# VLR HW4

**Akshay Antony**
**akshayan@andrew.cmu.edu**

*To run the codes:*
1. *Simple baseline: python3 main.py –model=simple*
2. *Coattention: python3 main.py –model=coattention*

## Task 1: Data Loader:

1.1: **getQuesIds(),** Total number of questions=**248349**

1.2: Content of question: **What position is the man squatting with a glove on playing?**
Image ID: **40938**

1.3 Most Voted answer: **Catcher**

1.7 We cannot use one-hot embedding for words, because we do not know the length of each answer to output answer words that form a sentence. An answer can have any number of words while in question maximum number of words is fixed for a sentence. And we are modeling the solution as a single label classification, so a single label should contain the whole answer sentence.

1.8 The size of the dataset should be the number of questions because a single image has multiple questions and each question has 10 answers. So we should choose the question number as the index in the data loader from which we can access the corresponding single image and 10 answers. This leads to all the data having the same dimensions which are impossible if we choose the answer number or number of images.

1.9

3. There is _max_question_length to a question, so that will be the dimension of the one-hot question vector. All the questions with less than _max_question_length will be effectively padded with zeroes after the sentence ends. Dim of question tensor: **_max_question_length * question_word_list_length**
4. Dim of answer tensor **10 * answer_list_length**

## Task 2: Simple Baseline

2.1 **Bag of words Advantages:** It is a very compact representation and it is irrespective of the sentence length. Any length sentence can be converted to a bag of words representation of the **question_word_list_length** dimension. Unlike one-hot embedding, you don't need to define a max_question_lenght.

**Bag of words Disadvantages:** The representation does not take into account the position of words in the sentence. The bag of words only takes into account the frequency of each word in a sentence. For example, "**Is it good?**" and "**it is good**" have the same representation.

To convert one-hot encoding to bow, You can sum across the first dimension of one-hot vector of size **_max_question_length * question_word_list_length** to obtain the bow tensor of dim **question_word_list_length.**

**2.2** The major components are **image_feature_extractor, word_feature_extractor, softmax classifier.**
1. Input and output dimension of image_feature_extractor: **input size=B*3*224*224, output size = B*1024,** B=Batch size

2. Input and output dimensions of word_feature_extractor: **B*question_word_list_length, output size= B*512**

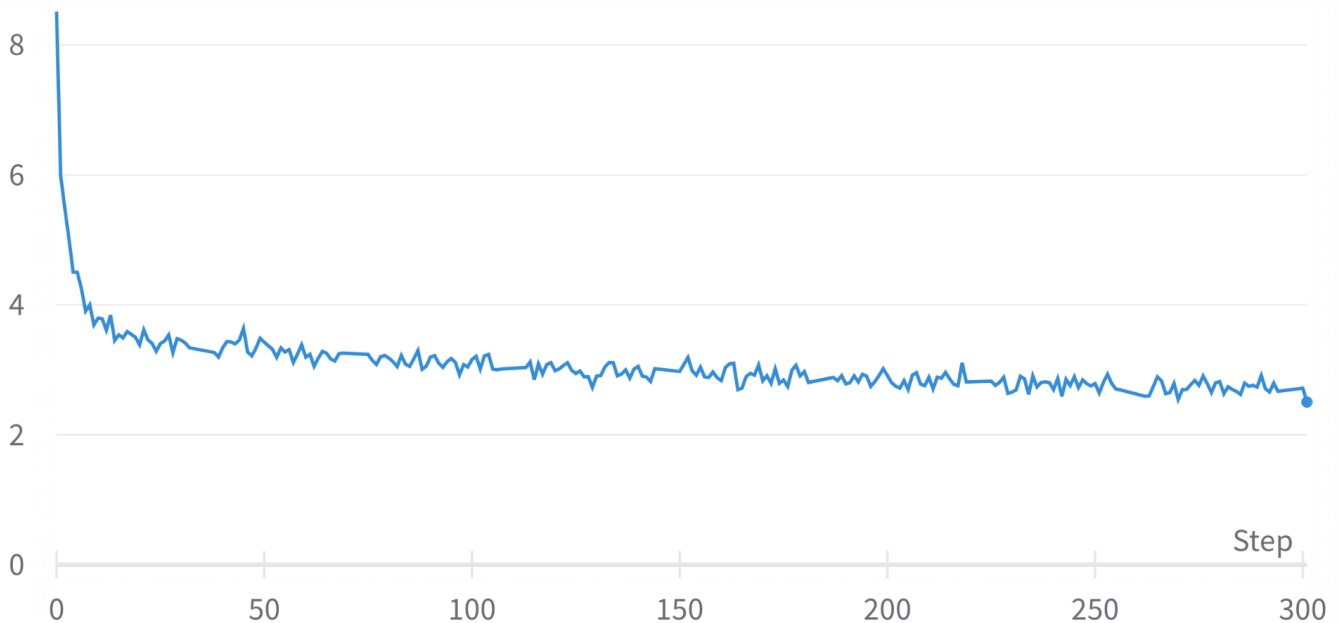3. Input and output dimensions of softmax classifier: **input size = B*1536, output size=B*answer_list_length**

2.4 As the network input and output dimensions are made w.r.t to training question and answer dictionary. We use the same **question_word_to_id_map** and **answer_to_id_map** for validation set too. So the predictions on the test set will also be done on the training vocabulary. If we produce a new dictionary for the validation set, the network design should also change to accommodate the new question and answer dimensions.

2.5 In  I defined two SGD optimizers, one of which takes the parameters of the question_embedding module of the model and the other optimizer takes the parameters of the softmax classifier module.

2.10 I blended the probabilities of all the 10 answers while training while during validation I used majority voting. I found the word_embediing layer output should not be very small so that the information is neglected when it is concatenated with the image embeddings.
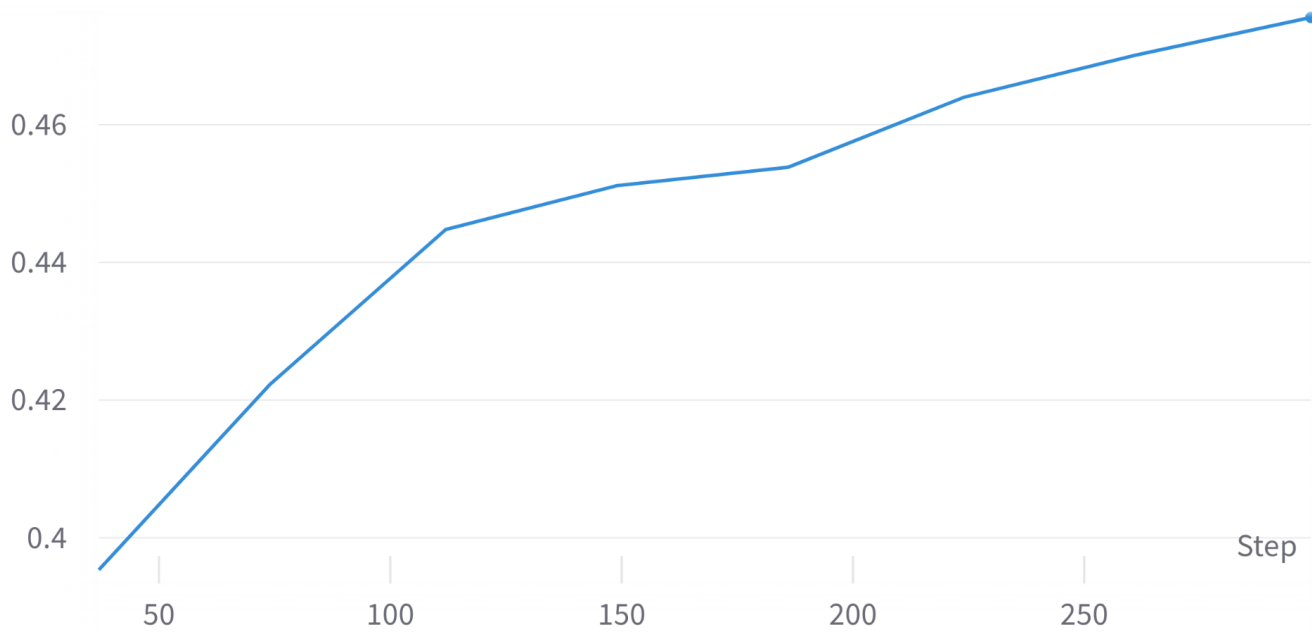
**Trained for 8 epochs with batch size 768, final accuracy obtained = 0.4756, final loss = 2.504**



Train loss vs step number

## SimpleBaseline/traintime_val_accuracy



Accuracy vs step number (after each epoch validation accuracy is calculated once)



Epoch 0
- Question: what color is the sink
  - Predicted answer: white
  - Ground truth: white

Epoch 6

- Question: how many busses are there
  - Pred Answer: 2
  - Ground truth:  3



*Epoch 8:*

- Question: how many people are in the image
  - Predicted answer: 1
  - Ground truth: 1

Link to wandb report

https://wandb.ai/akshayantony12/SimpleBaselines/reports/VLR-task-2--VmlldzoxODc3MDEx?accessToken=y2bekbn5bxi7fq1wpzh4y6iwph6fw1qlfttyty97v4truw27ik4y55euo9qktoim

**Task 3**

3.3

1. The three hierarchies of question representation used are
   **1. Word level 2. Phrase level 3. Sentence level.**

   1. **Word Level:** The word-level features are extracted using a multi-layer perceptron, which maps input word vector of size $B*max\_sentence\_length*number\_of\_words\_in\_question\_dictionary$ to **B*max_sentence_length*emb_dims,** where emb_dims is one of our choices. The word-level features consider each word separately and through the mlp layer, it projects each word to a higher-order latent feature. The second dimension is actually the words of the sentence which we treat in an mlp layer separately i.e we are aggregating the word embeddings along this axis. So this extracts the word-level features.

   2. **Phrase Level:** But as phrases would have more meaning to a question we need to extract the phrase level features. So we need something like convolution which aggregated the features around a neighborhood defined by kernel size. But we do not know how many words we should include in a phrase. So we use multiple kernel sizes to extract features. We use three kernel sizes of 1, 2, and 3 here. The input to this layer is the word-level features thus making it hierarchical. The three features after the 1d convolution are concatenated and aggregated using max operation. Thus the output is the same as word-level features i.e **B*max_sentence_length*emb_dims**

   3. **Sentence level:** Sometimes the question fully might give us the contextual information. In that cases, we need the whole sentence in a higher-dimensional space. Also, phrase level is limited to the number of words we choose for a phrase i.e. the kernel size, which can never cover the whole sentence. So sentence-level features are important.  For sentence-level features, a single lstm layer is used, which takes the phrase-level features as input. LSTM is used because we need a sequential feature. Thus the output is the same as word-level features and phrase-level features i.e **B*max_sentence_length*emb_dims**

**2.**

**a.   Attention** can be seen as concentrating only on the most important features that help our task. Without attention, we are treating all the words in the sentence as equal. Suppose after the word level feature extraction for a question we have the size of features as  **B*max_sentence_length*emb_dims,** we finally need a single feature of size **B*emb_dims,** we need to aggregate along the first dim i.e. **max_sentence_length**, so we cant just use sum because sum just takes all the words as equal. We need a weighted sum, these weights are obtained using the attention mechanism. If using the word features to get the wordwise weights that would be self-attention. Programming wise is obtained by using an MLP that projects emb_dims to some other dim and is projected back to 1, which gives a tensor of shape **B*max_sentence_length,** which is passed through a softmax to make the coefficients of weights sum to 1. This is multiplied by the initial word-level features of size B*max_sentence_length*emb_dims and added across the dim=1 to cause a tensor of size **B*emb_dims.**

**Coattention:** In the previous paragraph self-attention was explained, where the weights we calculated only using word features. But in the case of coattention, we use the image features as well to calculate the attention weights. So effectively both the image and question feature operates to give attention

weight or more formally guide the model to concentrate on certain parts of the sentence and certain regions of the image rather than focussing equally on the whole sentence and image.

**Help to vqa:** So in the vqa task, there is a high chance that the answer can be extracted for some regions and some parts of the sentence. In that cases looking at the whole sentence and image would not be a good idea because some regions might not contribute to the answer and may even hamper the performance by giving unnecessary information. Using coattention the model attends only to the critical parts of the questions and image giving higher accuracy of answer prediction.

3. As the process of task 3 involves various steps we will need to modularize all the independent steps like hierarchical feature extraction, feature aggregation using alternating co-attention, and a final module for the public to use. The final module should accept the input sentence and image and give out the output. But the final module should accept some arguments like
   a. word_input_size=None,
   b. embedding_size=512,
   c. dropout=0.5,
   d. num_answers=None,

The word input size should be the number of unique words in the question dictionary embed size should be what the features need to be set to. Users should be able to set this according to their requirements. Dropout should be users' choice depending on the need for regularization. num_answers is the number of unique answers in the answer dictionary. So by developing appropriate modules for question feature extraction and parallel/or alternative co-attention, we can modularize the attention code.

3.5

Something special is I used Adam optimizer with lr=0.001 for this task, I also used a batch size of 500 and trained for 10 epochs the final accuracy obtained was **57.93.** Also, I changed the dimension of self.Ww, self.Wp, self.Ws to 1024 (i.e linear layer to output a tensor of 1024 dimension). I also tried 512 but the accuracies were not high enough. The paper also recommends using 1024 for vqa dataset.
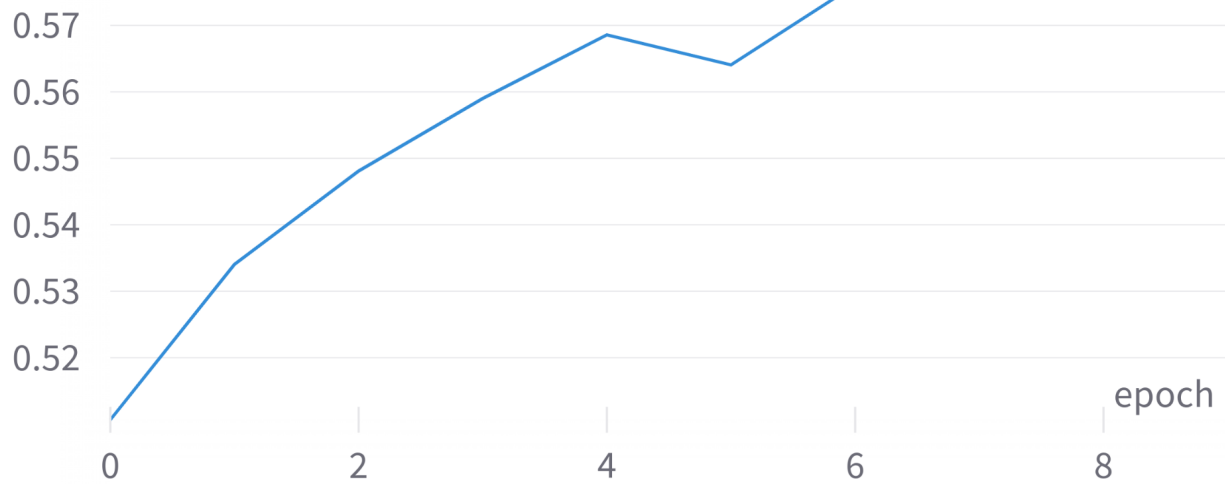
# Coattention/train_loss

# Coattention_val_accuracy

**—** stellar-cherry-123



Epoch 3

[OBJ][OBJ]question: what color is the sand
predicted answer: blue
ground truth: brown

Epoch 5
question: is the dog hungry
predicted answer: yes
ground truth: yes



Epoch 9
question: what system are they playing,
predicted answer: wii,
ground truth: wii

**Wandb link**

https://wandb.ai/akshayantony12/SimpleBaselines/reports/VLR-TASK-3--VmlldzoxODc2OTE2?accessToken=tadmmyqoas1ecsh7t86v3ni28wlpen0macs3m9zc9h7nri7recsoq4bnm1id1ubf

4.1

1. Self-attention in aggregating phrase-level features. In task 3 phrase-level features are aggregated using a mean pooling of different kernel sizes. But we do not know which kernel size namely 1, 2, or 3 is more important. If we use self-attention to find some weights to aggregate the three phrase-level features we might get better results.

2. In simple baselines instead of using MLPs, we can use LSTM to project the bag of words input to a higher level feature and concatenate it with image features. This would result in higher accuracy than the simple baseline.

3. Use a better sentence feature extraction like BERT.