# SPEAKER ROLE CONTEXTUAL MODELING FOR LANGUAGE UNDERSTANDING AND DIALOGUE POLICY LEARNING

National Taiwan University

*Ta-Chung Chi*[*], *Po-Chun Chen*[*], *Shang-Yu Su*[*], *Yun-Nung (Vivian) Chen*   *Three authors contribute this work equally.*

https://github.com/MiuLab/Spk-Dialogue

## Summary

- Task Definition
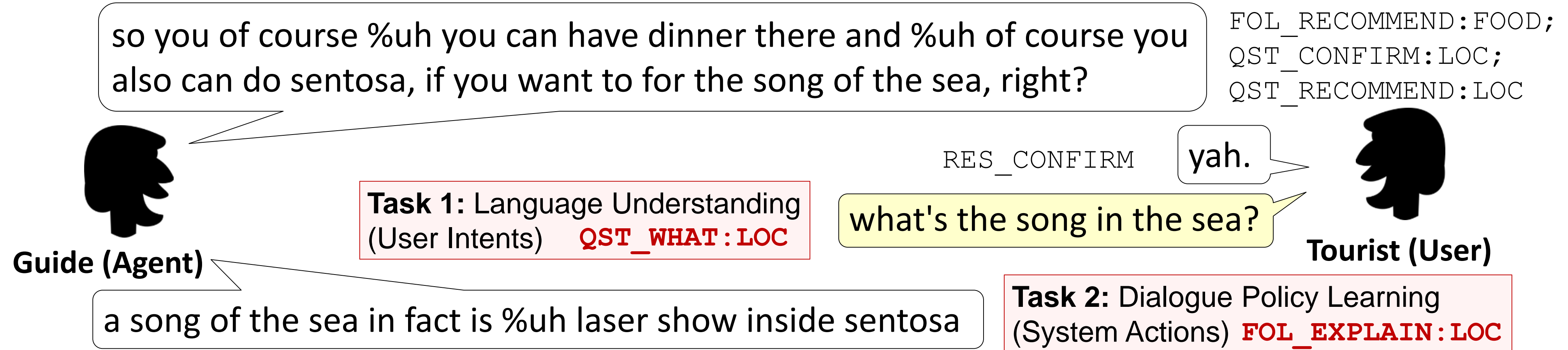  - DSTC4: human-human dialogues between tourists and guides
- Motivation
  - Human-human dialogues contain **rich and complex** human behaviors
  - **Different speaker roles** behave differently and cause notable variance in speaking habits
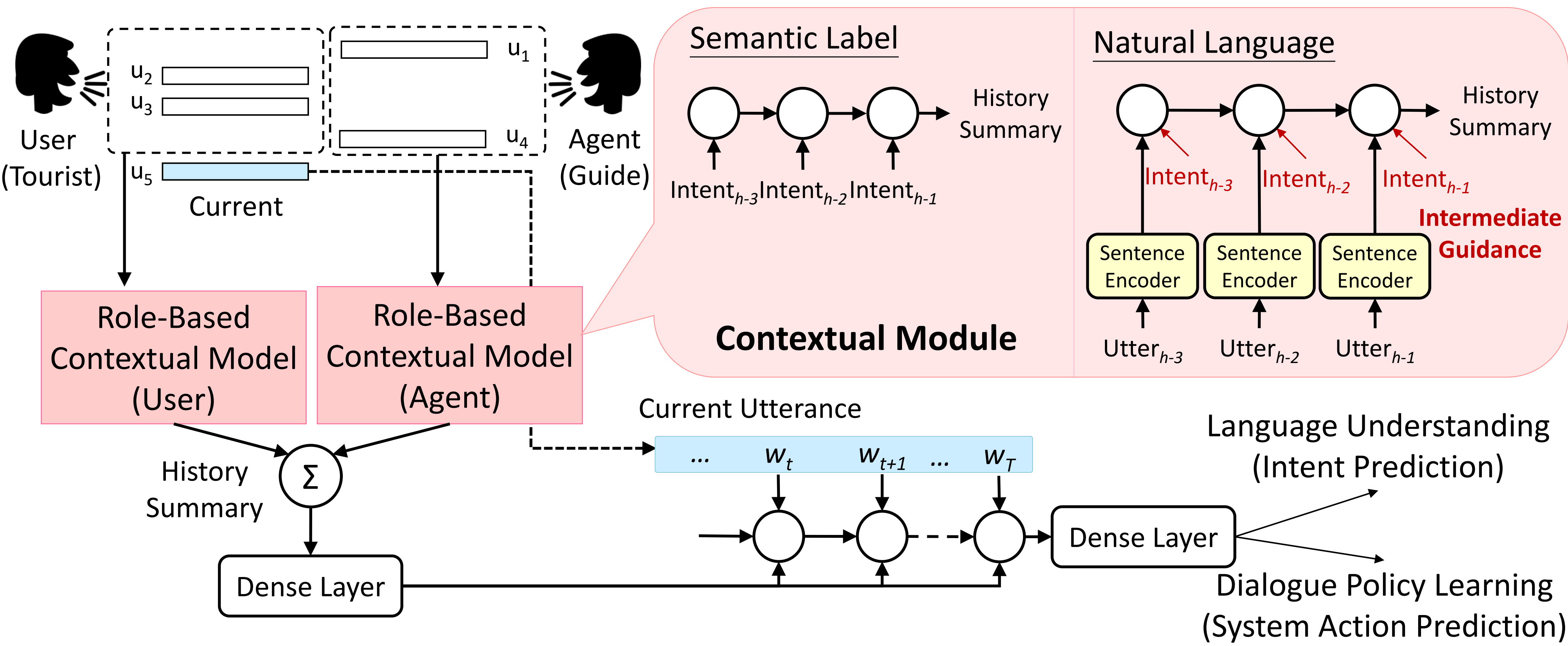- Method: **Role-Based Contextual Model for LU & PL**
  - Introduce two separate models to represent two speaker roles

- Result
  - The model achieves impressive improvement on the DSTC4 dataset

> so you of course %uh you can have dinner there and %uh of course you also can do sentosa, if you want to for the song of the sea, right?

`FOL_RECOMMEND:FOOD;`
`QST_CONFIRM:LOC;`
`QST_RECOMMEND:LOC`

`RES_CONFIRM`   yah.

**Guide (Agent)**

**Task 1:** Language Understanding (User Intents) `QST_WHAT:LOC`

what's the song in the sea?   **Tourist (User)**

a song of the sea in fact is %uh laser show inside sentosa

**Task 2:** Dialogue Policy Learning (System Actions) `FOL_EXPLAIN:LOC`

## The Proposed Approach: Role-Based Model for LU & PL



**Contextual Module**

### ❖ Contextual Model

- **encoding contexts as a history vector $v_{his}$**
  - ✓ **Semantic Label:** ground-truth intent tags are encoded as the 1-hot sentence semantics
  $$\vec{v}_{his} = \text{BLSTM}(\vec{\text{intent}}_t)$$
  - ✓ **Natural Language:** CNN-encoded sentence vector for practical situations
  $$\vec{v}_{his} = \text{BLSTM}(\text{CNN}(\vec{\text{utt}}_t))$$
  - ✓ **NL w/ Intermediate Guidance:** semantic labels act as middle supervision signal for guiding the sentence encoding module to project from input utterances to a more meaningful feature space

> Leverage contextual information for better understanding

### ❖ Speaker Role Modeling

- **train two role-specific models independently, $\text{BLSTM}_{role_a}$ and $\text{BLSTM}_{role_b}$**

$$\vec{v}_{his} = \text{BLSTM}_{role_a}(\vec{\text{intent}}_{t,role_a}) + \text{BLSTM}_{role_b}(\vec{\text{intent}}_{t,role_b})$$

$$\vec{v}_{his} = \text{BLSTM}_{role_a}(\text{CNN}(\vec{\text{utt}}_{t,role_a})) + \text{BLSTM}_{role_b}(\text{CNN}(\vec{\text{utt}}_{t,role_b}))$$

> User usually pay attention to **self history (reasoning)** and **others' utterances (listening)**
> Two speaker roles behave differently

> **End-to-End Training Objective**
> - BLSTM-encoded current utterance concatenated with the history vector for *multi-label intent prediction and system action prediction*
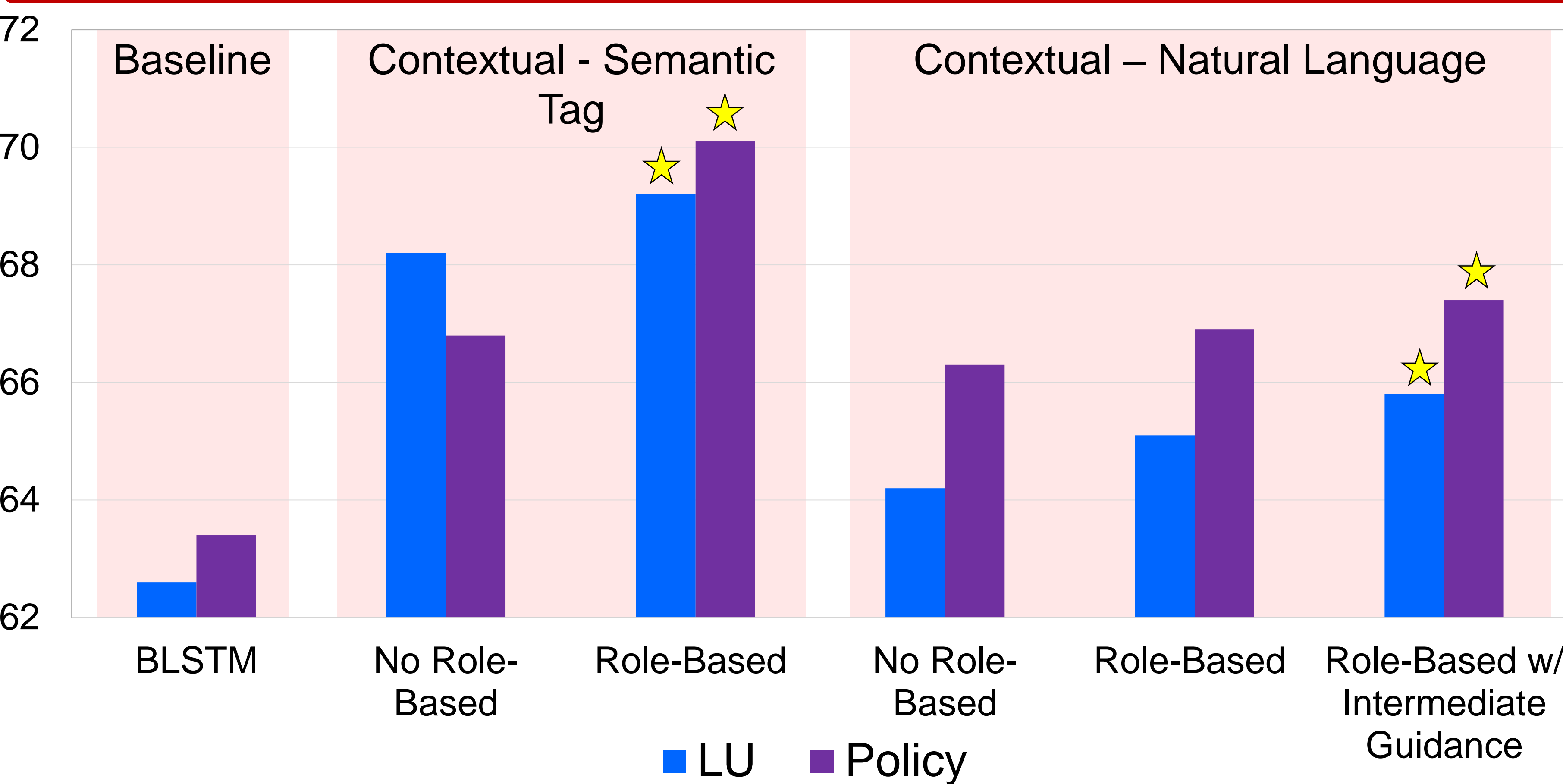
$$\vec{y} = \text{BLSTM}(\vec{v}_{his}, \vec{x})$$
$$p(\vec{y} \mid \vec{x}) = \prod_i p(y_i \mid w_1, \cdots, w_i)$$

## Experiments and Discussions



> Setup
> - Dataset: DSTC4 35 human-human dialogues
> - Evaluation metrics: F1 for multi-label classification
> Experimental Results
> - **Contextual** models significantly improve the baselines
> - The **role-based** models outperform the one without the role information for both tasks
> - **Intermediate guidance** improves semantic modeling
> Discussion
> - Most LU results are worse than dialogue policy learning results
> - The reason may be that the *guide has similar behavior patterns* (e.g. providing information and confirming questions) while *the user has more diverse interactions*
> - The idea about modeling speaker role information can be further extended to various research topics

> The proposed speaker role contextual model obtains the state-of-the-art results.

## Conclusions

- **Approach**: an end-to-end role-based contextual model that automatically learns speaker specific contextual encoding
- **Experiment**: impressive improvement on a benchmark multi-domain dialogue dataset
- **Result**: demonstrating that different speaker roles behave differently and focus on different goals

**Code Available:**
https://github.com/MiuLab/Spk-Dialogue

b02902019@ntu.edu.tw
r06922028@ntu.edu.tw
r05921117@ntu.edu.tw
y.v.chen@ieee.org