



# Empathetic Dialogue Generation via Sensitive Emotion Recognition and Sensible Knowledge Selection

**Lanrui Wang<sup>1,2</sup>, Jiangnan Li<sup>1,2</sup>, Zheng Lin<sup>1,2\*</sup>, Fandong Meng<sup>3</sup>,  
Chenxu Yang<sup>1,2</sup>, Weiping Wang<sup>1</sup>, Jie Zhou<sup>3</sup>**

<sup>1</sup>Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

<sup>2</sup>School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China

<sup>3</sup>Pattern Recognition Center, WeChat AI, Tencent Inc, China

{wanglanrui, lijiangnan, linzheng, yangchenxu, wangweiping}@iie.ac.cn  
{fandongmeng, withtomzhou}@tencent.com

code: <https://github.com/wlr737/EMNLP2022-SEEK>

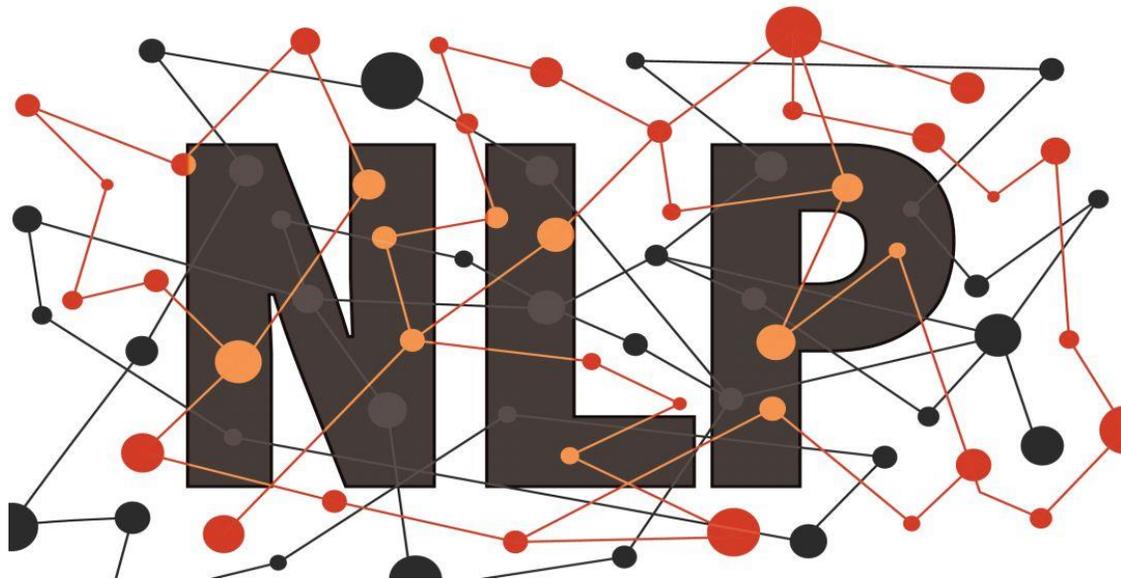
2022. 11. 27 • ChongQing

**2022\_EMNLP**





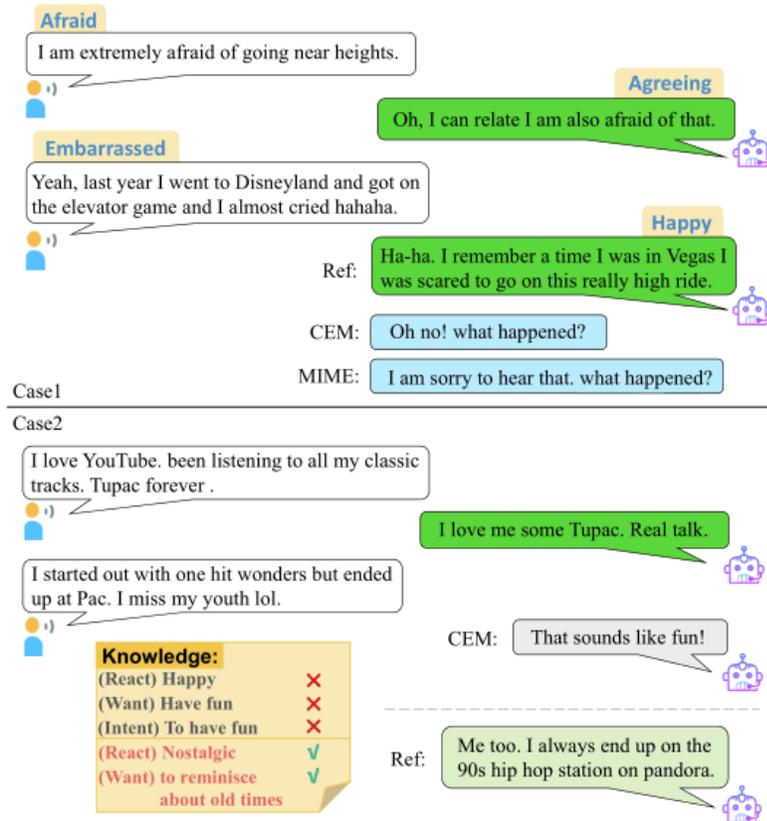
## NATURAL LANGUAGE PROCESSING



- 1. Introduction**
- 2. Method**
- 3. Experiments**



# Introduction



current approaches to empathetic response generation focus on capturing implicit emotion within dialogue context, where the **emotions are treated as a static variable** throughout the conversations.

simply importing commonsense knowledge without harmonization may trigger the **conflicts between knowledge and emotion**.

Figure 1: Two cases of multi-turn Empathetic Dialogues. The first case shows the speaker's emotion went from fear at the beginning of the conversation to an embarrassed self-deprecation, ending with a happy mood. And the second case shows that CEM chooses the wrong knowledge leading to inappropriate response.

# Method

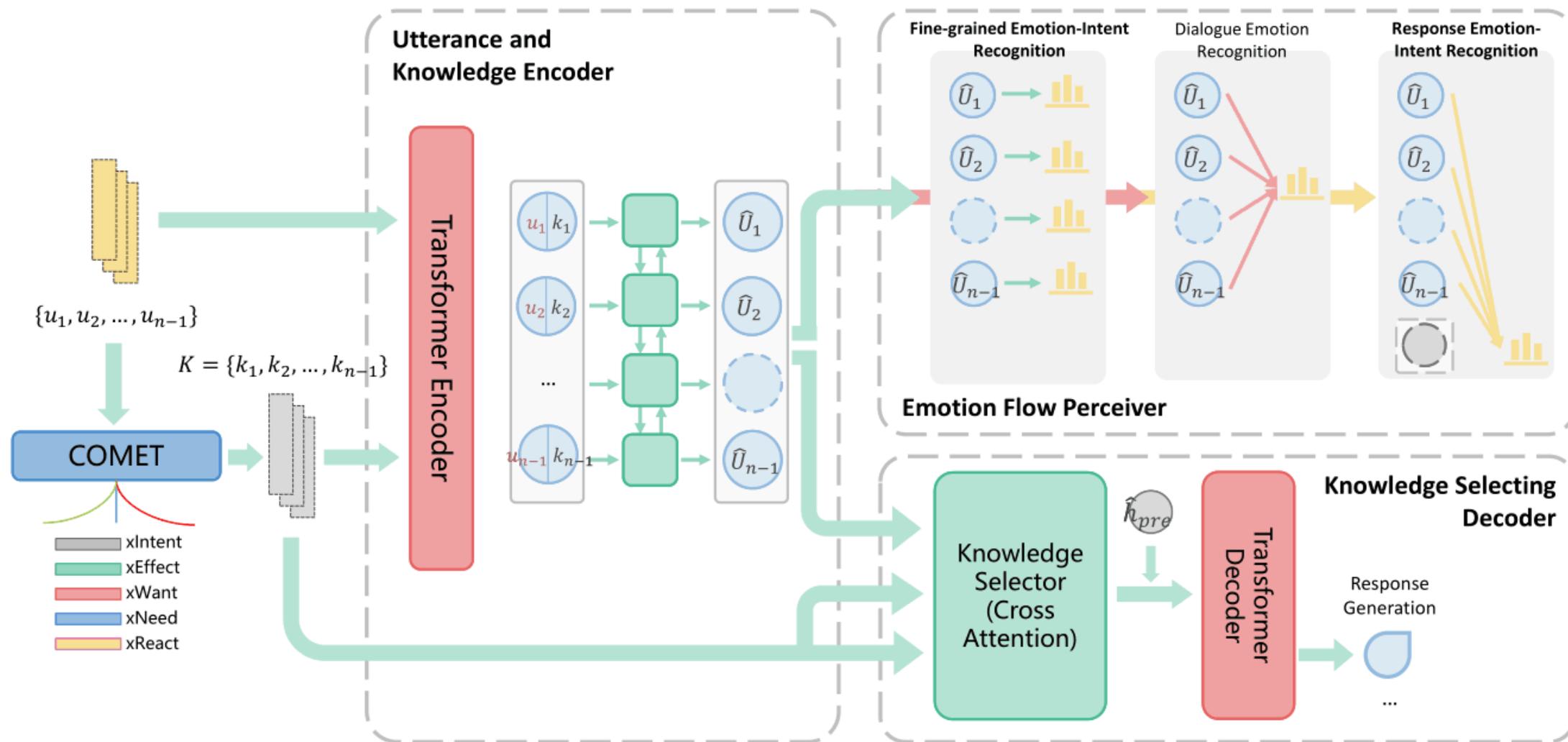


Figure 2: An overall architecture of our proposed model.

# Method

Task Formulation

$$C = [C_1, \dots, C_{N-1}]$$

$$EI = [ei_1, \dots, ei_{N-1}, ei_Y]$$

Utterance and Knowledge Encoder

Utterance Encoding

$$C_i = [w_{CLS}, w_1, w_2, \dots, w_{L_i}]$$

$$H_{U_i} = \mathbf{TRS}_{Enc}(EMB_{C_i}), \quad (1)$$

$$U_i = H_{U_i}[0]. \quad (2)$$

Knowledge Encoding

$$H_{K_i} = \mathbf{TRS}_{Enc}(\mathcal{K}_i) \quad (3)$$

$$K_i = \mathbf{Mean}(H_{K_i})$$

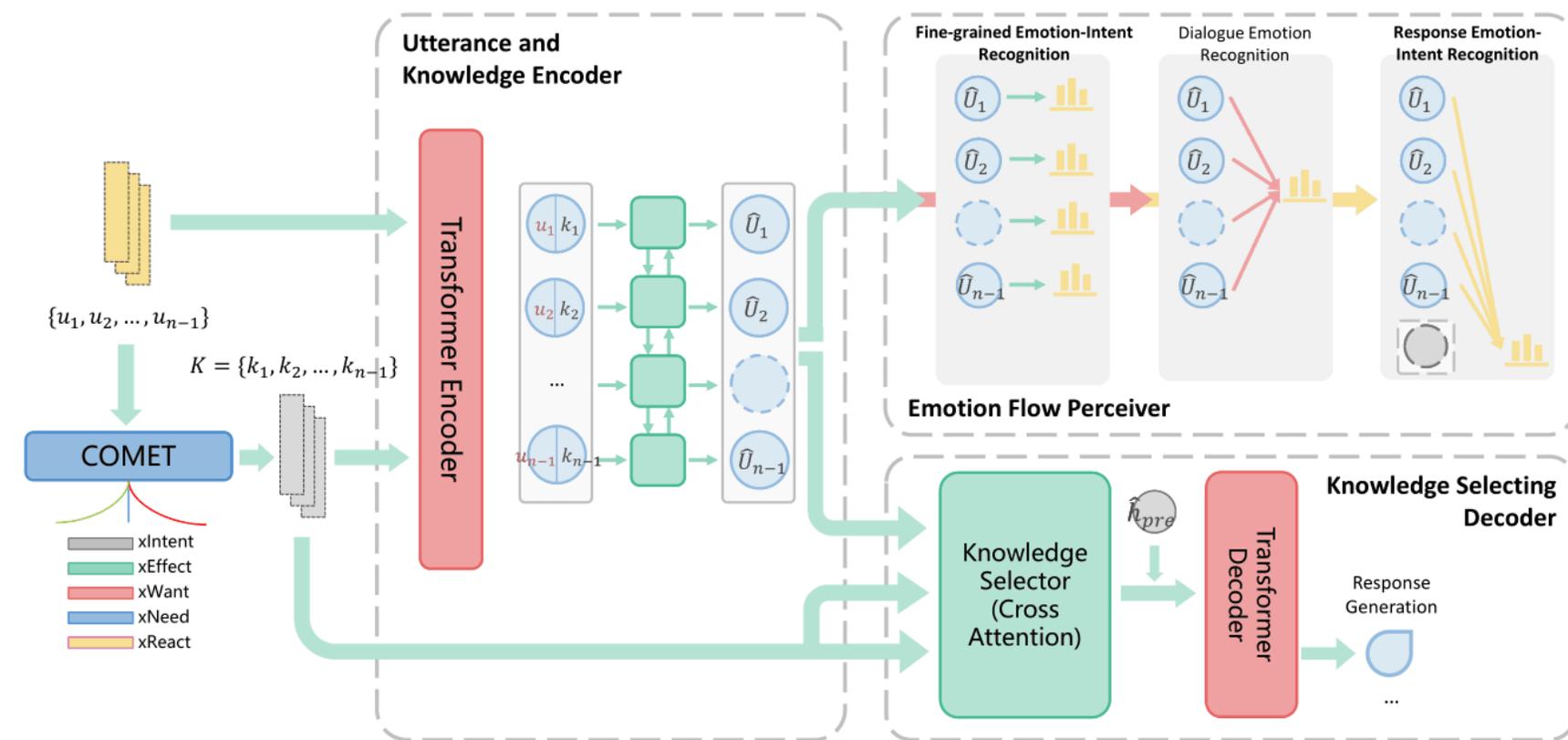


Figure 2: An overall architecture of our proposed model.

# Method

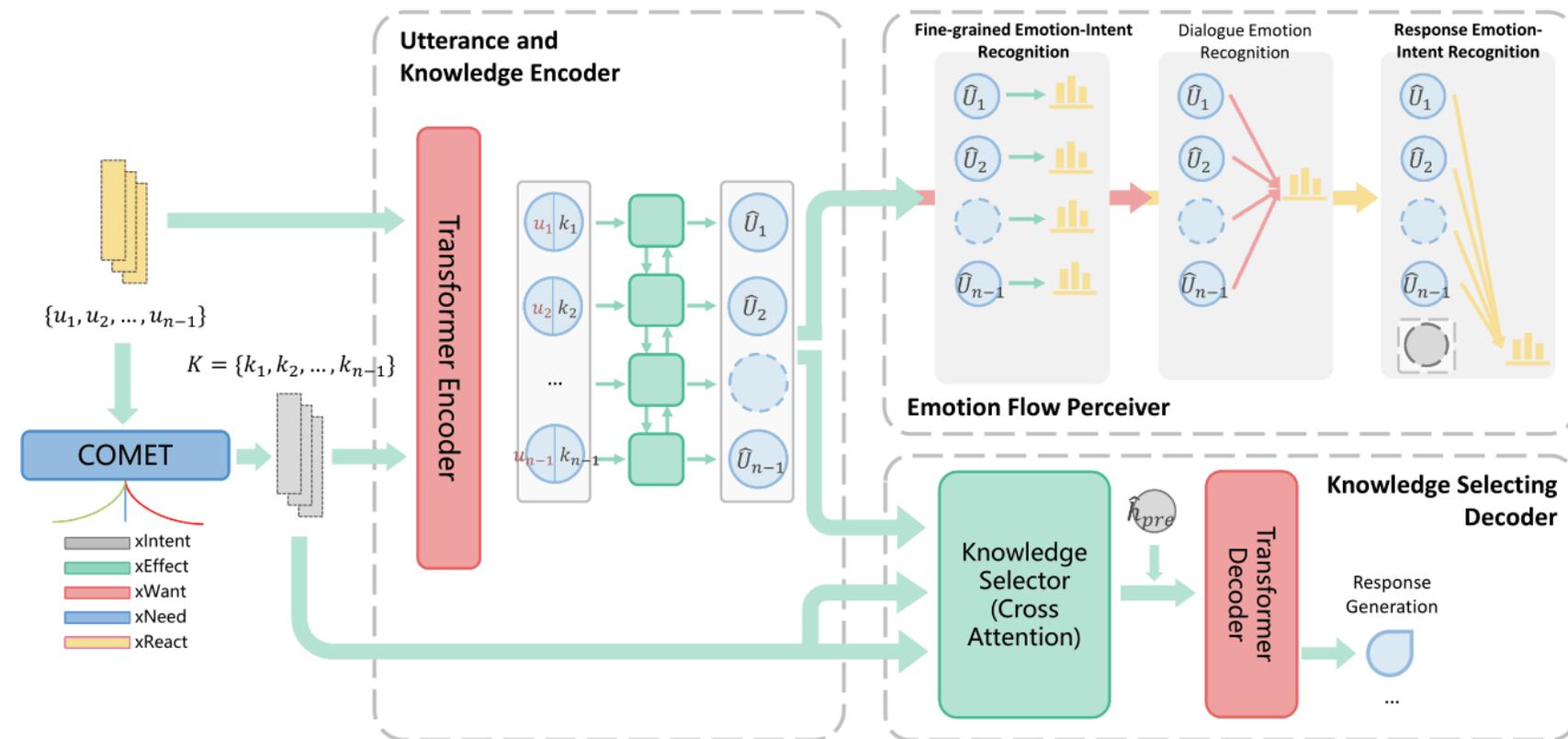


Figure 2: An overall architecture of our proposed model.

## Emotion Flow Perceiver

$$\mathbf{a}_i = [\mathbf{U}_i; \mathbf{K}_i], \quad (4)$$

$$\hat{U}_i = \text{BiLSTM}(\mathbf{W}_a \mathbf{a}_i),$$

Fine-grained Emotion Recognition

$$P_{tag}(\mathbf{e}i_i) = \text{Softmax}(\mathbf{W}_e \hat{U}_i) \quad (5)$$

$$\mathcal{L}_{emo} = - \sum_{i=1}^{N-1} \log(P_{tag}(\mathbf{e}i_i)). \quad (6)$$

Response Emotion-Intent Prediction

$$\hat{\mathbf{h}}_{pre} = \text{attention}([\hat{U}_1, \hat{U}_2, \dots, \hat{U}_{N-1}]), \quad (7)$$

$$P_{pre} = \text{Softmax}(\mathbf{W}_p \hat{\mathbf{h}}_{pre}),$$

$$\mathcal{L}_{pre} = -\log(P_{pre}(\mathbf{e}i_N)). \quad (8)$$

# Method

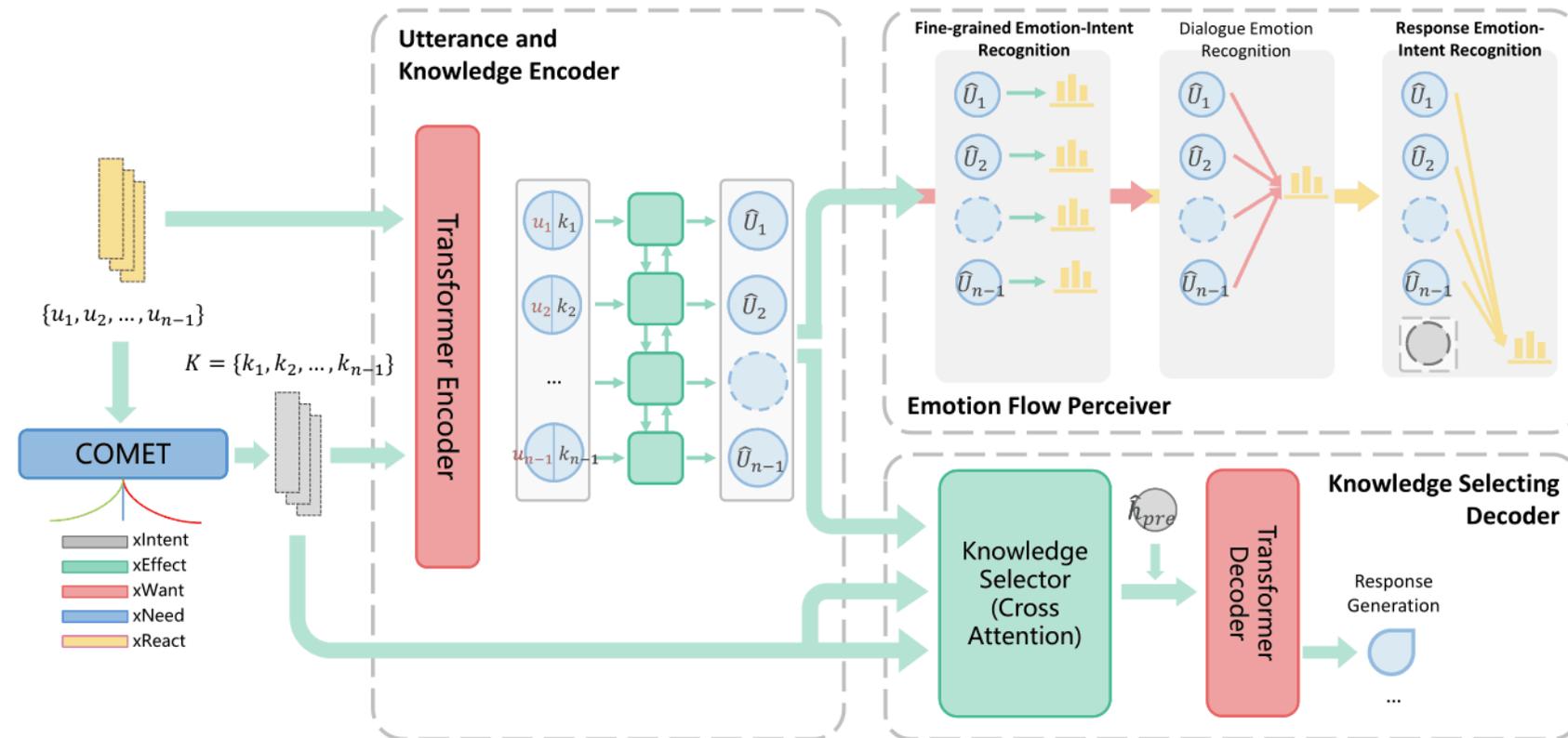


Figure 2: An overall architecture of our proposed model.

Dialogue Emotion Recognition

$$\hat{\mathbf{h}}_{dia} = \text{attention}([\hat{U}_1, \hat{U}_2, \dots, \hat{U}_{N-1}]), \quad (9)$$

$$P_{dia} = \text{Softmax}(\mathbf{W}_d \hat{\mathbf{h}}_{dia}),$$

$$\mathcal{L}_{dia} = -\log(P_{dia}(e^*)). \quad (10)$$

Knowledge Selecting Decoder

$$\mathcal{S} = \text{Cross-Attention}(\hat{U}, \mathcal{K}, \mathcal{K}), \quad (11)$$

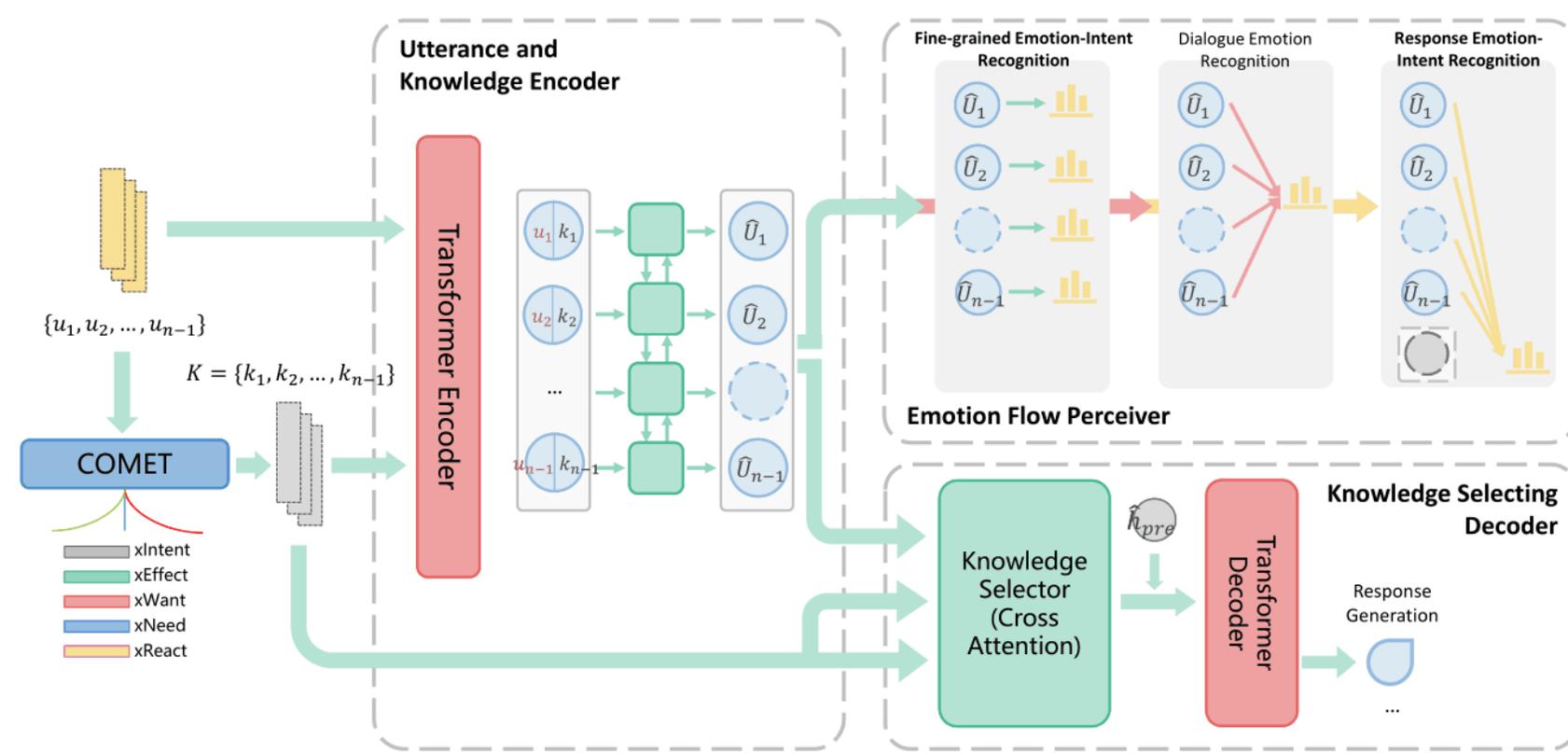
$$\mathcal{S} = \text{pooling}(\mathcal{S}). \quad (12)$$

$$[\text{SOS}] = \mathbf{W}_k([\mathcal{S}; \hat{\mathbf{h}}_{pre}]) \quad (13)$$

$$Y = [[\text{SOS}], y_1, \dots, y_T]$$

$$\mathcal{L}_{nll} = -\sum_{t=1}^T \log(P(y_t | C, y_{<t})). \quad (14)$$

# Method



## Training Objectives

$$\mathcal{L}_{cls} = \mathcal{L}_{tag_{emo}} + \mathcal{L}_{pre} + \mathcal{L}_{dia}. \quad (15)$$

$$\mathcal{L}_{div} = - \sum_{t=1}^T \sum_{i=1}^V w_i \delta_t(c_i) \log(P(y_t | C, y_{<t})), \quad (16)$$

$$\mathcal{L} = \alpha \mathcal{L}_{nll} + \beta \mathcal{L}_{cls} + \gamma \mathcal{L}_{div}, \quad (17)$$

Figure 2: An overall architecture of our proposed model.



# Experiment

Models	PPL	Dist-1	Dist-2	DE Acc.	UEI Acc.	REI Acc.
MIME	37.08	0.31	1.03	29.38	-	-
EmpDG	37.77	0.59	2.48	30.03	-	-
KEMP	<b>36.89</b>	0.61	2.65	37.58	-	-
CEM	37.03	0.66	2.99	36.44	-	-
SEEK	37.09	<b>0.73</b>	<b>3.23</b>	<b>41.85</b>	<b>34.08</b>	<b>25.67</b>

Table 1: Automatic Evaluation results of baselines and our model. The improvement of SEEK to four strong baselines is statistically significant (paired t-tests with p-values  $< 0.05$ ).



# Experiment

Models	PPL	Dist-1	Dist-2	DE Acc.	UEI Acc.	REI Acc.
<b>SEEK</b>	<b>37.09</b>	<b>0.73</b>	<b>3.23</b>	<b>41.85</b>	34.08	25.67
w/o Utter	37.37	0.70	3.13	38.9	-	<b>30.41</b>
w/o Res	37.97	0.63	2.74	40.82	<b>50.48</b>	-
w/o Utter & Res	38.48	0.60	2.70	39.7	-	-
w/o Emo	37.67	0.61	2.66	41.27	35.88	23.37
w/o Know	37.35	0.31	1.19	41.07	33.53	25.58
+ Others know	37.50	6.90	2.88	38.25	34.43	24.32
+ Context Enc	38.68	0.67	2.60	41.81	32.86	24.45

Table 2: Ablation study of our proposed model SEEK. The best results are marked with bold.

# Experiment

Models	Coh.	Emp.	Flu.
<b>MIME</b>	2.84	2.97	2.87
<b>EmpDG</b>	2.85	2.78	2.76
<b>KEMP</b>	2.73	2.80	2.80
<b>CEM</b>	2.82	2.99	2.75
<b>SEEK</b>	<b>2.91</b>	<b>3.02</b>	<b>3.07</b>

Table 3: Human evaluation results. We apply Fleiss’s Kappa, denoted as  $\kappa$ , to measure inter-annotator agreement, where  $0.4 < \kappa < 0.6$  indicates moderate agreement.

Comparisons	Aspects	Win	Lose	Tie
SEEK vs. MIME	Coh.	<b>24.3</b>	17.1	58.6
	Emp.	<b>31.4</b>	22.2	46.4
	Flu.	<b>28.6</b>	25.9	45.5
SEEK vs. EmpDG	Coh.	<b>32.1</b>	26.3	41.6
	Emp.	<b>35.5</b>	27.4	37.1
	Flu.	<b>26.9</b>	22.3	50.8
SEEK vs. KEMP	Coh.	<b>29.2</b>	25.2	45.6
	Emp.	<b>28.8</b>	19.9	51.3
	Flu.	<b>38.7</b>	15.6	45.7
SEEK vs. CEM	Coh.	<b>27.3</b>	24.8	47.9
	Emp.	<b>33.4</b>	27.5	39.1
	Flu.	<b>35.7</b>	21.6	42.7

Table 4: Human A/B test (%) on the three aspects: coherence, empathy, and fluency. The comparison results directly shows SEEK outperforms the four baselines models.

# Experiment

<b>Context</b>	<i>Speaker:</i> I love YouTube. I've been listening to all my <b>classic</b> tracks. Tupac forever. ( <b>Nostalgic</b> )
	<i>Listener:</i> I love me some Tupac. Real talk. ( <b>Acknowledging</b> )
	<i>Speaker:</i> I started out with <i>One Hit Wonders</i> but ended up at Pac. <b>I miss my youth lol.</b> ( <b>Nostalgic</b> )
<b>Speaker Emotion</b>	Nostalgic
<b>MIME</b>	I am sure you will be able to get a new one.
<b>EmpDG</b>	I am sure you will be fine.
<b>KEMP</b>	I am sure you will be fine.
<b>CEM</b>	That sounds like fun!
<b>SEEK</b>	I do not like that, but i think they will be able to <b>get the old times</b> . ( <b>Nostalgic</b> )
<b>Golden</b>	Me too. I always end up on the 90s hip hop station on Pandora. ( <b>Nostalgic</b> )
<b>Context</b>	<i>Speaker:</i> Yeah about 10 years ago I had a horrifying experience. It was 100% their fault, but they hit the water barrels and survived. They had no injuries, but they almost ran me off the road. ( <b>Guilty</b> )
	<i>Listener:</i> Did you suffer any injuries? ( <b>Questioning</b> )
	<i>Speaker:</i> No, I was not hit. It turned out they were drunk. I felt guilty, but realized it was his fault. ( <b>Guilty</b> )
<b>Speaker Emotion</b>	Guilty
<b>MIME</b>	I would be too. I would be so scared.
<b>EmpDG</b>	I am sorry to hear that.
<b>KEMP</b>	I am sorry to hear that.
<b>CEM</b>	That is so sad.
<b>SEEK</b>	That is pretty scary! I am <b>glad</b> that you were able to <b>get in out!</b> ( <b>Sympathizing</b> )
<b>Golden</b>	Why did you feel guilty? People really shouldn't drive drunk. ( <b>Questioning</b> )

Table 5: Two cases of generated responses by SEEK and the baselines. We annotated each turn with the emotional or intentional labels at the end of the utterances. The words relevant to the predicted labels in SEEK's response are highlighted in red.

# Experiment

Type	x_intent	x_need	x_want	x_effect	x_react
Knowledge	to see the <b>baby</b>	to have an ultrasound	to see what the <b>baby</b> is	to see the <b>baby</b>	<b>happy</b>
	to know the gender	to see the ultrasound	to show it to their friends	to see the <b>gender</b>	excited
	to know the sex	to have the ultrasound	to show it to everyone	to see the ultrasound	surprised
	to be informed	to have a <b>baby</b>	to show it to others	to be happy	joyful
	none	to get the ultrasound	to see the <b>baby</b>	we get excited	relieved
Context	We asked the doc to put the ultrasound in an envelope so we could record our reaction to the gender reveal. I was very happy when I finally saw it! ( <b>Excited</b> )				
SEEK	<b>Congratulations!</b>				
Gold	<b>Congrats!</b> what gender did your child end up being?				

Table 6: The visualization of the cross-attention weights of selecting knowledge in SEEK.



# Thank you!



gesis  
Leibniz-Institut  
für Sozialwissenschaften

