# 6 Supplementary Material for TAB-VCR: Tags and Attributes based Visual Commonsense Reasoning Baselines

We structure the supplementary into two subsections.

1. Details about implementation and training routine, including hyperparamters and design choices.
2. Additional qualitative results including error modes
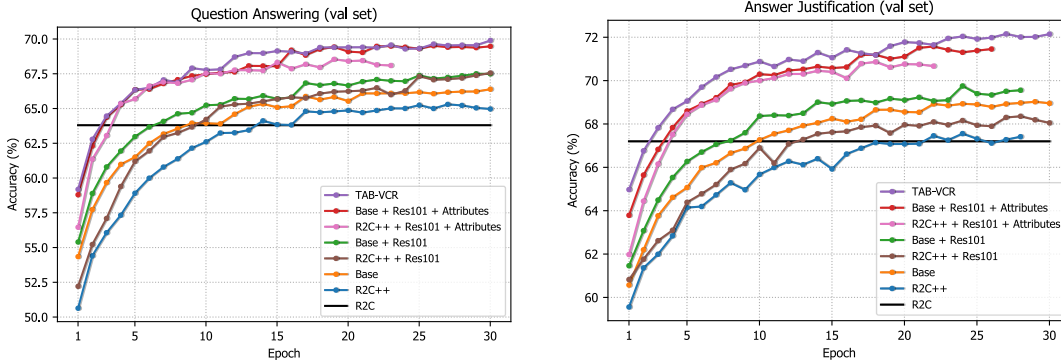
## 6.1 Implementation and training details



Figure 6: **Accuracy on validation set.** Performance for $Q{\to}A$ (left) and $QA{\to}R$ (right) tasks.

As explained in Sec. 3.1, our approach is composed of three components. Here, we provide implementation details for each: (1) BERT: Operates over query and response under consideration. The features of the penultimate layer are extracted for each word. Zellers et al. [103] release these embeddings with the VCR dataset and we use them as is. (2) Joint encoder: As detailed in Sec. 4.3, we assess different variants over the baseline model using two CNN models. The output dimension of each is 2048. The downsample net is a single fully connected layer with input dimension of 2048 (from the image CNN) and an output dimension of 512. We use a bidirectional LSTM with a hidden state dimension of $2 \cdot 256 = 512$. The outputs of which are average pooled. (3) MLP: Our MLP is much slimmer than the one from the R2C model. The pooled query and response representations are concatenated to give a $512 + 512 = 1024$ dimensional input. The MLP has a 512 dimensional hidden layer and a final output (score) of dimension 1. The threshold for Wu Palmer similarity $k$ is set to $0.95$.

We used the cross-entropy loss function for end-to-end training, Adam optimizer with learning rate $2\mathrm{e}{-}4$, and LR scheduler that reduce the learning rate by half after two consecutive epochs without improvement. We train our model for 30 epochs. We also employ early stopping, *i.e.*, we stop training after 4 consecutive epochs without validation set improvement. Fig. 6 shows validation accuracy for both the subtasks of VCR over the training epochs. We observe the proposed approach to very quickly exceed the results reported by previous state-of-the-art (marked via a solid horizontal black line).

## 6.2 Additional qualitative results

Examples of TAB-VCR performance on the VCR dataset are included in Fig. 7. They supplement the qualitative evaluation in the main paper (Sec. 4.4 & Fig. 3). Our model correctly predicts for each of these examples. Note how our model can ground important words. These are highlighted in **bold**. For instance, for Fig. 7(a), the correct rationale prediction is based on the expression of the **lamp**, which we ground. The lamp wasn't grounded in the original VCR dataset. Similarly grounding the **tag**, and **face** helps answer and reason for the image in Fig. 7(b) and Fig. 7(c). As illustrated via the **couch** in Fig. 7(d), it is interesting that the same noun is present in detections yet not grounded to words in the VCR dataset. This could be accounted to the data collection methodology, as explained in Sec. 4.4 ('explanation of missed tags') of the main paper.

In Fig. 8(a), we provide additional examples to supplement the discussion of error modes in the main paper (Sec. 4.4 & Fig. 5). TAB-VCR gets the question answering subtask (left) incorrect, which we detail next. Once the model knows the correct answer it can correctly reason about it, as evidenced by being correct on the answer justification subtask (right). In Fig. 8(a) both the responses 'Yes, she does like [1]' and 'Yes, she likes him a lot' are very similar, and our model misses the 'correct' response. Since the VCR dataset is composed by an automated adversarial matching, these options could end up being very overlapping and cause these errors. In Fig. 8(b) it is difficult to infer that the the audience are watching a live band play. This could be due to the missing context as video captions aren't available to our models, but were available to workers during dataset collection. In Fig. 8(c) multiple stories could follow the current observation, and TAB-VCR makes errors in examples with ambiguity regarding the future.
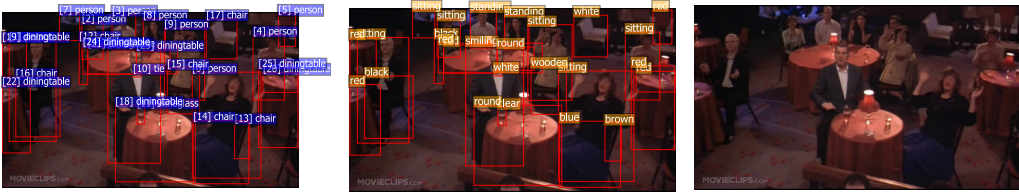
Figure 7: **Qualitative results.** More examples of the proposed **TAB-VCR** model, which incorporates attributes and augments image-text grounding. The image on the left shows the object detections provided by VCR. The image in the middle shows the attributes predicted by our model and thereby captured in visual features. The image on the right shows *new tags* detected by our proposed method. Below the images are the question answering and answer justification subtasks. The *new tags* are highlighted in **bold**.

Figure 8: **Qualitative analysis of error modes.** Responses with (a) similar meaning, (b) lack of context and (c) ambiguity in future actions. Correct answers are marked with ticks and our models incorrect prediction is outlined in red.