

Amortized Noisy Channel Neural Machine Translation



Richard Yuanzhe Pang¹
¹ New York University

He He¹
² Genentech

Kyunghyun Cho^{1,2,3}
³ CIFAR Fellow

1 Background and Goal

Naïve decoding based on the forward translator

Training: train p_f using (\mathbf{X}, \mathbf{Y})

Inference: greedy decoding or beam search with small beam size (e.g., $b=5$)

One way of noisy channel decoding: **beam search and rerank (BSR)**

Training: train p_f and p_r using (\mathbf{X}, \mathbf{Y})

p_f : forward translator

models $p(\text{target-lang sentence} \mid \text{source-lang sentence})$

p_r : reverse translator

models $p(\text{source-lang sentence} \mid \text{target-lang sentence})$

Inference: For each source sentence \mathbf{x} , (1) do beam search with beam size 50–100 (**SLOW!**); (2) rerank using the following objective and pick the top-ranked translation

$$\log p_f(\mathbf{y} \mid \mathbf{x}) + \gamma \log p_r(\mathbf{x} \mid \mathbf{y}) + \gamma' \log p_{\text{lm}}(\mathbf{y})$$

Used in many top/winning models in WMT competitions

Can we train a **new network** such that if we do **greedy decoding** using the new network, the translations will maximize $R(\mathbf{y}) = \log p_f(\mathbf{y} \mid \mathbf{x}) + \gamma \log p_r(\mathbf{x} \mid \mathbf{y})$?

Criteria for successful amortization

Inference speed

Successful if the inference is faster than BSR. Guaranteed!

Translation reward

Successful if the forward rewards of the generated sentences are comparable to the forward rewards by BSR, and the reverse rewards are comparable to the reverse rewards by BSR.

Translation quality (approximated by BLEU/BLEURT)

Successful if the BLEURT of our translations are similar to the BLEURT by BSR.

3 Discussion

- “BSR → high BLEURT” doesn’t imply “higher reward → higher BLEURT”
- KD/IL-generated translations are similar (in terms of corpus-level BLEU); they are different from Q-generated translations, possibly due to how reverse reward is presented to KD/IL vs. Q
- Q learning also applies to text generation (we trained Q from scratch!) – rarely used in NLG; but Q learning doesn’t do well when the source sentence is long (> 80 tokens) possibly due to the optimization difficulty given by the sparse reverse reward

2 Methods and Results

Approach 1: knowledge distillation (KD)

Step 1: train p_f using (\mathbf{X}, \mathbf{Y})

Step 2: generate pseudo-corpus Y_{pseudo} by BSR

Step 3: train p_{KD} using $(\mathbf{X}, Y_{\text{pseudo}})$

Effectively minimizing the KL-div between the distribution induced by the pseudo-corpus obtained through BSR and our model distribution

Approach 2: one-step deviation imitation learning (IM)

Call our new network A . To train A , intuitively:

Use cross entropy to...

- match the t -th step distribution of A and the t -th step distribution of p_f
- match $\text{onehot}(\mathbf{x}_t)$ and the t -th step distribution of p_r (p_r is a function of A)

Approach 3: Q learning *adapted from DQN used to train Atari games*

Want: Q (“future return” – higher is better);

Define: $s_t = (\mathbf{y}_{<t}, \mathbf{x})$, $a_t = y_t$

$$r_t = \log p_f(y_t \mid \mathbf{y}_{<t}, \mathbf{x}), \text{ if } t < T$$

$$= \log p_f(y_T \mid \mathbf{y}_{<T}, \mathbf{x}) + \gamma \log p_r(\mathbf{x} \mid \mathbf{y}), \text{ if } t = T$$

Given p_f, p_r , translation dataset D . Initialize Q_ϕ and Q'_ϕ by p_f .

while not converged do

- | Collect training trajectories, and sample a minibatch B
- | Compute target R_t :
 - | if $t < T$, then $R_t = r_t + \max_{a_{t+1}} Q'_\phi(s_{t+1}, a_{t+1})$
 - | if $t = T$, then $R_t = r_T$
- | Update ϕ (using gradient descent) by the objective
 - | $\text{argmin}_\phi [Q_\phi(s_t, a_t) - R_t]^2$
- | Update Q'_ϕ : $Q'_\phi \leftarrow Q_\phi$ every K steps

Results intuitively...

