

## BACKGROUND: BEAM SEARCH ALGORITHM

While decoding with a sequence-to-sequence model, we could not afford to search globally for optimal output sequence, so researchers often resort to beam search algorithm to approximate exact search. The beam search algorithm expands  $B_{t-1}$  to  $B_t$  as follows:

$$B_0 = [\langle s \rangle, p(\langle s \rangle | \mathbf{x})]$$

$$B_t = \text{top}^b \{ \langle y' \circ y_t, s \cdot p(y_t | \mathbf{x}, y) \rangle \mid \langle y', s \rangle \in B_{t-1} \}$$

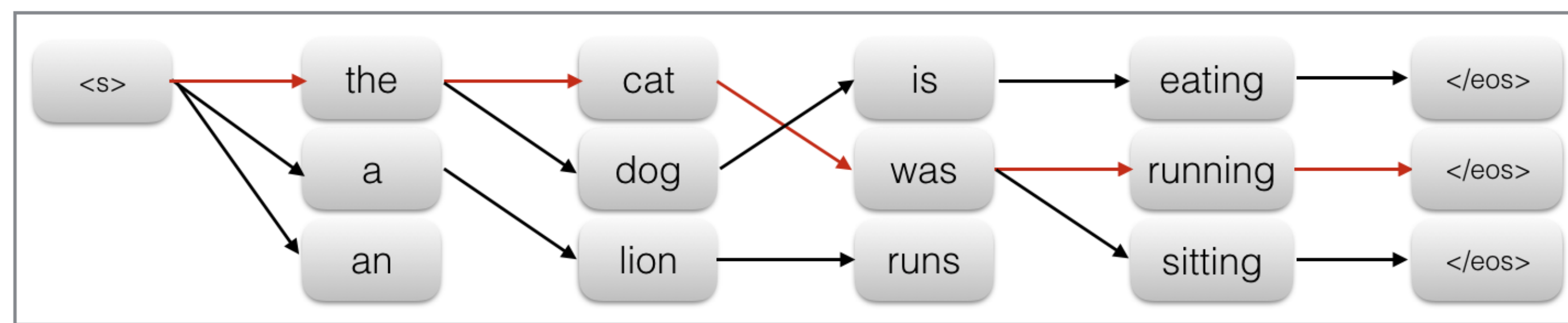


Figure 1: Examples of beam search algorithm with beam size 3. Red arrows denote greedy search (beam size 1).

In the end, the algorithm chooses the candidate with highest log-probability:

$$\mathbf{y}^* = \underset{y: \text{comp}(y)}{\text{argmax}} sc(\mathbf{x}, \mathbf{y}) = \underset{y: \text{comp}(y)}{\text{argmax}} \sum_{t \leq |y|} \log p(y_t | \mathbf{x}, y_{<t})$$

where  $\text{comp}(y) \triangleq \{y | y| = \langle /eos \rangle\}$  returns the completeness of a hypothesis.

## BEAM SEARCH CURSE

It's widely observed that as beam size increases after 5, the performance of sequence-to-sequence models, as quantified by the BLEU score, drops greatly. Since the models could not leverage the computational power from wider beams, we call this phenomenon the *Beam Search Curse*.

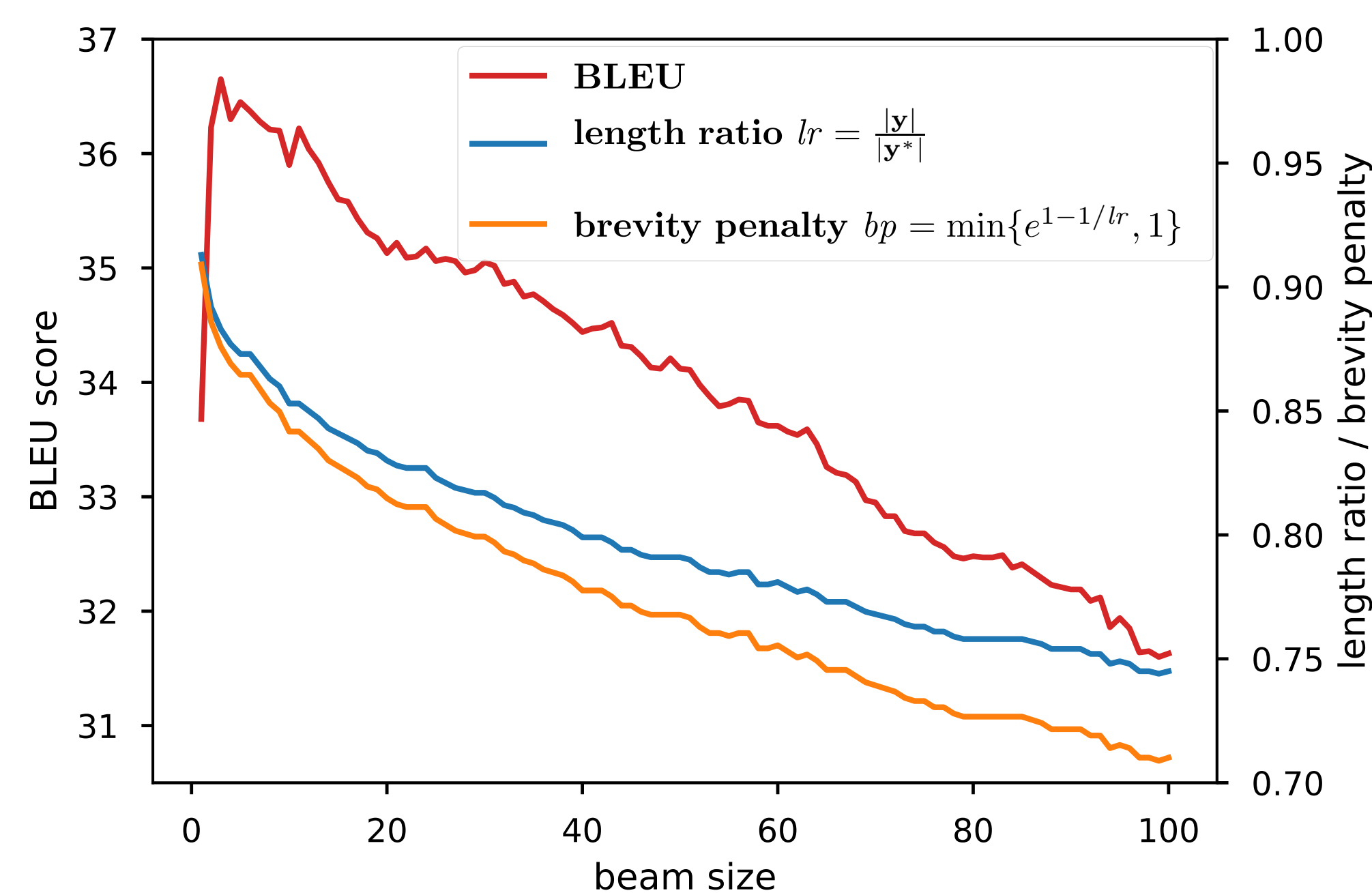


Figure 2: While the BLEU score drops with an increasing beam size (after 5), the brevity penalty drops with a similar curve.

## EXPERIMENTAL SETUP

- Based on OpenNMT-py, a PyTorch reimplementation of Torch-based OpenNMT (Klein et al., 2017).
- 2M Chinese-English sentence pairs for training.
- Used byte-pair encoding (BPE) (Senrich et al., 2015) to reduce vocabulary sizes down to 18k/10k respectively.
- Chinese to English: NIST 06 newswire portion (616 sentences) for dev; NIST 08 newswire portion (691 sentences) for test.

## WHY THE CURSE EXISTS

- As beam size increases, the more candidates it would explore. Therefore, it becomes easier to find the  $\langle /eos \rangle$  symbol and terminate. Left figure shows that the  $\langle /eos \rangle$  indices decrease steadily with wider beams.
- Then, because of the internal property of log-probability, shorter candidates have clear advantages *w.r.t.* model score.

As a conclusion, the search algorithm would find shorter candidates, and prefer even shorter ones among them.

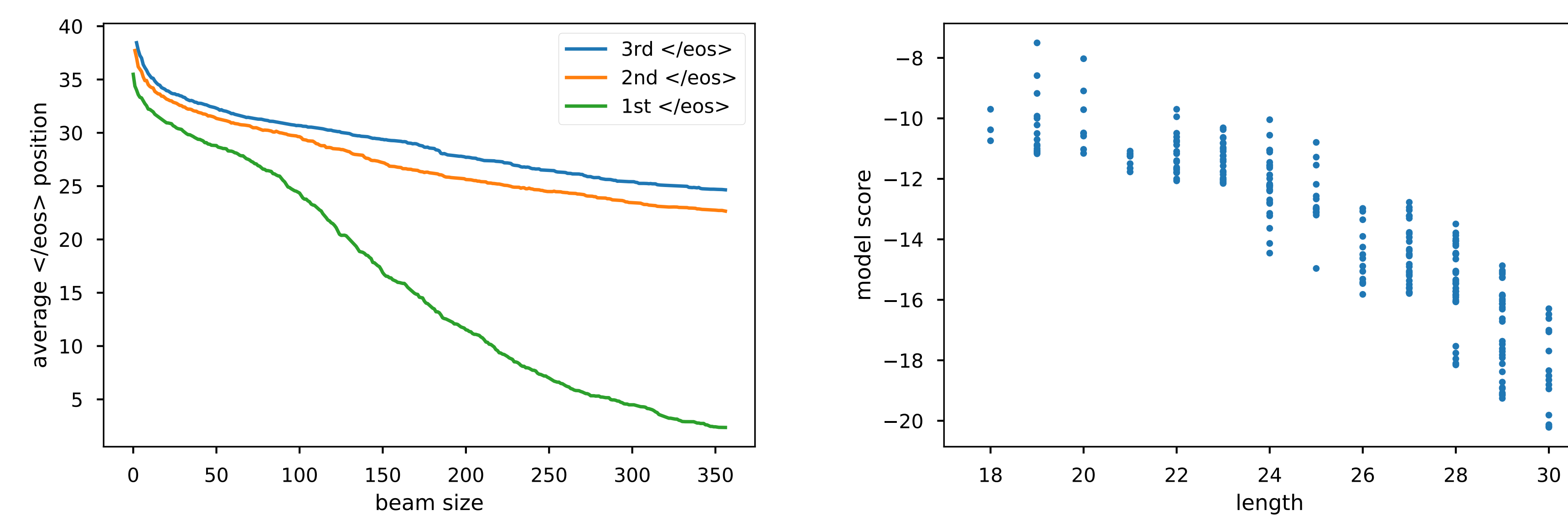


Figure 3: Left: Searching algorithm with wider beams generates  $\langle /eos \rangle$  earlier. Right: The model score (log-probability) strongly prefers shorter candidates.

## HOW TO BREAK THE CURSE

### Previous Methods

- Length Normalization (Bahdanau et al., 2014): normalize the score by its length.
- Word-Reward (He et al., 2016): add reward  $r$  to each word.
- Bounded Word-Reward (Liang et al., 2017): add reward  $r$  to each word up to a bound.

**Rescoring with Length Prediction** We use a 2-layer MLP, which takes the mean of source hidden states as input, to predict the generation ratio  $gr(\mathbf{x})$ . Then we can get our predicted length  $L_{pred}(\mathbf{x}) = gr(\mathbf{x}) \cdot |\mathbf{x}|$ .

**Bounded Word-Reward w/ Predicted Length** To favor longer generation, we add rewards  $r$  to each word up to its predicted length.

$$L(\mathbf{x}, \mathbf{y}) = \min\{|\mathbf{y}|, L_{pred}(\mathbf{x})\} \quad \hat{sc}(\mathbf{x}, \mathbf{y}) = sc(\mathbf{x}, \mathbf{y}) + r \cdot L(\mathbf{x}, \mathbf{y})$$

where  $sc(\mathbf{x}, \mathbf{y})$  is the original model score (log-probability).

**Bounded Adaptive-Reward** Instead of a tuned reward  $r$ , we add an adaptive reward to each step based off local beam information. With beam size  $b$ , the reward for time step  $t$  is the average negative log-probability of the words in the current beam.

$$r_t = -(1/b) \sum_{i=1}^b \log p(\text{word}_i) \quad \hat{sc}(\mathbf{x}, \mathbf{y}) = sc(\mathbf{x}, \mathbf{y}) + \sum_{t=1}^{L(\mathbf{x}, \mathbf{y})} r_t$$

**BP-Norm** Instead of adding rewards, we apply brevity penalty to the length-normalized model score.

$$bp = \min\{e^{1-1/lr}, 1\} \quad \hat{sc}(\mathbf{x}, \mathbf{y}) = \log bp + sc(\mathbf{x}, \mathbf{y})/|\mathbf{y}|$$

## DISCUSSION

Among all methods, we recommend **BP-Norm** for the following reasons:

- BP-Norm works equally well with others, while doesn't contain any hyper-parameters.
- BP-Norm is intuitive and **in the same form as BLEU**. Both of their exponential forms are products of brevity penalty term and geometric mean of probabilities (BP-Norm) or accuracies (BLEU).

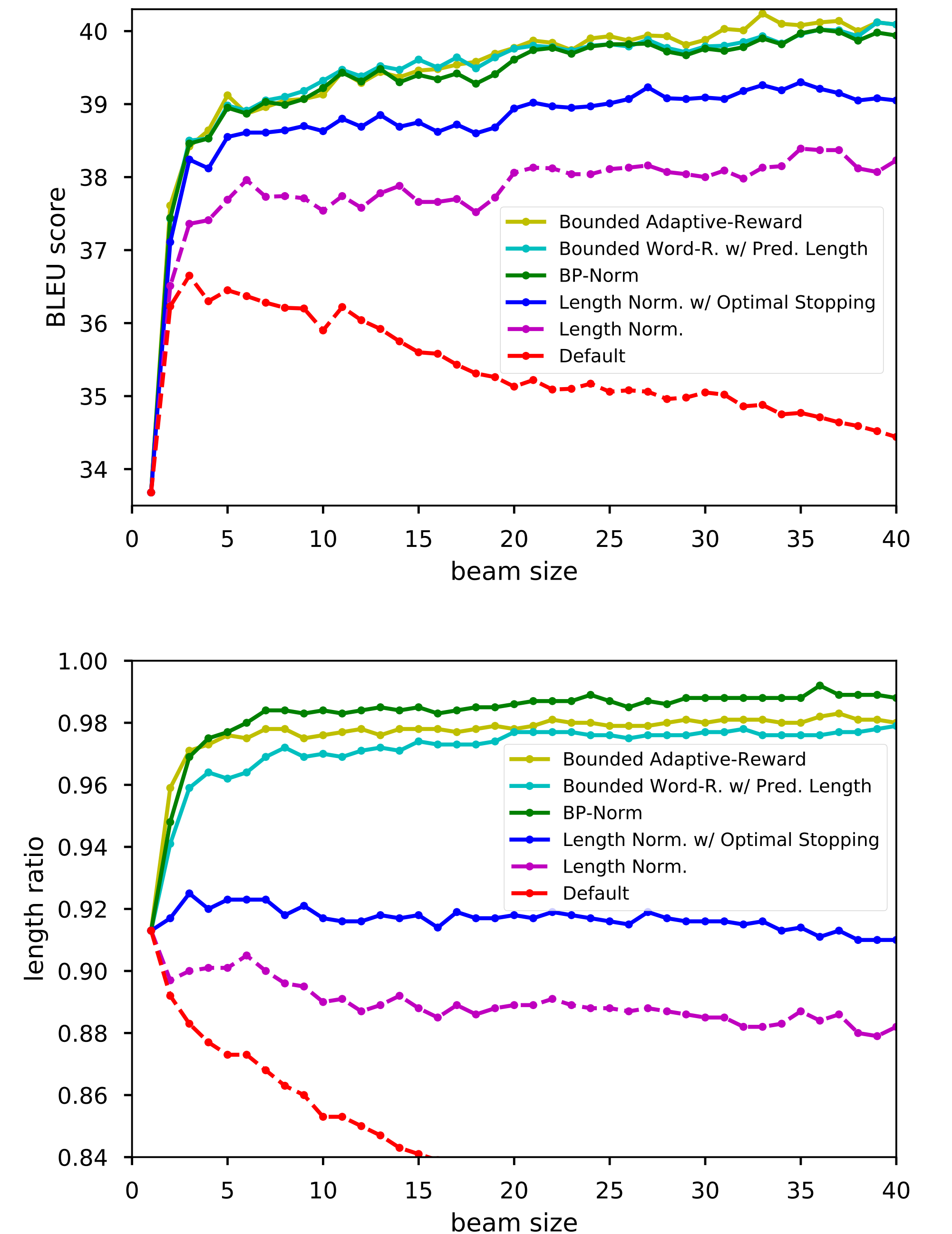


Figure 4: BLEU and length ratios of various rescoring methods.

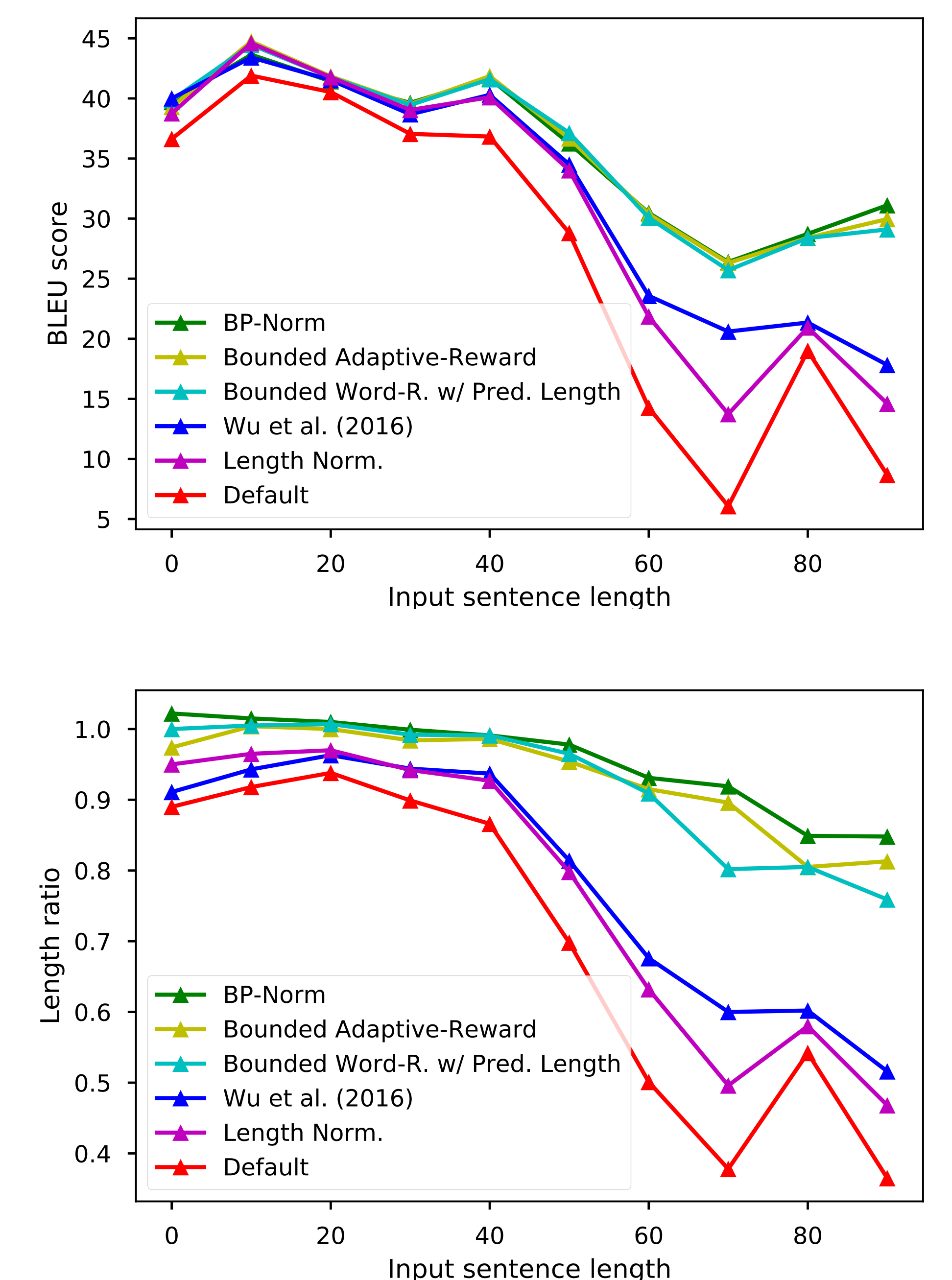


Figure 5: BLEU and length ratios over various input sentence lengths.