



# Posterior-GAN: Towards Informative and Coherent Response Generation with Posterior Generative Adversarial Network

Shaoxiong Feng<sup>1</sup>, Hongshen Chen<sup>2</sup>, Kan Li<sup>1</sup>, Dawei Yin<sup>2</sup>

<sup>1</sup> School of Computer Science & Technology, Beijing Institute of Technology

<sup>2</sup> Data Science Lab, JD.com

## Introduction

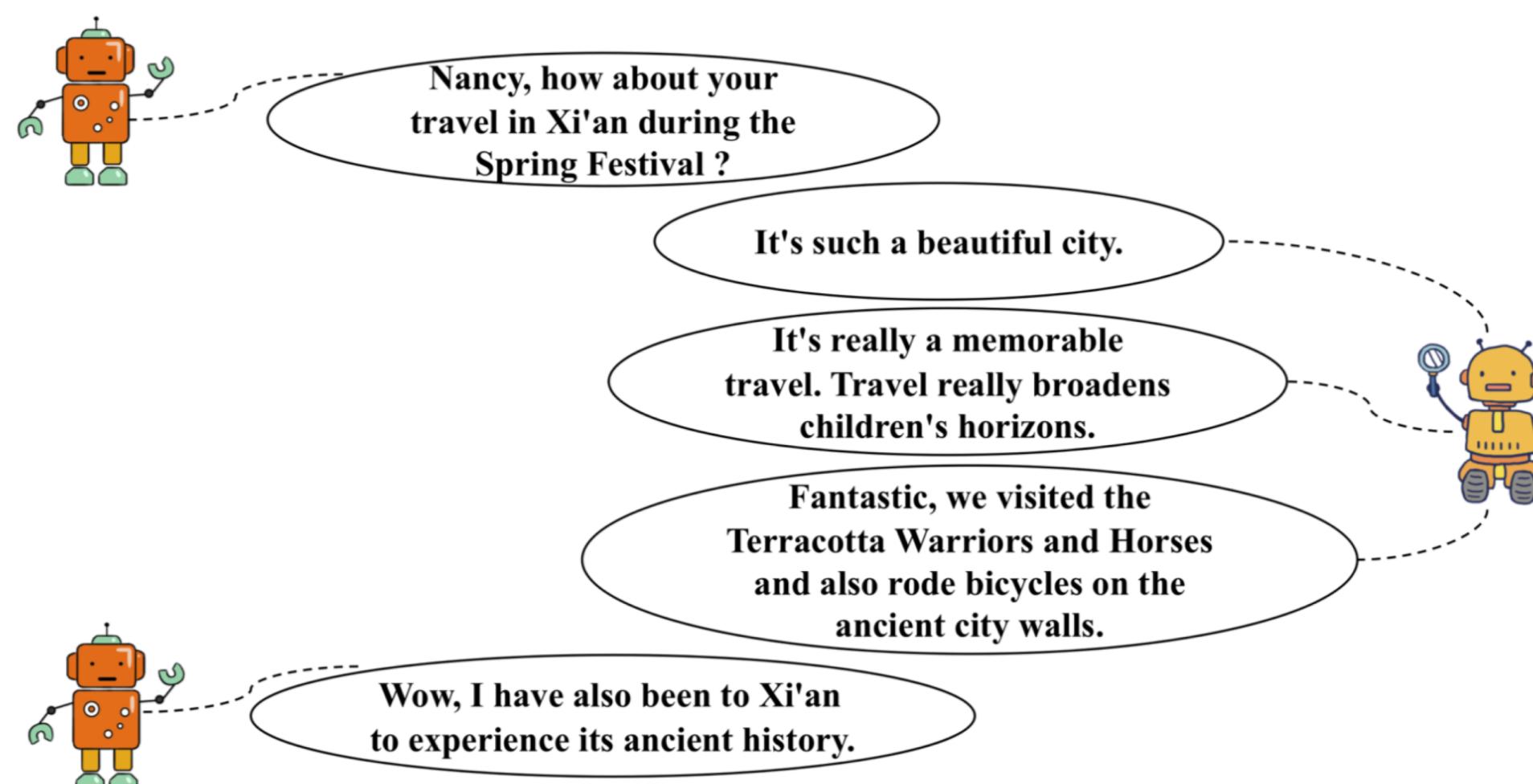


Figure 1: Dialogue examples with various responses regarding informativeness and coherence.

- We identify an unexplored type of metadata, *query-response-future turn triples*, for response generation. Compared to general query-response tuples, the triples help the model use bidirectional information to learn the response generation in training.

## Method

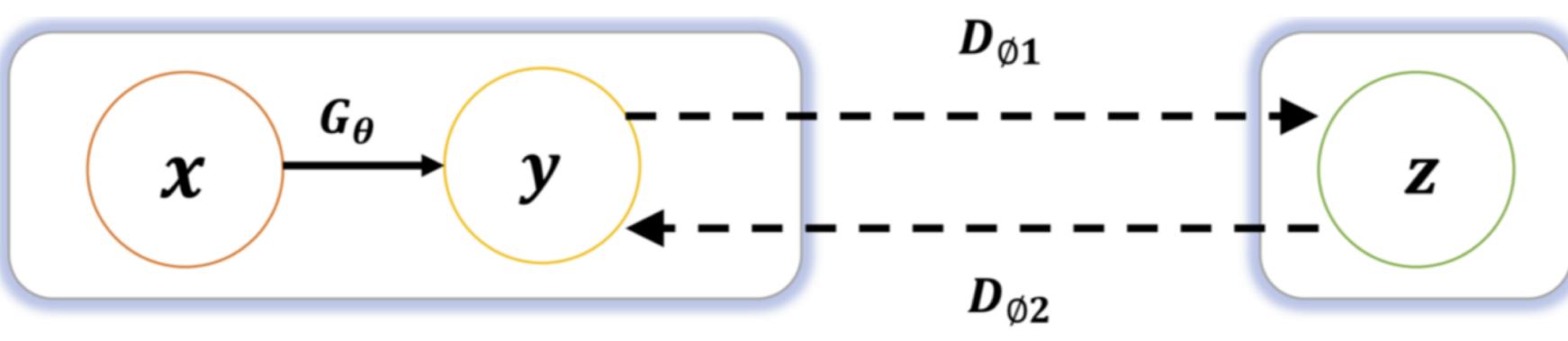


Figure 2: Illustration of Posterior-GAN. Brown for the query, yellow for the current response, and green for the future turn.  $G_\theta$  represent generator.  $D_{\phi 1}$  and  $D_{\phi 2}$  represent forward and backward generative discriminator respectively.

- We propose a novel encoder-decoder based generative adversarial learning framework, **Posterior-GAN**, to facilitate the query-response-future turn modeling, which induces the generated response to be informative and coherent by constructing **two generative discriminators, a forward one and a backward one respectively**.

### Forward Generative Discriminator ( $D_{\phi 1}$ )

$$R_1(y) = \frac{1}{K} \sum_{k=1}^K \log D_{\phi 1}(z_k | y, z_{<k})$$

$$J(\phi 1) = - \left( E_{y_{true} \sim p_{data}} [R_1(y_{true})] - E_{y_{G_\theta} \sim G_\theta} [R_1(y_{G_\theta})] \right)$$

### Backward Generative Discriminator ( $D_{\phi 2}$ )

$$R_2(y_m) = - \log D_{\phi 2}(y_m | z, y_{<m})$$

$$J(\phi 2) = - \left( E_{y_{true} \sim p_{data}} [R_2(y_{true})] - E_{y_{G_\theta} \sim G_\theta} [R_2(y_{G_\theta})] \right)$$

### Training Objective ( $J(\theta)$ )

$$J(\theta) = E_{y_{G_\theta} \sim G_\theta} \left( Q_{D_{\phi 1}, D_{\phi 2}}^{G_\theta}(x, y_{G_\theta}, z) \right)$$

$$\nabla_\theta J(\theta) \simeq$$

$$\sum_{n=1}^N \sum_{m=1}^L R_m^n \nabla_\theta \log G_\theta(y_m^n | x, y_{<m}^n)$$

$$R_m^n = \sum_{i=m}^L \lambda^i (R_1(y_i^n) - \text{MIN}(R_1)) R_2(y_i^n)$$

## Experiments

- We perform detailed experiments to demonstrate the effectiveness of the proposed framework and verifies the ability of bidirectional generative discriminators on assessing the quality of response.
- Datasets:** DailyDialog and Opensubtitles
- Baselines:** Seq2Seq-att, Adver-REGS, and DP-GAN
- Metrics:** BLEU, Embedding-based Metrics (Embedding Average, Embedding Greedy and Embedding Extrema), and Distinct (Dist-{1,2,3})

## Automatic Evaluation

Models	DailyDialog						
	Dist-1	Dist-2	Dist-3	BLEU	Greedy	Average	Extrema
Seq2Seq-att	0.0277	0.1625	0.3868	0.1878	0.4825	0.5993	0.3080
Adver-REGS	0.0541	0.2877	0.5542	0.2116	0.4857	0.6215	0.3542
DP-GAN	0.0656	0.3088	0.5630	0.1992	0.4749	0.6144	0.3409
Posterior-GAN(F)	0.0659	0.2995	0.5578	0.2067	0.4754	0.6218	0.3343
Posterior-GAN(B)	0.0578	0.2950	0.5545	0.2130	0.4818	0.6220	0.3442
Posterior-GAN(A)	<b>0.0678</b>	<b>0.3549</b>	<b>0.6006</b>	<b>0.2183</b>	<b>0.4916</b>	<b>0.6260</b>	<b>0.3544</b>

Models	OpenSubtitles (OSDb)						
	Dist-1	Dist-2	Dist-3	BLEU	Greedy	Average	Extrema
Seq2Seq-att	0.0016	0.0064	0.0150	0.1405	0.3900	0.4527	0.2243
Adver-REGS	0.0041	0.0136	0.0248	0.1609	0.4655	0.5523	0.2645
DP-GAN	0.0044	0.0143	0.0262	0.1484	0.4600	0.5509	0.2589
Posterior-GAN(F)	0.0049	0.0170	0.0322	0.1733	0.4753	0.5708	0.2573
Posterior-GAN(B)	0.0045	0.0148	0.0321	0.1756	0.4947	0.6217	0.2690
Posterior-GAN(A)	<b>0.0049</b>	<b>0.0180</b>	<b>0.0330</b>	<b>0.1955</b>	<b>0.4973</b>	<b>0.6346</b>	<b>0.2778</b>

Table 1: The automatic metrics evaluation results. Higher is better. "(F)", "(B)" and "(A)" represent Posterior-GAN with a forward generative discriminator, a backward generative discriminator and both two discriminators, respectively.

## Human Evaluation

Models	DailyDialog	
	Coherence	Informativeness
Seq2Seq-att	3.8550	3.8933
Adver-REGS	3.4717	3.3683
DP-GAN	3.4400	3.2350
Posterior-GAN	<b>3.2883</b>	<b>3.2250</b>

Models	OpenSubtitles (OSDb)	
	Coherence	Informativeness
Seq2Seq-att	3.8549	3.6952
Adver-REGS	4.0365	4.0432
DP-GAN	3.7638	3.8088
Posterior-GAN	<b>3.4806</b>	<b>3.4567</b>

Table 2: The human evaluation results. We calculate each score by averaging the rank of each model in corresponding metrics. Lower is better.

## Automatic Analysis

Models	Averaged Greedy Matching for $y$ and $(x/z)$	Frequency-based Similarity	
		for $y$ and $(x/z)$	Similarity
Seq2Seq-att	0.5942	0.5478	
Adver-REGS	0.6403	0.7165	
DP-GAN	0.6514	0.6739	
Posterior-GAN	<b>0.7276</b>	<b>0.7460</b>	

Table 4: The results of Embedding-based Averaged Greedy Matching for response  $y$  with the given query  $x$  and future conversations  $z$ , which reflects the coherence of the response, and Frequency-based Similarity, which illustrates the informativeness of the response.

## Visualization

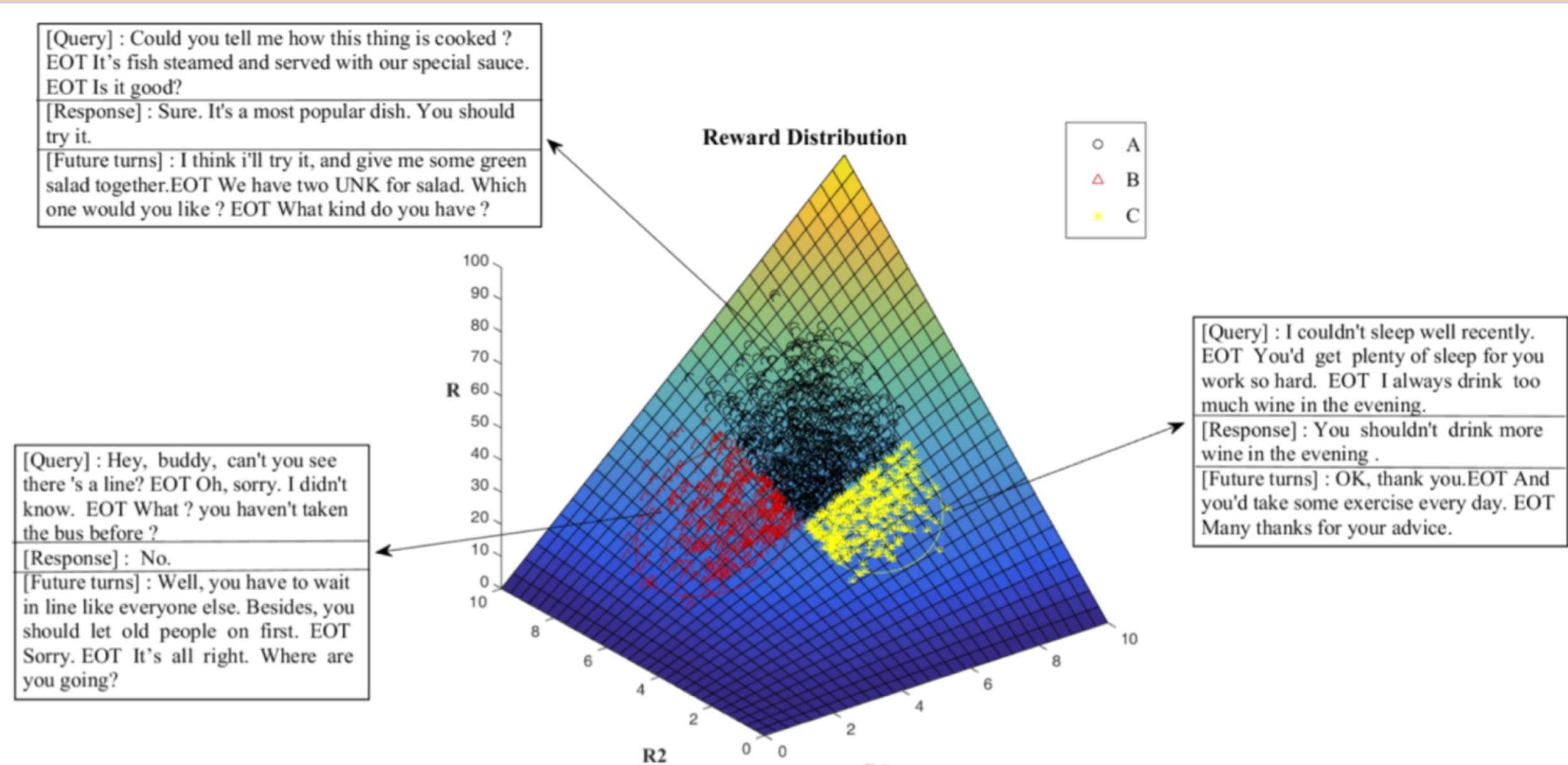


Figure 4: The distribution of sample rewards calculated by the forward generative discriminator  $R1$  and the backward generative discriminator  $R2$  on DailyDialog.  $R$  is the combination of  $R1$  and  $R2$ . We use three regions A, B, and C to represent three types of samples. Samples in region A gain high rewards in both discriminators. Samples in region B achieve higher reward in the backward generative discriminator than in the forward one, while Samples in region C obtain higher reward in the forward generative discriminator than in the backward one.

## Conclusion

- We propose the *query-response-future turn triples* instead of the conventional query-response pairs for neural dialog response generation.
- To facilitate the triple modeling and alleviate the overproducing of generic and repetitive responses problem, **Posterior-GAN** that consists of a forward and a backward encoder-decoder based generative discriminator is further introduced.
- Augmented with future conversations and Posterior-GAN in training, detailed experiments and analysis demonstrate that **the model effectively generates more informative and coherent responses**.