

Document Level Hierarchical Transformer

Najam Zaidi (Monash University), Trevor Cohn (University of Melbourne),
Gholamreza Haffari (Monash University)

Abstract

Generating long and coherent text is an important and challenging task encompassing many application areas such as summarization, document level machine translation and story generation. Despite the success in modeling intra-sentence coherence, existing long text generation models (e.g., BART and GPT-3) still struggle to maintain a coherent event sequence throughout the generated text. We conjecture that this is because of the difficulty for the model to revise, replace or revoke any part that has been generated by the model.

In this paper, we present a novel semi-autoregressive document generation model capable of revising and editing the generated text. Building on some recent semi-autoregressive models, we propose document generation as a hierarchical Markov decision process with a two level hierarchy, where the high and low level editing programs generate and refine the document. We train our model using imitation learning and introduce roll-in policy such that each policy learns on the output of applying the previous action. Experiments applying the proposed approach convey various insights on the problems of long text generation using our model. We suggest various remedies such as using distilled dataset, designing better attention mechanisms and using autoregressive models as a low level program.

Problem Formulation

We cast document generation and refinement as a hierarchical Markov decision process (HMDP) with a two level hierarchy. The high level program is defined by the tuple $(\mathcal{D}, \mathcal{A}_H, \mathcal{E}, \mathcal{R}, \mathbf{d}_0)$ where a state $\mathbf{d} \in \mathcal{D}$ corresponds to a set of sequences $\mathbf{d} = (\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_L)$ up to length L , and $\mathbf{d}_0 \in \mathcal{D}$ is the initial document. The low level program corresponds to the tuple $(\mathcal{S}, \mathcal{A}_L, \mathcal{E}, \mathcal{R}, \mathbf{s}_0)$ where a state $\mathbf{s} \in \mathcal{S}$ corresponds to a sequence of tokens $\mathbf{s} = (w_1, w_2, \dots, w_n)$ from the vocabulary V up to length n , and $\mathbf{s}_0 \in \mathcal{S}$ is the initial sequence.

Policies

High Level Actions work at the sentence level:

- **Reposition:** The reposition policy reads in the document consisting of set of sequences. For every sentence, the reposition policy makes a categorical decision between 0 and L+1 where L is the number of sequences in the document.
- **Insert:** The insertion policy reads the input document consisting of set of sequences and for every possible slot, the insertion policy makes a binary decision which is 1 (insert here) or 0 (do not insert).
- **Update:** The update policy reads the input document, consisting of set of sequences, and for every sequence position, the update policy makes a binary decision which is 1 (update this sentence) or 0 (do not update)

Low Level Actions work at the word level:

- **Reposition and Insert:** Similar to sentence level but works at word level

Training and Generation

Algorithm 1 Generation in HMDP

```
Require: Initial document  $\mathbf{d}_0$ , policy:  $\pi_{\theta_H}$ 
1:  $\mathbf{d} \leftarrow \mathbf{d}_0$ 
2: while Termination condition is not met do
3:    $\text{rep\_index} \leftarrow \arg\max_r \sum_{s_i \in \mathbf{d}} \log \pi_{\theta_H}^{rep}(r_i | \mathbf{s}_i, \mathbf{d})$  ▷ Do reposition
4:    $\mathbf{d} \leftarrow \mathcal{E}(\mathbf{d}, \text{rep\_index})$ 
5:    $\text{ins\_index} \leftarrow \arg\max_p \sum_{s_i, s_{i+1} \in \mathbf{d}} \log \pi_{\theta_H}^{ins}(p_i | \mathbf{s}_i, \mathbf{s}_{i+1}, \mathbf{d})$  ▷ Do insertion
6:    $\mathbf{d} \leftarrow \mathcal{E}(\mathbf{d}, \text{ins\_index})$  ▷ Call to Low level MDP
7:    $\text{upd\_index} \leftarrow \arg\max_u \sum_{s_i \in \mathbf{d}} \log \pi_{\theta_H}^{upd}(u_i | \mathbf{s}_i, \mathbf{d})$  ▷ Do update
8:    $\mathbf{d} \leftarrow \mathcal{E}(\mathbf{d}, \text{upd\_index})$  ▷ Call to Low level MDP
9: end while
```

Algorithm 2 Low Level MDP

```
Require: Document  $\mathbf{d}$ , policy:  $\pi_{\theta_L}$ , Hi Level MDP action:  $\mathbf{H}$ 
1: while Termination condition is not met do
2:    $\mathbf{s}_0 \leftarrow \text{buildFrame}(\mathbf{d}, \mathbf{H})$ 
3:   if  $\mathbf{s}_0$  is empty then
4:      $\mathbf{s} \leftarrow \mathbf{s}_0$  ▷ Skip reposition
5:   else
6:      $\text{rep\_index} \leftarrow \arg\max_r \sum_{w_i \in \mathbf{s}} \log \pi_{\theta_L}^{rep}(r_i | w_i, \mathbf{s}, \mathbf{d})$  ▷ Do reposition
7:      $\mathbf{d} \leftarrow \mathcal{E}(\mathbf{s}, \text{rep\_index})$ 
8:   end if
9:    $\text{plh\_index} \leftarrow \arg\max_p \sum_{w_i, w_{i+1} \in \mathbf{s}} \log \pi_{\theta_L}^{ins}(p_i | w_i, w_{i+1}, \mathbf{s}, \mathbf{d})$  ▷ Insert placeholders
10:   $\mathbf{s} \leftarrow \mathcal{E}(\mathbf{s}, \text{plh\_index})$ 
11:   $\text{tok\_index} \leftarrow \arg\max_t \sum_{w_i \in \mathbf{s}, w_i = \langle \text{mask} \rangle} \log \pi_{\theta_L}^{tok}(t_i | w_i, \mathbf{s}, \mathbf{d})$  ▷ Fill placeholders
12:   $\mathbf{s} \leftarrow \mathcal{E}(\mathbf{s}, \text{tok\_index})$ 
13: end while
14:  $\mathbf{d} \leftarrow \text{documentUpdate}(\mathbf{d}, \mathbf{s})$ 
```

Algorithm 3 Training for Hierarchical Levenshtein Transformer

```
Require: Training data  $\mathcal{T}$ , Model policy:  $\pi_{\theta}$ , Expert policy:  $\pi_*$ 
1: while Maximum training steps reached do
2:    $(\mathbf{d}, \mathbf{d}_*) \sim \mathcal{T}$  ▷ Sample a training pair
3:    $\text{repH}^*, \text{insH}^*, \text{updH}^* \leftarrow \pi_{\theta}^H(\mathbf{d}, \mathbf{d}_*)$  ▷ Get oracle actions
4:    $\text{repL1}^*, \text{insL1}^*, \text{tokL1}^*, \text{repL2}^*, \text{insL2}^*, \text{tokL2}^* \leftarrow \pi_{\theta}^L(\mathbf{d}, \mathbf{d}_*)$ 
5:    $\mathcal{L}_{\theta_L}^{rep} \leftarrow - \sum_{s_i \in \mathbf{d}} \log \pi_{\theta_L}^{rep}(\text{repH}_i^* | \mathbf{s}_i, \mathbf{d})$ 
6:    $\mathbf{d} \leftarrow \text{applyAction}(\mathbf{d}, \text{repH}^*)$ 
7:    $\mathcal{L}_{\theta_L}^{ins} \leftarrow - \sum_{s_i, s_{i+1} \in \mathbf{d}} \log \pi_{\theta_L}^{ins}(\text{insH}_i^* | \mathbf{s}_i, \mathbf{s}_{i+1}, \mathbf{d})$ 
8:    $\mathbf{s} \leftarrow \text{buildFrame}(\text{insH}^*, \mathbf{d})$ 
9:    $\mathcal{L}_{\theta_L}^{rep1} \leftarrow - \sum_{w_i \in \mathbf{s}} \log \pi_{\theta_L}^{rep}(\text{repL1}_i^* | w_i, \mathbf{s}, \mathbf{d})$  ▷ Low Level
10:   $\mathbf{s} \leftarrow \text{applyAction}(\mathbf{s}, \text{repL1}^*)$ 
11:   $\mathcal{L}_{\theta_L}^{ins1} \leftarrow - \sum_{w_i, w_{i+1} \in \mathbf{s}} \log \pi_{\theta_L}^{ins}(\text{insL1}_i^* | w_i, w_{i+1}, \mathbf{s}, \mathbf{d})$ 
12:   $\mathbf{s} \leftarrow \text{applyAction}(\mathbf{s}, \text{insL1}^*)$ 
13:   $\mathcal{L}_{\theta_L}^{tok1} \leftarrow - \sum_{w_i \in \mathbf{s}, w_i = \langle \text{mask} \rangle} \log \pi_{\theta_L}^{tok}(\text{tokL1}_i^* | w_i, \mathbf{s}, \mathbf{d})$ 
14:   $\mathbf{d} \leftarrow \text{applyAction}(\mathbf{d}, \text{insH}^*)$ 
15:   $\mathcal{L}_{\theta_L}^{upd} \leftarrow - \sum_{s_i \in \mathbf{d}} \log \pi_{\theta_L}^{upd}(\text{updH}_i^* | \mathbf{s}_i, \mathbf{d})$ 
16:   $\mathbf{s} \leftarrow \text{buildFrame}(\text{updH}^*, \mathbf{d})$ 
17:   $\mathcal{L}_{\theta_L}^{rep2} \leftarrow - \sum_{w_i \in \mathbf{s}} \log \pi_{\theta_L}^{rep}(\text{repL2}_i^* | w_i, \mathbf{s}, \mathbf{d})$  ▷ Low Level
18:   $\mathbf{s} \leftarrow \text{applyAction}(\mathbf{s}, \text{repL2}^*)$ 
19:   $\mathcal{L}_{\theta_L}^{ins2} \leftarrow - \sum_{w_i, w_{i+1} \in \mathbf{s}} \log \pi_{\theta_L}^{ins}(\text{insL2}_i^* | w_i, w_{i+1}, \mathbf{s}, \mathbf{d})$ 
20:   $\mathbf{s} \leftarrow \text{applyAction}(\mathbf{s}, \text{insL2}^*)$ 
21:   $\mathcal{L}_{\theta_L}^{tok2} \leftarrow - \sum_{w_i \in \mathbf{s}, w_i = \langle \text{mask} \rangle} \log \pi_{\theta_L}^{tok}(\text{tokL2}_i^* | w_i, \mathbf{s}, \mathbf{d})$ 
22:   $\theta \leftarrow \theta - \lambda \nabla [\mathcal{L}_{\theta_H}^{rep} + \mathcal{L}_{\theta_H}^{ins} + \mathcal{L}_{\theta_H}^{upd} + \mathcal{L}_{\theta_L}^{rep1} + \mathcal{L}_{\theta_L}^{ins1} + \mathcal{L}_{\theta_L}^{tok1} + \mathcal{L}_{\theta_L}^{rep2} + \mathcal{L}_{\theta_L}^{ins2} + \mathcal{L}_{\theta_L}^{tok2}]$ 
23: end while
```

Results

	Multi-News			DUC-2004		
	R-1	R-2	R-L	R-1	R-2	R-L
Copy	42.32	13.28	37.86	36.30	8.47	32.52
Transformer	40.62	12.42	36.37	35.4	7.78	31.71
LevT	25.93	8.59	28.95	23.45	4.89	25.12
Editor	25.56	8.13	28.33	23.17	4.21	25.01
Ours	21.67	5.89	24.03	18.22	2.17	20.87

Table 1: Experiment Results on Multi-News and DUC2004 dataset

	Synthetic ROC-Stories	
	Copy	Transformer
Copy	23.59	28.82
Transformer	30.17	35.72
LevT	22.42	25.29
Editor	22.78	25.89
Ours	20.63	23.10

Table 2: Experiment Results on Synthetic and ROC-stories dataset. We report the BLEU score in the table.

The main results for summarization are shown in table1. The best result is obtained by copy across both dataset indicating that post editing of long sequences may hurt its quality. Copy consist of output from SummPip system. SummPip uses graph clustering to find relevant sentences which are then used to generate the summary. Among other models, the Vanilla transformer performed better showing a strong bias present in the languages for autoregressive monotone generation. Levenshtein and the Editor transformer performed comparably whereas as our model showed no improvement over the baselines. We see similar performance in Synthetic and ROC-stories dataset in table \ref{tab:results2} with Vanilla transformer performing better than the other models.

Conclusion and future directions

We present a hierarchical document generation model, that is capable of revising and editing its generated text thus bringing it closer to human-level intelligence. Although results showed that our approach lags behind the baselines, it did shed light into various problems present in semi-autoregressive models and long document generation. In the future, we will be incorporating

- Distilled Dataset
- Better training regimes
- Use of autoregressive model as low level programmer
- Attention mechanism

Acknowledgements

The Authors would like to acknowledge reviewers for their constructive comments.