

1 **To All Reviewers** We wholeheartedly thank all reviewers for all valuable and insightful feedback, along with taking  
2 the time to review our submission. We have provided detailed comments to each reviewer below.

3 **To Reviewer 1** Thanks for the valuable and insightful feedback!

4 Regarding the suggestion about more experiments on more datasets, we full wholeheartedly agree that more experiments  
5 will definitely improve the comprehensiveness of the paper. We are running experiments on these tasks right now and  
6 include results on these datasets (SST, SNLI, WMT En-De etc.) by the camera ready version as supplementary material.

7 Pertaining to Q3, weights are extracted from a trained model (on NLI) on sampled  $1K$  datapoints from the dev set.  
8 Similar graphs appear with repeated sampling. Regarding deletion being prevalent, our hypothesis is that only tokens  
9 that are particularly significant will be compositionally (add/subtracted). For most words, the sigmoid function provides  
10 flexibility of deleting tokens.

11 Regarding the value of  $\alpha, \beta$  in equation 1 and equation 2, we set them to 1 in the experiments. The higher value of  $\alpha, \beta$   
12 is, the ‘harder’ the compositional pooling becomes. The hardness of CoDA required is possibly related and analogous  
13 to hard vs soft attention and can be domain dependent. For most language tasks, we find that not biasing CoDA towards  
14 being hard is quite sufficient. Not all tokens are important, so CoDA maintains the flexibility of standard attention while  
15 enabling arithmetic compositionality. We will include more supplementary distribution visualisations on different tasks  
16 in the revised version.

17 Regarding the form of Equation 9 used in the experiments, we used one layer non-linear projection network to compute  
18 the pairwise affinity. We thank the reviewer for pointing it out and we will mention it clearly in the revision.

19 For other comments such as references or typos, we will correct and add them in the revision. We will also be sure to  
20 include a discussion about sparsemax and softgen.

21 **To Reviewer 2** Thanks for the insightful and valuable feedback!

22 Regarding the evaluation on Tensor2Tensor, IWSLT En-Vi was chosen because of the size and resource limitations.  
23 Please be assured that the tasks were not cherry-picked and we have not experienced any failure cases on T2T tasks yet.  
24 We also had success on the arithmetic T2T and subject-verb agreement tasks but did not report due to lack of space.  
25 We will prepare detailed supplementary material to cater to extra experiments. Moreover, we are currently running  
26 WMT En-De and WMT En-Ro. This may take sometime due to our limited hardware but will definitely be ready by the  
27 camera-ready version.

28 Regarding the form of CoDA, we use the following version our experiments, i.e. subtract the mean before applying  
29 the sigmoid or tanh, and they are scaled before applying self-attention. We will make this absolutely clear in the final  
30 version. We left several variations open in our technical exposition since certain hyperparameters (e.g.,  $\alpha$ ) could allow  
31 practitioners to control properties (i.e., hardness) of CoDA. We apologize for any confusion. In the revision, we will  
32 also include more ablation studies of different alternations of CoDA form, providing better guidance for usage of CoDA  
33 configurations. We will also include supplementary visualisation as requested.

34 Regarding the confusion of using notions and references in the papers, as pointed out in the detailed comments, we  
35 thank the reviewer for pointing them out. We will correct and clean them in the revision.

36 **To Reviewer 3** Thanks for the insightful and valuable feedback!

37 Our intuition is that tanh saturates at  $\{-1, 1\}$ , which doesn’t allow the model the flexibility to *delete* (erase) tokens.  
38 Sigmoid provides this flexibility to our model. Some early ablation results on retrieval tasks are reported in Table 1. We  
39 will include more comprehensive ablations in the final version.

40 Thanks for the suggestion about adding a visualization of the learned attention weights on examples for interpretability.  
41 We will definitely include it in the revision. Pertaining aesthetic comments, thanks for pointing them out and we will be  
42 sure to correct them in the revision.

Method	TrecQA	WikiQA
<i>tanh</i> only	66.78 / 73.49	67.13 / 67.81
CoDA	79.84 / 84.78	68.07 / 68.28

Table 1: Dev set results of *tanh* only versus CoDA (*tanh* \* *sigmoid*).