

1 **General Comments:** We thank the reviewers for their thoughtful and useful comments.

2 Density of Presentation: Several reviewers commented on the density of the presentation. We will add more high-level
3 explanations particularly to make cascaded decoding more intuitive (R2, R3, R4). We agree with the general spirit of
4 simplicity and note that the main methods in this work are actually quite basic. The pseudocode presented for cascade
5 decoding is nearly identical to the implementation used, and the Markov transformer can be implemented with only a
6 few modifications to any transformer architecture.

7 **Specific Comments:**

8 **R1** -Long-Term Dependencies: It is a good point that there is a tradeoff between better inference and longer term
9 dependencies among outputs. We note that this method is a compromise in this regard between NAT and fully
10 autoregressive. It is true that we lack long-term dependencies of very long-term models (GPT-3); however it is an
11 open-question of whether all these dependencies are required for conditional generation. For example: one might
12 imagine introducing latent variables to model global topics, or use cascaded decoding to prune the search space for a
13 fully autoregressive model. We think these are promising directions for future work.

14 -Length Prediction: Like most NAT methods, our work needs to predict length beforehand, which we can do with linear
15 or non-linear regression. However, unlike most of those works, our length constraint only needs to be a *maximum length*.
16 For MT we predicted this value, but for other settings with more variable length it could be just set to a large constant.

17 **R2** -More powerful CRF modeling: For training speed, we use a locally normalized LM (a special case of a CRF). It
18 is a very interesting suggestion to instead train a globally normalized CRF. This would require running the cascade
19 decoding at training (a difficult but perhaps feasible approach). We think this is interesting future work.

20 -Relationship to Sun et. al. 2019: We will elaborate on the relationship in our next version. Sun et al introduce a
21 1st-order CRF on top of a non-autoregressive model. Our work goes beyond a 1st-order CRF and approximates argmax
22 from higher-order CRFs. In addition, we use a single Markov transformer to parameterize all log potentials, instead
23 of using additional side-parameters for pairwise potentials. Lastly, we propose tree decoding to make the parallel
24 complexity sublinear, whereas Sun et al’s work is linear.

25 -Why max marginals and not ngram scores: The problem with using ngram scores is that they do not consider
26 compatibility with other positions. Max-marginal fixes this issue with negligible extra time. On WMT14 En-De val,
27 using ngram scores would get BLEU of 28.42 at 123.48ms, while using max marginals reaches 29.24 at 128.58ms.

28 **R3** -Performance of approach: NAT research methods inevitably sacrifice accuracy for speedups vs fully autoregressive
29 models with serial search. Compared to these methods, we present a new point on the Pareto frontier of accuracy vs
30 speed: many of our results are very close in accuracy while achieving major speedups. Within this space, we believe
31 these numbers are valuable, particularly given that this is a novel and clean approach: for example, while Levenshtein
32 gets a bit better performance on some datasets, it requires hand-crafted policies for training, and additional tricks at
33 inference. On the other hand, our approach requires much less tuning: at training time we simply used the same set of
34 hyper-parameters as training language models, while tuning K and iters only happens at inference (like beam size in
35 beam search).

36 -How to generate beyond length 3: We run a parallel version of the forward-backward algorithm, each time pruning
37 low-scoring ngram choices at each position, and we gradually increase n from 1 to M . Finally, to generate an entire
38 sentence, we use dynamic programming (Viterbi) to find the sequence with the highest score under the pruned search
39 space. We presented a few examples in appendix A which might be more intuitive.

40 -Whether $x_{i:j}$ in $\mathcal{X}(x_{i:j})$ should appear at the same position or not: It needs to appear at the same position.

41 -Tree dynamic programming notations: We will clarify our notations in our next version. C_{j,k_1,k_2}^i is the max possible
42 score of all spans from the $(j * 2^i + 1)$ -th position to the $(j + 1) * 2^i$ -th position (the length of any span is 2^i), with the
43 constraint of the left end being word k_1 and the right end being word k_2 . We compute C^i bottom-up, starting from
44 $i = 0$ and merge adjacent spans in C^i to get C^{i+1} (such that span length becomes $2 * 2^i = 2^{i+1}$). The *prefix* score
45 P_{j,k_1,k_2}^i stores the max possible score of all spans from the *1st* position to the $(j * 2^i + 1)$ -th position (also with end
46 constraints), which we compute iteratively top-down using P^{i+1} and C^{i+1} . Conversely, the *suffix* score S_{j,k_1,k_2}^i is the
47 max score of all spans from the $(j + 1) * 2^i$ -th position to the *last* position, also computed top-down. Eventually, we
48 use the prefix scores P^0 , suffix scores S^0 , and log potentials C^0 to calculate max marginals of any edge.

49 **R4** -Performance: Please see comment above for R3. We will also include more recent papers in the next version.

50 -MT Baselines: We followed the exact settings in fairseq translation examples. The 0.3 BLEU score difference might be
51 due to randomness. Our models use the same machine and settings as baselines, so the comparisons aim to be fair. Our
52 IWSLT14 performance baseline is already much better than existing NAT works.