

# Structuring Latent Spaces for Stylized Response Generation

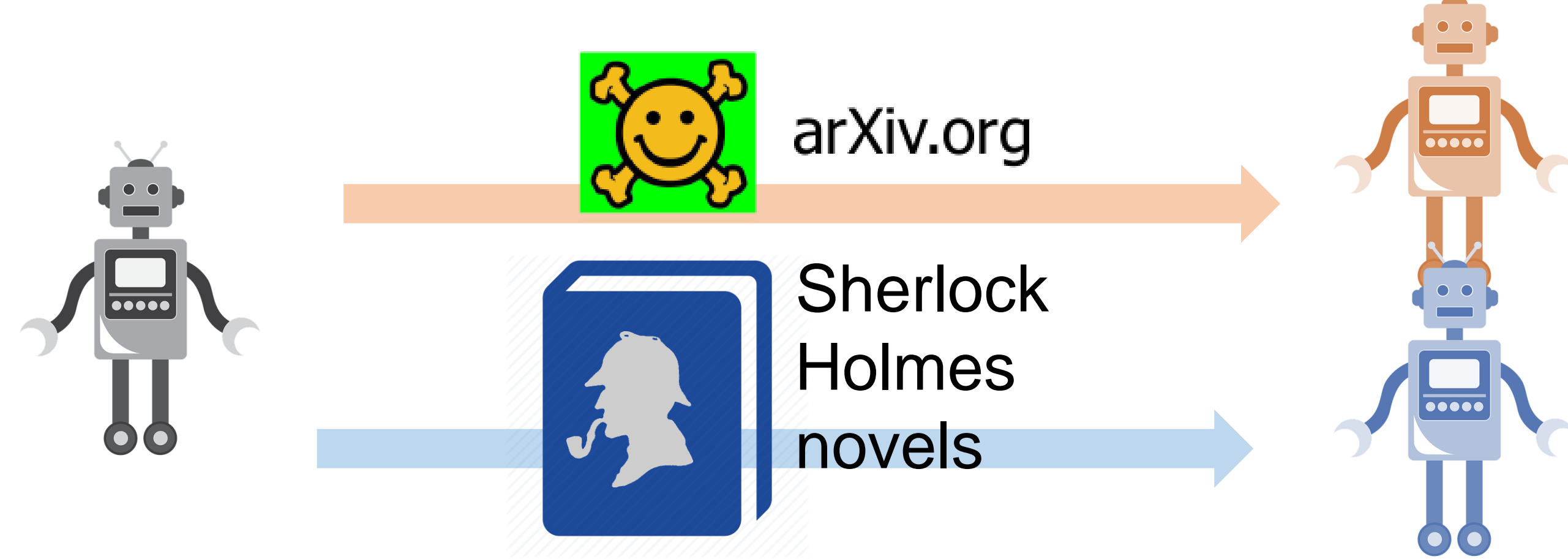
Xiang Gao, Yizhe Zhang, Sungjin Lee\*, Michel Galley, Chris Brockett, Jianfeng Gao, Bill Dolan



\* Now at Amazon Alexa AI

## Motivation & Task

- The master of response style is an important step towards humanlike chatbots.
- However, challenging if no parallel conversational data in different styles

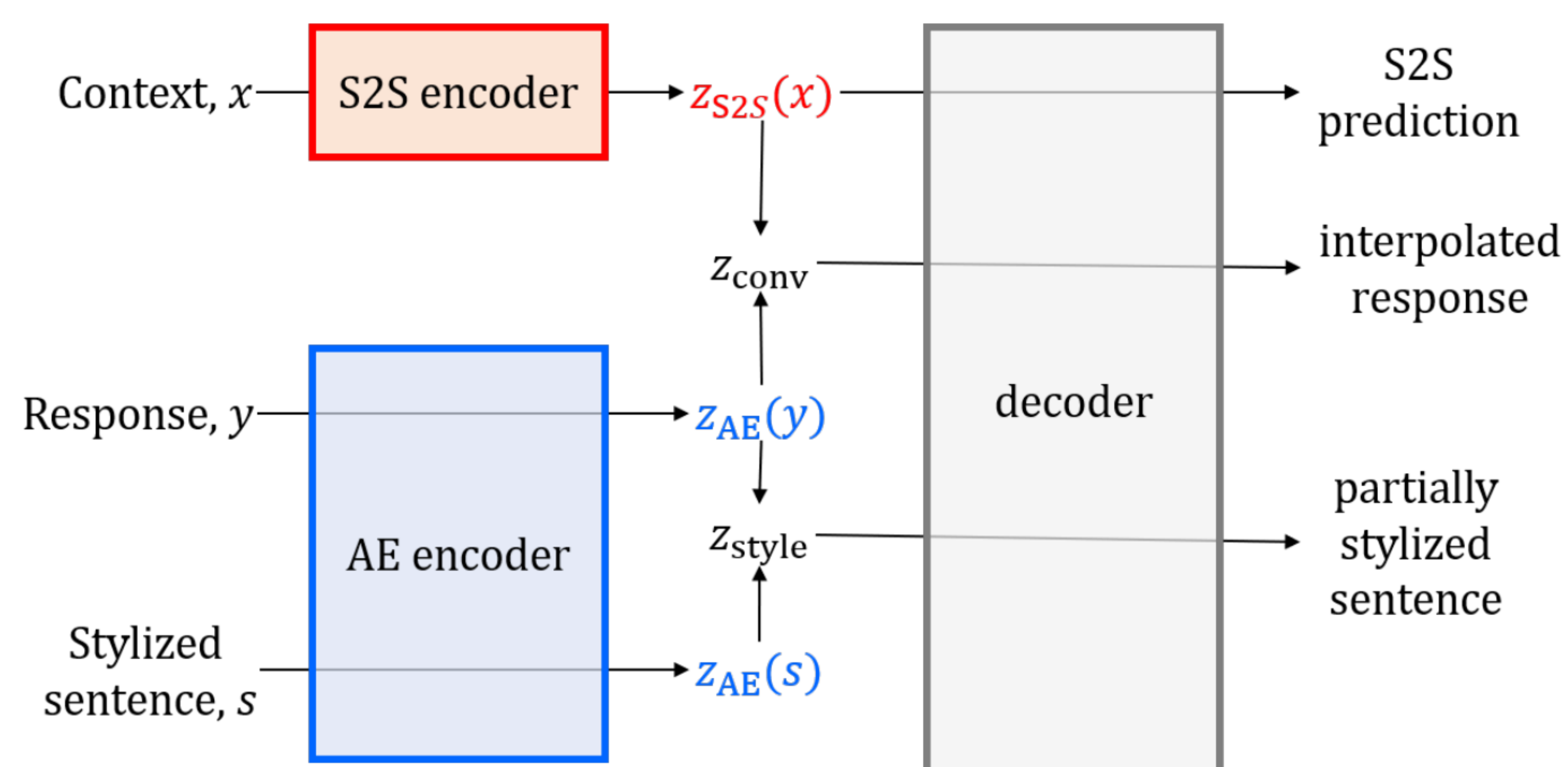


### Our task

- Train an agent on conversational dataset  $D_{conv}$
- which learns style from non-conversational non-parallel text dataset of target style  $D_{style}$
- to generate response in target style and appropriate to context

## Approach

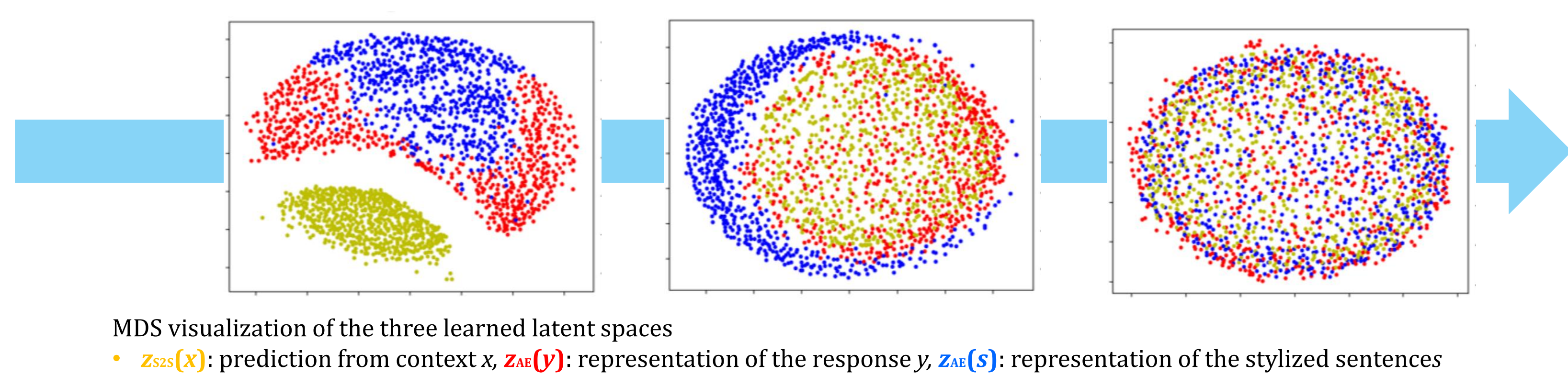
- We propose a regularized multi-task framework to align latent spaces of conversational and stylized datasets



- The loss function considers the generation probability as well as the structure of the latent space, as an extension of our previous SpaceFusion work (Gao et al., NAACL'19)

$$\mathcal{L} = -\frac{1}{|y|} \log p(y|z_{S2S}) - \frac{1}{|y|} \log p(y|z_{AE}) + \mathcal{L}_{smooth,conv} + \mathcal{L}_{fuse,conv} + \mathcal{L}_{smooth,style} + \mathcal{L}_{fuse,style}$$

Vanilla S2S (Luan et al., IJCNLP'17)    Vanilla MTask (Luan et al., IJCNLP'17)    SpaceFusion (Gao et al., NAACL'19)    "StyleFusion"



- $\mathcal{L}_{smooth}$  encourages smooth semantic interpolation on latent space  
 $z_{conv} = (1-u)z_{AE}(y) + uz_{S2S}(x) + \epsilon$      $z_{style} = (1-u)z_{AE}(x) + uz_{AE}(s) + \epsilon$   
 $\mathcal{L}_{smooth,conv} = -\frac{1}{|y|} \log p(y|z_{conv})$      $\mathcal{L}_{smooth,style} = -\frac{1}{|x|} \log p(x|z_{style}) - u \frac{1}{|s|} \log p(s|z_{style})$
- $\mathcal{L}_{fuse}$  encourages different latent spaces to overlap with each other  
 $\mathcal{L}_{fuse,conv} = d_{conv} - d_{spread-out}$      $\mathcal{L}_{fuse,style} = d_{style} - d_{spread-out}$

- $d_{conv}$  measures the distance between a pair of  $z_{S2S}(x)$  and  $z_{AE}(y)$
- $d_{style}$  measures the distance between a  $z_{S2S}(x)$  point and its nearest neighbor from  $z_{AE}(s)$
- $d_{spread-out}$  measures the average distance between a point and its nearest neighbor from the same latent space

## Inference

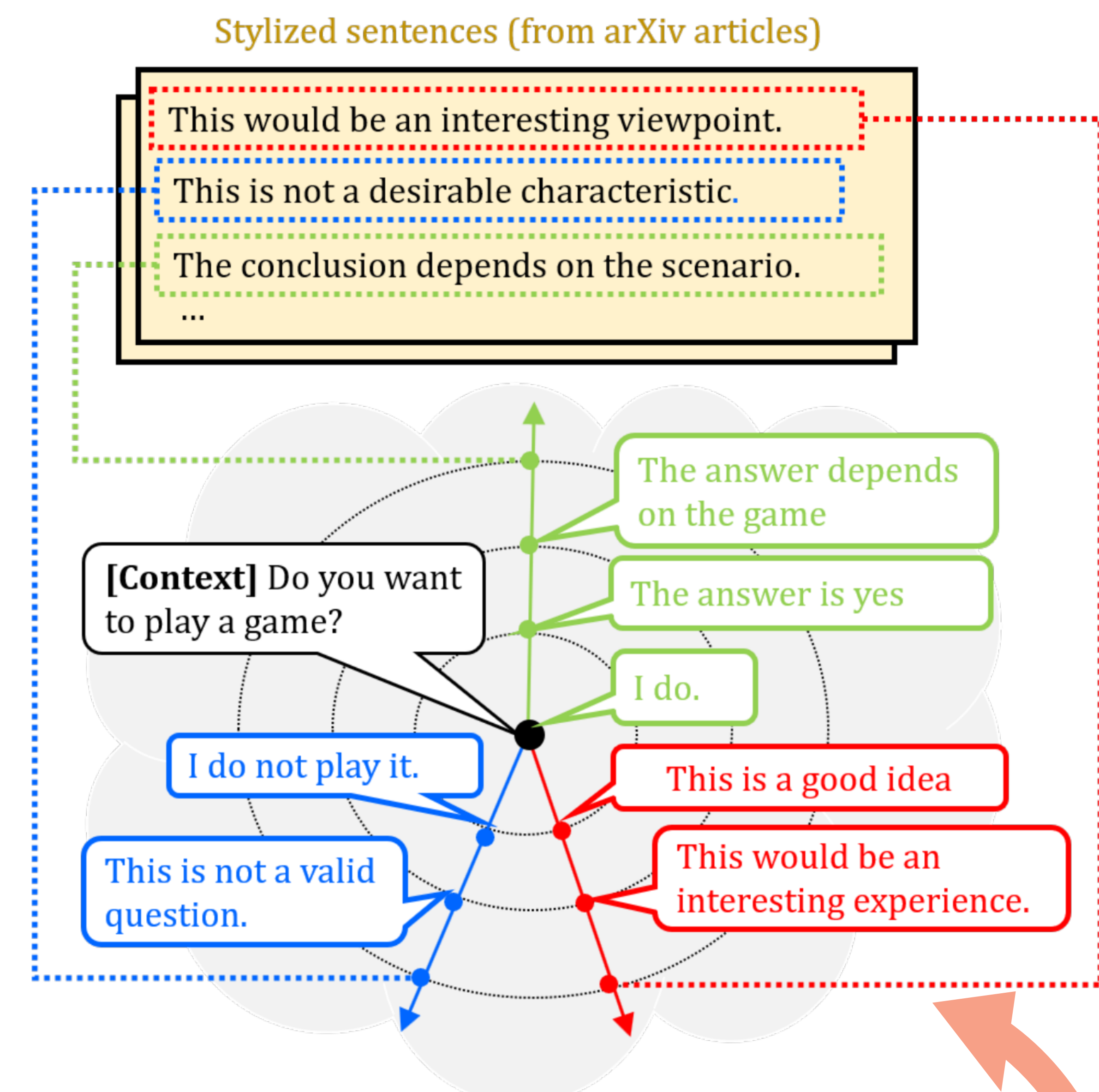
we sample in the neighborhood of  $z_{S2S}(x)$  by adding a noise  $r$  of a given length  $|r|$  towards a direction randomly drawn from the uniform distribution

$$z = z_{S2S}(x) + r \quad \rho = |r|/(\sigma\sqrt{l})$$

$$\text{score}(h_i) = (1-\lambda)P(h_i|z_{S2S}(x)) + \lambda P_{style}(h_i)$$

## Intuition

Structure a latent space such that the stylized texts are mapped around the related conversation data

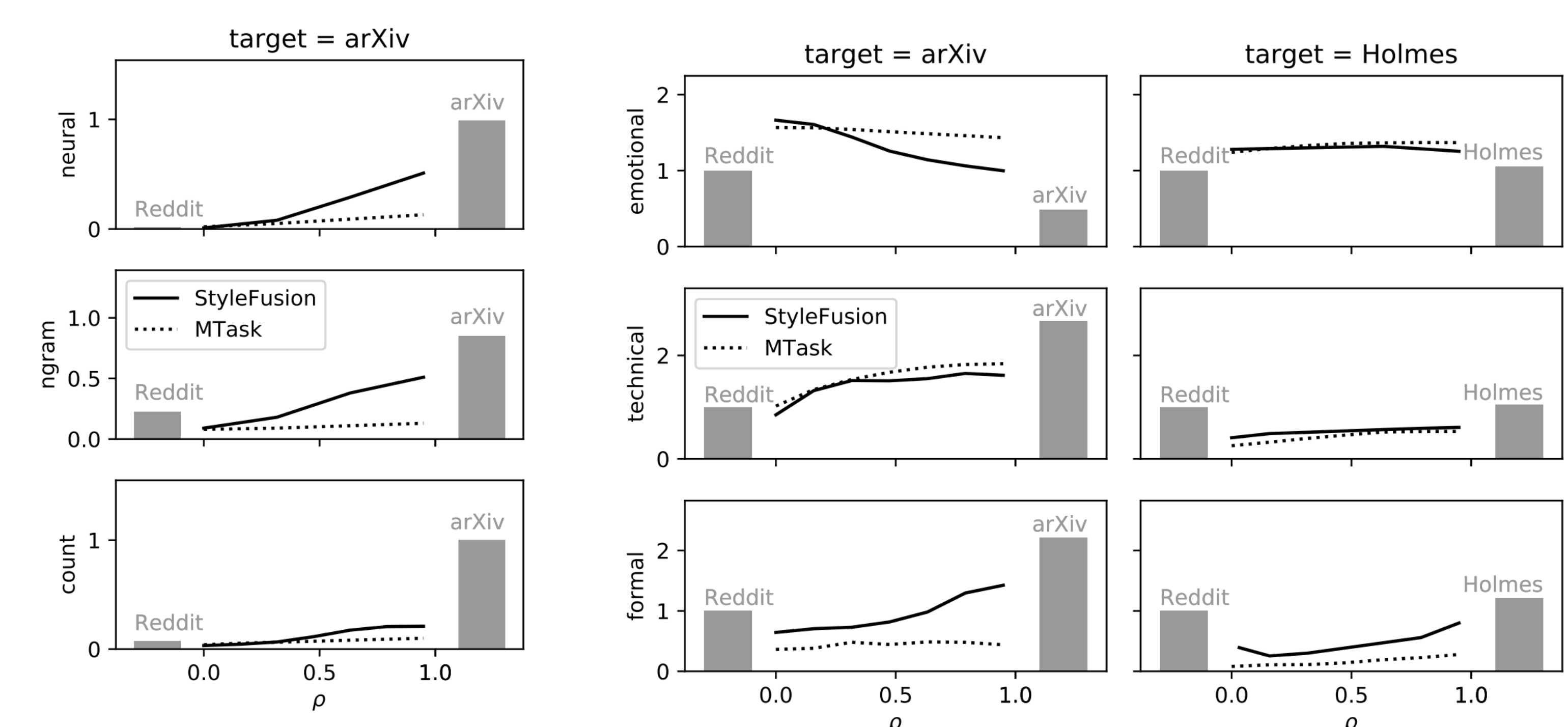


## Results

- Direction controls the content, radius  $\rho$  controls style intensity

context	Do you want to play a game?
towards	The conclusion depends on the scenario .
$\rho = 0.0$	I do.
$\rho = 0.5$	The answer is yes.
$\rho = 1.0$	The answer depends on the game.
towards	This would be an interesting viewpoint.
$\rho = 0.4$	This is a good idea.
$\rho = 0.9$	This would be an interesting experience
towards	This is not a desirable characteristic.
$\rho = 0.5$	I don't play it.
$\rho = 1.0$	This is not a valid question.

- As radius  $\rho$  increases, fine-granularity styles (emotional, formal, technical) also become similar to the target style



- Human evaluation shows that, our proposed approach, StyleFusion, jointly improve response appropriateness and style intensity, Compared with competitive baselines,

target	model	appropriateness	style intensity	harmonic mean
arXiv	STYLEFUSION	<b>0.17</b>	0.26	<b>0.20</b>
	MTask	<b>0.17</b>	0.11	0.14
	S2S+LM	0.09	0.51	0.15
	Retrieval	0.07	0.71	0.14
	Rand	0.04	<b>0.96</b>	0.07
	Human	0.51	0.28	0.36
Holmes	STYLEFUSION	<b>0.22</b>	0.28	<b>0.25</b>
	MTask	<b>0.23</b>	0.15	0.18
	S2S+LM	0.03	0.55	0.05
	Retrieval	0.14	0.30	0.19
	Rand	0.08	<b>0.69</b>	0.14
	Human	0.63	0.26	0.37

### Notes

- S2S+LM refers to the method proposed by Niu and Bansal (2018), which uses the weighted average of a S2S model, trained on conversational dataset, and a LM model, trained on style dataset, as the token probability distribution at inference time.
- MTask refers to the vanilla multi-task learning model proposed in (Luan et al., 2017) trained on both  $D_{conv}$  and  $D_{style}$ .



EMNLP-IJCNLP 2019

Paper: [arxiv.org/abs/1909.05361](https://arxiv.org/abs/1909.05361)

Code/Data: [github.com/golsun/StyleFusion](https://github.com/golsun/StyleFusion)