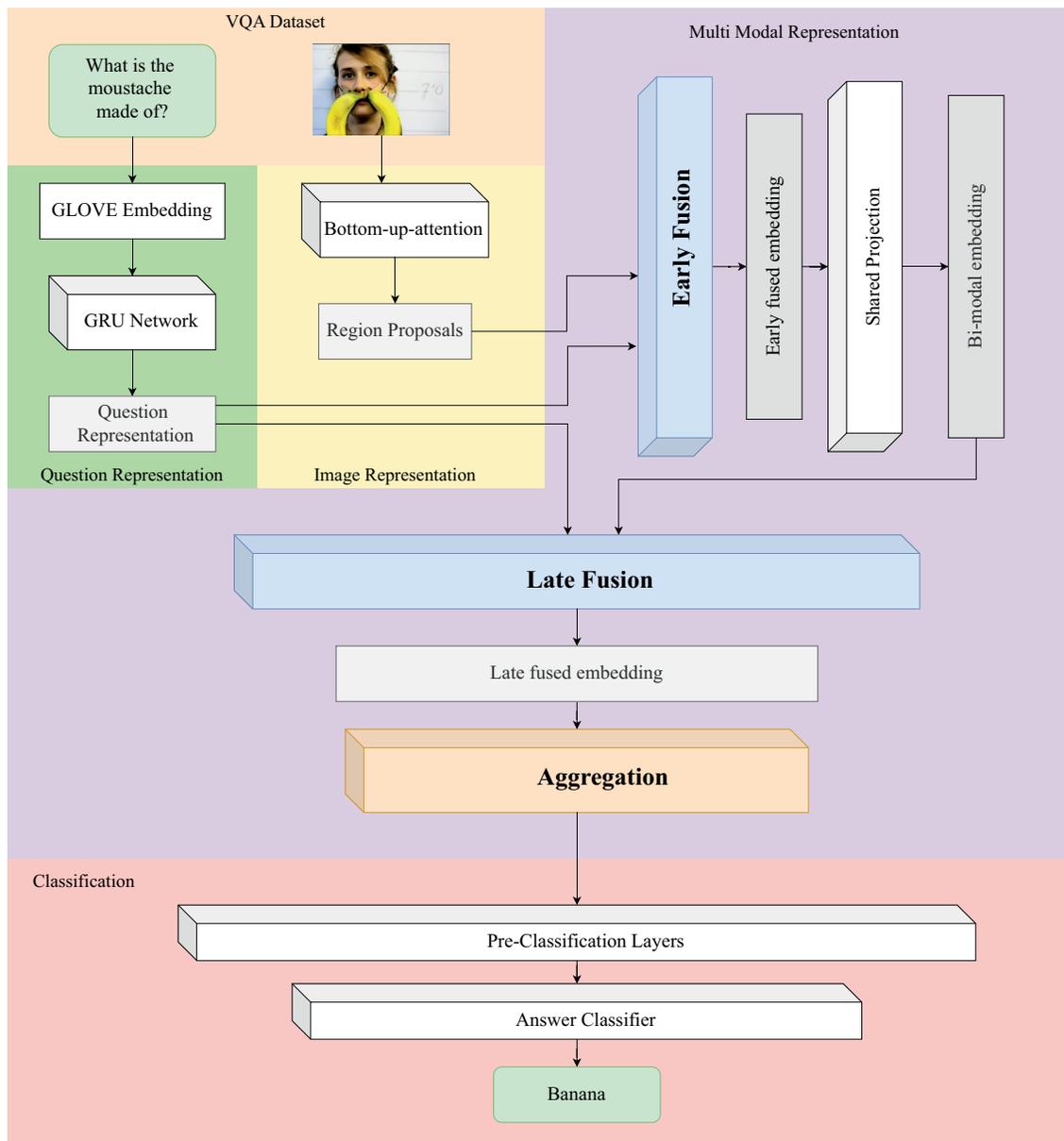


Fig. 1: High-level architecture of the RAMEN model

of the model will be affected. Many datasets were introduced, with each dataset focusing on either solving biases [6]–[8], [10], [11] or solving a specific type of question domain [12]–[15].

This study focuses on the datasets used by [1], where the datasets were divided into two groups; VQA datasets for natural image understanding and VQA datasets to test reasoning. Table I summarizes the VQA datasets used by [1] to test for generalization.

1) VQA Datasets for Natural Image Understanding:

These datasets aims to provide answers by identifying objects in the image. This can be through colour, count or other visual cues. All the datasets in this group uses the MSCOCO dataset [16] as the base image dataset except for TDIUC which adds

extra images.

a) *VQAv1* [2]: One of the most widely known datasets with the current SOTA accuracy of 75.26%. This dataset mainly focuses on detection questions such as *Is there food on the table* [2] and lacks reasoning questions such as *What is behind the computer in the corner of the table?* [10]. However, it consists of inherent question-answer biases where questions such as *Where is the giraffe standing on?* always contains the answer *grass*.

b) *VQAv2* [4]: The successor to *VQAv1*, was able to reduce the question-answer biases by introducing complementary questions. However, even though this allowed the *VQAv2* dataset to be more balanced, the bias of having more detection questions is still prevalent in this dataset; which makes the

TABLE I: Summary of VQA datasets used in this study. [1]

Dataset	Image Type	Question Type	Images	Q&A Pairs	Download links
VQAv1 [2]	Natural	Human	204K	614K	https://visualqa.org/vqa_v1_download.html
VQAv2 [4]	Natural	Human	204K	1.1M	https://visualqa.org/download.html
TDIUC [5]	Natural	Both	167K	1.6M	https://kushalkafle.com/projects/tdiuc.html
C-VQA [6]	Natural	Human	123K	369K	https://computing.ece.vt.edu/~aish/cvqa/
VQACpv2 [7]	Natural	Human	219K	603K	https://computing.ece.vt.edu/~aish/vqacp/
CLEVR [8]	Synthetic	Synthetic	100K	999K	https://cs.stanford.edu/people/jcjohns/clevr/
CLEVR-Humans [9]	Synthetic	Human	32K	32K	https://cs.stanford.edu/people/jcjohns/iep/
CLEVR-CoGenT-A [8]	Synthetic	Synthetic	100K	999K	https://cs.stanford.edu/people/jcjohns/clevr/
CLEVR-CoGenT-B [8]	Synthetic	Synthetic	30K	299K	https://cs.stanford.edu/people/jcjohns/clevr/

models trained on VQAv2 datasets inherently weaker when answering questions with reasoning.

c) TDIUC [5]: This dataset was created with the primary aim of evaluating the performance of models on 12 distinct types of VQA tasks. Color attributes, positional reasoning and object presence are some of the types of tasks. A new metric called Mean-per-type was also introduced as shown in Equation 7 in Section 7. Therefore, it is evident that a model needs to perform well across all the question types to get a good performance score.

d) C-VQA [6]: Aims to re-split the VQAv1 dataset to introduce novel combinations for the question-answer pairs when testing. During testing the models will come across new combinations of question-answer pairs. Therefore, the models need to be able to generalize on the task and not the question and answer.

e) VQACpv2 [7]: Overcomes the question and language bias by splitting the VQAv1 and VQAv2 dataset. A completely different answer distribution is present in the test split compared to the training split. This allows the models to test their ability to generalize by not over-fitting on the training set.

2) VQA Datasets to Test Reasoning: These datasets aim to test the ability of models to answer reasoning based questions by using synthetic images. These synthetic computer generated images allow this dataset to generate complex reasoning questions automatically. All the datasets in this group use the images from the CLEVR dataset, with each dataset having different question-answer pairs.

a) CLEVR [8]: The main goal of this dataset is to test the reasoning capability of models on geometric shapes. Similar to TDIUC, this dataset is classified into five categories.

b) CLEVR-Humans [9]: The main downside of CLEVR dataset is that the questions are computer generated, thus being very structured. The CLEVR-Humans dataset addresses this issue by using free form human generated question-answer pairs. It still uses the same images from the CLEVR dataset.

c) CLEVR-CoGenT [8]: This dataset was introduced with the CLEVR dataset having two splits with mutually exclusive color and shapes, namely, CLEVR-CoGenTA and CLEVR-CoGenTB. This dataset aims to study the model’s ability to recognize novel combinations of attributes such as color and shapes at test time. For example, CLEVR-CoGenTA

contains red colour cylinders in the training set, in contrast, CLEVR-CoGenTB does not contain red colour cylinders.

B. RAMEN

The VQA pipeline consists of five main components; VQA dataset, Image representation, Question representation, Multi-modal representation and Answer classification. Many studies have been done in the field of VQA, with each focusing on improving different sections of the VQA pipeline [2], [17]–[28].

[1] proposed a framework to compare the performance of VQA algorithms across different domains. They standardize the image representation and the question representation across the VQA datasets of multiple domains. With this, they were able to compare the performance of multiple algorithms [12], [24], [26], [29], [30] across domains and assess the generalization ability of the model architectures.

They also proposed a model named RAMEN with a conceptually simple architecture that was able to generalize across multiple domains. Figure 1 shows the high-level architecture of the RAMEN model, where the five main components of the VQA pipeline can be identified in this model.

a) Image Representation: The image representation module focuses on extracting features from the image and converting them into visual features. Various methods exist that focus on extracting features from images using different techniques such as VGG-Net [2], [18], [31]–[33], ResNet [19], [22], [23], [34], [35] and Faster-RCNN [24], [36]–[38]. The RAMEN model uses a Faster-RCNN based technique where the image is passed through the bottom-up-top-down network [24], which uses attention on the object level to return visual features as a set of regions. These regions correspond to the main object regions in the image which are used to answer the questions. For VQAv1, VQAv2, CVQA, VQACpv2 and TDIUC datasets the bottom-up attention module returns 36 regions and for the CLEVR family of datasets it returns 15 regions.

b) Question Representation: The question representation module converts the question into a vector representation. This vector representation encodes all the words while maintaining the flow and positional information of the question. Studies have proposed multiple ways to extract these features by using CNN [21], [39], LSTM [19], [23], [31] and GRU [24],

[34], [37] networks. The RAMEN model uses a GRU based approach by first splitting the questions into multiple word tokens. Each token is instantiated with the GLOVE embedding [40]. Then the embeddings are passed through a GRU based RNN [41] to obtain the question representation.

c) Multi Modal representation: Once the image and question representations are passed into this module, the two vectors are fused together using concatenation. This step is also known as early fusion in the model architecture. Next, the RAMEN model uses a Multi Layer Perceptron (MLP) to create a bimodal embedding. This allows the model to learn the relationship between the image and question representation. Then the bimodal embedding is concatenated with the question representation, with the Late Fusion module. The fused vector is then passed to the aggregation module where it is passed through a bi-directional GRU network. This step captures the relationships between the bimodal embeddings.

d) Classification: In the final module, the output of the multi-modal representation is pass through a series of linear layers to perform the pre-classification step. This is then followed by a single linear layer for classification.

C. Transformer

The introduction of the transformer architecture [42] has been a pivotal moment in the NLP community. The main use case of the transformer network is for machine translation. The main advantage of using transformer over traditional RNN networks is that the sequences are processed as a whole compared to one by one. Moreover, the transformer uses multi-head attention and positional encoding to obtain more information about the relationships between the features. This allows the transformer architecture to be parallelizable compared to sequential RNN networks.

Many studies have been done in the field of VQA that incorporated transformers into the architecture [43]–[45]. Most of them focus on encoding the question using the transformer architecture. The Bidirectional Encoder Representations from Transformers (BERT) architecture [46], which is derived from the transformer architecture is commonly used to encode the question [47], [48]. However, limited research has been done where the transformer architecture is used to capture the relationship between visual and question features. [47] showed that the transformer architecture is able to capture intra-modality and cross-modality relationships on the VQA and GQA [49] datasets. Therefore, this study aims to investigate the effect of using transformer as an encoder in VQA.

III. METHODOLOGY

The key focus of the improvements are on the multi modal representation section of the RAMEN model. Based on the ablation study done by [1], the early fusion module has a significant effect on the performance of the model. The aggregation technique also has an effect on the performance of the model, whereas, the late fusion module has the minimum effect on the performance of the model. Therefore, the experiments are done on the Early Fusion, Late Fusion and the Aggregation modules of the model as shown in Figure 1.

A. Fusion Strategies

[3] performed a survey on the fusion strategies in ImageVQA and VideoVQA studies. They classified the fusion strategies into three main types; Vector operations, Neural Networks (NN) and Bilinear pooling.

The RAMEN model uses a mix of vector operations and neural networks to perform the multi-modal fusion. The early and late sub-modules uses simple concatenation of the features and the shared projection and aggregation sub-modules uses neural networks as the strategy. First in the early fusion module, the regional visual features are fused using concatenation with the question embedding to obtain the early fused embedding. This is then passed through the neural network based shared projection and the output bi-modal embedding is obtained. In the late fusion module, the bi-modal embedding is again fused using concatenation with the question embedding to obtain the late fused embedding. After the fusion operation, both early and late fusion embeddings are passed through a Batch Normalization [50]. Finally, the vector is passed through the aggregation module, which is a Recurrent Neural Network based fusion strategy to obtain the fused vector for classification.

In the survey, [3] also categorized the vector operations into three main sections; concatenation, addition and multiplication. In this study, these three main vector operations will be experimented on the RAMEN model to observe the performance effect on the NN. Figure 2 shows the overview of the fusion strategies tested in this study.

1) Concat Fusion: This is the baseline strategy used by the RAMEN model. A question embedding size of 1024 and visual feature size of 2048 is used to obtain a final embedding size of 3072. In this approach the output embedding passes all the information from both embeddings to the NN to identify the relationships. No information is lost in this approach and all the feature points are given similar weights. In order to perform vector operations, the question embedding is repeated to match the size of the visual features and bi-modal embedding. To do this, the question embedding is repeated 36 times for the VQA family of datasets and 15 times for the CLEVR family of datasets. This is done for all the fusion strategies experimented in this study.

Equation 1 is used to obtain the final embedding (c_i), where q_i is the question embedding and v_i is the regional visual features or the bi-modal embedding.

$$c_i = \text{BatchNorm}([q_i, v_i]) \quad (1)$$

2) Additive Fusion: For the additive fusion, the question embedding is matched to the same size as the visual features. Therefore, the embedding size is changed from 1024 to 2048 to obtain the final embedding size of 2048. This approach emphasizes on the different feature points which allows the model to update the question embeddings to focus on them. This approach has information loss due to the addition operation, however, it is compensated by the increase in the

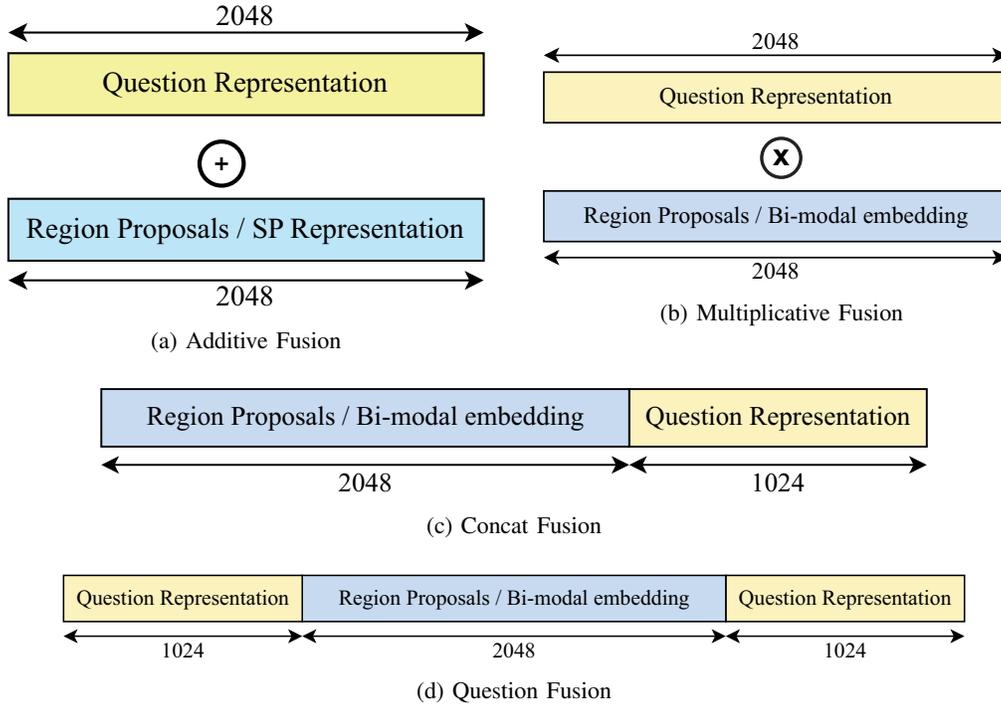


Fig. 2: Fusion Strategies for Early and Late fusion

question embedding size. Equation 2 is used to obtain the final embedding (c_i). q_i and v_i are as described as in Equation 1.

$$c_i = \text{BatchNorm}(v_i \oplus q_i) \quad (2)$$

3) **Multiplicative Fusion:** Similar to additive, the question embedding size of 2048 is used for this strategy. In this approach, the emphasis on different feature points is greater than the additive fusion. Equation 3 is used to obtain the final embedding (c_i). q_i and v_i are as described as in Equation 1.

$$c_i = \text{BatchNorm}(q_i \odot v_i) \quad (3)$$

4) **Question Fusion:** Question fusion uses a dual concatenation strategy. The question embedding is concatenated before and after the visual features. An embedding size of 1024 is used for the question embedding, with the final embedding size of 4096. The main emphasis in this strategy is to provide more feature points for the question embedding. Datasets such as the CLEVR family that are used to test reasoning contains longer questions compared to other datasets. Therefore, limiting the question embedding to a single vector of size 1024 or 2048 can effect the emphasis of the question on the model. Equation 4 is used to obtain the final embedding (c_i). q_i and v_i are as described as in Equation 1.

$$c_i = \text{BatchNorm}([q_i, v_i, q_i]) \quad (4)$$

B. Aggregation Strategies

This module is used to calculate the relationship between the question and bi-modal embeddings. The bi-modal embeddings

contains the relationship between the question and each regional visual feature. Thereby this module aims to identify the relationships between the visual regions. Higher performance on this module will lead to better results on questions that require multi object or localized information to answer.

1) **bi-GRU network:** The baseline aggregation strategy in the RAMEN model uses a bidirectional GRU based RNN to calculate the feature vector. The main downside to this approach is that the model goes through each region sequentially in both directions. Therefore, to obtain a relationship between two regions of the image, the model needs to pass through the other regions which can lead to information loss. This approach is best used when all the regions are equally important for the question.

2) **Transformer network:** The transformer architecture is stronger in identifying the relationship between the multiple regions/vectors, because the network processes all the regions at once and not sequentially. This is the reason why the transformer model performs well on the machine translation tasks [36]. This allows it to capture relationships among regions better than RNNs.

However, the positional encoder in the traditional transformer network masks out half of the regions. This is to ensure the model is not able to see the next word in machine translation. For the RAMEN model, this is not an issue. So the mask is removed and the transformer is able to view all the regions.

The output of the original transformer model is a set of decoders for the translated sentence. But in this case, the main aim is to obtain a representation to be passed to the

classification module. Therefore, the decoder is replaced with a fully connected NN that returns a vector representation instead of the transformer decoder module.

One of the main downside of the transformer network is the slow convergence. Typically the transformer network might take upto sixty hours to fully convergence on translation tasks.

IV. EXPERIMENT

A. Dataset specification

In the baseline paper, the accuracy of the CLEVR-CoGenTB dataset was obtained on a sub split of the test set. But in this study, the accuracy is obtained on the complete test set. Similarly, the original paper fine-tuned the model trained on the CLEVR-dataset with the CLEVR-Humans dataset to obtain the accuracy. However, this study, trains the CLEVR-Humans dataset from scratch. All other training and testing splits of the datasets are identical to the baseline paper.

B. Model specification

Due to the changes in the datasets, the baseline accuracies are all re-calculated to ensure consistency. All the model hyper-parameters are maintained as mentioned in the baseline paper.

The model with the transformer as the aggregation strategy is named as the TransformerNet and for baseline model with the bi-GRU network the name RAMEN model is used. With each model, the four different fusion strategies are experimented for both the early and late fusion modules. Therefore, in total the nine datasets are trained on eight versions of the models.

C. Evaluation metrics

Three types of evaluations metrics are used in this study to compare the results between the datasets. These are the same metrics used in the baseline study.

a) *10-choose-3*: Equation 5 shows the evaluation metric used by VQAv1, VQAv2, CVQA and VQACpv2. These datasets provide multiple answers for each question from multiple human annotators. Thus using this metric reduces the inter-human variability [2].

$$Acc(answer) = \min\left\{\frac{\# \text{ of annotators provided answer}}{3}, 1\right\} \quad (5)$$

b) *Simple Accuracy*: CLEVR, CLEVR-Humans, CLEVR-CoGenT-A and CLEVR-CoGenT-B uses the simple accuracy shown in Equation 6 as the evaluation metric [8].

$$Acc(answer) = \frac{\# \text{ correct answer}}{\# \text{ questions}} \quad (6)$$

c) *Mean-per-type*: The TDIUC dataset uses the mean-per-type evaluation metric as shown in Equation 7 [5]. This ensures that the model is able to perform well on each category, even though the number of test instances of each category are different.

$$Acc(answer) = \frac{\sum \left\{ \frac{\# \text{ correct answer per type}}{\# \text{ of questions per type}} \right\}}{\# \text{ of types}} \quad (7)$$

D. Training specifications

All the experiments were done on a PC running *Ubuntu 18.04.1 LTS* with an *Intel® Xeon(R) W-2145 CPU @ 3.70GHz* with 16 logical cores and 64GB RAM. A single *Quadro P5000* GPU was used to perform the NN training with a *7200RPM Seagate* hard drive to store the data. The gradual learning rate warm up is used similar to [1], [26]. The mini-batch size of 256 is used for all the experiments. The models are trained until 25 epochs with some exceptions in the TransformerNet experiments; mainly due to the slower convergence rate. As shown in Appendix A, an average training time of 46 minutes per epoch was elapsed for all experiments.

V. RESULTS & DISCUSSION

This section focuses on the key observations from the experiments and the effect of changing the aggregation and fusion strategies. Table II show the scores obtained by each experiment. The top three models are highlighted with darker colours indicating better performance. Appendix B contains the full results table with the training scores and number of epochs.

When comparing the baseline results (*Ramen-Concat*) with the results from the RAMEN study [1], it is evident that there exist a minor difference in the scores. The training was done with the same hyper-parameters as stated in the baseline paper even though there exists a deficiency in the scores. However, since all the experiments are done with the fixed set of hyper-parameters, the scores in this study are consistent.

A. Overall observations

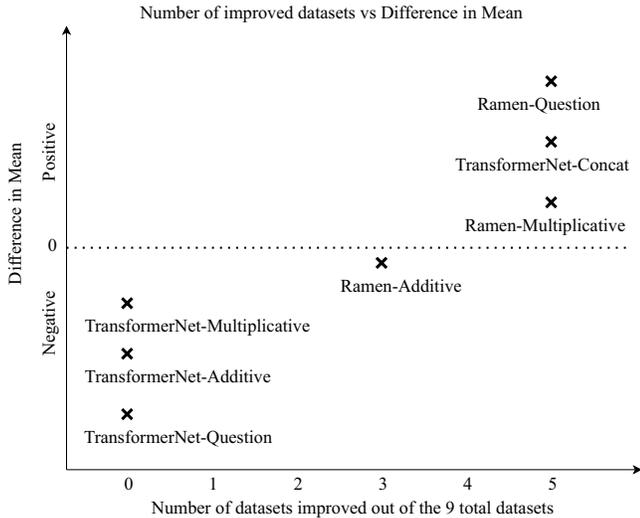
First, considering the model with the highest mean score across all the datasets, the *Ramen-Question* model has a score of 68.76. With a percentage difference of about 1%, it is evident that the improvement of using different fusion and aggregation strategies is minor. However, many other patterns in the results can be observed which can help improve the performance of future models. It is also noted that *Ramen-Multiplicative* and *TransformerNet-Concat* was also able to improve the performance by about 0.5% and 0.65% respectively.

Next, it is evident that TransformerNet model did not perform at all on most of the datasets, with differences in accuracies of more than 25% on the CLEVR datasets. This issue is addressed in Section V-D. However, the *TransformerNet-Concat* model performed well on most of the datasets.

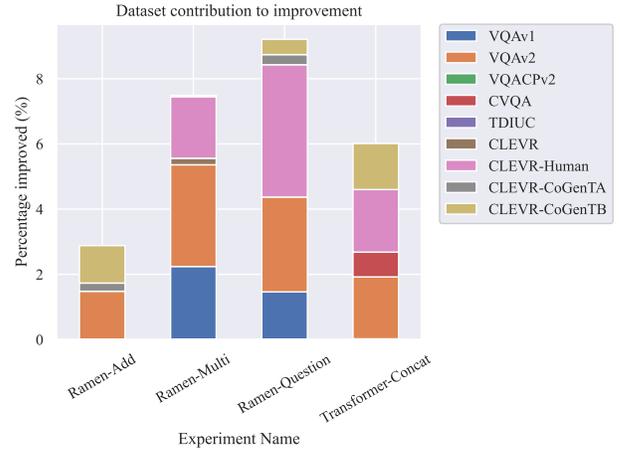
When considering the model with the highest number of top scores on the nine datasets, the *Ramen-Multiplicative* model has achieved the highest score on three of the main datasets. Therefore, it is evident that this model can perform well on both natural and synthetic types of VQA datasets. However, the model is not able to generalize to question and attribute biases well. This is identified from the CVQA and VQACpv2 datasets, due to the lower performance on them. When comparing the performance of the model between CLEVR-CoGenTA and CLEVR-CoGenTB datasets, it is found to be evident that the model is seeing a dip in performance.

TABLE II: Results from all eight model with the nine VQA datasets.

Dataset	Ramen				TransformerNet			
	Baseline Concat [1]	Additive	Multiplicative	Question	Concat	Additive	Multiplicative	Question
VQAv1	63.30	63.21	65.54	64.76	63.32	55.88	60.91	55.08
VQAv2	62.16	63.64	65.28	65.07	64.06	59.79	56.32	50.14
VQACpv2	37.61	36.73	36.28	37.03	37.47	27.60	27.60	26.89
CVQA	56.98	55.82	56.58	56.81	57.74	53.49	54.20	48.47
TDIUC	66.48	64.90	65.69	64.81	65.47	58.03	56.34	53.90
CLEVR	96.52	96.26	96.72	96.06	95.79	50.52	57.31	50.07
CLEVR-Humans	44.57	40.21	46.46	48.63	46.49	38.45	40.07	37.99
CLEVR-CoGenTA	96.59	96.84	96.63	96.90	96.43	64.26	72.07	60.24
CLEVR-CoGenTB	88.27	89.42	86.22	88.74	89.68	55.80	60.18	55.19
Mean	68.05	67.45	68.38	68.76	68.49	51.53	53.89	48.66



(a) The difference in mean compared to the overall number of improved datasets.



(b) Improvement of the datasets for each model compared to the baseline.

Fig. 3: Summaries of the results.

Next, comparing the performance of the model based on the scores in the top-3 rank of the experiments, it is evident that *Ramen-Question* has the overall best performance. In seven out of the nine dataset, this model was able to achieve the top three results. This indicates the ability of the model to generalize across multiple datasets. Also the *TransformerNet-Concat* model was able to achieve top three in seven out of nine datasets even though it only have the highest scores for the CVQA dataset.

With reference to Figure 3a, the *Ramen-Additive* model was able to improve on three datasets. However, due to the lower scores on the other datasets, especially CLEVR-Humans, the model is not able to achieve a positive mean. This indicates that the model is not able to perform well on free-form questions.

In the following section the performance of the models on specific datasets are analysed.

B. Dataset observations

Figure 3b demonstrates the datasets impact on the improvement of the models. It only showcase the experiments with at least one improvement, hence *TransformerNet-Additive*, *TransformerNet-Multiplicative* and *TransformerNet-Question* are ignored. All of the remaining models were able to improve on the VQAv2 dataset. This suggests that all the fusion strategies and the transformer aggregation strategy can improve the performance of localizing and detecting items.

However, only the Multiplicative and Question fusion were able to improve the VQAv1 dataset. Since, VQAv1 has inherent bias in the questions, this implies that the these two fusion strategies are more prone to over-fitting on the VQA based datasets. This is also made clear by the poor performance on the CVQA and VQACpv2 datasets.

However, the most performance gain is by the *Ramen-Question* on the CLEVR-Humans dataset. It points to an

indication that the pre and post concatenation of the question embedding has an improvement on the free form questions. This is also true for the models performance on VQAv1 and VQAv2 datasets.

Compared to the baseline, the *TransformerNet-Concat* is the only model to have an improvement in the score for CVQA. This highlights that the transformer aggregation module provides the ability for the module to generalize. This is further emphasized since the CLEVR-CoGenTB and VQAv2 datasets also show improvement in the scores.

C. Fusion Strategies

Overall it is clear that the different fusion strategies favor various datasets due to the unique characteristics in them.

1) *Concat Fusion*: The baseline concatenation fusion approach was able to get the highest score for VQACPv2 and TDIUC dataset using the RAMEN model and CLEVR-CoGenTB using the TransformerNet model. The VQACPv2 dataset aims to test the answer bias in the models. Therefore, concatenation based fusion is able to generalize well in terms of answer biases as both *Ramen-concat* and *TransformerNet-Concat* was able to achieve high scores.

Next, the TDIUC MPT metric measures the performance of the model on multiple question types. Considering that both the *Ramen-Concat* and *TransformerNet-Concat* have high scores on the TDIUC dataset, it is clear that concatenation based fusion allows for much more question type based generalization.

The *TransformerNet-Concat* model is trained on CLEVR-CoGenTA and tested on CLEVR-CoGenTB. Therefore, the model will not learn any details about the complementary attributes in the dataset. This indicates that the model is able to generalization well onto unseen combinations of attributes. However, the relationship between concatenation based fusion and attribute based generalization cannot be established. This is due to the lower score in the *Ramen-Concat* model, which implies that the performance gain is due to the transformer aggregation strategy.

2) *Additive Fusion*: The additive fusion strategy was not able to get the highest score for any of the datasets. However, the *Ramen-Additive* model was able to improve on the CLEVR-CoGenTA and CLEVR-CoGenB datasets, which point towards the model's ability to generalize to new concept compositions. The main issue with the additive fusion strategy is the information loss and the lower emphasis on the vector operation.

3) *Multiplicative Fusion*: The multiplicative fusion strategy achieved the highest score for VQAv1, VQAv2 and CLEVR datasets. As mentioned in Section V-A, the model suffers with generalization. However, the emphasis on the vector operation is higher compared to additive fusion, therefore the most important details are passed on through the fusion module.

4) *Question Fusion*: The question fusion obtained the highest score for CLEVR-Humans and CLEVR-CoGenTA. It also achieved high scores for all the CLEVR datasets. This is an indication that for reasoning type datasets with higher

significance on the question, the double concatenation of the question has an effect.

The ability for the model to generalize on new concept compositions can be observed due to the performance on CLEVR-CoGenTA and CLEVR-CoGenTB.

D. Aggregation Strategies

The transformer module as the aggregation strategy does not perform well. With only decent performance using the concat fusion strategy, it may not be suitable to be part of the RAMEN model. Many issues were faced when training the transformer module such as slow convergence and longer training time. However, the slow convergence of the transformer module remains a significant drawback.

This is evident when considering the number of epochs used to train *TransformerNet-Concat* vs *TransformerNet-Question* on the VQAv2 dataset. *TransformerNet-Concat* was trained for 50 epochs where the highest score was obtained at epoch 46, whereas *TransformerNet-additive* was only trained for 25 epochs. Appendix B reports all the training details used for each dataset. Therefore, training the TransformerNet models for longer can provide better scores.

Additionally, due to the longer training times and the time constraint, hyper-parameter tuning was not an option. With an average time of 58 minutes per epoch for the *TransformerNet-Question* model, the training would take more than 48 hours on a single GPU. However, as observed by *TransformerNet-Concat*, in an ideal scenario the model is able to convergence.

VI. CONCLUSION

The proposed improvements of this study resulted in minor gains in performance of about 1%. However, in-terms of domain generalization the *Ramen-Multiplicative*, *Ramen-Question* and *TransformerNetwork-Concat* models were able to achieve improvements in five out of the nine datasets. Also, *Ramen-Question* and *TransformerNetwork-Concat* models were able to achieve the top three scores in seven out of the nine datasets.

Analyzing the fusion strategies, provided insights to the different characteristics that effect domain generalization. For example, Question fusion performed well on reasoning questions due to the increase in the number of question embedding data points, which resulted in an increase in the amount of information passed into the aggregation module. This knowledge can be used to improve the performance of model and domain generalization.

When studying the effects of the Transformer module as the aggregation strategy, it is clear that selecting the correct hyper-parameters and providing the necessary amount of training time to converge are two main requirements when training the transformer module. This is one of the main limitations of this study. The time constraint and high computational cost of the experiments led to some of the experiments not converging to the higher scores. Therefore, since VQA datasets tend to be larger in size, more powerful hardware is needed

to perform a broader hyper-parameter search to optimize the models performance.

Focusing only on the vector operation based fusion strategies was also another limitation of the study. Bilinear Pooling techniques for fusion has proven to be effective when it comes to identifying relationships. However, the added computational cost of pooling on top of the transformer module will result in poor performance if training is not done till convergence. Nonetheless, this is a probable path that can be explored in the future.

Another limitation of the study is that the RAMEN architecture itself may be causing the bottle neck when trying to improve the generalization. Since both improvements were done to sub-modules of the multi-modal section of the RAMEN model, inherent limitations may exist in the architecture. Using the knowledge gained from the analysis of the fusion and aggregation module, a new architecture can be developed to well suite domain generalization.

In conclusion, this study paved the path to understanding the characteristics required for domain generalization as well as improving the performance of the RAMEN model.

VII. ACKNOWLEDGMENT

A special thanks to Robik Shrestha from the Rochester Institute of Technology who was one of the main researchers of the RAMEN study [1] for the guidance and support in providing the links to the datasets and setting up the baseline for this study.

REFERENCES

- [1] R. Shrestha, K. Kafle, and C. Kanan, "Answer them all! toward universal visual question answering models," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10472–10481.
- [2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "VQA: Visual question answering," *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2015 Inter, pp. 2425–2433, 2015.
- [3] D. Zhang, R. Cao, and S. Wu, "Information fusion in visual question answering: A survey," *Information Fusion*, vol. 52, pp. 268–280, 2019.
- [4] Y. Goyal, T. Khot, A. Agrawal, D. Summers-Stay, D. Batra, and D. Parikh, "Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering," *International Journal of Computer Vision*, vol. 127, no. 4, pp. 398–414, 2019.
- [5] K. Kafle and C. Kanan, "Visual question answering: Datasets, algorithms, and future challenges," *Computer Vision and Image Understanding*, vol. 163, pp. 3–20, 2017.
- [6] A. Agrawal, A. Kembhavi, D. Batra, and D. Parikh, "C-VQA: A Compositional Split of the Visual Question Answering (VQA) v1.0 Dataset," pp. 1–10, 2017. [Online]. Available: <http://arxiv.org/abs/1704.08243>
- [7] A. Agrawal, D. Batra, D. Parikh, and A. Kembhavi, "Don't just assume; look and answer: Overcoming priors for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4971–4980.
- [8] J. Johnson, L. Fei-Fei, B. Hariharan, C. L. Zitnick, L. Van Der Maaten, and R. Girshick, "CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning," *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, vol. 2017-Janua, pp. 1988–1997, 2017.
- [9] J. Johnson, J. Hoffman, L. Fei-fei, C. L. Zitnick, and R. Girshick, "Inferring and Executing Programs for Visual Reasoning," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2989–2998, 2017. [Online]. Available: http://openaccess.thecvf.com/content_ICCV_2017/papers/Johnson_Inferring_and_Executing_Programs_for_Visual_Reasoning_ICCV_2017_paper.pdf
- [10] M. Malinowski and M. Fritz, "A multi-world approach to question answering about real-world scenes based on uncertain input," in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 1682–1690.
- [11] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. Bernstein, and L. Fei-Fei, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," 2016. [Online]. Available: <https://arxiv.org/abs/1602.07332>
- [12] M. Malinowski, C. Doersch, A. Santoro, and P. Battaglia, "Learning visual question answering by bootstrapping hard attention," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11210 LNCS, pp. 3–20, 2018.
- [13] S. A. Hasan, Y. Ling, O. Farri, J. Liu, H. Müller, and M. Lungren, "Overview of ImageCLEF 2018 medical domain Visual Question Answering task," *CEUR Workshop Proceedings*, vol. 2125, 2018.
- [14] S. Shah, A. Mishra, N. Yadati, and P. P. Talukdar, "KVQA: Knowledge-Aware Visual Question Answering," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 8876–8884, 2019.
- [15] K. Kafle, B. Price, S. Cohen, and C. Kanan, "Supplementary Materials for DVQA : Understanding Data Visualizations via Question Answering," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5648–5656, 2018.
- [16] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, "Microsoft coco: Common objects in context," 2014.
- [17] H. Xu and K. Saenko, "Ask, Attend and Answer: Exploring Question-Guided Spatial Attention for Visual Question Answering," ser. Lecture Notes in Computer Science, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, vol. 9911, pp. 451–466. [Online]. Available: <http://link.springer.com/10.1007/978-3-319-46448-0http://link.springer.com/10.1007/978-3-319-46448-0http://link.springer.com/10.1007/978-3-319-46448-0>
- [18] J. Lu, J. Yang, D. Batra, and D. Parikh, "Hierarchical question-image co-attention for visual question answering," in *Advances in neural information processing systems*, 2016, pp. 289–297.
- [19] I. Ilievski, S. Yan, and J. Feng, "A Focused Dynamic Attention Model for Visual Question Answering," 2016. [Online]. Available: <http://arxiv.org/abs/1604.01485>
- [20] H. Noh and B. Han, "Training Recurrent Answering Units with Joint Loss Minimization for VQA," 2016. [Online]. Available: <http://arxiv.org/abs/1606.03647>
- [21] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, "Stacked attention networks for image question answering," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2016-Decem, no. 1, 2016, pp. 21–29.
- [22] H. Nam, J.-W. Ha, and J. Kim, "Dual Attention Networks for Multimodal Reasoning and Matching," pp. 299–307, nov 2016. [Online]. Available: <http://arxiv.org/abs/1611.00471>
- [23] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach, "Multimodal compact bilinear pooling for visual question answering and visual grounding," *arXiv preprint arXiv:1606.01847*, 2016.
- [24] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pp. 6077–6086, 2018.
- [25] Y. Jiang, V. Natarajan, X. Chen, M. Rohrbach, D. Batra, and D. Parikh, "Pythia v0.1: the Winning Entry to the VQA Challenge 2018," pp. 4–6, 2018. [Online]. Available: <http://arxiv.org/abs/1807.09956>
- [26] J. H. Kim, J. Jun, and B. T. Zhang, "Bilinear attention networks," *Advances in Neural Information Processing Systems*, vol. 2018-Decem, no. NeurIPS, pp. 1564–1574, 2018.
- [27] B. Liu, Z. Huang, Z. Zeng, Z. Chen, and J. Fu, "Learning Rich Image Region Representation for Visual Question Answering," pp. 14–16, 2019. [Online]. Available: <http://arxiv.org/abs/1910.13077>
- [28] Z. Yu, J. Yu, Y. Cui, D. Tao, and Q. Tian, "Deep modular co-attention networks for visual question answering," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR 2017*, pp. 6274–6283, 2017.

- [29] W. Norcliffe-brown and S. Parisot, "Learning Conditioned Graph Structures for Interpretable Visual Question Answering arXiv : 1806 . 07243v6 [cs . CV] 1 Nov 2018," no. Nips, 2018.
- [30] D. A. Hudson and C. D. Manning, "Compositional attention networks for machine reasoning," *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, pp. 1–20, 2018.
- [31] K. Chen, J. Wang, L.-C. Chen, H. Gao, W. Xu, and R. Nevatia, "Abccnn: An attention based convolutional neural network for visual question answering," *arXiv preprint arXiv:1511.05960*, 2015.
- [32] C. Xiong, S. Merity, and R. Socher, "Dynamic memory networks for visual and textual question answering," in *International conference on machine learning*, 2016, pp. 2397–2406.
- [33] M. Ren, R. Kiros, and R. Zemel, "Exploring models and data for image question answering," in *Advances in neural information processing systems*, 2015, pp. 2953–2961.
- [34] H. Ben-Younes, R. Cadene, M. Cord, and N. Thome, "Mutan: Multi-modal tucker fusion for visual question answering," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2612–2620.
- [35] Z. Su, C. Zhu, Y. Dong, D. Cai, Y. Chen, and J. Li, "Learning visual knowledge memory networks for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7736–7745.
- [36] J. Singh, V. Ying, and A. Nutkiewicz, "Attention on attention: Architectures for visual question answering (vqa)," *arXiv preprint arXiv:1803.07724*, 2018.
- [37] J. Song, P. Zeng, L. Gao, and H. T. Shen, "From pixels to objects: Cubic visual attention for visual question answering," in *IJCAI*, 2018, pp. 906–912.
- [38] P. Lu, H. Li, W. Zhang, J. Wang, and X. Wang, "Co-attending free-form regions and detections with multi-modal multiplicative feature embedding for visual question answering," *arXiv preprint arXiv:1711.06794*, 2017.
- [39] L. Ma, Z. Lu, and H. Li, "Learning to answer questions from image using convolutional neural network," *arXiv preprint arXiv:1506.00333*, 2015.
- [40] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014, pp. 1532–1543.
- [41] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [42] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [43] G. Li, L. Zhu, P. Liu, and Y. Yang, "Entangled transformer for image captioning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8928–8937.
- [44] C. Sur, "Self-segregating and coordinated-segregating transformer for focused deep multi-modular network for visual question answering," *arXiv preprint arXiv:2006.14264*, 2020.
- [45] Y. Kant, D. Batra, P. Anderson, A. Schwing, D. Parikh, J. Lu, and H. Agrawal, "Spatially aware multimodal transformers for textvqa," *arXiv preprint arXiv:2007.12146*, 2020.
- [46] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, "Vi-bert: Pre-training of generic visual-linguistic representations," *arXiv preprint arXiv:1908.08530*, 2019.
- [47] H. Tan and M. Bansal, "Lxmert: Learning cross-modality encoder representations from transformers," *arXiv preprint arXiv:1908.07490*, 2019.
- [48] A. U. Khan, A. Mazaheri, N. d. V. Lobo, and M. Shah, "Mmft-bert: Multimodal fusion transformer with bert encodings for visual question answering," *arXiv preprint arXiv:2010.14095*, 2020.
- [49] D. A. Hudson and C. 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*, 2019, pp. 6700–6709.
- [50] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

APPENDIX A
TOTAL TRAINING TIME

Dataset	Training Time per Epoch (minutes)							
	RAMEN				TransformerNet			
	Baseline Concat [1]	Additive	Multiplicative	Question Fusion	Concatenation	Additive	Multiplicative	Question Fusion
VQAv1	17.00	27.51	27.67	41.52	32.75	29.04	29.10	58.12
VQAv2	27.15	50.80	50.94	76.59	58.45	77.88	78.30	110.40
VQACP2	16.29	26.27	32.74	40.08	25.00	45.39	45.67	77.67
CVQA	11.09	14.87	15.04	26.31	21.67	42.94	28.11	39.85
TDIUC	48.79	85.32	86.13	129.87	91.65	89.73	90.00	183.57
CLEVR	28.51	49.31	51.49	47.26	28.18	56.28	56.82	63.57
CLEVR-Human	0.80	1.01	1.39	1.28	0.82	1.52	92.38	2.19
CLEVR-CoGenTA	16.40	47.47	49.28	44.95	28.11	53.52	53.74	60.49
CLEVR-CoGenTB	-	-	-	-	-	-	-	-
Average	20.75	37.82	39.33	50.98	35.83	49.54	59.26	74.48
Total Average Time	46.00							

APPENDIX B
TRAINING RESULTS

Dataset	Ramen Concat				TransformerNet Concat				Ramen Additive				TransformerNet Additive			
	Training Score		Test Score		Training Score		Test Score		Training Score		Test Score		Training Score		Test Score	
	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score
VQAv1	25	88.44	25	63.3	50	8388.85	50	63.32	25	84.65	25	63.21	50	58.33	50	55.88
VQAv2	25	84.95	25	62.16	46	82.56	46	64.06	25	82.43	25	63.64	25	62.99	25	59.79
VQACP2	18	80.81	18	37.61	16	80.92	16	37.47	14	86.53	14	36.73	93	53.93	93	27.6
CVQA	8	68.36	8	56.98	16	73.61	16	57.74	8	68.23	8	55.82	20	79.66	20	53.49
TDIUC	-	-	9	66.48	-	-	16	65.47	-	-	14	64.9	-	-	7	58.03
CLEVR	20	99.75	20	96.52	68	96.94	68	95.79	18	99.63	18	96.26	16	50.5	16	50.52
CLEVR-Human	31	100	31	44.57	77	99.2	77	46.49	10	87.67	10	40.21	83	38.96	83	38.45
CLEVR-CoGenTA	15	99.76	15	96.59	83	98.86	83	96.43	24	99.82	24	96.84	32	64.48	32	64.26
CLEVR-CoGenTB	-	-	15	88.27	-	-	83	89.68	-	-	20	89.42	-	-	32	55.8
Average	68.05				68.49				67.45				51.53			
Dataset	Ramen Question Fusion				TransformerNet Question Fusion				Ramen Multiplicative				TransformerNet Multiplicative			
	Training Score		Test Score		Training Score		Test Score		Training Score		Test Score		Training Score		Test Score	
	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score	Epoch	Score
VQAv1	25	76.9	25	64.76	25	54.61	25	55.08	25	70.32	25	65.54	25	62.83	25	60.91
VQAv2	25	79.53	25	65.07	25	49.31	25	50.14	25	68.98	25	65.28	25	56.1	25	56.32
VQACP2	30	87.09	30	37.03	26	55.75	26	26.89	25	71.76	25	36.28	93	53.93	93	27.6
CVQA	8	75.38	8	56.81	25	56.16	25	48.47	38	90.41	38	56.58	17	54.2	17	54.2
TDIUC	-	-	11	64.81	-	-	11	53.9	-	-	15	65.69	-	-	8	56.34
CLEVR	14	97.83	14	96.06	25	49.84	25	50.07	21	99.86	21	96.72	25	57.24	25	57.31
CLEVR-Human	17	99.76	17	48.63	89	39.36	89	37.99	25	100	25	46.46	25	54.64	25	40.07
CLEVR-CoGenTA	21	99.68	21	96.9	23	60.35	23	60.24	16	99.81	16	96.63	25	72.42	25	72.07
CLEVR-CoGenTB	-	-	21	88.74	-	-	23	55.19	-	-	16	86.22	-	-	16	60.18
Average	68.76				48.66				68.38				53.89			