

# Constrained Beam Search

## 1 Idea

We first obtain the list of all english tokens by decoding the tokens into string and checking the ASCII values of the characters. We maintain two separate beams of  $b$  length. The first beam is like the normal beam in vanilla Beam search Algorithm call this beam "**vanilla\_beam**" while in the other beam all the sequence have the invariant that they contain atleast one english token we call this beam "**english\_beam**". In every update step, we compute 3 list of sequences,  $L1$ ,  $L2$ ,  $L3$ .

- **L1**: The sequences obtained by adding the top  $(b+1)$  tokens into each vanilla sequence.
- **L2**: The sequences obtained by adding the top  $(b+1)$  **english** tokens into each vanilla sequence.
- **L3**: The sequences obtained by adding the top  $(b+1)$  tokens into each english sequence.

We update the beams as:

- **english\_beam**: Best  $b$  sequences in  $(L2 \cup L3)$ .
- **vanilla\_beam**: Best  $b$  sequences in  $L1$ .

Whenever a  $< eot >$  gets appended into a sequence of the **english\_beam** we add it to the finished list. Also in the finalize function we only consider the english beams.

## 2 Algorithm

### 2.1 Update

---

**Algorithm 1** Update Function for Constrained Beam Search

---

```
1: Input: tokens, logits, sum_logprobs
2: // Initialization
3: Initialize lists L1, L2, L3 as empty
4: for ray in vanilla_beam do
5:   l1  $\leftarrow$  best  $k + 1$  tokens for the ray
6:   for each token in l1 add ray.concat(token) to L1
7:   l2  $\leftarrow$  best  $k + 1$  English tokens for the ray
8:   for each token in l2 add ray.concat(token) to L2
9: end for
10: for ray in english_beam do
11:   l3  $\leftarrow$  best  $k + 1$  tokens for the ray
12:   for each token in l3 add ray.concat(token) to L3
13: end for
14: Leng  $\leftarrow$  top  $k + 1$  sequences of (L2  $\cup$  L3)
15: Lall  $\leftarrow$  top  $k + 1$  sequences of L1
16: if any ray in Lall is completed (token is eot) then
17:   pop the ray from Lall
18: end if
19: if any ray in Leng is completed (token is eot) then
20:   pop the ray from Leng and add to finished_sequences
21: end if
22: Set vanilla_beam to top  $k$  elements from Lall
23: Set english_beam to top  $k$  elements from Leng
24: if size of finished is max_candidates then
25:   completed  $\leftarrow$  true
26: else
27:   completed  $\leftarrow$  false
28: end if
29: return vanilla_beam, english_beam, completed
```

---