



# DARER: Dual-task Temporal Relational Recurrent Reasoning Network for Joint Dialog Sentiment Classification and Act Recognition

Bowen Xing<sup>1</sup> and Ivor W. Tsang<sup>2,1</sup>

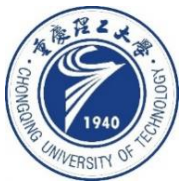
<sup>1</sup>Australian Artificial Intelligence Institute, University of Technology Sydney, Australia

<sup>2</sup>Centre for Frontier Artificial Intelligence Research, A\*STAR, Singapore

[bxing714@gmail.com](mailto:bxing714@gmail.com), [ivor\\_tsang@ihpc.a-star.edu.sg](mailto:ivor_tsang@ihpc.a-star.edu.sg)

2022. 3. 17 • ChongQing

— ACL2022



gesis  
Leibniz-Institut  
für Sozialwissenschaften



Reported by Sijin Liu



# 1.Introduction

## 2.Method

### 3.Experiments



# Introduction

Utterances	Act	Sentiment
$u_a$ : I highly recommend it. Really awesome progression and added difficulty	Statement	Positive
$u_b$ : I never have.	Disagreement	Negative

Table 1: A dialog snippet from the Mastodon dataset.

- Previous works only consider the **parameter sharing** and **semantics-level interactions**, while the **label information** is not integrated into the dual-task interactions.
- On the other hand, previous works do not consider the **temporal relations between utterances** in dual-task reasoning, while in which they play a key role.

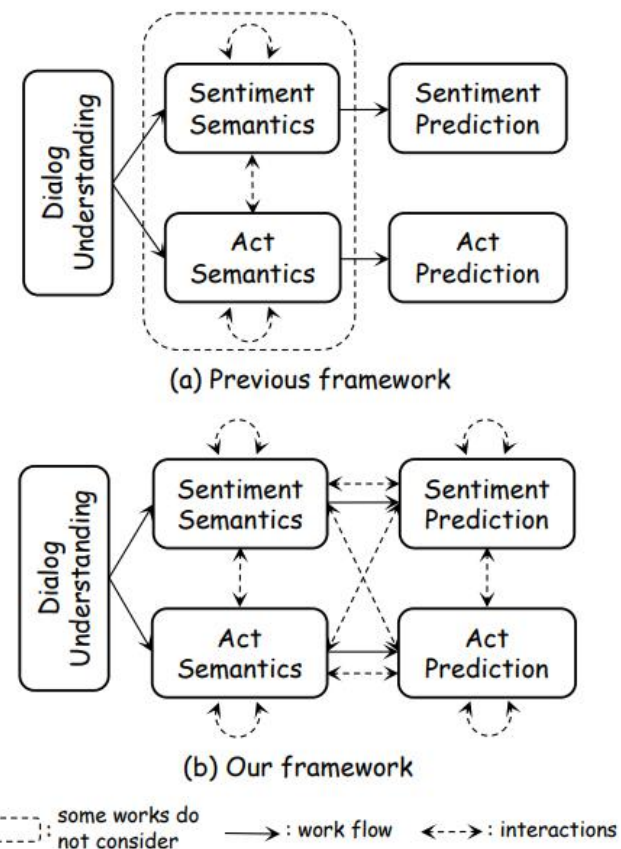


Figure 1: Illustration of previous framework and ours.



# Method

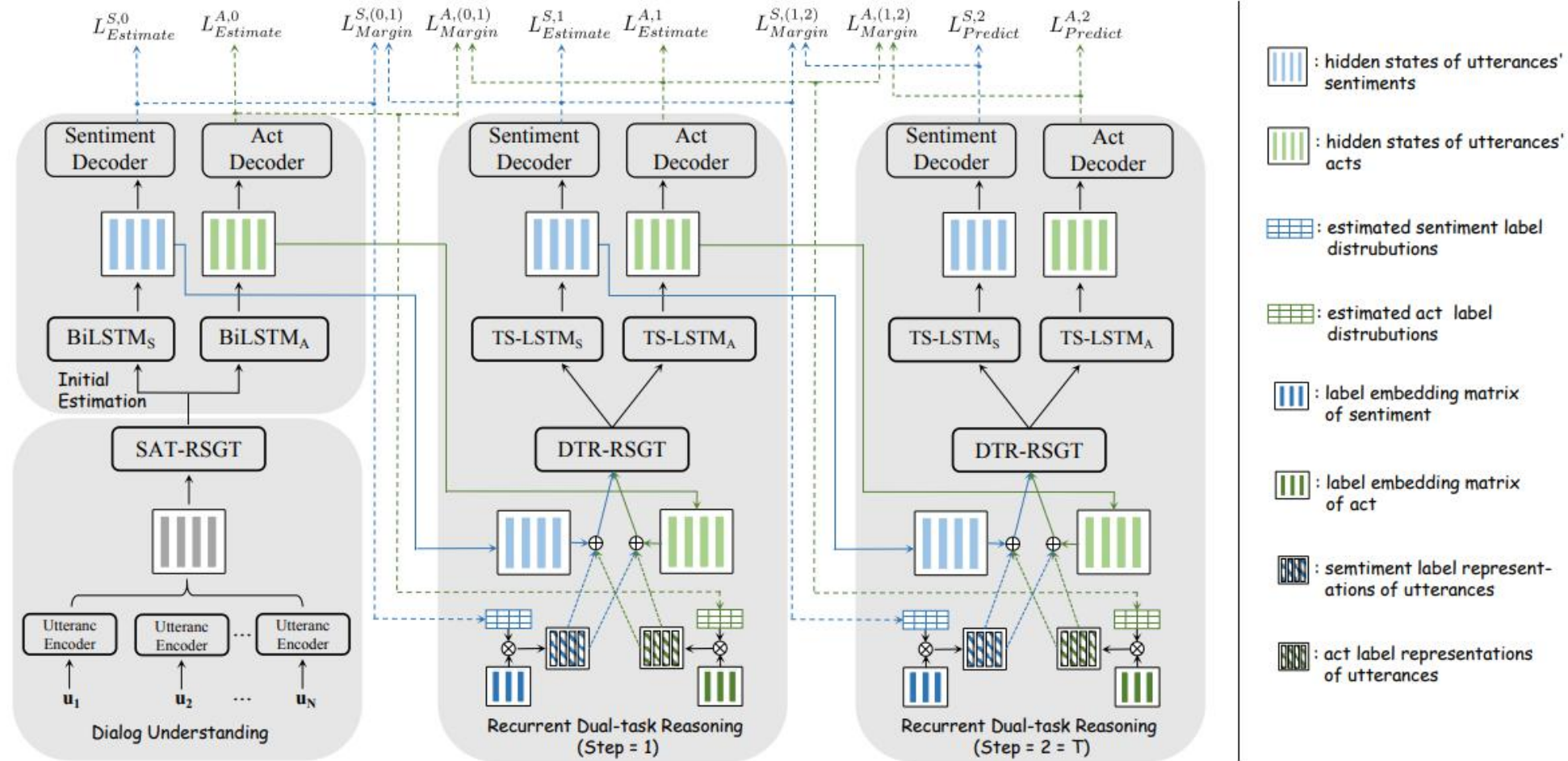
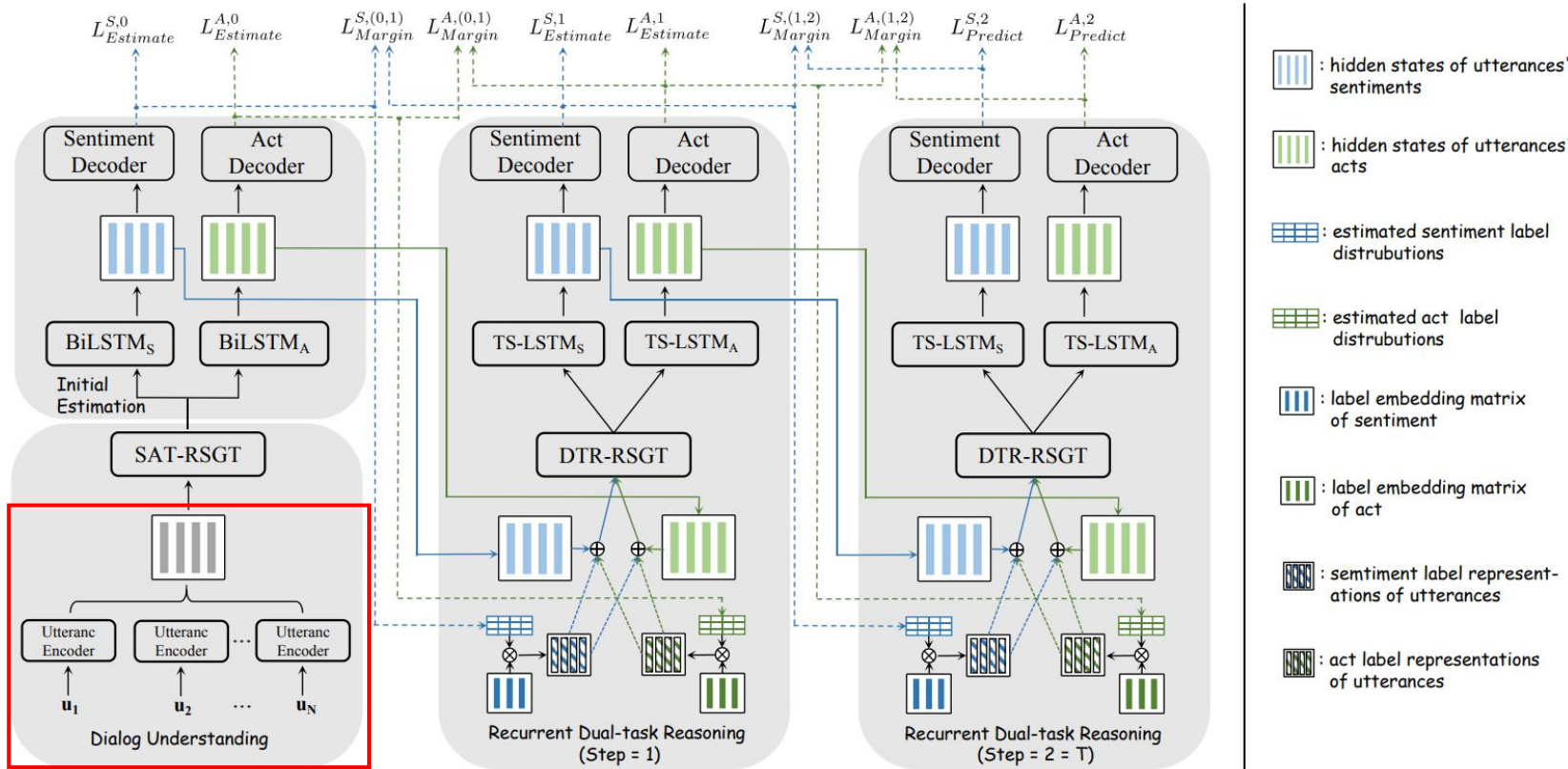


Figure 4: The network architecture of our proposed DARER model. Without loss of generality, the step number  $T$  in this illustration is set 2.

# Method



## Problem Definition

$$\mathcal{D} = \{u_1, u_2, \dots, u_N\}$$

dialog sentiment labels  $Y^S = y_1^s, \dots, y_N^s$

dialog act labels  $Y^A = y_1^a, \dots, y_N^a$

## Utterance Encoder

$$H = (h_0, \dots, h_N)$$

$$H_{u,i} = (h_{u,i}^0, \dots, h_{u,i}^{l_i})$$

# Method

$$M_i^l = \sum_{j \in \mathcal{N}_i} \alpha_{ij} W_{r_{ij}}^l H_j^l, \quad (5)$$

$r_{ij}$	1	2	3	4	5	6	7	8
$I_s(i)$	1	1	1	1	2	2	2	2
$I_s(j)$	1	1	2	2	1	1	2	2
$pos(i, j)$	>	≤	>	≤	>	≤	>	≤

Table 2: All relation types in SATG (assume there are two speakers).  $I_s(i)$  indicates the speaker node  $i$  is from.  $pos(i, j)$  indicates the relative position of node  $i$  and  $j$ .

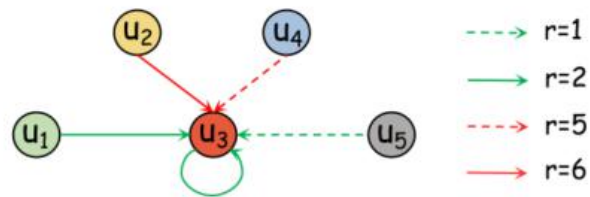


Figure 2: An example of SATG.  $u_1, u_3$  and  $u_5$  are from speaker 1 while  $u_2$  and  $u_4$  are from speaker 2. w.l.o.g, only the edges directed into  $u_3$  node are illustrated.

Speaker-aware Temporal relation-specific graph transformations

## Speaker-aware Temporal RSGT

$$\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{R})$$

$$\hat{h}_i = W_1 h_i^0 + \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}_i^r} \frac{1}{|\mathcal{N}_i^r|} W_1^r h_j^0 \quad (1)$$

$$\hat{H} = (\hat{h}_0, \dots, \hat{h}_N)$$

## Initial Estimation

$$H_s^0 = \text{BiLSTM}_S(\hat{H})$$

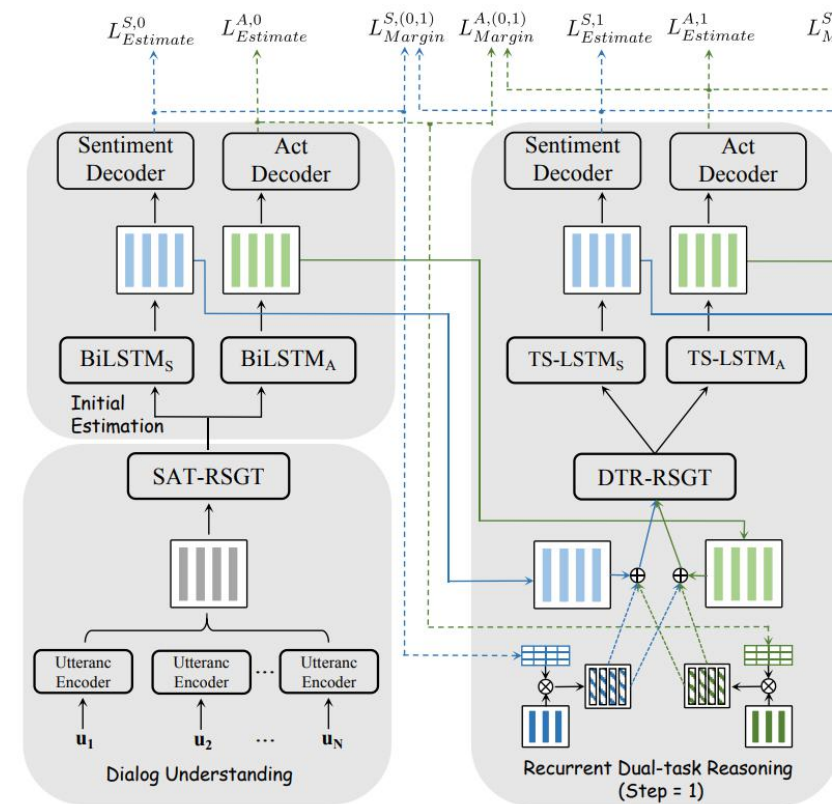
$$H_a^0 = \text{BiLSTM}_A(\hat{H})$$

where  $H_s^0 = \{h_{s,i}^0\}_{i=1}^N$  and  $H_a^0 = \{h_{a,i}^0\}_{i=1}^N$

$$\begin{aligned} P_S^0 &= \{P_{S,i}^0\}_{i=1}^N, P_A^0 = \{P_{A,i}^0\}_{i=1}^N \\ P_{S,i}^0 &= \text{softmax}(W_d^s h_{a,i}^0 + b_d^s) \\ &= [p_{s,i}^0[0], \dots, p_{s,i}^0[k], \dots, p_{s,i}^0(|\mathcal{C}_s| - 1)] \quad (2) \end{aligned}$$

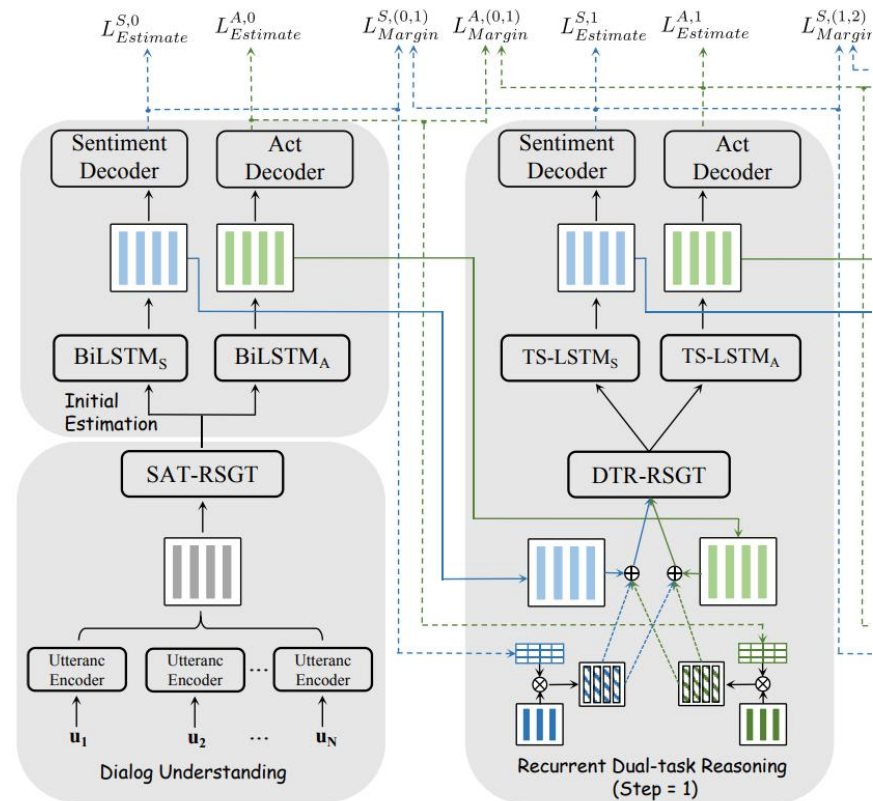
$$\begin{aligned} P_{A,i}^0 &= \text{softmax}(W_d^a h_{s,i}^0 + b_d^a) \\ &= [p_{a,i}^0[0], \dots, p_{a,i}^0[k], \dots, p_{a,i}^0(|\mathcal{C}_a| - 1)] \end{aligned}$$

where  $W_d^*$  and  $b_d^*$  are weight matrices and biases,  $\mathcal{C}_s$  and  $\mathcal{C}_a$  are sentiment class set and act class set.





# Method



$r'_{ij}$	1	2	3	4	5	6	7	8	9	10	11	12
$I_t(i)$	S	S	S	S	S	S	A	A	A	A	A	A
$I_t(j)$	S	S	S	A	A	A	S	S	S	A	A	A
$pos(i, j)$	<	=	>	<	=	>	<	=	>	<	=	>

Table 3: All relation types in DRTG.  $I_t(i)$  indicates that node  $i$  is a sentiment (S) node or act (A) node.

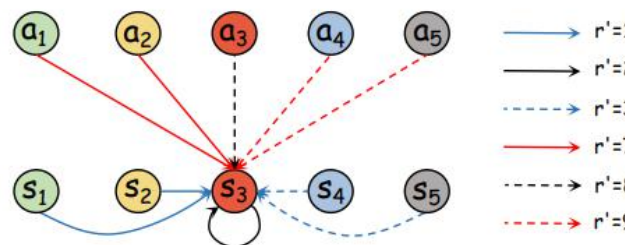


Figure 3: An example of DRTG.  $s_i$  and  $a_i$  respectively denote the node of DAC task and DAR task. w.l.o.g, only the edges directed into  $s_3$  are illustrated.

## Recurrent Dual-task Reasoning

- Projection of Label Distribution

$$e_{s,i}^t = \sum_{k=0}^{|\mathcal{C}_s|-1} p_{s,i}^{t-1}[k] \cdot v_s^k \quad (3)$$

$$e_{a,i}^t = \sum_{k'=0}^{|\mathcal{C}_a|-1} p_{a,i}^{t-1}[k'] \cdot v_a^{k'}$$

where  $v_s^k$  and  $v_a^{k'}$  are the label embeddings of sentiment class  $k$  and act class  $k'$ , respectively.

$$\hat{h}_{s,i}^t = h_{s,i}^{t-1} + e_{s,i}^t + e_{a,i}^t \quad (4)$$

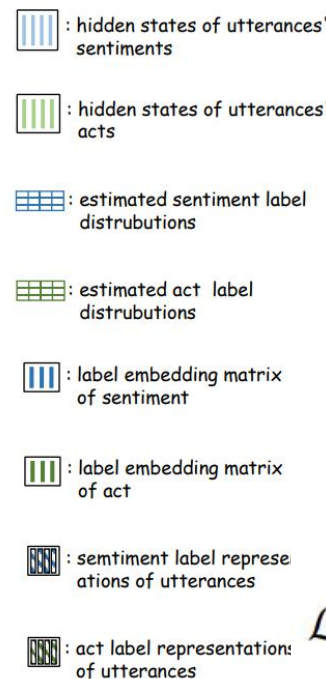
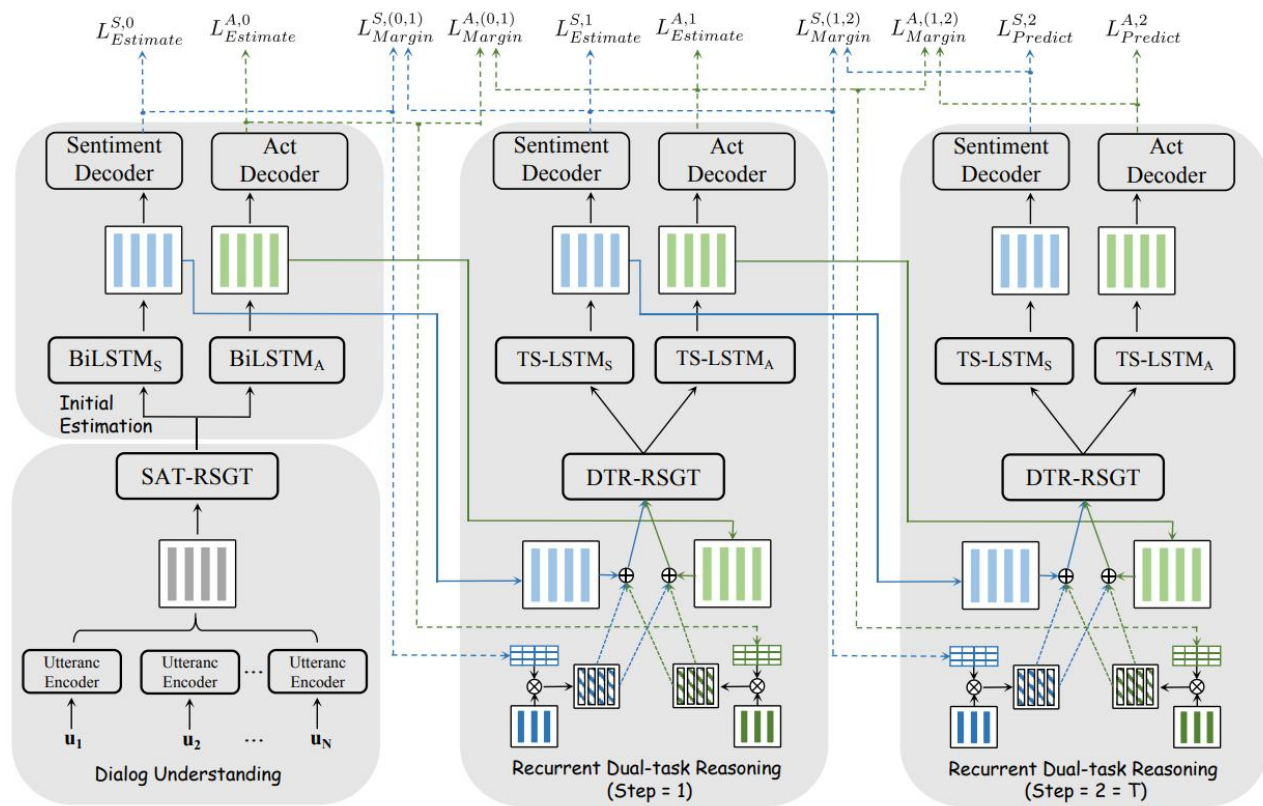
$$\hat{h}_{a,i}^t = h_{a,i}^{t-1} + e_{s,i}^t + e_{a,i}^t$$

$$\bar{h}_i^t = W_2 \hat{h}_i^t + \sum_{r \in \mathcal{R}'} \sum_{j \in \mathcal{N}_i^{r'}} \frac{1}{|\mathcal{N}_i^{r'}|} W_2^r \hat{h}_j^t \quad (5)$$

$$H_s^t = \text{TS-BiLSTM}_S(\bar{H}_s^t) \quad (6)$$

$$H_a^t = \text{TS-BiLSTM}_A(\bar{H}_a^t)$$

# Method



## Estimate Loss

$$\mathcal{L}_{Estimate}^{S,t} = \sum_{i=1}^N \sum_{k=0}^{|\mathcal{C}_s|-1} y_{s,i}^k \log(p_{s,i}^t[k]) \quad (7)$$

## Margin Loss

$$\mathcal{L}_{Margin}^{S,(t,t-1)} = \sum_{i=1}^N \sum_{k=0}^{|\mathcal{C}_s|-1} y_{s,i}^k \max(0, p_{s,i}^{t-1}[k] - p_{s,i}^t[k]) \quad (8)$$

## Constraint loss

$$\mathcal{L}_{Constraint}^S = \sum_{t=0}^{T-1} \mathcal{L}_{Estimate}^{S,t} + \gamma * \sum_{t=1}^T \mathcal{L}_{margin}^{S,(t,t-1)} \quad (9)$$

## Prediction loss

$$\mathcal{L}_{Prediction}^S = \sum_{i=1}^N \sum_{k=0}^{|\mathcal{C}_s|-1} y_{s,i}^k \log(p_{s,i}^T[k]) \quad (11)$$

## Final Training Objective

$$\mathcal{L}^S = \mathcal{L}_{Prediction}^S + \mathcal{L}_{Constraint}^S \quad (10)$$

$$\mathcal{L} = \mathcal{L}^S + \mathcal{L}^A \quad (12)$$



# Experiments

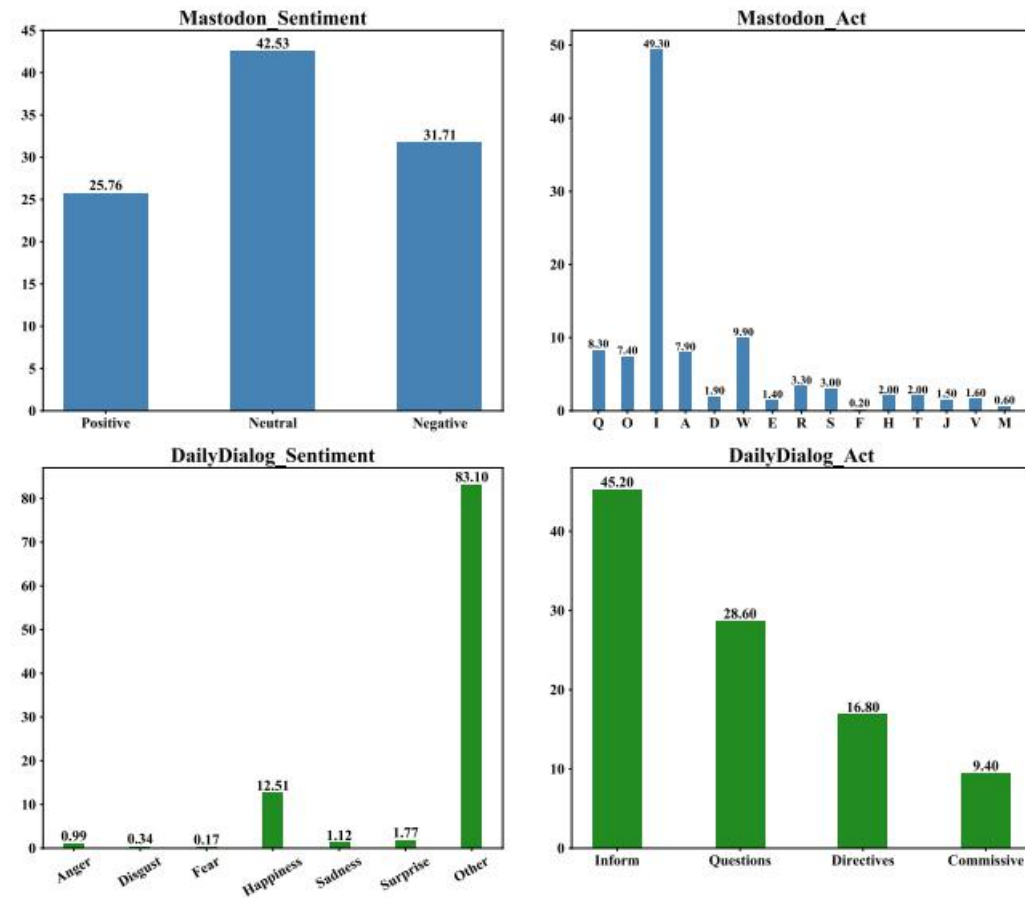


Figure 5: Illustration of class distributions.

# Experiments

Models	Mastodon						DailyDialog					
	DSC			DAR			DSC			DAR		
	P(%)	R(%)	F1(%)	P(%)	R(%)	F1(%)	P(%)	R(%)	F1(%)	P(%)	R(%)	F1(%)
JointDAS	36.1	41.6	37.6	55.6	51.9	53.2	35.4	28.8	31.2	76.2	74.5	75.1
IIM	38.7	40.1	39.4	56.3	52.2	54.3	38.9	28.5	33.0	76.5	74.9	75.7
DCR-Net	43.2	47.3	45.1	60.3	56.9	58.6	56.0	40.1	45.4	79.1	79.0	79.1
BCDCN	38.2	62.0	45.9	57.3	61.7	59.4	55.2	45.7	48.6	80.0	80.6	80.3
Co-GAT	44.0	53.2	48.1	60.4	60.6	60.5	65.9	45.3	51.0	81.0	78.1	79.4
Co-GAT*	45.40 $\pm 2.31$	48.11 $\pm 2.91$	46.47 $\pm 0.37$	62.55 $\pm 0.46$	58.66 $\pm 1.71$	60.54 $\pm 1.10$	58.04 $\pm 0.84$	44.65 $\pm 0.36$	48.82 $\pm 0.22$	79.14 $\pm 0.40$	79.71 $\pm 0.16$	79.39 $\pm 0.14$
DARER	<b>56.04</b> <sup>†</sup> $\pm 0.85$	<b>63.33</b> <sup>†</sup> $\pm 0.30$	<b>59.59</b> <sup>†</sup> $\pm 0.70$	<b>65.08</b> <sup>‡</sup> $\pm 1.25$	<b>61.88</b> <sup>†</sup> $\pm 0.37$	<b>63.43</b> <sup>†</sup> $\pm 0.85$	<b>59.96</b> <sup>‡</sup> $\pm 1.25$	<b>49.51</b> <sup>†</sup> $\pm 1.33$	<b>53.42</b> <sup>†</sup> $\pm 0.18$	<b>81.39</b> <sup>†</sup> $\pm 0.55$	<b>80.80</b> <sup>‡</sup> $\pm 0.43$	<b>81.06</b> <sup>†</sup> $\pm 0.04$

Table 4: Experiment results. \* denotes we reproduce the results using official code.  $\pm$  denotes standard deviation. <sup>†</sup> denotes that our DARER significantly outperforms Co-GAT with  $p < 0.01$  under t-test and <sup>‡</sup> denotes  $p < 0.05$ .



# Experiments

Variants	Mastodon		DailyDialog	
	DSC	DAR	DSC	DAR
DARER	<b>59.59</b>	<b>63.43</b>	<b>53.42</b>	<b>81.39</b>
w/o Label Embeddings	56.76	62.15	50.64	79.87
w/o Harness Loss	56.22	61.99	49.94	79.76
w/o SAT-RSGT	57.37	62.96	50.25	80.52
w/o DTR-RSGT	56.69	61.69	50.11	79.76
w/o TS-LSTMs	56.30	61.49	51.61	80.33
w/o Tpl Rels in SATG	58.23	62.21	50.99	80.70
w/o Tpl Rels in DRTG	57.22	62.15	50.52	80.28

Table 5: Results of ablation experiments on F1 score.



# Experiments

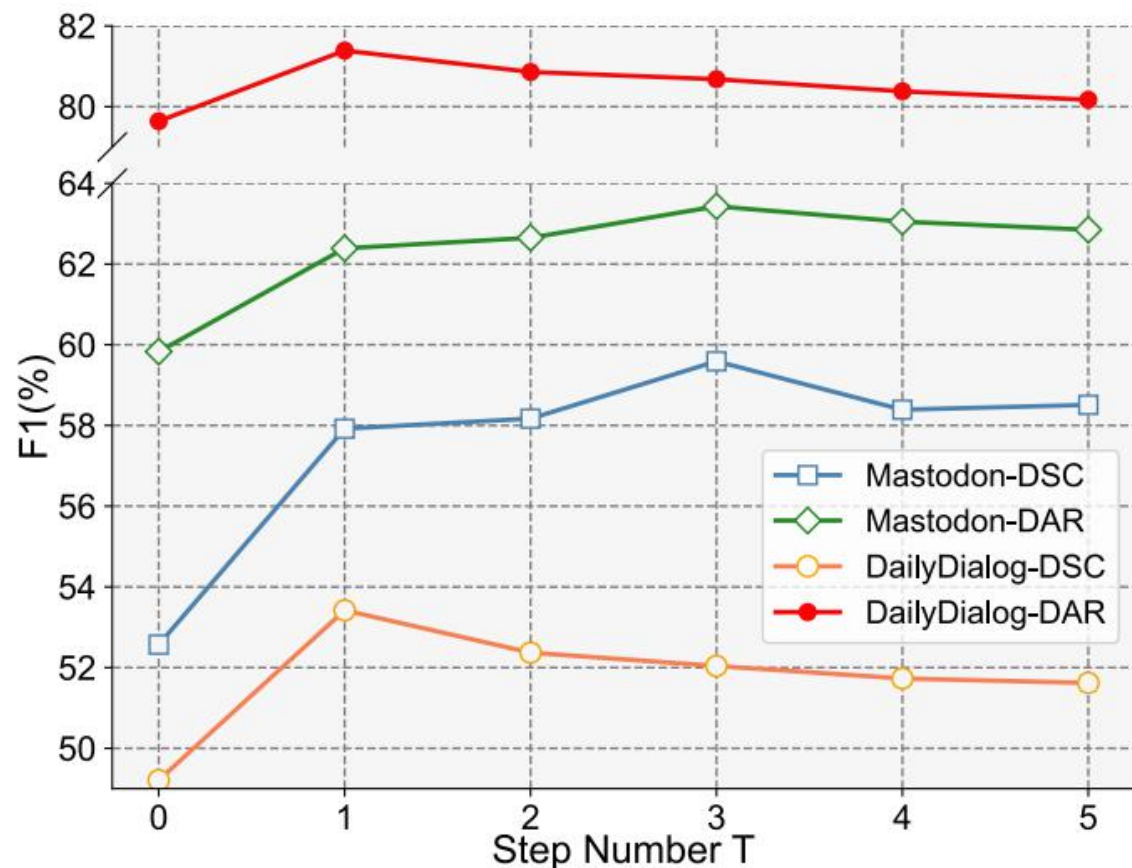


Figure 6: Performances of DARER over different  $T$ .

Models	Mastodon						
	DSC			DAR			
	P(%)	R(%)	F1(%)	P(%)	R(%)	F1(%)	
BERT	+ Linear	61.79	61.09	60.60	70.20	67.49	68.82
	+ Co-GAT	66.03	58.13	61.56	70.66	67.62	69.08
	+ DARER	65.98	67.39	<b>66.42</b>	73.82	71.67	<b>72.73</b>
RoBERTa	+ Linear	57.83	60.54	57.83	62.49	61.93	62.20
	+ Co-GAT	61.28	57.25	58.26	66.46	64.01	65.21
	+ DARER	61.36	67.27	<b>63.66</b>	70.87	68.68	<b>69.75</b>
XLNet	+ Linear	61.42	67.80	63.35	67.31	63.04	65.09
	+ Co-GAT	64.01	65.30	63.71	67.19	64.09	65.60
	+ DARER	68.05	69.47	<b>68.66</b>	72.04	69.63	<b>70.81</b>

Table 6: Results based on different PTLM encoders.



# Experiments

Models	Number of Parameters	Training Time per Epoch	GPU Memory	Avg. F1
Co-GAT	6.93M	2.35s	2007MB	53.66%
DARER	2.50M	2.20s	1167MB	61.51%
<b>Improve</b>	<b>-63.92%</b>	<b>-6.38%</b>	<b>-41.85%</b>	<b>14.63%</b>

Table 7: Comparison with SOTA on different aspects.