

291K

Deep Learning for Machine Translation Decoding

Lei Li

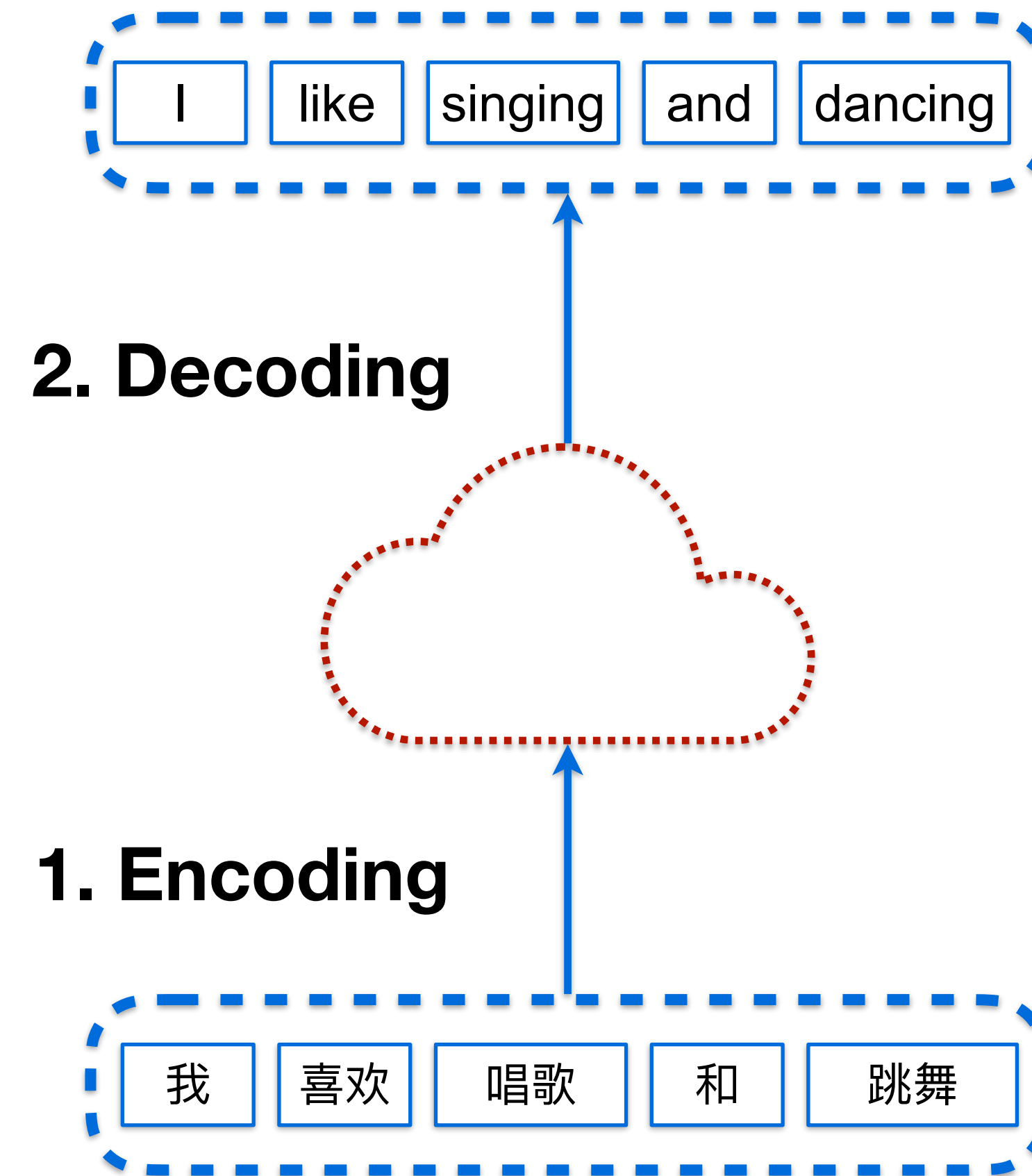
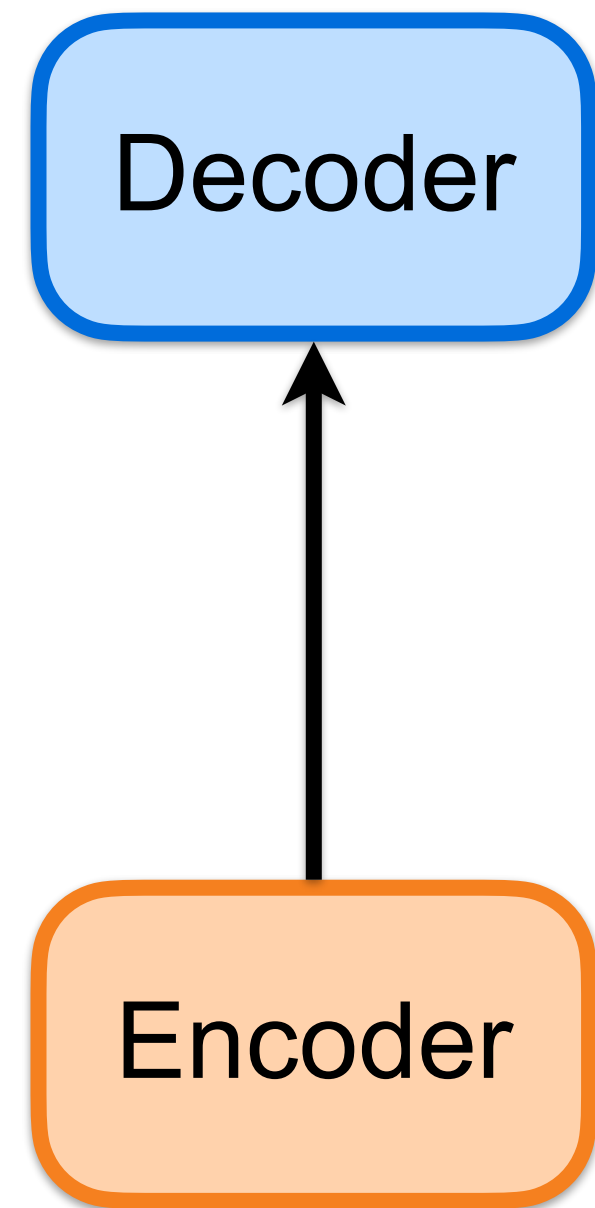
UCSB

10/18/2021

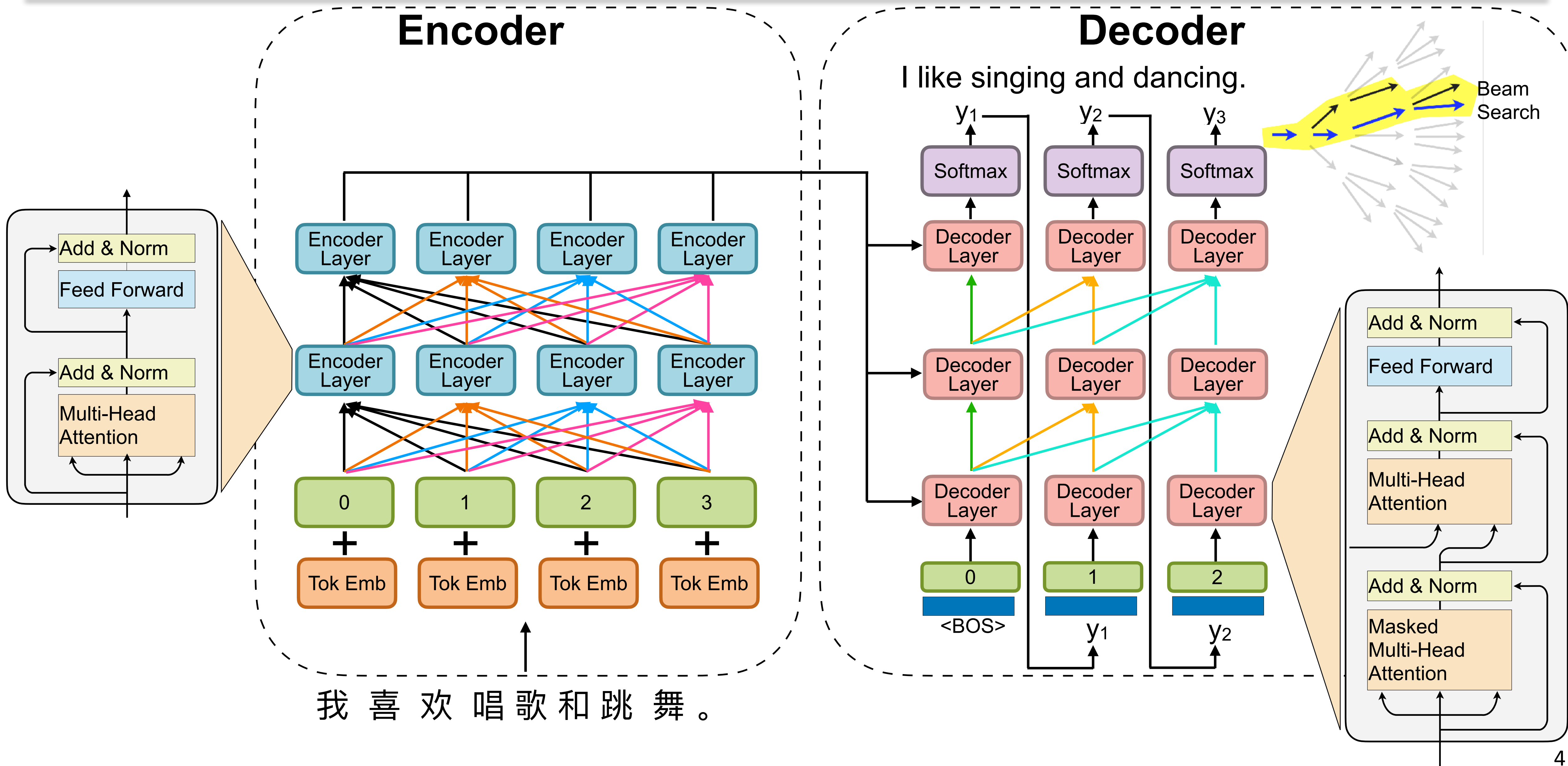
Outline

- Beam Search
- Diverse Beam Search
- Reranking
- Sampling
- Constrained decoding
- Model Average
- Model Ensemble
- Minimum Bayes Risk Decoding

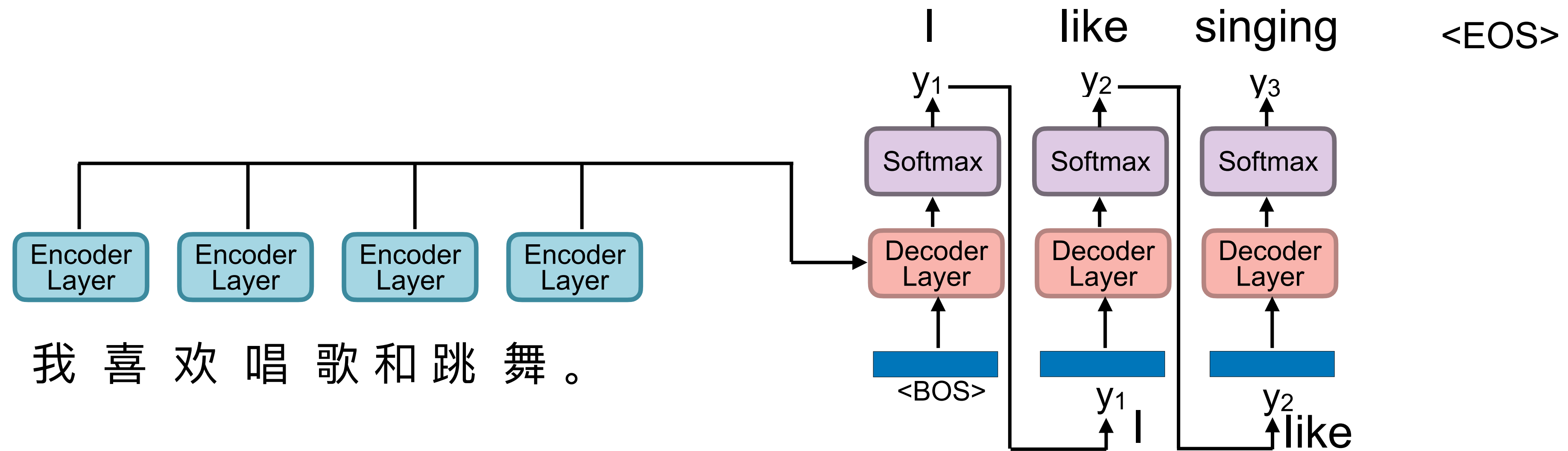
Encoder-Decoder Paradigm



Transformer



Autoregressive Generation



But, this is not necessary the best

Inference

- Now already trained a model θ
- Decoding/Generation: Given an input sentence x , to generate the target sentence y that maximize the probability $P(y | x; \theta)$
- $\operatorname{argmax}_y P(y | x) = f_{\theta}(x, y)$
- Two types of error
 - the most probable translation is bad \rightarrow fix the model
 - search does not find the most probably translation \rightarrow fix the search
- Most probable translation is not necessary the highest BLEU one!

Decoding

- $\operatorname{argmax}_y P(y | x) = f_{\theta}(x, y)$
- naive solution: exhaustive search
 - too expensive
- Beam search
 - (approximate) dynamic programming

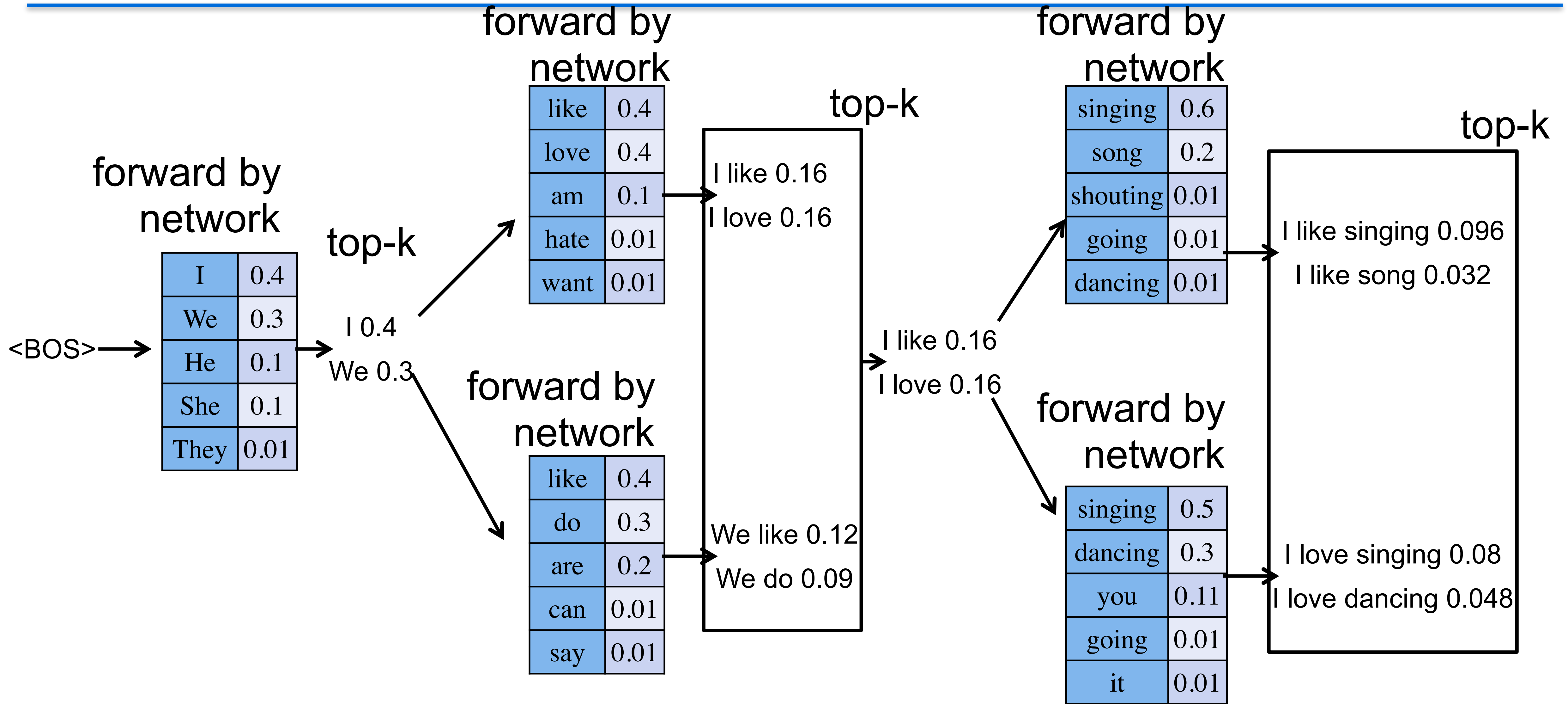
Beam Search

- start with empty S
- at each step, keep k best partial sequences
- expand them with one more forward generation
- collect new partial results and keep top- k

Beam Search (pseudocode)

```
best_scores = []
add {[0], 0.0} to best_scores # 0 is for beginning of sentence token
for i in 1 to max_length:
    new_seqs = PriorityQueue()
    for (candidate, s) in best_scores:
        if candidate[-1] is EOS:
            prob = all -inf
            prob[EOS] = 0
        else:
            prob = using model to take candidate and compute next token probabilities (logp)
    pick top k scores from prob, and their index
    for each score, index in the top-k of prob:
        new_candidate = candidate.append(index)
        new_score = s + score
        if not new_seqs.full():
            add (new_candidate, new_score) to new_seqs
        else:
            if new_seqs.queue[0][1] < new_score:
                new_seqs.get() # pop the one with lowest score
                add (new_candidate, new_score) to new_seqs
```

Beam Search

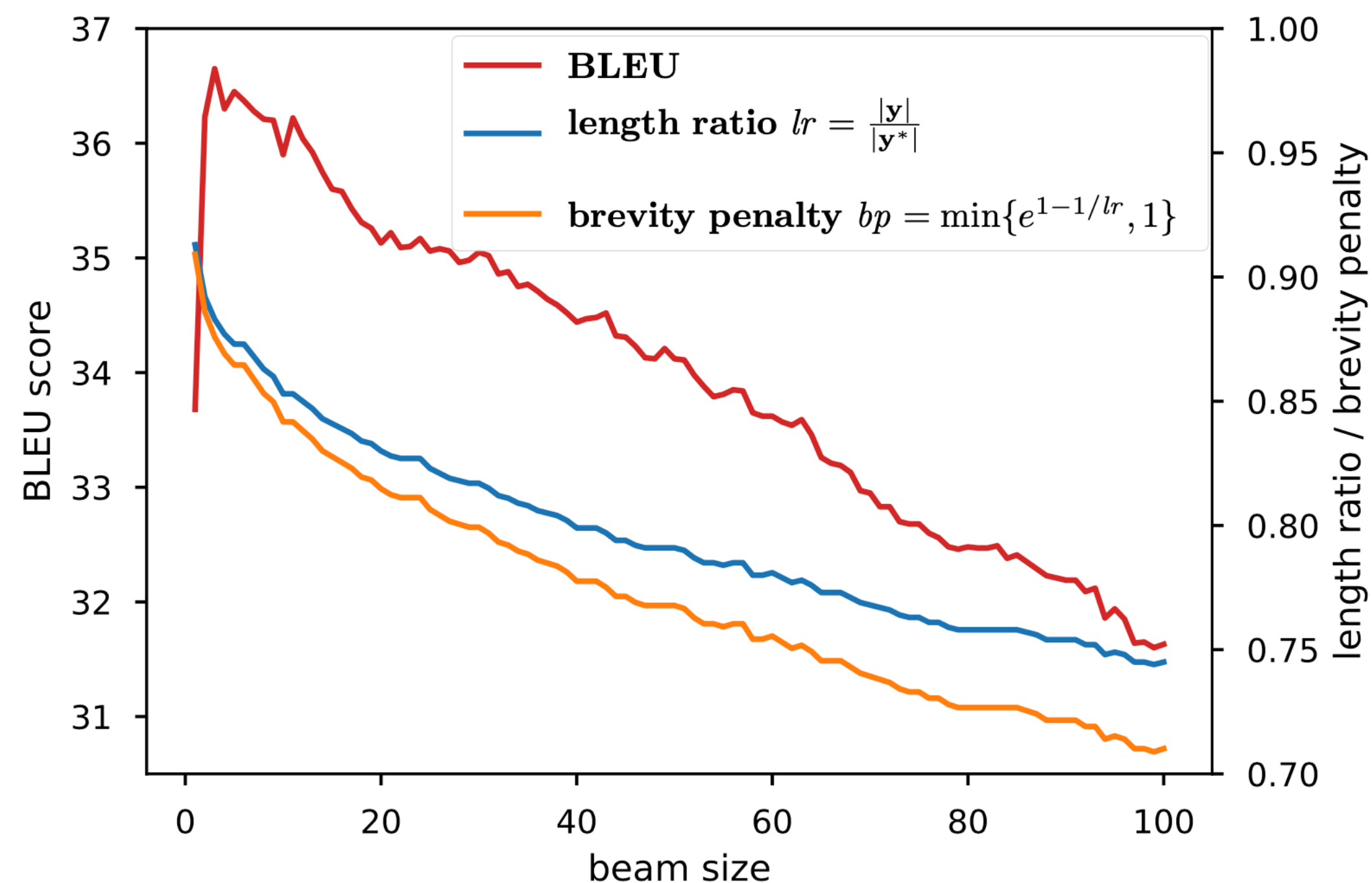


Pruning for Beam Search

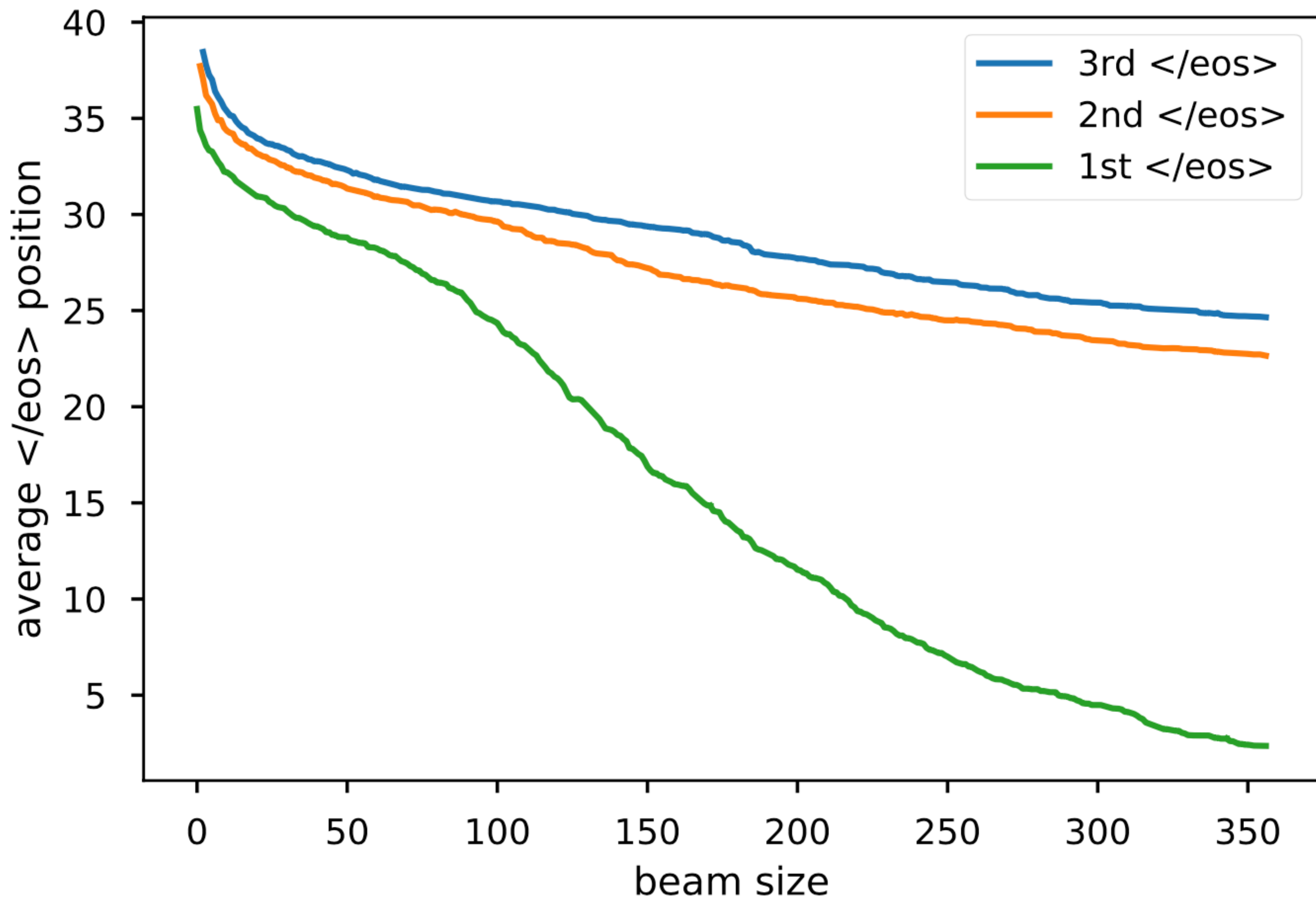
- Relative threshold pruning
 - prune candidates with too low score from the top one
 - Given a pruning threshold rp and an active candidate list C , a candidate $cand \in C$ is discarded if: $score(cand) \leq rp * \max\{score(c)\}$
- Absolute threshold pruning:
 - $score(cand) \leq \max\{score(c)\} - ap$
- Relative local threshold pruning

What is Beam size?

- 3 to 5
- Why not larger?
 - larger does not necessarily produce higher BLEU



Larger Beam -> Shorter Translation



Normalization of Score

- Length normalization: $\hat{S}_{\text{length_norm}}(\mathbf{x}, \mathbf{y}) = S(\mathbf{x}, \mathbf{y})/|\mathbf{y}|$
- Word-reward: promoting longer sentences

— $\hat{S}_{\text{WR}}(\mathbf{x}, \mathbf{y}) = S(\mathbf{x}, \mathbf{y}) + r \cdot |\mathbf{y}|$

- Bounded word reward with length prediction

— $L_{\text{pred}}(\mathbf{x}) = gr^*(\mathbf{x}) \cdot |\mathbf{x}|$

$$L^*(\mathbf{x}, \mathbf{y}) = \min\{|\mathbf{y}|, L_{\text{pred}}(\mathbf{x})\}$$

— $\hat{S}_{\text{BWR}^*}(\mathbf{x}, \mathbf{y}) = S(\mathbf{x}, \mathbf{y}) + r \cdot L^*(\mathbf{x}, \mathbf{y})$

Diverse Beam Search

- Top k results from NMT decoding are very similar
- Same for other text generation tasks
- Need more diversity?
 - e.g. in image-captioning, diverse candidates are desired



Beam Search	Diverse Beam Search
<p>A tree diagram for Beam Search. The root node is 'A'. The first level is 'table that has a'. The second level is 'bunch of'. From 'bunch of', there are five branches: 'bowls on it and some bottles', 'food on it', 'bottles on it', 'plates on it', and 'flowers on it'.</p>	<p>A tree diagram for Diverse Beam Search. The root node is 'A'. The first level is 'table with a vase of flowers', 'table that has some', and 'an empty kitchen table with a'. From 'table with a vase of flowers', there are two branches: 'on it' and 'and a window'. From 'table that has some', there are two branches: 'food on it' and 'pots and pans on it'. From 'an empty kitchen table with a', there are two branches: 'vase of flowers' and 'bowl of fruit on it'.</p>
<p>A table that has a bunch of bowls on it and some bottles A table that has a bunch of bowls on it A table that has a bunch of food on it A table that has a bunch of bottles on it A table that has a bunch of plates on it A table that has a bunch of flowers on it</p>	<p>A table with a vase of flowers on it A table with a vase of flowers on it and a window A table that has some food on it A table that has some pots and pans on it An empty kitchen table with a vase of flowers An empty kitchen table with a bowl of fruit on it</p>

How

- Two approaches
 - MMI: maximizing mutual information of $MI(X, Y)$ instead of $P(Y|X)$
 - Maximize the penalized score: $\log P(Y|X) + \text{distance}(Y \text{ and existing candidates})$

Maximize mutual information (MMI)

- Mutual Information

$$\text{MI}(X, Y) = \frac{p(X, Y)}{p(X)p(Y)}$$

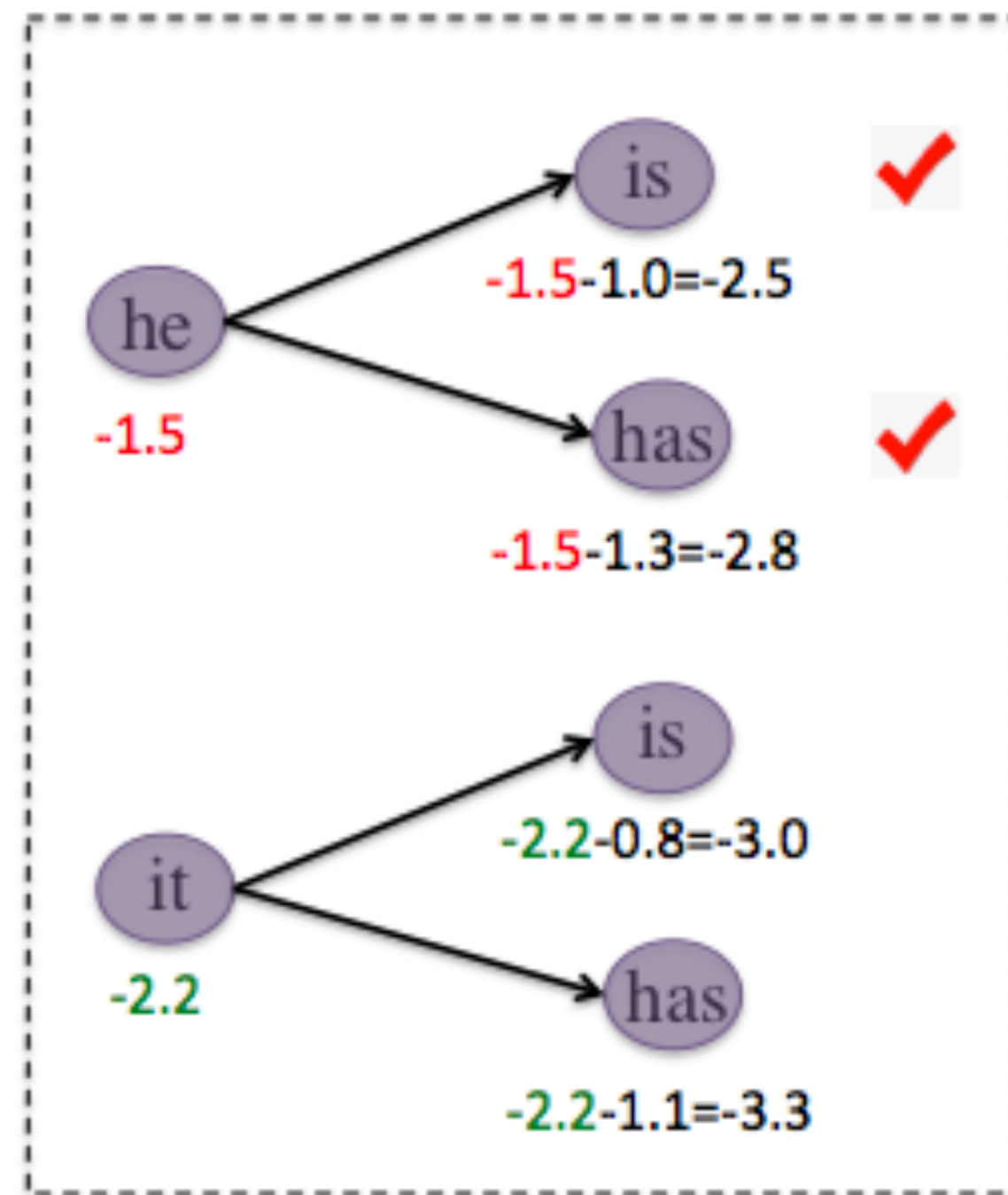
- $\arg \max \log p(Y|X) - \lambda \log p(Y)$
 - need a separate Language model $p(Y)$ for target language

Maximizing Mutual Information

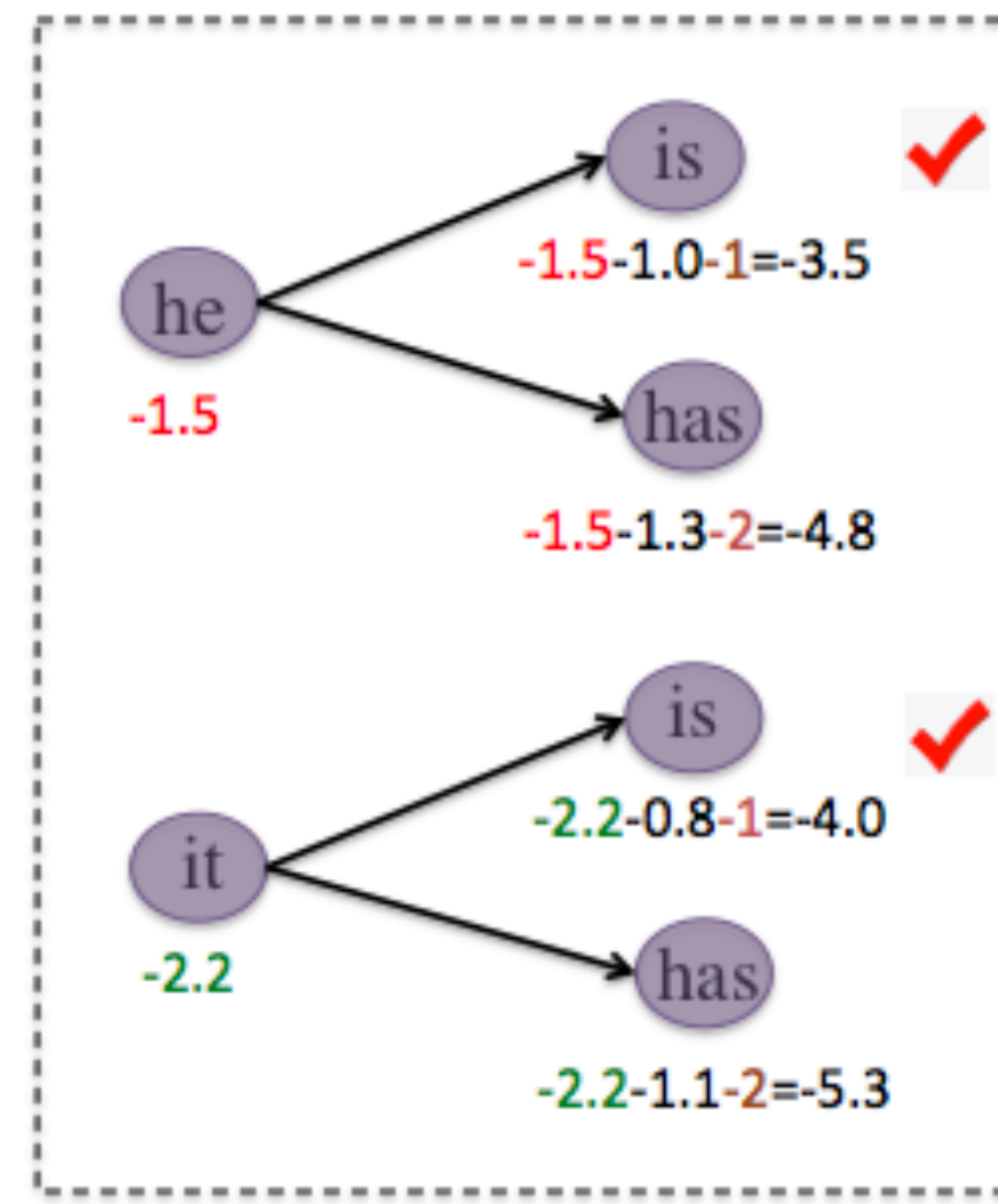
- $\arg \max (1 - \lambda) \log p(Y|X) + \lambda \log p(X|Y)$
- penalized forward decoding

– $p(Y|X) - \gamma \text{rank}_y$

$$\hat{S}(Y_{t-1}^k, y_t^{k,k'} | x) = S(Y_{t-1}^k, y_t^{k,k'} | x) - \gamma k'$$



Standard Beam Search



Diversity Promoting Beam Search (γ set to 1)

Reranking

- Obtain N-best from beam search
- Rerank based on:

$$\text{Score}(y) = \log p(y|x) + \lambda \log p(x|y) + \gamma \log p(y) + \eta LT$$

- **Alternative: learned reranking**

- Lee et al. Discriminative Reranking for Neural Machine Translation. 2021

Sampling

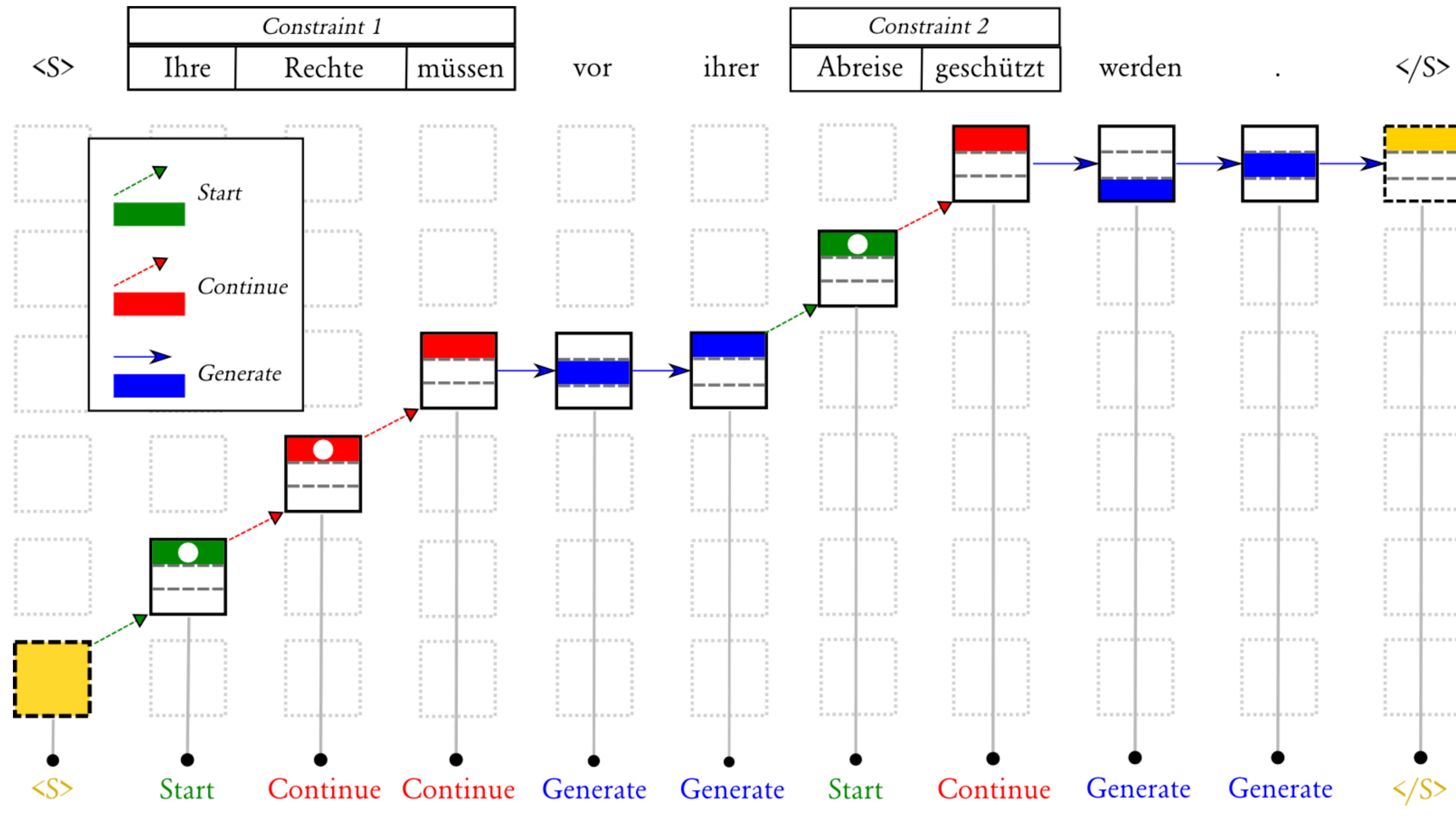
- Instead of $\operatorname{argmax}_y P(y | x) = f_{\theta}(x, y)$
- Generate samples of translation Y from the distribution $P(Y|X)$
- Q: how to generate samples from a discrete distribution?

Combine Sample and Beam Search

- Sample the first tokens
- continue beam search for the later

Lexical Constrained Decoding

- The generated sentence must contain given keywords
- To generate from
 - Vocabulary
 - Keywords



Input: Rights protection should begin before their departure .

Order-agnostic Constraints

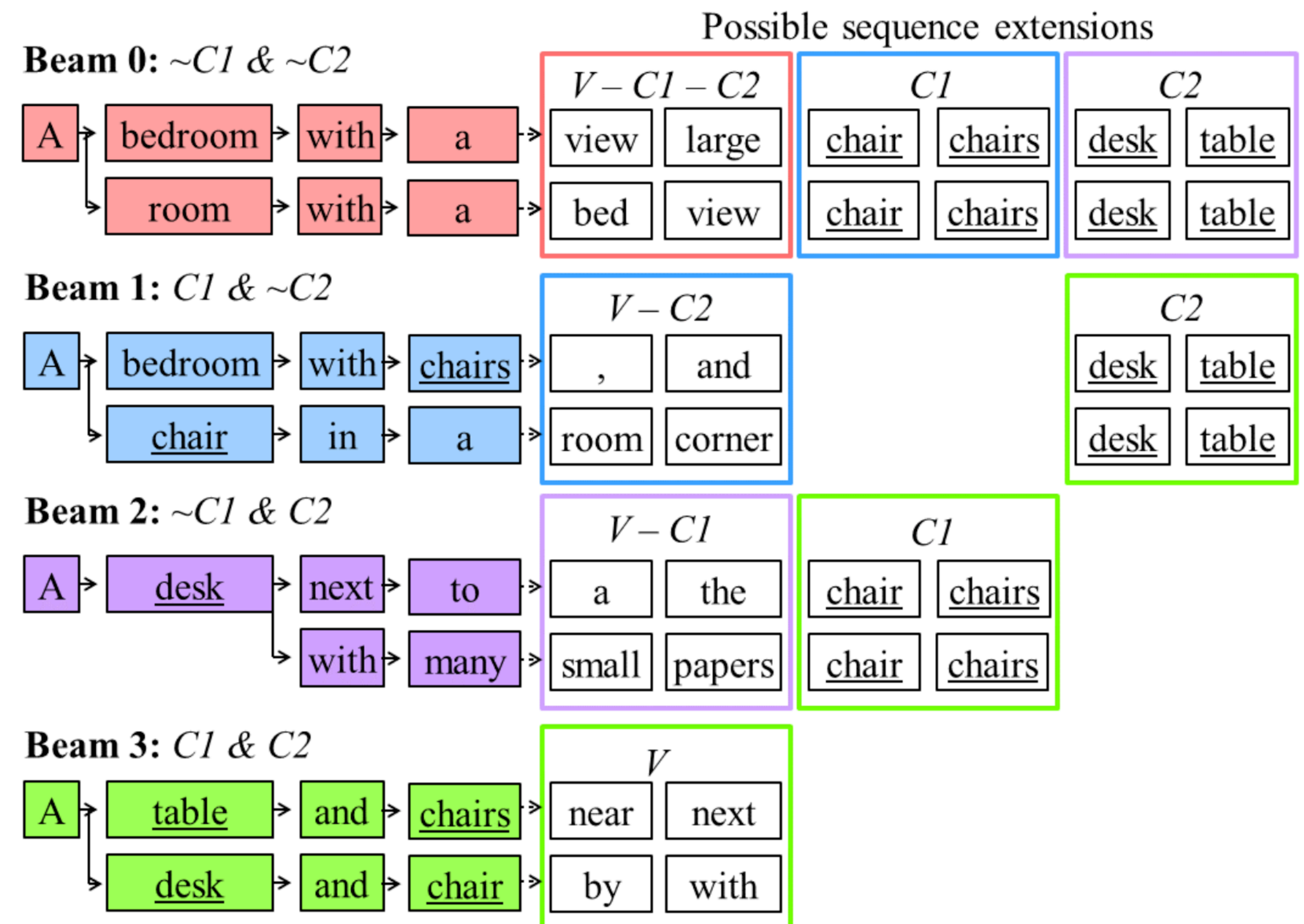
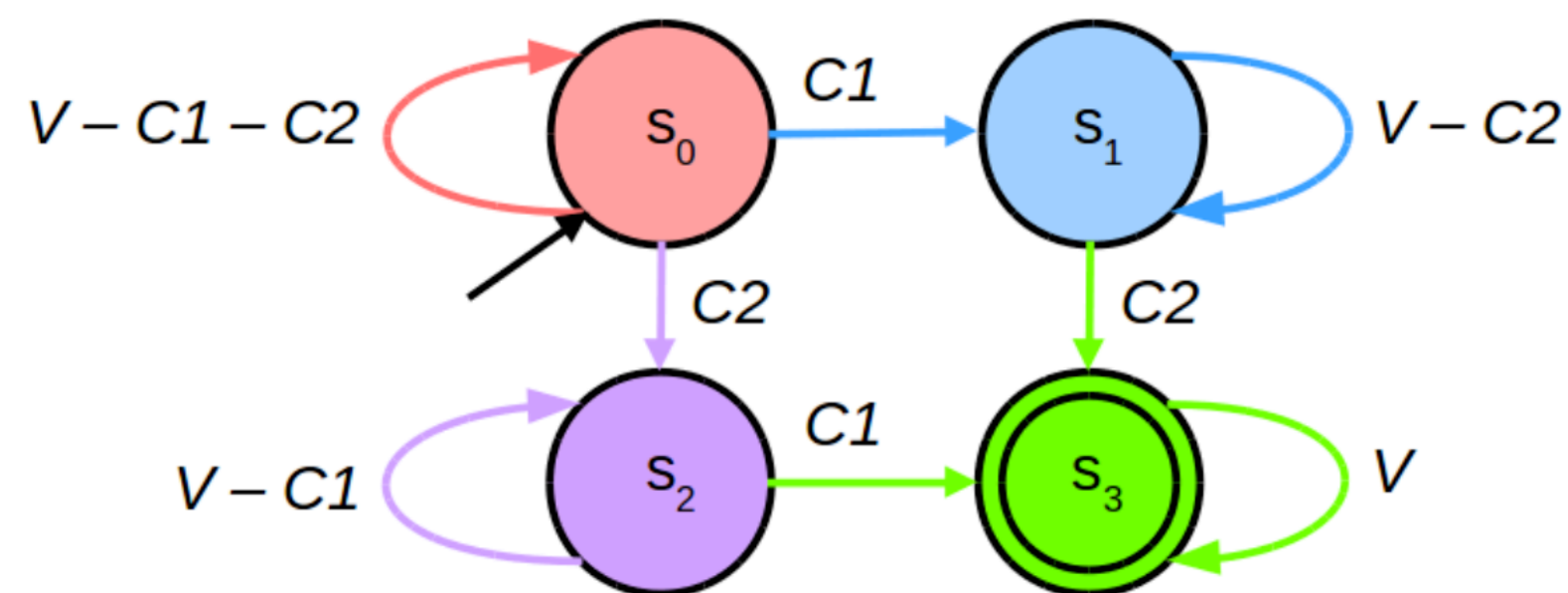
- The generated sentence must contain given keywords
- Using finite state machine to represent constraint state.

- Expand with

- Vocabulary
- Constraint keywords

Finite-state machine

$C1 = \{chair, chairs\}, C2 = \{desk, table\}$



Post-training Processing: Model Average

- Pick the model when converges
- Model average:
 - instead, using the last 5-10 epoch's models, and average the parameters to get one model
 - This turns out to generalize better than the last one.
 - Why? (over-fit)

Model Ensemble

- Train several separate MT model
- decode with

$$\arg \max_{y_t} \sum_k \log P(y_t | y_{<t}, x; M_k)$$

Distillation with Ensemble

- In order to obtain a single model with good performance.
- Use ensemble model to create pseudo-parallel data
- Train a single MT model using both original training data and pseudo-parallel data.

Minimum Bayes Risk Decoding

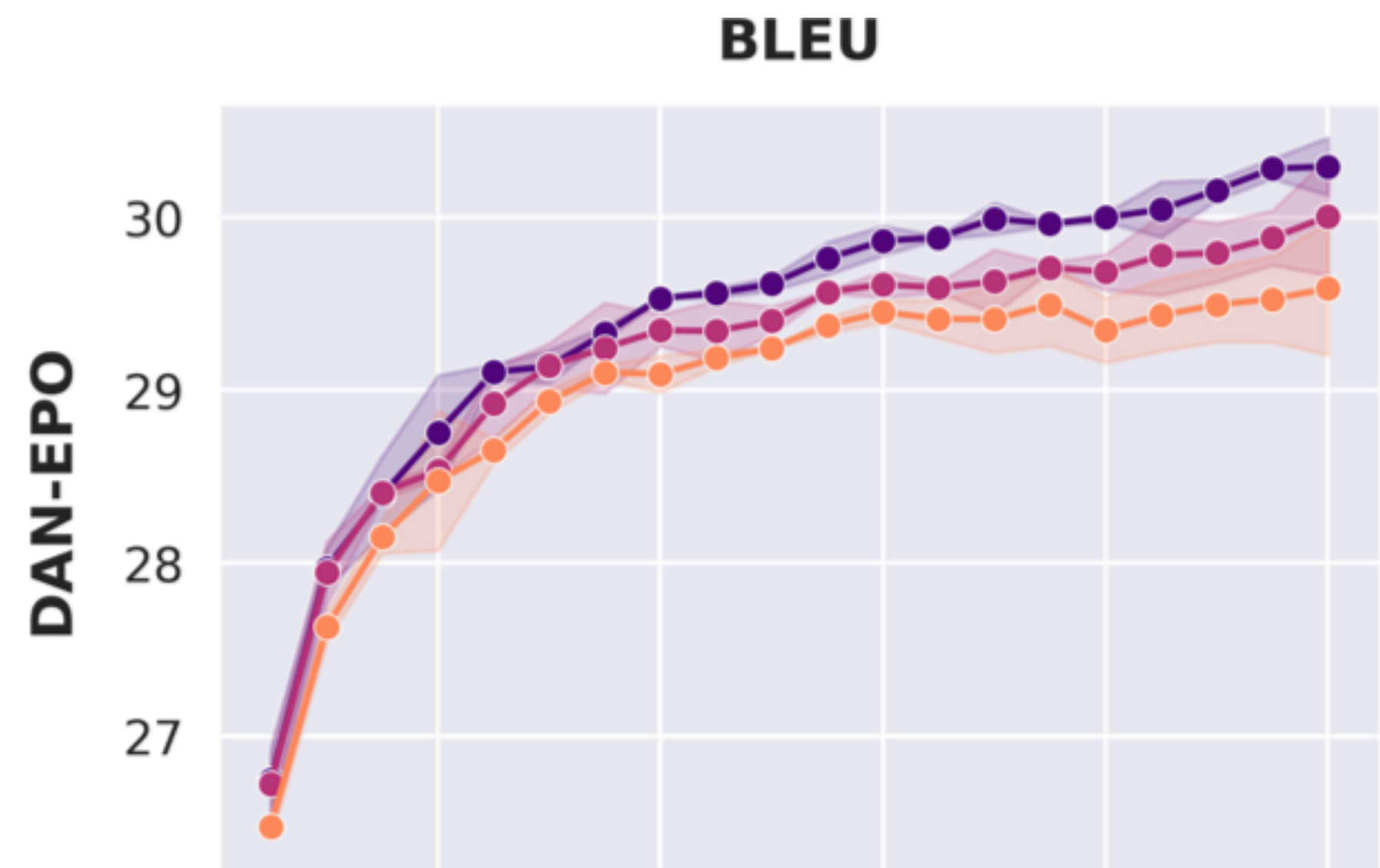
- Bias in decoding:
 - length bias
 - word frequency
 - beam search curse
 - copy noise
 - low domain
- Decoding with Mode vs. with most “common” one

Minimum Bayes Risk Decoding

- Minimize risk = maximize average utility
- Utility: similarity among samples.

- $S_1, S_2, \dots, S_n \sim P(y | x, \theta)$

- $\hat{y} \arg \max_{s_i} \frac{1}{n} \sum_j u(s_i, s_j)$



Language Presentation

Reading

- Freitag & Al-Onaizan. Beam Search Strategies for Neural Machine Translation. 2017.
- Muller and Sennrich. Understanding the Properties of Minimum Bayes Risk Decoding in Neural Machine Translation. 2021.