

1 **General Response** We thank all reviewers for their insightful comments! We are sorry that Figure 1 in the submission
2 version of this paper is not the latest so the descriptions in Section 4.3 are ambiguous. We will correct it in the updated
3 version. We have revised typos and removed duplicated sentences.

4 **Regarding computation efficiency and memory cost.** For memory, SAC consumes smaller memory since the
5 connections are sparse. For computation efficiency, at training time, SAC is only slightly faster than vanilla transformers,
6 since learning which node should be linked to which node using RL is time-consuming. But at test time, SAC is
7 significantly faster due to significantly smaller cost in self-attention computations. We will add more details in the
8 updated version.

9 **To Reviewer #1** Yes, we agree that the improvements compared with state-of-the-art models are marginal. But
10 the main goal of this paper is to **reduce the memory cost of vanilla self-attention while achieving slightly better**
11 **performances.** We do not attempt to improve the results, and instead, we show that with less attention connections, the
12 model is also strong and consistent over tasks.

13 **To Reviewer #2** Thanks for your careful and insightful feedback! We will correct all the typos and incorrect references
14 in the next version. For the order, it is indeed what you think, i.e. $\{a_{11}, a_{12}, a_{21}, a_{22}, \dots\}$, where a_{i1} is the start node
15 for the i -th edge and a_{i2} is the end node.

16 **Correctness:** Yes, we cannot directly say the other four methods are "special cases" of SAC, but by imposing extra
17 constraints when training the LSTM edge predictor, we can actually induce each of them. For example, for vanilla
18 self-attention (take the encoder side as an example), we can feed each node N times into the predictor so that it recovers
19 vanilla self-attention with the help of distant encodings. For Transformer-XL, we can still segment the sequence and
20 apply the above operations.

21 In terms of decoding in MT, all models (including baselines) use beam search. We are sorry for the confusion and will
22 make this point clearer in the updated version.

23 **Weakness:** See **General Response**.

24 **To Reviewer #3** Thanks for you helpful and insightful comments!

25 **Weakness 1&2:** See **General Response**.

26 **Weakness 3:** We omit the parameters due to the limited space of the page. In fact, the numbers of parameters for these
27 methods including SAC are very close, as you can see from Table 2 and Table 4 that the most parameters come from the
28 main model, i.e. Φ rather than Θ .

29 **Weakness 4:** Thank you for the sensible comment. We used a simple version of the edge predictor, in which all layers
30 share the same structure and each node has to be connected to some other nodes for each layer. For the head adaptive
31 strategy, we reported results for both (head adaptive or not adaptive) for different tasks. We will make these points
32 clearer in the updated version.

33 **To Reviewer #4** We thank you for your insightful comments!

34 **Weakness 1:** See **General Response**.

35 **Weakness 2:** Yes, at first glance, the inference speed of SAC is slower than vanillan Transformer since it introduces
36 the extra process of link prediction. But in fact, SAC does not need to do full self-attention, which makes a significant
37 remedy for test speed. We will show these in the next version. As for baselines, connecting to nearest nodes is
38 actually what CNNs do, for which many recent works have discovered for sparse self-attention. We will compare the
39 performance and the speed of these methods.

40 **Correctness:** Thanks for your suggestions! The intuition of using LSTM edge predictor is to learn different attention
41 patterns for different downstream tasks. In NMT, we find that the learned attention prefers more semantics-related
42 words. We will give a deeper analysis and plot more figures to show how SAC construct attentions for different tasks.

43 **Clarity:** The LSTM predictor is initialized randomly (uniform and gaussian distributions are both possible). We will
44 clarify this in the next version.

45 **Reference:** We are sorry for the confusion and will make it clearer in the updated version.