

# Reasoning before Responding: Integrating Commonsense-based Causality Explanation for Empathetic Response Generation

Yahui Fu, Koji Inoue, Chenhui Chu, Tatsuya Kawahara

KYOTO UNIVERSITY

京都大学



# Research Background

## What's empathy?

Empathy is a desirable capacity of humans to place themselves in another's position to **show understanding** of his/her **experience and feelings**.

## Why empathy?

An empathetic dialogue system can serve as chit-chat friends for **companion**, psychologists for **health care**, etc.



Figure 1: Empathetic dialogue system by multi-modality avatar Gene<sup>[1]</sup>

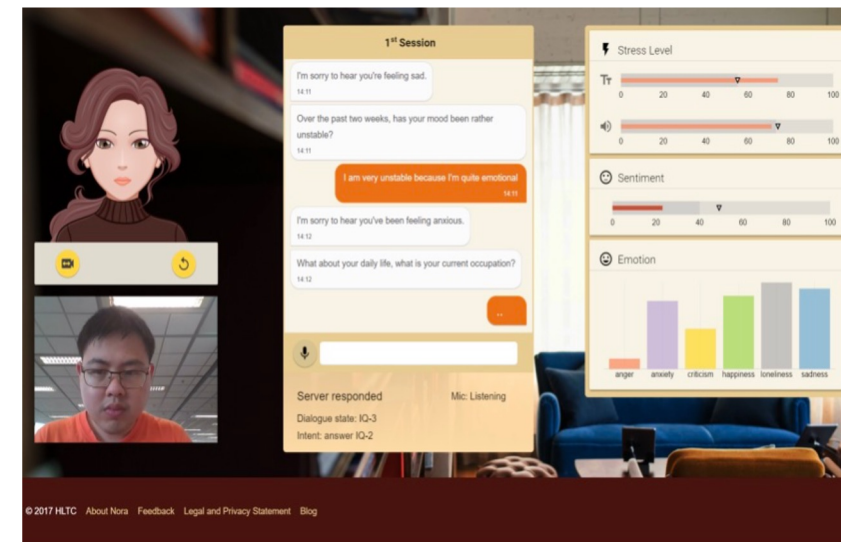


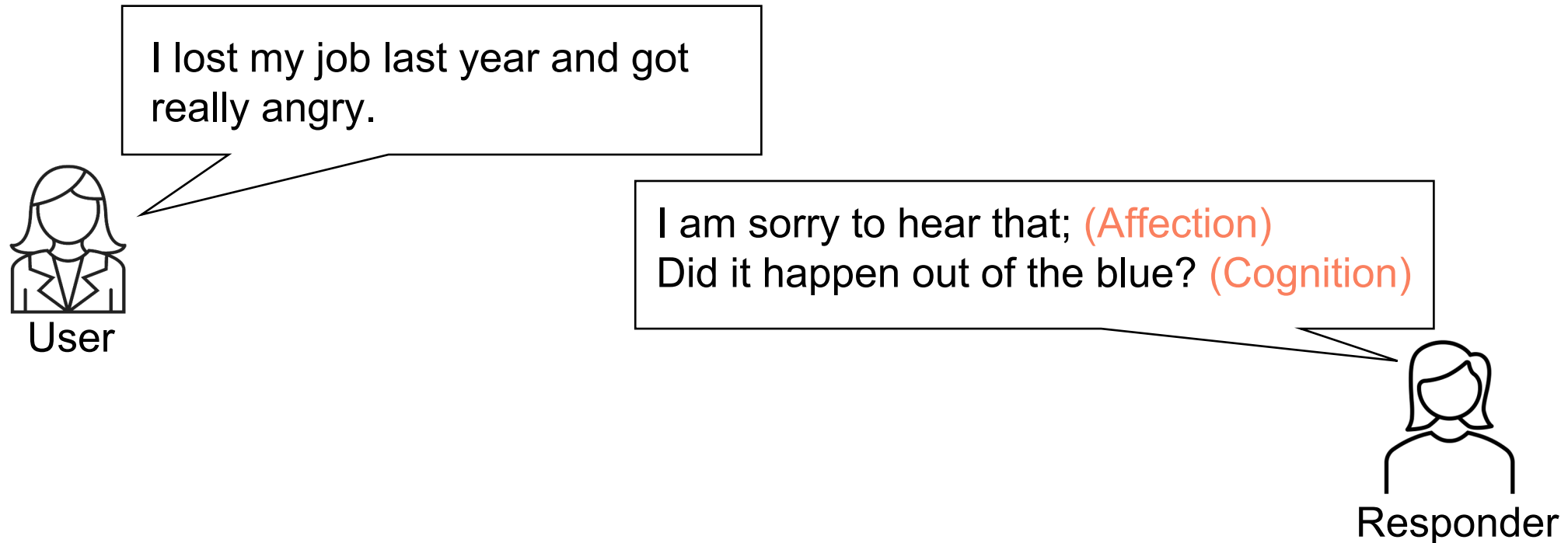
Figure 2: Nora, the empathetic psychologist<sup>[2]</sup>

# Research Background

## *How to express empathy?*

Empathy includes two aspects: Cognition and Affection<sup>[3]</sup>.

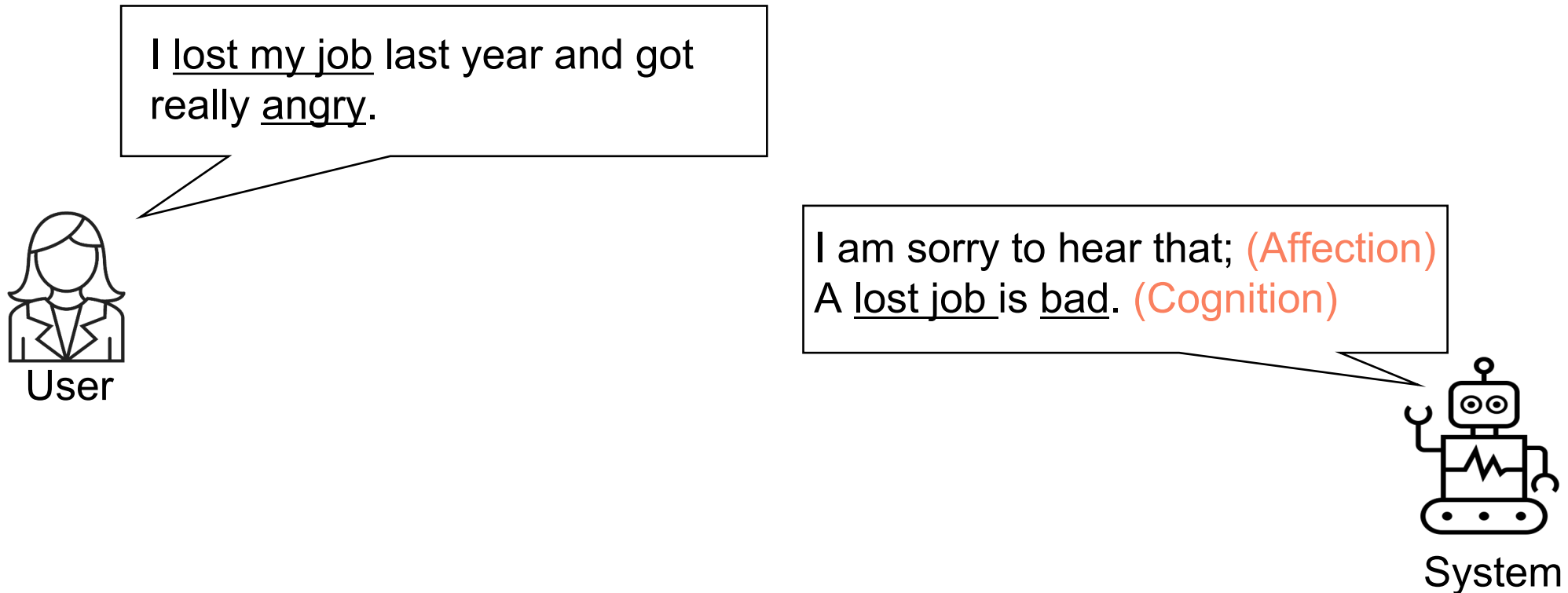
- ❖ **Cognition**: understand the other person's perspective and situation.
- ❖ **Affection**: express suitable emotion



# Research Background

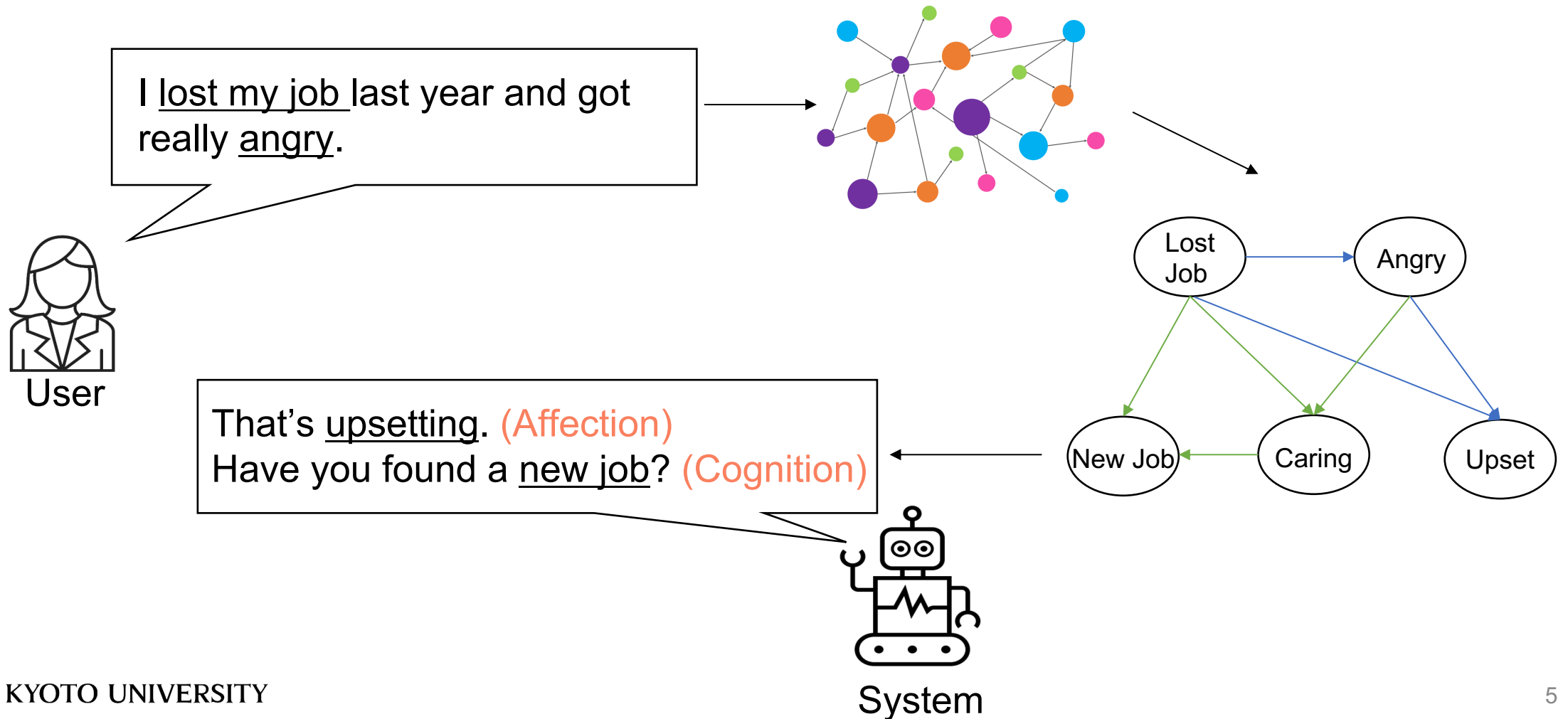
## *How to generate empathetic response?*

*A case with no causality explanation, generating an empathetic response based on context information.*



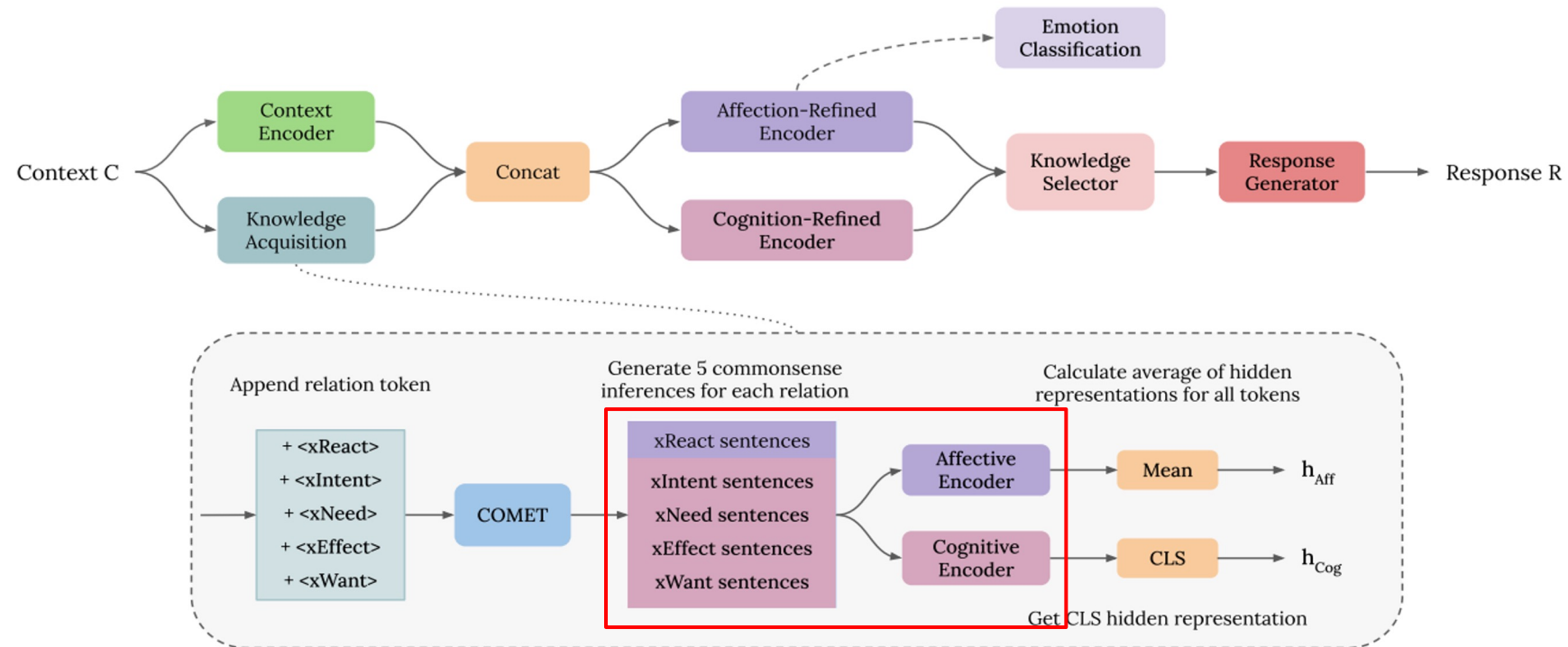
# Research Background

*A case with causality explanation, generating an empathetic response based on knowledge reasoning.*



# Related Work

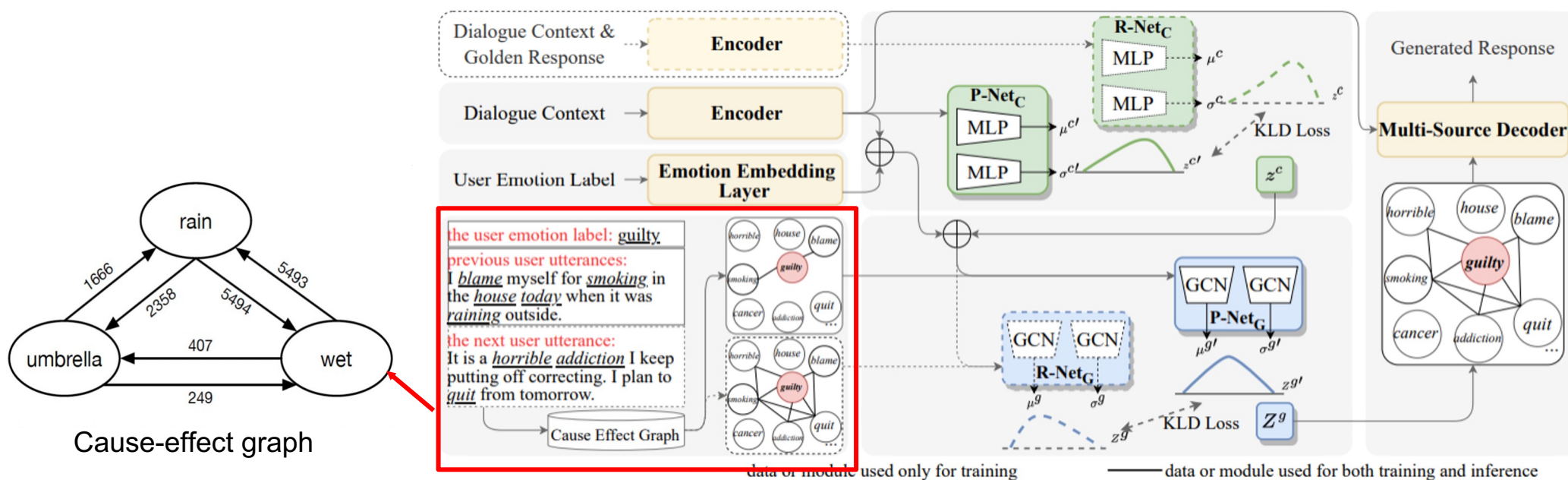
[Sabour et al AAI 2022] used a knowledge model COMET to obtain the user's *react and situation* for *affective and cognitive* encoding.



**Weakness:** concatenated related knowledges, no reasoning process.

# Related Work

[Wang et al EMNLP 2022] used **cause-effect graph** to build the causality interdependence between user's emotion to user's context, and user's emotion to system's response.



**Weakness:** It only reasoned casualties to the user's emotion, did not reason more fine-grained user's want and system's intent.

# Motivation

Exploring user's perspective:

**Affection:** angry

**Desire:** to get a new job.

System's intention is aligned  
with user's desire:

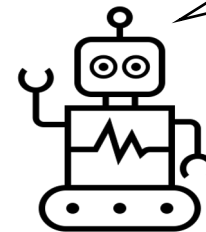
**Affection:** sad

**Intention:** to give a new job.



User

I lost my job last year and got really angry.



System

I am sorry to hear that (**Affection**);  
I wish can give you a new job!  
(**Cognition**)



# Motivation

**Reasoning user's perspective:**

**Affection:** angry

**Desire:** want to complain.



User

I lost my job last year and got really angry.

**Reasoning responder's perspective to mimic humans:**

**Affection:** sad

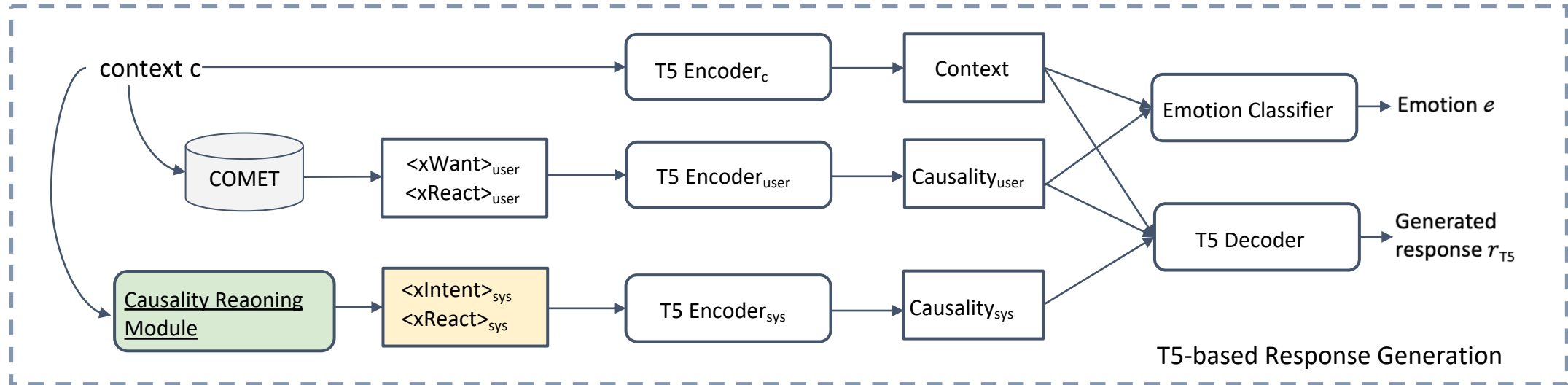
**Intention:** to know what happened.



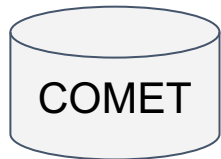
Responder

I am sorry to hear that (**Affection**);  
Did it happen out of the blue?  
(**Cognition**)

# Proposed Method

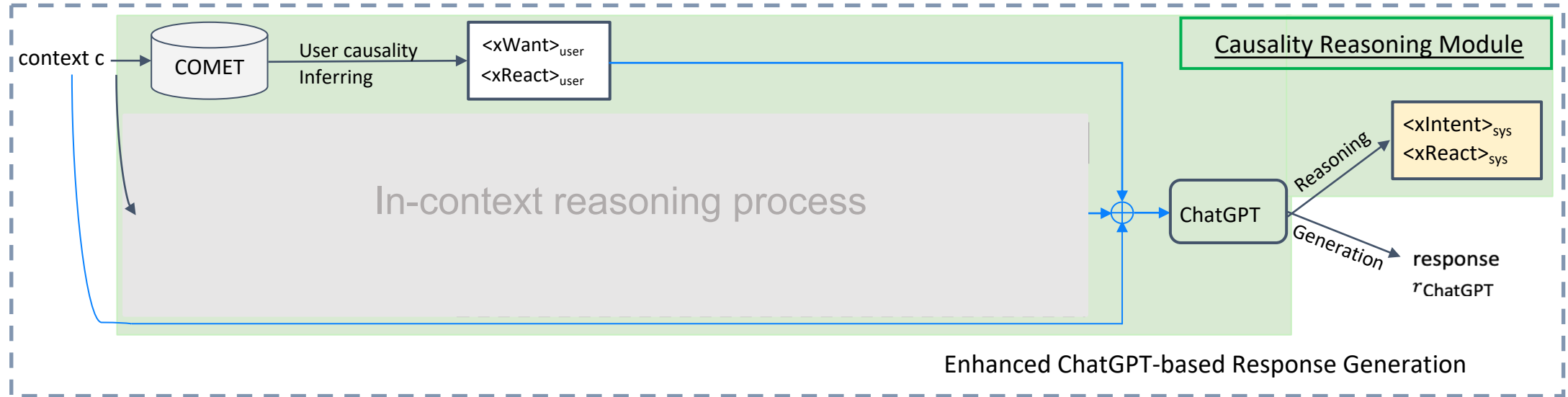


- For the context  $c$ , using COMET to predict user's want/react, and causality reasoning module to predict system's intent/react.



- It is a BART-based model which is fine-tuned on the **cause-effect** graph from ATOMIC-2020 dataset.

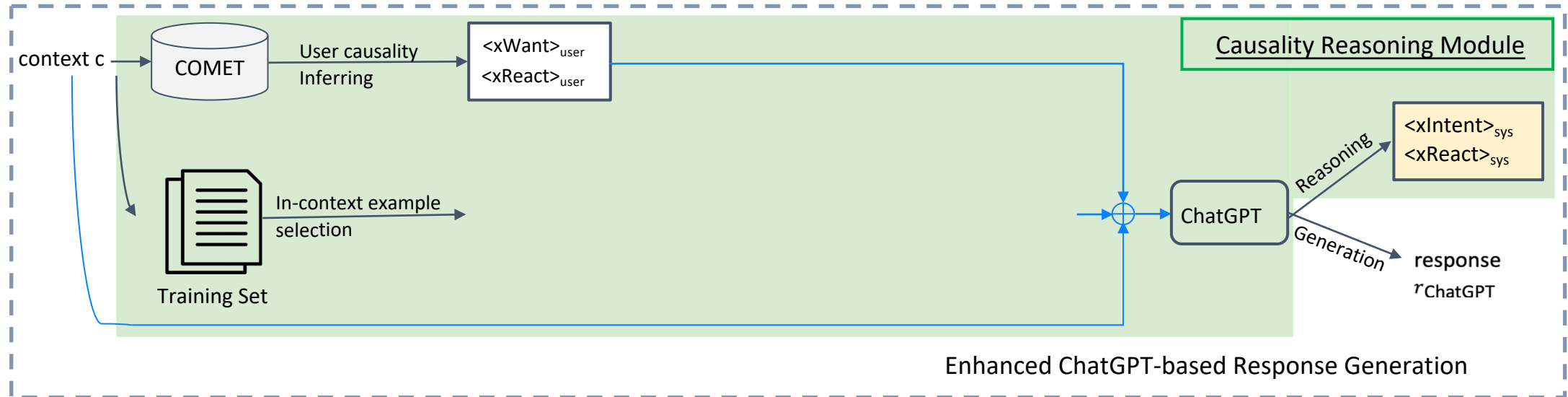
# Proposed Method



Input: *context c; user's want/react, outputs of in-context reasoning process*

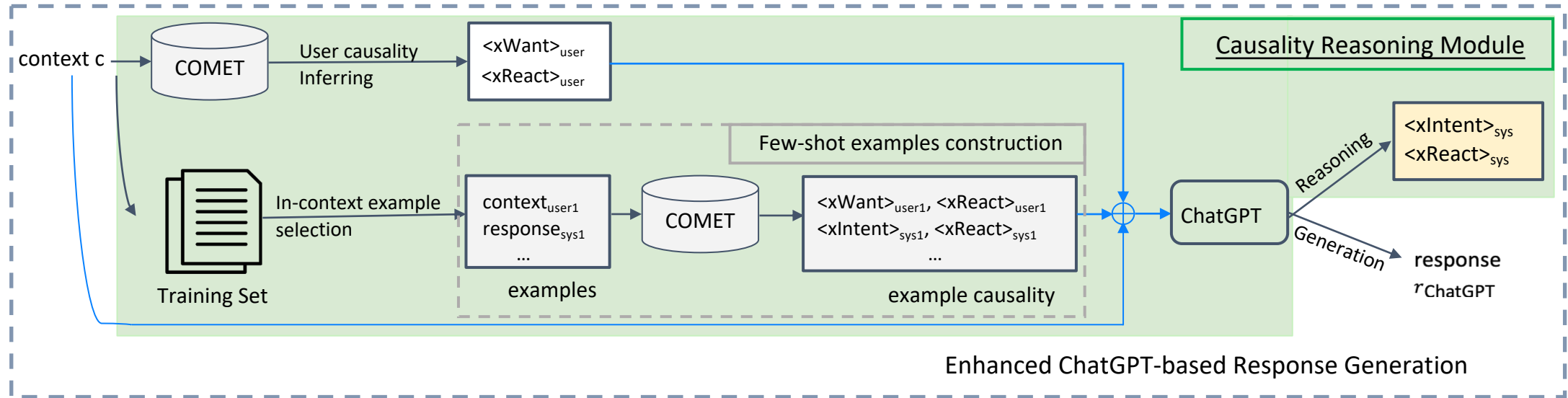
Output: *system's intention and reaction; response*

# Proposed Method



- Select top- $k$  conversations from training set based on cosine similarity as the in-context example.

# Proposed Method



- For the <context, response> in each selected example, predict user's want/react and system's intention/reaction as the example causality.

# In-context reasoning

Test input user: I'm so excited because I'm finally going to visit my parents next month! I didn't see them for 3 years.

Predictions user wants: to spend time with family; to have fun with them; to see them again.  
user reacts to: excited; happy; nostalgic; anxious; joyful.

*User causality referring*

Few-shot1	context1	user1: Someone is visiting me soon and I can't wait! sys1: Who is it? user1: My mom, she is amazing.
	example causality	user1 wants: to have a good time; to talk to their mom; to have fun with Mom. user1 reacts to: excited; happy; satisfied; good; loved. sys1's intent: to be with her; to be loved; to be nice; happy. sys1 reacts to: happy; excited; proud; good; loving.
	response1	sys1: I bet she is! I am so glad you get to see her. Mom's are awesome!
Few-shot2	context2	user2: My family is coming to visit! sys2: Awesome. When are they coming and for how long? user2: They are coming next year from Africa!
	example causality	user2 wants: to have a good time; to go to the airport; to have fun with the family. user2 reacts to: happy; excited; happy; excited; loved. sys2's intent: to see the sights; to be with family; to be with them; to have fun. sys2 reacts to: happy; excited; satisfied; tired; relieved.
	response2	sys2: That's a long trip. I hope they have a good time.

*In-context reasoning process*

Reasoning sys's intent: to be supportive; to be happy for them; to ask about the visit.  
sys reacts to: happy; excited; curious; supportive; interested.  
response: That's great news! I'm so happy for you. What are you planning to do when you visit them?

*Resoned results*

# Experiment: Dataset

**Dataset:** EmpatheticDialogue [6]

25k empathetic conversations with 32 emotion labels.

The ratio for training/validation/test is 8:1:1.

# Experiment: Number of few-shots

	EMOACC	IP	EX	ER
$k=2$	0.24	0.08	<b>0.57</b>	<b>1.10</b>
$k=3$	0.25	0.09	0.48	1.05
$k=4$	<b>0.27</b>	0.09	0.40	1.04
$k=5$	0.25	<b>0.10</b>	0.33	1.00
$k=6$	0.25	0.08	0.32	1.01

- EMOACC = Emotion accuracy, measured by a fine-tuned BERT-base model on the EmpatheticDialogue dataset.
- IP, EX, ER is measured by separately fine-tune pre-trained empathy identification models for each metric<sup>[7]</sup>.
- IP = Interpretation
- EX= Exploration
- ER= Emotion reaction



# Experiment: Results on ChatGPT

Results of automatic evaluations for single-turn.

Method	Empathy				Coherence		
	EMOACC	ER	IP	EX	PBERT	RBERT	FBERT
k=2 ChatGPT	0.060	0.923	0.073	0.341	0.877	0.872	0.875
ChatGPT+Causality <sub>user,sys</sub>	<b>0.280</b>	<b>1.116</b>	<b>0.104</b>	<b>0.768</b>	<b>0.886</b>	<b>0.878</b>	<b>0.882</b>

Results of automatic evaluations for multi-turn.

Method	Empathy				Coherence		
	EMOACC	ER	IP	EX	PBERT	RBERT	FBERT
k=2 ChatGPT	0.083	0.917	<b>0.065</b>	0.318	0.891	0.902	0.894
ChatGPT+Causality <sub>user,sys</sub>	<b>0.199</b>	<b>1.094</b>	0.058	<b>0.397</b>	<b>0.899</b>	<b>0.907</b>	<b>0.901</b>

Emotion expression

Cognition

- Compared with ChatGPT, ChatGPT with causality explanation can generate response with appropriate emotion and contents.

# Experiment: Results on ChatGPT

Results of *human A/B test* evaluations.

Emp., Coh., Inf. refer to **Empathy**, **Coherence**, and **Informativeness**

Comparisons	Aspects	Win	Loss	Tie
ChatGPT+Causality <sub>user,sys</sub> vs. ChatGPT ( $k=2$ )	Emp.	<b>50.7</b>	36.0	13.3
	Coh.	<b>42.7</b>	42.0	15.3
	Inf.	<b>51.3</b>	37.3	11.3

# Experiment: Results on T5

Results of automatic evaluations

	Methods	PPL ↓	BLEU-2	BLEU-3	BLEU-4	D1	D2	PBERT	RBERT	FBERT
Baselines	MOEL	37.63	8.63	4.25	2.43	0.38	1.74	86.19	85.67	85.91
	MIME	36.84	8.37	4.31	2.51	0.28	0.95	86.27	85.59	85.92
	EmpDG	38.08	7.74	4.09	2.49	0.46	1.90	86.09	85.49	85.78
	CEM	36.36	6.35	3.55	2.26	0.54	2.38	86.61	85.39	85.98
	LEMPEx	30.42	2.1	0.8	0.35	1.02	<b>10.81</b>	83.60	83.09	83.34
Ours	T5	46.13	3.59	1.94	1.15	0.49	2.82	86.69	84.07	85.35
	T5+Causality <sub>user</sub>	15.26	4.84	1.97	0.89	<b>1.08</b>	10.75	90.16	89.48	89.80
	T5+Causality <sub>user,sys</sub>	<b>13.07</b>	<b>10.53</b>	<b>6.34</b>	<b>4.06</b>	0.75	5.52	<b>92.24</b>	<b>90.76</b>	<b>91.48</b>

# Experiment: Results on T5

Results of *human A/B test* evaluations.

Emp., Coh., Inf. refer to **Empathy**, **Coherence**, and **Informativeness**

Comparisons	Aspects	Win	Loss	Tie
T5+Causality <sub>user,sys</sub> vs. CEM	Emp.	<b>42.0</b>	40.0	18.0
	Coh.	<b>38.7</b>	33.3	28.0
	Inf.	38.3	<b>44.3</b>	17.3
T5+Causality <sub>user,sys</sub> vs. LEMPE <sub>x</sub>	Emp.	<b>53.0</b>	35.0	12.0
	Coh.	<b>39.0</b>	33.3	27.7
	Inf.	<b>50.0</b>	38.0	12.0

# Comparison between ChatGPT and T5

Model	Empathy				Diversity		BLEU		
	EMOACC	IP	EX	ER	D1	D2	BLEU-2	BLEU-3	BLEU-4
T5+Causality <sub>user,sys</sub>	0.125	<b>0.271</b>	0.498	0.751	0.75	5.52	<b>10.53</b>	<b>6.34</b>	<b>4.06</b>
ChatGPT+Causality <sub>user,sys</sub>	<b>0.235</b>	0.046	<b>0.668</b>	<b>1.109</b>	<b>2.91</b>	<b>16.44</b>	3.95	2.17	1.32



- Stronger interpretation (IP), which involves *understanding and empathizing through shared experiences* is more frequently observed in the trained T5 model, while ChatGPT tends to respond from the view of a machine.
- Indicating a potential need of balancing the response diversity and the accuracy in generating empathetic response.

# Case studies

	Context	user: I'm so sad because I've read an article about a newborn girl who died because her parents didn't believe in medication and doctors.
Baselines	MOEL	sys: I am so sorry to hear that. I hope you are able to find her.
	MIME	sys: That is sad. I am sorry to hear that.
	EmpDG	sys: I am sorry to hear that. I hope you find out of it.
	CEM	sys: Oh no, I am so sorry to hear that.
	LEMPEX	sys: I'm sorry to hear that. I'm so sorry to hear that.
	ChatGPT	sys: That's heartbreaking. It's important to trust medical professionals and seek help when needed.
Ours	Causality <sub>user</sub>	user wants: to tell their friends; to cry; to find out what happened. user reacts to: sad; sorrowful; depressed; upset.
	Causality <sub>sys</sub>	sys's intent: to comfort user; to show empathy; to be supportive sys reacts to: sad; upset; empathetic; sympathetic; concerned.
	T5+Causality <sub>user,sys</sub>	sys: I know what you mean, I would have been so mad!
	ChatGPT+Causality <sub>user,sys</sub>	sys: That's really heartbreaking. I'm here for you if you need to talk or just need someone to listen.

# Conclusion

- ❖ A *commonsense-based causality explanation* approach that reasons not only the user's desires/reaction but also the system's proper intention/reaction.
- ❖ Integration of T5 with ChatGPT's reasoning capability realizes more empathetic responses that result in better evaluations.
- ❖ They are *more accurate and empathetic* than the responses by ChatGPT while not so diverse.

**Thanks for you attention!**



# Q&A