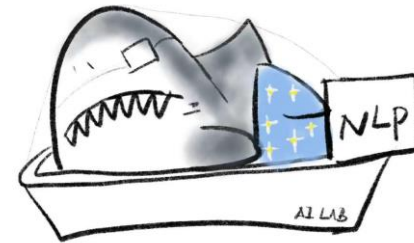


DiffuSeq: Sequence to Sequence Text Generation With Diffusion Models

Shark-NLP
Shanghai, China
hisansas@gmail.com



Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, Lingpeng Kong

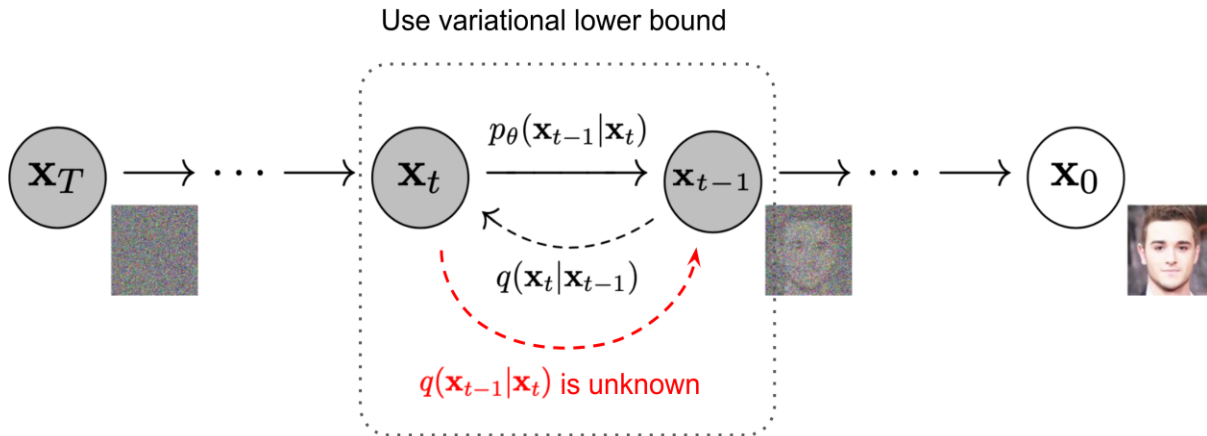
<https://github.com/Shark-NLP/DiffuSeq>

Contents

- Preliminary and motivation
- DiffuSeq and beyond
- Experiments and analysis
- Conclusion and future work

1.1 Preliminary

Diffusion process in continuous space:
(applied in vision, audio, time series and etc....)



1. Noise-conditioned score network (NCSN; Yang & Ermon, 2019)
2. Denoising diffusion probabilistic models (DDPM; Ho et al. 2020)

△ Forward process:

- $\mathbf{x}_0 \sim q(\mathbf{x}) \rightarrow \mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$
- $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$

△ Reverse process:

- $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_\theta(\mathbf{x}_t, t))$
- $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}$

△ Training loss:

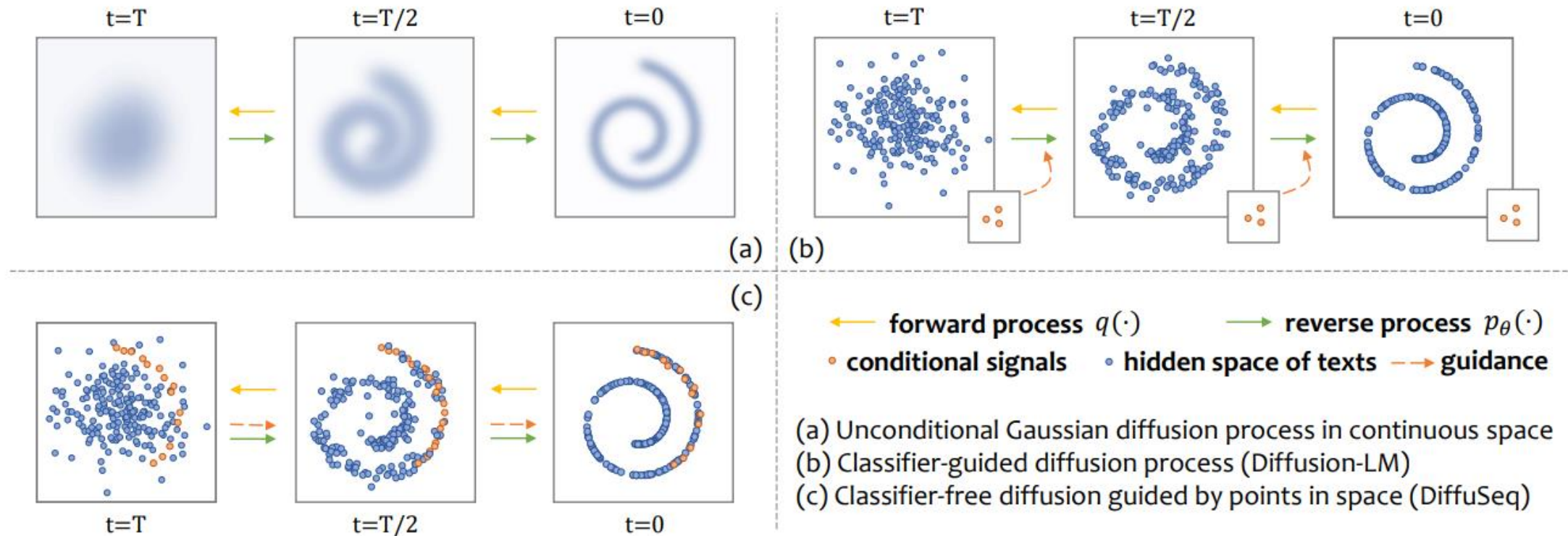
- $L_t = D_{KL}(q||p_\theta)$
- Parameterization of $L_t =$

$$\mathbb{E}_{\mathbf{x}_0}(\|\mathbf{x}_0 - f_\theta(\mathbf{x}_t, t)\|^2)$$

1.2 Motivation

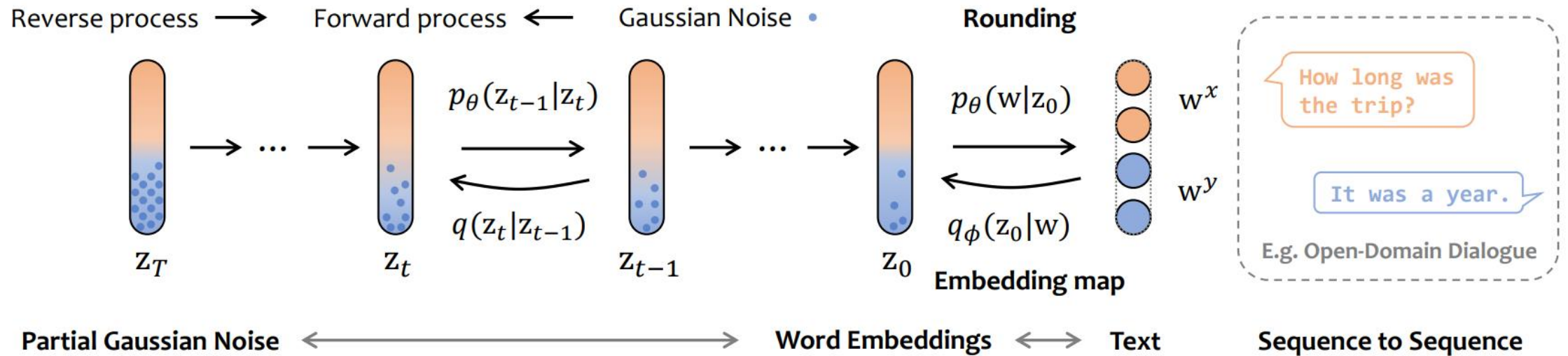
From unconditional models to conditional models:

Diffusion-LM (classifier-guided) v.s. DiffuSeq (classifier-free)



Seq2Seq tasks: $\mathbf{x} \rightarrow \mathbf{y}$

2.1 DiffuSeq



△ Forward Process with Partial Noising:

- $q(\mathbf{z}_0|\mathbf{w}^{x\oplus y}) = \mathcal{N}(\text{EMB}(\mathbf{w}^{x\oplus y}), \beta_0 \mathbf{I}); \mathbf{z}_t = \mathbf{x}_t \oplus \mathbf{y}_t$

△ Reverse Process with Conditional Denoising:

- $L_t = \mathbb{E}_{\mathbf{x}_0, \mathbf{y}_0} (\|\mathbf{y}_0 - f_{\tilde{\theta}}(\mathbf{z}_t, t)\|^2)$

△ Training:

- importance sampling

△ Inference:

- Rounding to embeddings
- Anchoring input signals

2.2 Connections of different models

AR/iter-NAR/DiffuSeq: Generation process is along with different dimensions:

$$p_{\text{AR}}(\mathbf{w}_{1:n}^y | \mathbf{w}^x) = \underbrace{p(w_1^y | \mathbf{w}^x)}_{\text{initial prediction}} \underbrace{\prod_{i=1, \dots, n-1} p(w_{i+1}^y | \mathbf{w}_{1:i}^y, \mathbf{w}^x)}_{\text{progressive left-context prediction}},$$

$$p_{\text{iter-NAR}}(\mathbf{w}_{1:n}^y | \mathbf{w}^x) = \sum_{\mathbf{w}_1^y, \dots, \mathbf{w}_{K-1}^y} \underbrace{\prod_{i=1 \dots n} p(w_{1,i}^y | \mathbf{w}^x)}_{\text{initial prediction}} \underbrace{\prod_{k=1 \dots K-1} \prod_{i=1 \dots n} p(w_{k+1,i}^y | \mathbf{w}_{k,1:n}^y, \mathbf{w}^x)}_{\text{progressive full-context prediction}}.$$

$$p_{\text{DIFFUSEQ}}(\mathbf{w}^y | \mathbf{w}^x) = \sum_{\substack{\mathbf{w}_T^y, \dots, \mathbf{w}_1^y \\ \mathbf{y}_T, \dots, \mathbf{y}_0}} p(\mathbf{w}^y | \mathbf{y}_0, \mathbf{w}^x) \prod_{t=T, \dots, 1} p(\mathbf{w}_t^y | \mathbf{y}_t, \mathbf{w}^x) p(\mathbf{y}_{t-1} | \mathbf{w}_t^y)$$

3.1 Experiments

Four tasks: Dialogue, QG, Text Simplification, Paraphrase

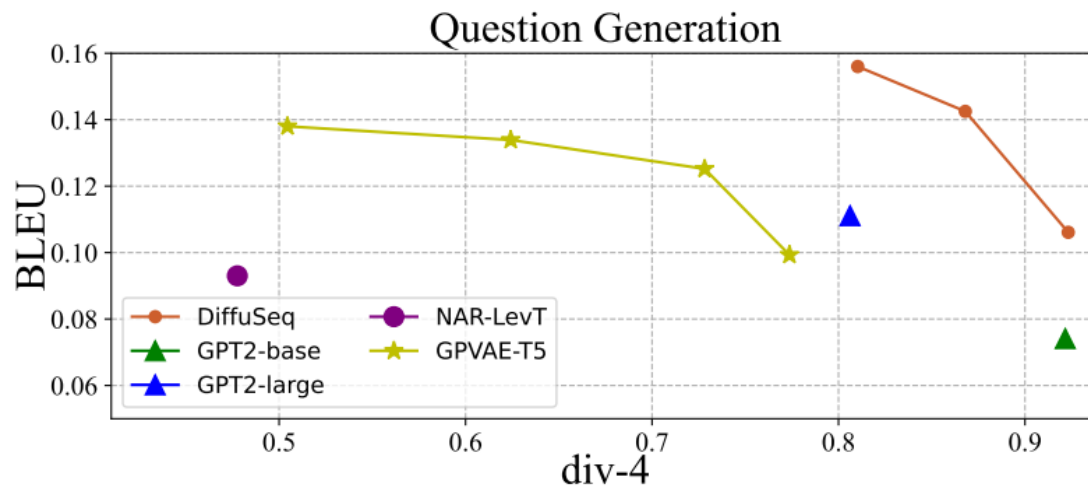
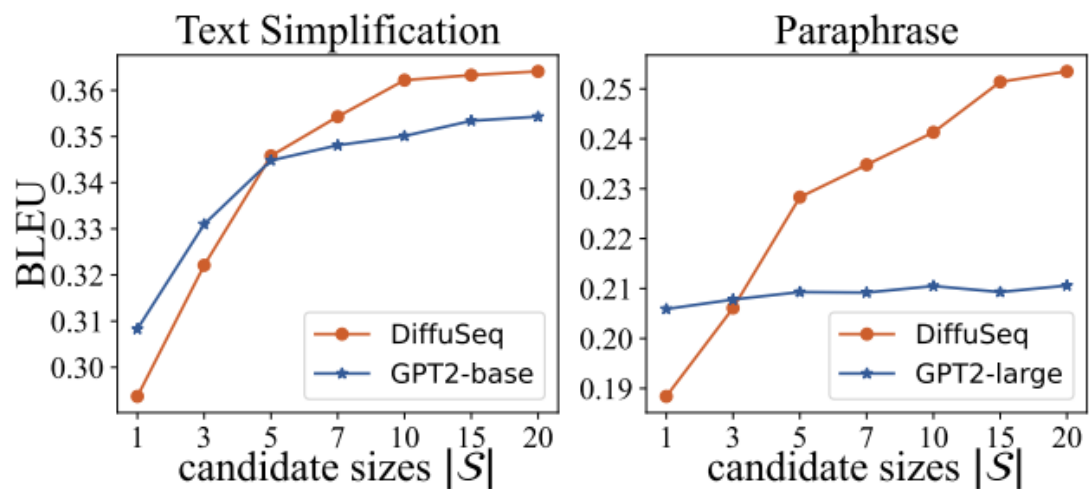
Three groups of baselines: Plain encoder-decoder, PLMs, NAR

Tasks	Methods	BLEU \uparrow	R-L \uparrow	Score \uparrow	dist-1 \uparrow	selfB \downarrow / div-4 \uparrow	Len
Paraphrase	GRU-attention \diamond	0.1894	0.5129	0.7763	0.9423	0.9958/0.3287	8.30
	Transformer-base \diamond	<u>0.2722</u>	0.5748	<u>0.8381</u>	0.9748	0.4483/0.7345	11.2
	GPT2-base FT \bullet	0.1980	0.5212	0.8246	0.9798	0.5480/0.6245	9.67
	GPT2-large FT \bullet	0.2059	0.5415	0.8363	0.9819	0.7325/0.5020	9.53
	GPVAE-T5 \bullet	0.2409	0.5886	0.8466	0.9688	0.5604/0.6169	9.60
	NAR-LevT \ddagger	0.2268	0.5795	0.8344	0.9790	0.9995/0.3329	8.85
	DIFFUSEQ (Ours) \ddagger	0.2413	<u>0.5880</u>	0.8365	<u>0.9807</u>	<u>0.2732/0.8641</u>	11.2

Comparable quality, better diversity

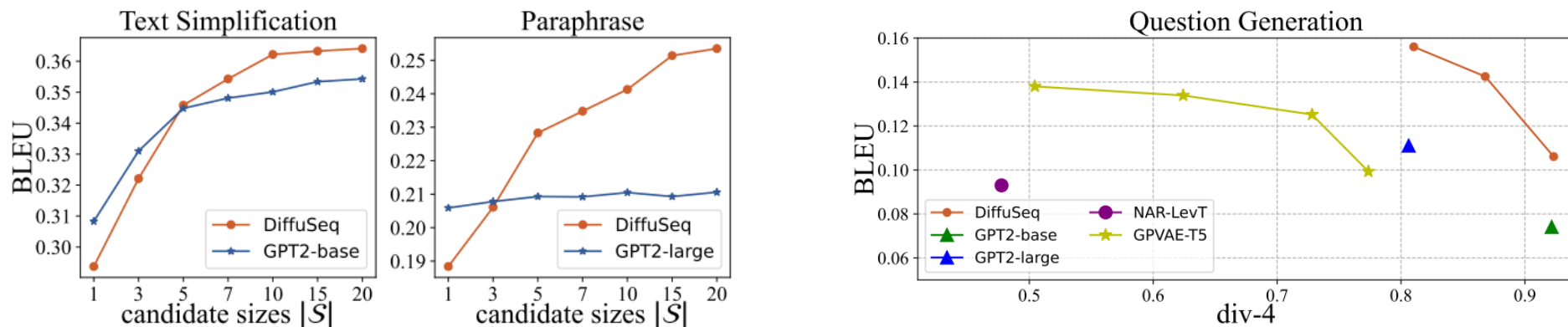
3.2 Analysis

Diversity Ensures Quality

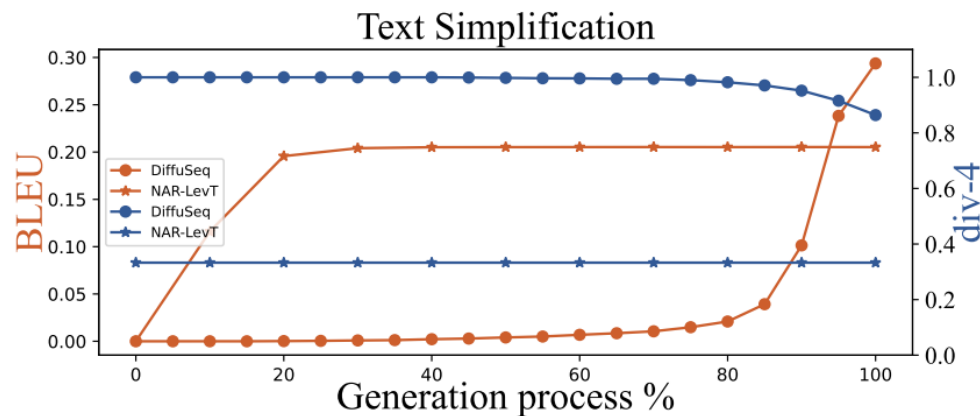


3.2 Analysis

Diversity Ensures Quality

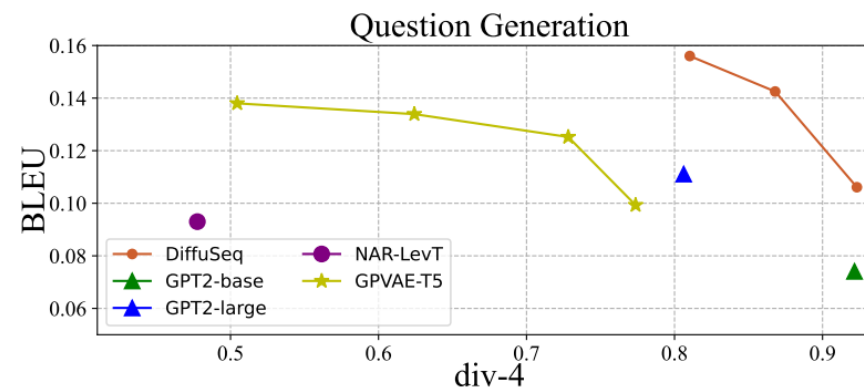
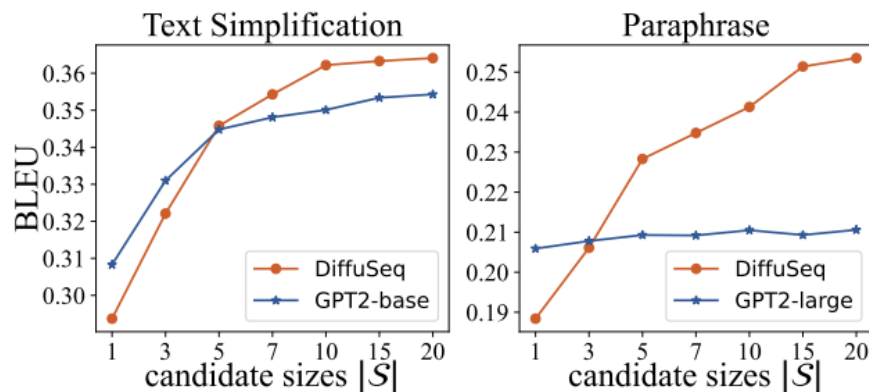


Step-wise Analysis against Iterative NAR

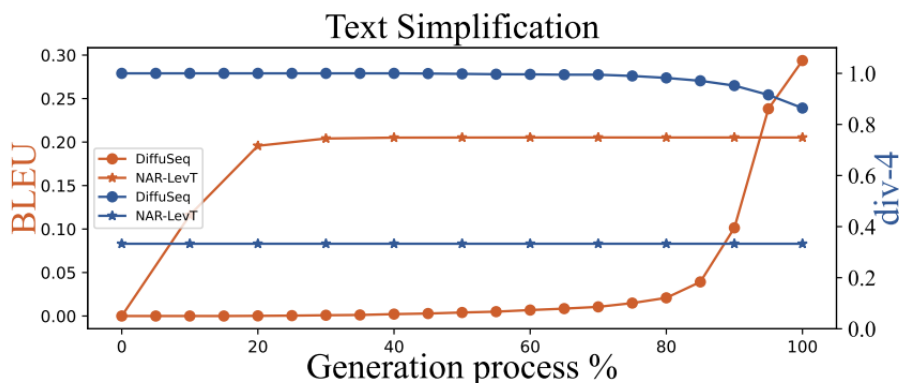


3.2 Analysis

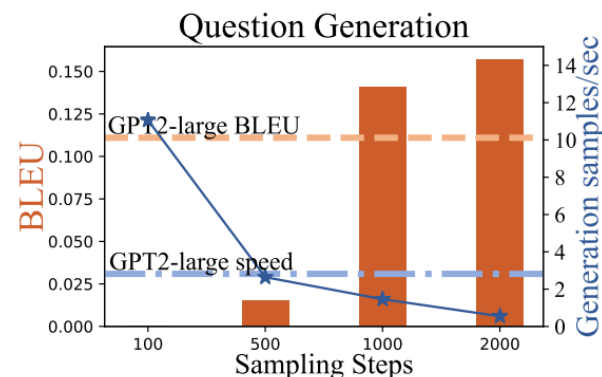
Diversity Ensures Quality



Step-wise Analysis against Iterative NAR



Inference Speed



4 Conclusion and future work

- DiffuSeq: as a new generation paradigm
 - Potential: competitive results on Seq2Seq tasks
 - Analysis: diversity DiffuSeq v.s. iter-NAR
- Future work: inference speed and sentence fluency

Thank you for watching!

Shark-NLP
Shanghai, China
hisansas@gmail.com



<https://github.com/Shark-NLP/DiffuSeq>