

1 We would like to thank the reviewers for their helpful comments and their thorough evaluation of our work. We are
2 encouraged that reviewers find our paper clear and well written (R1, R2, R3) and our method to be theoretically sound
3 (R1, R2), of high practical impact (R2, R3), and intuitive and easy to implement (R3). Below, we address the issues that
4 have been raised.

5 **[R1] Should the clusters merge in the next layer?:** The output of the attention layer involves (a) residual connections:
6 this means queries in the same cluster can have different values (b) multi-head attention: this allows for a query to
7 be clustered differently in two different heads. This is followed by a feed-forward projection, which results in new
8 clustering patterns in the subsequent layers.

9 **[R1] Is LSH-X a simpler version of the Reformer?:** LSH-X refers to Reformer as described in the paper except
10 for the reversible layers. Reversible layers is a technique introduced by Gomez et al. (2017) and is orthogonal and
11 independent to both clustered and LSH attention. As our focus is on the efficiency of the attention implementation, we
12 do not consider the reversible layers for both methods.

13 **[R1] Can Reformer be used for heterogeneous queries and keys?:** Using Reformer with heterogeneous queries
14 and keys, requires significant changes in its core component, namely sorting and chunking the hash values to create
15 query/key groups. In contrast, clustered attention places no such restriction.

16 **[R1] Analysis of clusters and Set Transformers:** Thank you for the suggestions, we will add further analysis of the
17 clustering in the supplementary material as suggested. We will also add Set Transformers to the related work section.

18 **[R1, R2] Experiments on other modalities. Is speech favorable to clustering?:** Thank you for your comments.
19 We would like to mention that our NLP approximation experiment for GLUE and SQuAD tasks in § 4.3 shows that
20 clustered attention can capture arbitrarily complex attention patterns. In the future, we will also explore applications to
21 NLP/vision tasks in the long context setting, as suggested.

22 **[R1, R2] Other baselines to investigate: Longformer, Sparse Transformer, Pay-less-attention, and Local atten-**
23 **tion:** We thank R1, R2 for the valuable suggestions. Sparse Transformer’s reliance on the blocksparse library and
24 the lack of documentation significantly hinder its evaluation on new tasks (evident by the lack of comparisons in the
25 community). Longformer is a very recent work that we will compare with in the future. Upon R1’s suggestion, we
26 evaluate local attention on the synthetic copy task (§ 3.2 in the supplementary). We observe that for a sequence of
27 length 128, local attention with context of 32 and 4 layers cannot solve the task whereas all clustered attention variants
28 can solve it perfectly with 30 clusters.

29 **[R3] Is sparse dot product used in improved-clustered attention practically efficient?:** We refer R3 to Fig. 3 of
30 the supplementary, where we present the time-memory benchmark. Notably, improved attention is faster than vanilla
31 for sequence lengths greater than 1024. Unlike general sparse dot products, our implementation exploits the fact that
32 all queries in a cluster use the same top-k keys. This allows us to cache the keys in CUDA shared memory for better
33 performance. We will add details about our implementation in the supplementary material.

34 **[R3] If C grows with N, then the method has quadratic complexity:** We agree with R3 that C should be small for
35 good performance. Note that, C is not necessarily a function of N as it depends heavily on the task. We would like to
36 refer R3 to an ablation on the relationship between clusters and sequence length in § 3.2 of the supplementary. For
37 the masked copy task, improved-clustered solves the task for every sequence length (up to 512) with just 15 clusters.
38 Similarly, for SQuAD approximation, using only 25 clusters for a sequence length of 380 leads to good performance.

39 **[R3] Training time (table 2) is only slightly faster with worse performance than the baseline?:** While for CTC
40 loss (table 2), the training time improvements were only small, for LF-MMI loss (table 3), the improvements were
41 significant. Also note that while the average sequence lengths for WSJ and Switchboard are 780 and 530, the maximum
42 lengths are 2500 and 3850 respectively. Finally, during inference under fixed computational budget, improved-clustered
43 is consistently better than the baselines (Fig. 1).

44 **[R3] Is clustering complexity $O(NCD)$ so bad? What is actual speed improvement on RoBERTa? Minor com-**
45 **ments:** Thank you for the suggestion, we agree that for long sequences, k-means with euclidean distance $O(NCD)$
46 could improve performance due to better clustering. For RoBERTa approximation, vanilla attention is faster due to
47 small sequences. The purpose of the experiment is to understand the ability of the clustered attention variants to capture
48 complex attention distributions using only a few clusters. We will make it clear in the camera-ready version that vanilla
49 attention is faster in this experiment. We also want to thank R3 for pointing out the mistakes, we will fix them.