# Adaptive Prior-Dependent Correction Enhanced Reinforcement Learning for Natural Language Generation

Background & Motivations

- > Neural natural language generation (NLG) aims to generate a piece of new text. NLG models have recently > Neural machine translation (NMT)
  - Image captioning
  - ➢ Text summarization
  - ▶ ...
- > Most of previous works use Maximum Likelihood Estimation (MLE) to train NLG models, but MLE-based training methods suffer from three issues:
- $\succ$  Exposure bias: The model is not exposed to the full range of errors during training.
- > Loss inconsistency: During training, we maximize the log-likelihood (token level), but during inference, the model is evaluated by a different metric such as BLEU or ROUGE (sentence level).
- outputs are treated equally during training.
- - to different incorrect model outputs.
- instability of training.

Methods

> Advantage weighted Policy Gradient (APG) training:

> Advantage function: We introduce the Generalized Advantage Estimation (GAE) method to reduce the variance of gradient.

$$\nabla_{\theta} L_{RL}(\theta) = \sum_{t=0}^{T} \mathbb{E}_{\pi} [A^{\pi}(s_t, a_t) \nabla_{\theta} log(\pi_{\theta}(a_t \mid s_t))] \qquad A^{\pi}(s_t, a_t) = \sum_{l=1}^{\infty} (\gamma \lambda)^l (r_t + \gamma V^{\pi}(s_{t+l+1}) - V^{\pi}(s_{t+l}))$$

- to estimate state values such as value-network. However, we can estimate state value indirectly.  $Q^{\pi}(s_t, a_t) = r_t + \gamma \sum P(s_{t+1} | s_t, a_t) V^{\pi}(s_{t+1})$
- $\succ$  In NLG task, the discount factor is set to 1 and state transition probability is equal to 1, thus we have  $A\pi(s_t, a_t) = Q^{\pi}(s_t, a_t) - Q^{\pi}(s_{t-1}, a_{t-1})$

## > Adaptive Prior-Dependent Correction. To alleviate deviation ignorance issue, we enhance the RL objective with a KL term with an adaptive factor:

cosine distance between tokens and label token:

$$p^*(w_t) = \sigma \left( cos\_sim \right)$$

> KL loss: KL loss measures the distance between output and prior distribution. We use the Adaptive Factor (negatively correlated to advantage function) to prevent the RL loss from being overcorrected by KL loss:



Figure1: Architecture of our approach: APG with APDC

Wei Cheng, Ziyan Luo<sup>1</sup>, Qiyue Yin<sup>2</sup> 1 Computer Science Department, University of California, San Diego 2 CRISE, Institute of Automation, Chinese Academy of Sciences, Beijing, China

shown remarkable progress in language fluency and coherence. We focus the training of Seq2Seq NLG model.

> Deviation ignorance: MLE fails to assign proper scores to different incorrect model outputs, which means that all incorrect

# > The main motivation of our research is that Reinforcement Learning (RL) based training method can solve the issues of exposure bias and loss inconsistency, but RL-based training methods also suffer from two issues:

> Deviation ignorance: RL-based training methods use the reward such as BLEU or ROUGE also fails to assign proper scores

> Large gradient variance: RL methods such as Policy Gradient have the issue of Large gradient variance, which leads to

> Estimation of state value: In order to estimate advantage, we need to estimate state value, but there have no dedicated module

> The prior distribution: We use fast-text to pretrain the word embedding and compute the prior distribution by measuring the

 $(emb(w_t^*),emb(w_t))$ 

Experiments & Results

### > Our method v.s. baselines

 $\succ$  Neural machine translation (NMT)

Method	En-De	En-Zh	Zh
Transformer+MLE (Vaswani et al. 2017)	27.30	34.12	24.
MIXER (Ranzato et al. 2016)	27.43	34.38	24.
RL4NMT (Wu et al. 2018)	27.52	34.46	24.
APG	27.63	34.54	24.
APG+PDC	27.70	34.56	24.
RL4NMT+PDC	27.81	34.62	24.
APG+APDC	28.03	34.91	25.

Table1: BLEU score of NMT for En-De, En-Zh, and Zh-En

#### Image captioning

Method	BLEU-1	BLEU-4	METEOR	ROUGE-L
SCST (Att2all) (Rennie et al. 2017)	-	34.2	26.7	55.7
OTRL (Chen et al. 2020a)	79.3	34.4	26.8	56.2
Top-Down+MLE (Anderson et al. 2018)	77.2	36.2	27.0	56.4
Top-Down+SCST (Anderson et al. 2018)	79.8	36.3	27.7	56.9
Top-Down+MIXER (Ranzato et al. 2016)	78.4	36.2	27.4	56.5
Top-Down+SCST+APDC	79.9	36.6	27.8	56.9
APG	80.1	36.4	28.0	56.9
APG+PDC	80.3	36.7	28.4	57.3
APG+APDC (METEOR)	80.1	36.5	29.2	57.1
APG+APDC	80.8	37.9	28.9	58.1

Table2: Performance of image captioning on the MSCOCO Karpathy test sp

### ► Abstractive text summarization

Method	ROUGE-1	ROUGE-2	ROUC
Pointer (See, Liu, and Manning 2017)	39.53	17.28	36.38
MIXER (Ranzato et al. 2016)	39.78	17.91	37.15
OTRL (Chen et al. 2020a)	41.40	18.22	38.86
APG	41.51	18.34	38.93
APG+PDC	41.60	18.48	39.02
APG+APDC	42.73	18.81	39.85

Table3: Results of abstractive text summarization on CNN/Daily Mail datas



Summary

> We propose a novel technique: adaptive prior-dependent correction (APDC) to further address the deviation ignorance issue that former RL-based approaches on NLG seldom study. > We utilize advantage-function-weighted policy gradient (APG) to work well with APDC, meanwhile alleviate the sparse reward issue. Enhancing APG with APDC can strike a balance between tokenlevel and sequence-level optimization.

> Further works: make the algorithm more efficient, ....

n-En 29 59 70 81 90 94 94 5.28	<ul> <li>Results shows that our method consistently outperforms the existing state-of-art methods on three NLG tasks.</li> <li>To evaluate the effectiveness of different components, we compare the results of applying APG, APG+PDC, and APG+APDC on the three tasks.</li> </ul>
CIDEr 114.0 111.8 113.5 120.1 115.6 120.2 120.2 120.2 121.2 119.7 <b>123.6</b>	<ul> <li>The results also show that APG+PDC outperform the APG (other RL), since PDC alleviates the deviation ignorance issue. However, the improvement of PDC is not significant.</li> <li>PDC is a token-level objective, but the advantage of RL is that the model can be trained at the sequence level, and PDC</li> </ul>
GE-L	<ul> <li>weakens this advantage. Therefore, PDC needs a sequence-level adaptive factor to adjust how much the token-level objective affects sequence-level training.</li> <li>The results that APG+APDC</li> </ul>
set.	significantly outperforms other methods also proves the importance of the adaptive mechanism.