

1 To reviewer #1:

2 **Broader applicability of the method.** Our method can be applied to a broader range of tasks that fall into visual
3 reasoning. To name a few, Newtonian physics problem solving task [1], geometric problem solving task [6] and CLEVR
4 as reviewer #3 suggested. We have observed a growing number of reasoning tasks in NLP such as multi-hop text
5 question answering, knowledge graph reasoning and conversational models. Our approach can potentially be applied to
6 strengthen the reasoning abilities of such tasks.

7 **Availability of the problem categories.** Category is widely available for abstract reasoning / visual reasoning which
8 state-of-the-art models leverages as type loss. When such categories are missing, one can either get it using an
9 unsupervised way such as the paper[5] that reviewer #2 suggested. Another way to deal with such a problem is to
10 assume a latent category variable and optimizes it together with teacher model which could be a promising future work.

11 **Differences to [7].** We compared methods used in [7] and as shown in table 1 and table 2, it is not performing well for
12 abstract reasoning. In [7], only the loss of the last step is used as opposed to the complete trajectory in our method.
13 Additionally, in [7] action is taken to be a 0/1 decision on sample id but our action is a proportion of the related problems.
14 All of these leads to the better performance of our teacher model. We will improve the paper to clearly outline the
15 differences.

16 **Confirmation of the source of performance gain.** We have illustrated in the empirical study that training with a
17 specific trajectory with different proportion of the distracting/reasoning features can dramatically improve model
18 performance. This is illustrated in Figure 1, which shows that the difficulty of abstract reasoning task lies in the
19 existence of distracting features. Table 1 shows that the appropriate training trajectory can greatly improve the model
20 in the presence of distracting features. At the same time, the visualization of Sec 5.4 also shows that our method
21 can distinguish distracting features better. We will put a visual training trajectory map in the final paper for better
22 illustration.

23 **Other issues.** The reviewer is right that χ_k (L151) is a set of embedding of all triple-panels. We will fix it.

24 To reviewer #2:

25 **Disentangled representations.** Disentangled representation separates information on a single input while our method
26 select inputs for a model. Our method is orthogonal to disentangled representation and can be applied on top of it. It is
27 also worth mentioning that [4] actually implemented the idea of disentangled representation but with little improvements
28 on abstract reasoning. One intuitive explanation is that distracting features live in a much more illusive manifold and
29 disentangled representation along is not capable of separating it from reasoning features.

30 **Related work.** We will add discussions of disentangled representation into our related work.

31 **Extrapolation experiments.** We actually did an extrapolation experiment on the PGM. We separated training and
32 testing in a way that they have non-overlapping values of “color” attribute and have achieved 8% improvement in
33 accuracy. We are happy to include this result along with additional experiments in the final paper.

34 **Performances of models other than LEN.** We actually have a complete set of comparisons of WReN with/without
35 teacher model in table 2. Performance of WReN is improved from 75.6% to 77.8% with type loss and from 62.8% to
36 68.9% without type loss. We will complete comparisons of other baselines (e.g., RN) in the final paper.

37 To reviewer #3:

38 **Performance improvements of the teacher network on LEN compares to WReN model on PGM.** The performance
39 improvements on WReN is as significant as the one on LEN. Since the codebase of WReN has never been released,
40 we implemented it on our own with an accuracy of 70.1% using type loss but we have never been able to reproduce
41 the reported benchmark (i.e., 75.6%). Nevertheless, we include the performance of the published results of WReN in
42 our paper for fair comparison. As a matter of fact, if we compare the improvements of teacher model against our own
43 baseline on WReN the improvement is 7.2%. Comparing to 11.1% of improvements with LEN model.

44 **Testing on additional visual reasoning tasks.** We actually did an experiment on the CLEVR dataset but we didn’t
45 include it into the paper. our LEN model achieves 1.7% accuracy increase(95.5% to 97.2%) compared to RN[2]. Please
46 note that CLEVR dataset is much easier than the datasets we used in the paper and the already high performance on
47 baseline method allow only a marginal improvements using our teacher model. Nevertheless, the performance gain
48 seems to be significant and consistent. We will explore the efficacy of the model on more widely used visual reasoning
49 tasks (E.g., CLEVR-CoGeNT task) in the final paper.

50 **Writing Issues.** We will fix these issues as the reviewer suggested.

51 [1] Sachan M, et al. Parsing to programs: A framework for situated qa. KDD 2018.

52 [2] Santoro A, et al. A simple neural network module for relational reasoning. NIPS 2017.

53 [4] Steenbrugge, X., et al. Improving generalization for abstract reasoning tasks using disentangled feature representa-
54 tions. In Workshop on NIPS, 2018.

55 [5] Hsu, Kyle, et al. Unsupervised learning via meta-learning. arXiv 2018.

56 [6] Seo M, et al. Solving geometry problems: Combining text and diagram interpretation. EMNLP 2015.

57 [7] Yang Fan, et al. Learning to teach. ICLR 2018.