

## 1 **General response**

2 We thank all reviewers for their valuable feedback and thoughtful suggestions.

3 > Comparison with Gu et al., Stern et al, Welleck et al

4 • > Direct comparison with Gu et al.

5 To the best of our knowledge, there is no official implementation for the paper by Gu et al. (no link to the code  
6 in the paper, no repository stating that it's an implementation for the paper). Thus, direct comparison using the  
7 original implementation was not possible.

8 However, in Section 5.1 we compare the lower bound on the objective we use with the one of Gu et  
9 al.(2019), when applied to our model. This is a more fair comparison of the two approaches since both model  
10 architectures can be used with these lower bounds and should be compared separately.

11 • > Comparison with Stern et al, Welleck et al

12 These works do not report significant improvements in BLEU scores against the autoregressive baselines.  
13 Stern et al.(2019) focus on parallel decoding (with the final result matching the vanilla Transformer). Welleck  
14 et al.(2019) find left-to-right models superior to their approach.

15 While we took into account these papers for the presentation of our work, we want to highlight that this was  
16 concurrent work.

17 > More experiments on standard data sets, e.g., WMT EN-DE

18 We agree that this is important for future research on this topic and would allow for easier comparison with our and  
19 previous work. While we can not provide these results during the author response period (due to large training time of  
20 NMT models for high-resource language pairs), we will add them should the paper get accepted.

## 21 **Response to reviewer 1**

22 > Why not consider tasks without language output? ... For example, you could consider music generation tasks, ...

23 Note that we consider not only natural language output, but also Image-to-Latex, where output is LaTeX formulas.

24 We believe applying our approach to music generation task could be an interesting direction for future work.

25 > Notation issue: ... Do you mean that you also compute an expectation over sampling from  $T^*(Y)$  first, and then use  
26 that to sample  $\tau$ ?

27 By  $\tau \sim p(\tau|X, T^*(Y), \theta)$  we actually mean  $\tau \sim p(\tau|X, \tau \in T^*(Y), \theta)$ . Thus, we sample the trajectory from our  
28 model, but only allow the correct insertion operations (i.e. the ones that lead to producing  $Y$ ).

29 We will make it more clear and expand the discussion.

30 > The footnote 8 on page 6 is important. You should have introduced [28] earlier and more explicitly.

31 > In Section 3, you write ... I would suggest moving some of this discussion up earlier; it will make the presentation of  
32 the algorithm clearer.

33 We agree this would improve the presentation. We will revise the paper accordingly.

## 34 **Response to reviewer 3**

35 > how to decide the output length

36 We describe this in Section 2.2. Specifically, we introduce the EOS element, which indicates the end of generation  
37 process. The model can produce this element instead of a pair of position and a token.