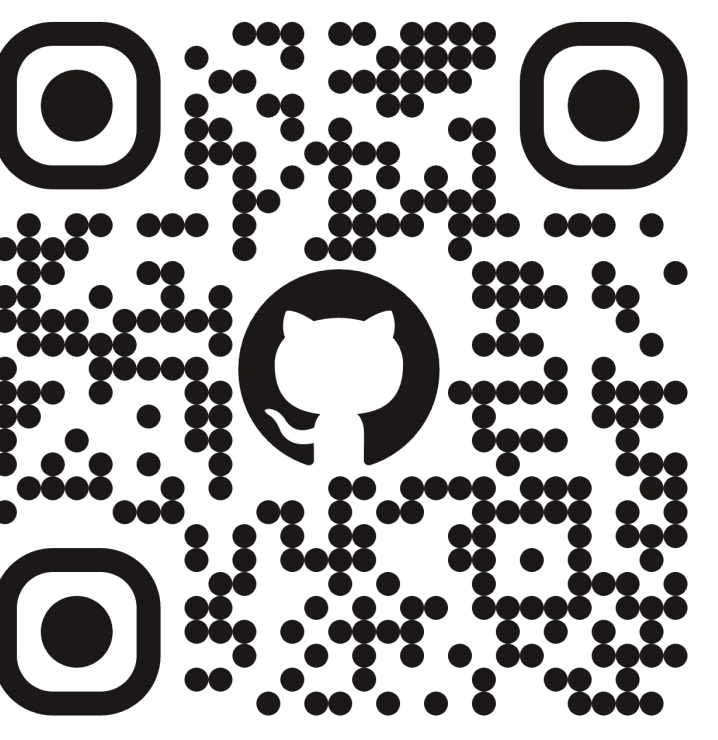


# Disambiguating Reference in Visually Grounded Dialogues through Joint Modeling of Textual and Multimodal Semantic Structures

Shun Inadumi<sup>1,2</sup>, Nobuhiro Ueda<sup>3</sup>, Koichiro Yoshino<sup>4,2,1</sup>

1. NAIST, 2. RIKEN GRP, 3. Kyoto University, 4. Science Tokyo



## Understanding References in Dialogues

### TRR: Textual Reference Resolution

Identifies **textual reference** relations **between phrases**

TRR partially consists of the following tasks:

- Coreference** Resolution: **this** = the coffee cup
- Predicate-argument structure** analysis: **take** → ACC the coffee cup

### MRR: Multimodal Reference Resolution [Ueda+, 2024]

Identifies **phrase-to-objects** **direct/indirect** references

- Phrase Grounding [Plummer+, 2015]: task limited to **direct** references. **the coffee cup** =

MRR enables systems to understand **dialogue events** — such as "**who does what to whom**" — linked to objects.

## Example of the system analyzes a two-person dialogue

P<sub>1</sub>: Person1

Would [you]  
**take** [me] **this** ?

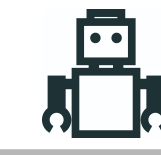
Thank you !

Japanese often omit  
subjects and objects.

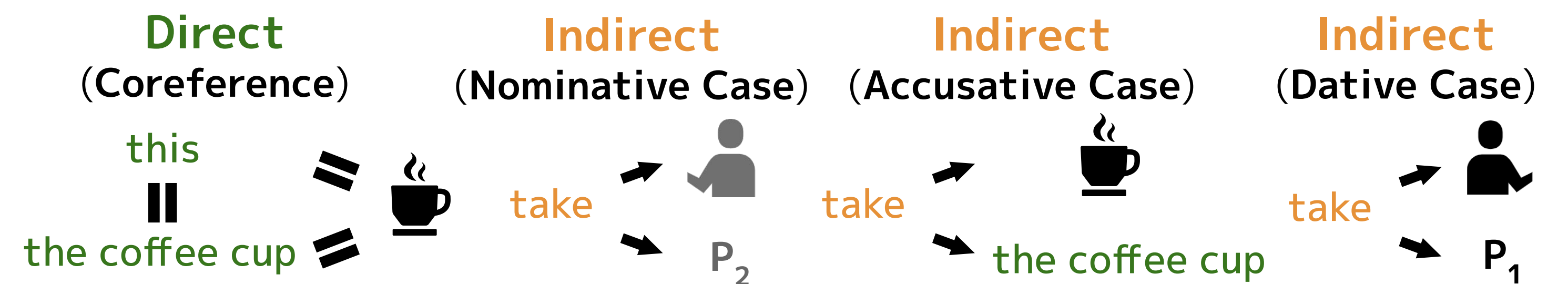
P<sub>2</sub>: Person2

Yes. **The coffee cup**, right ?

Not at all.



First person vision of the System

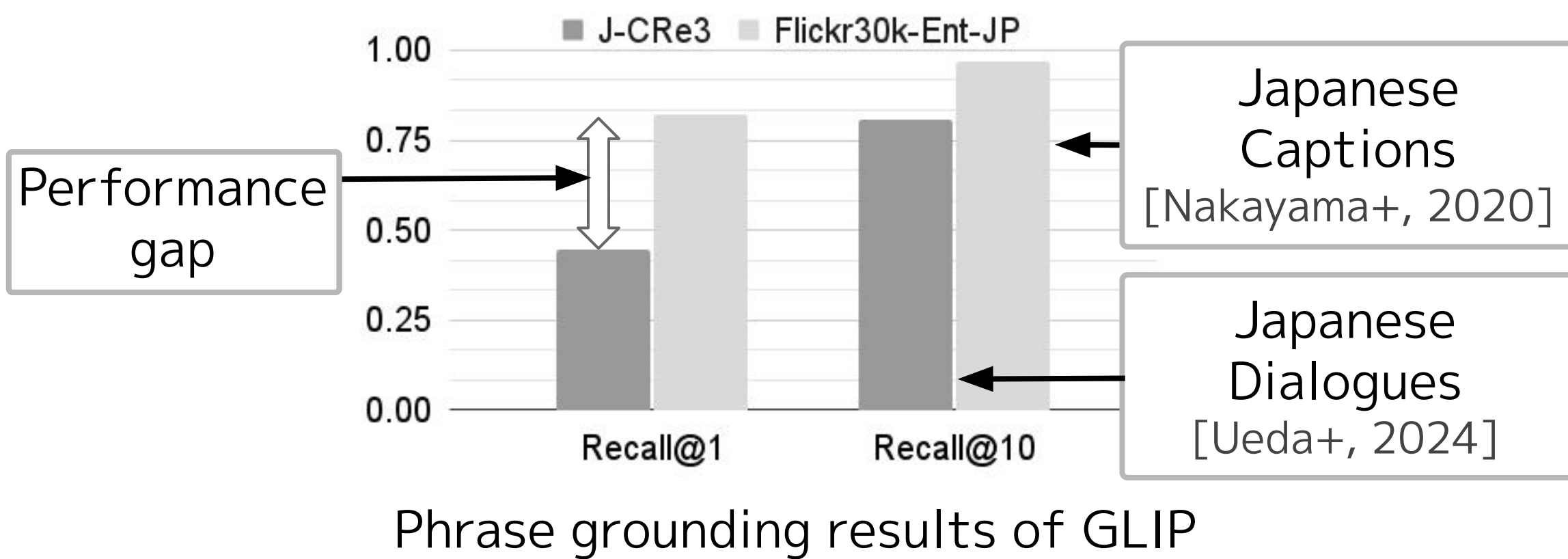


## Limitations of Existing Models

- GLIP [Li+, CVPR2022]: A phrase grounding model trained on large-scale **image-caption pairs** with **direct** references.

- In dialogue parsing, GLIP struggles with:

- Resolving **direct** references made via **pronouns**: **this** =
- Parsing **indirect** references involving **ellipses**: **take** → DAT, NOM



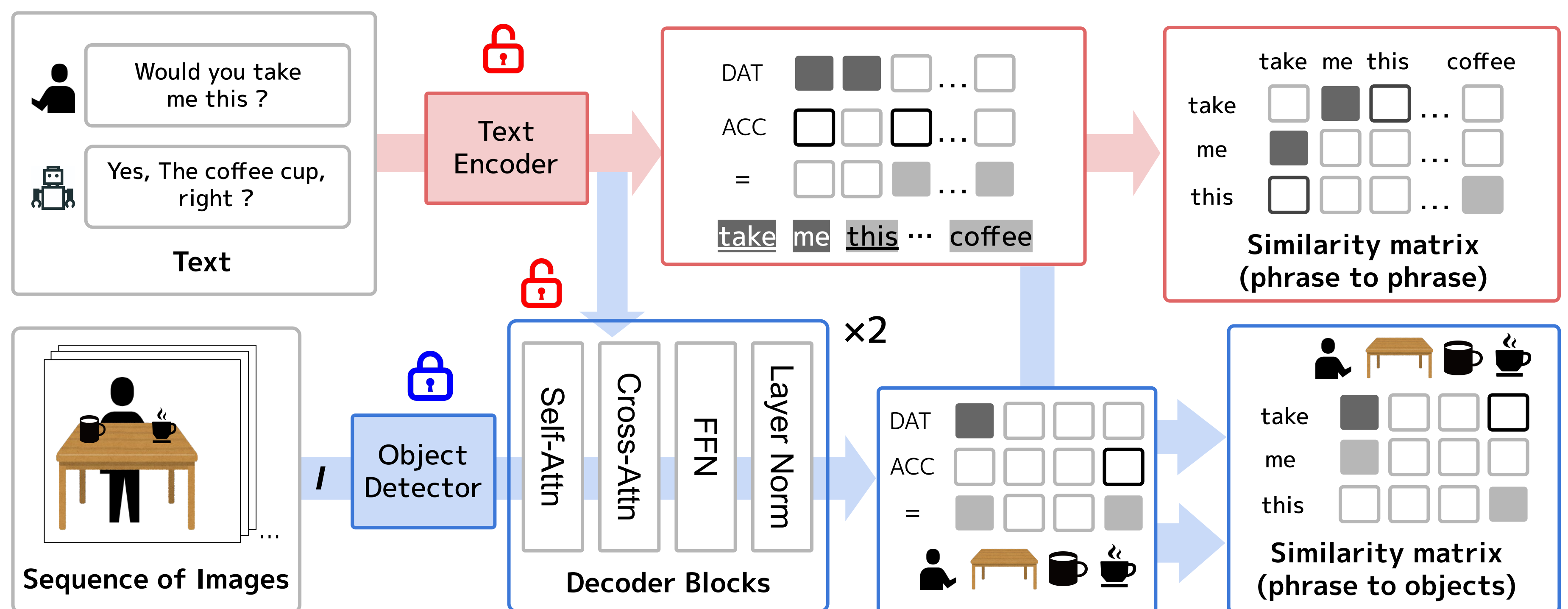
By resolving ambiguities in ① and ②, we aim to better understand real-world dialogues.

## Proposed Framework

- By incorporating **textual references**, we can improve MRR performance.

e.g.) If **the coffee cup** is known, **this** = can be uniquely identified.

- We propose a framework to **jointly** model **TRR** and **MRR**.

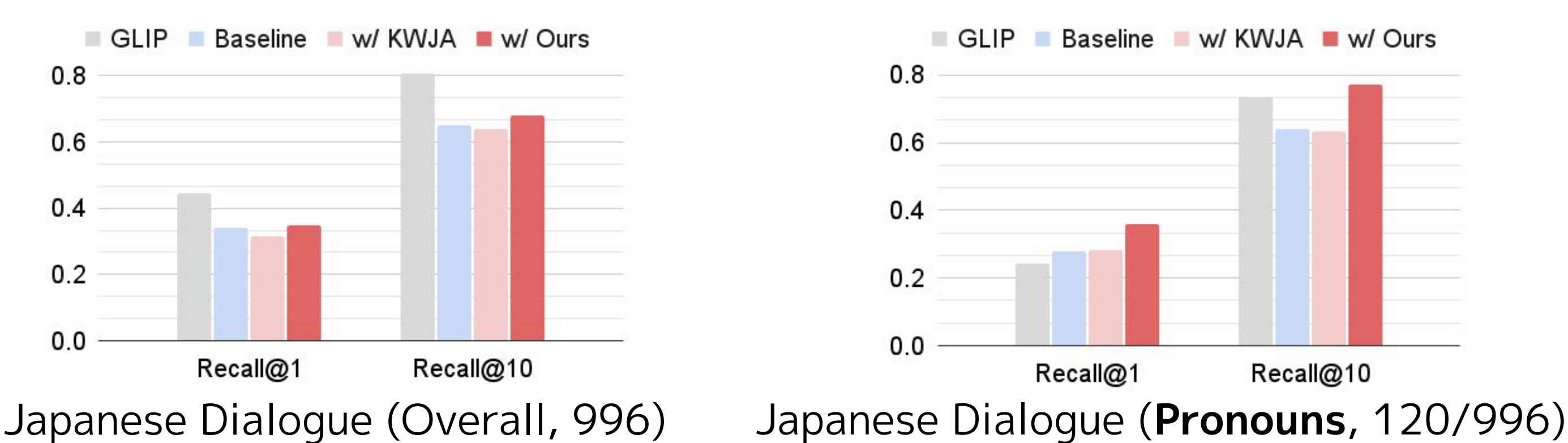


## Phrase Grounding Results

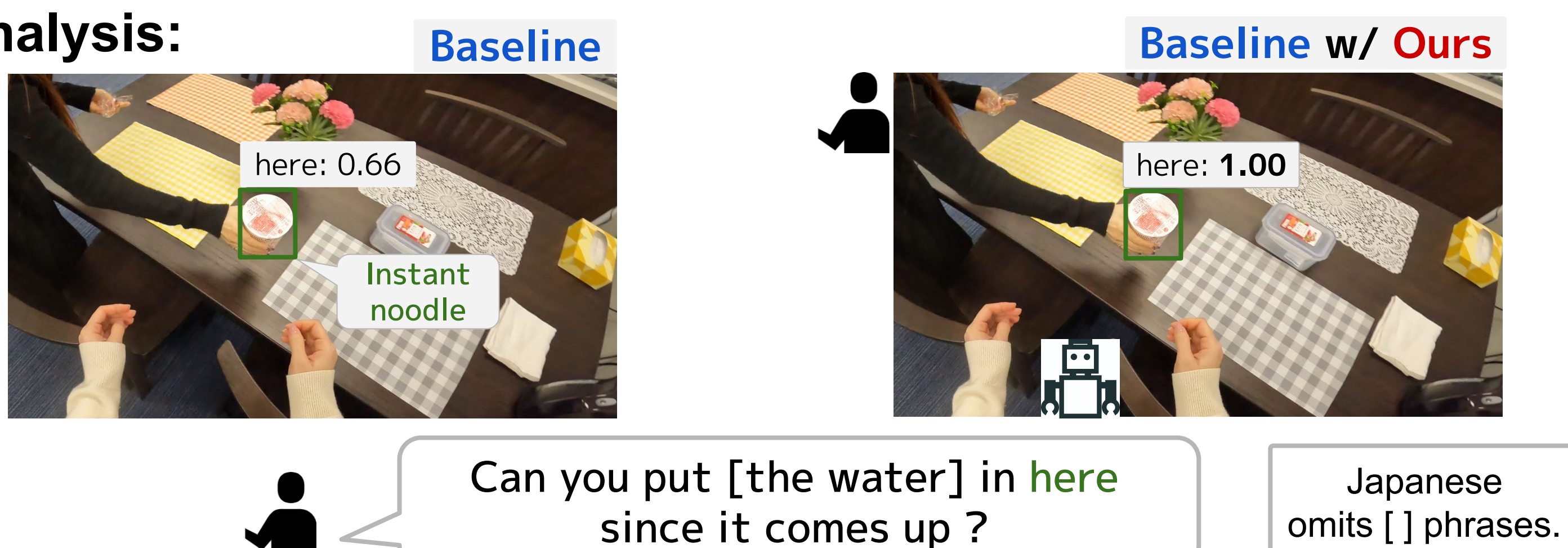
### Compared models:

- Baseline
  - Baseline w/ **Ours**
  - Baseline w/ **KWJA** [Ueda+, 2023]
  - GLIP [Li+, CVPR2022]
- Phrase grounding model with **coreference resolution** (fine-tuned on Japanese data)
- Pre-trained on English data [Krishna+, 2017, Hudson+, 2019]

### Main Results:



### Analysis:



Improvements from coreference resolution:

- Improved pronoun phrase grounding
- Increased confidence scores for **pronouns**

## MRR Results

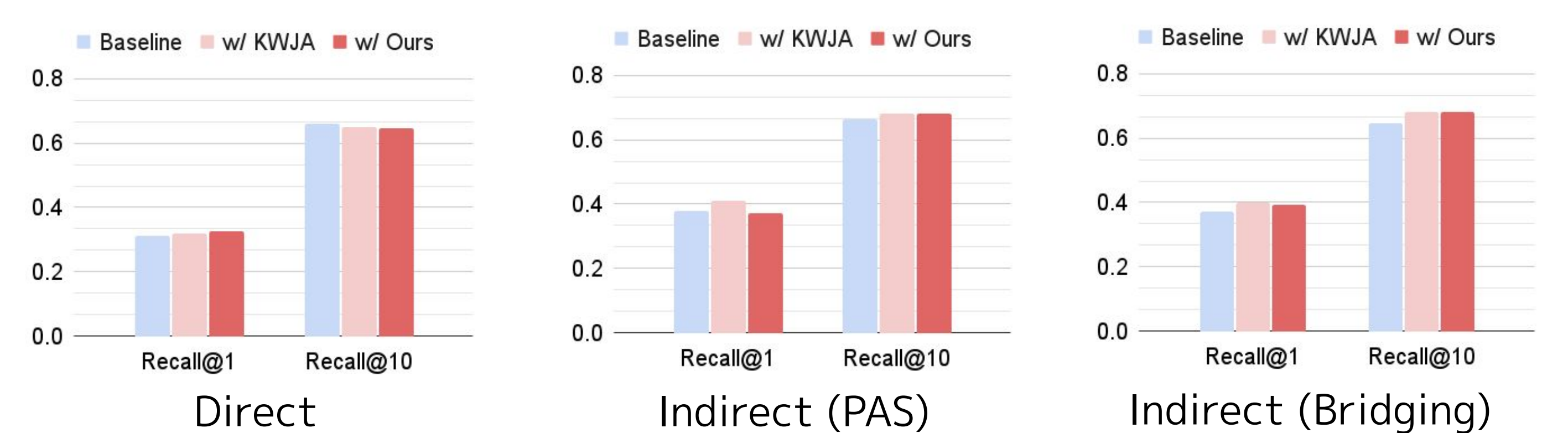
### Parsing Targets:

- Direct Reference
- Indirect References:
  - Predicate-argument structures (PAS)
  - Bridging Anaphora

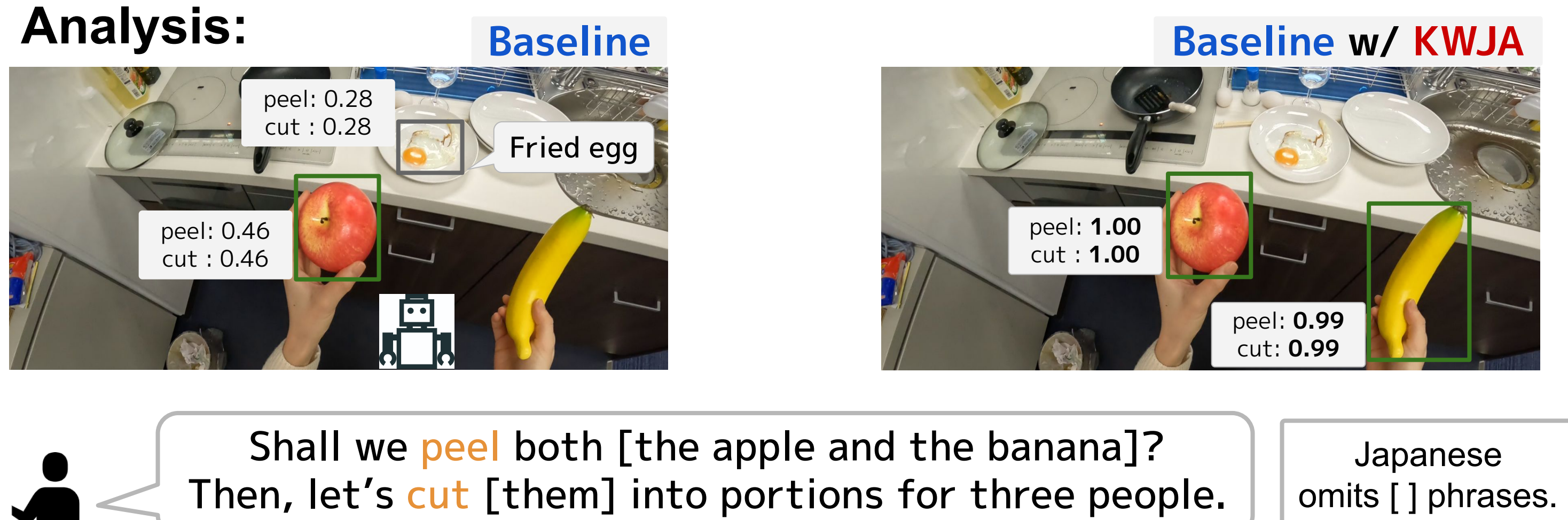
### Compared models:

- Baseline
  - Baseline w/ **Ours**
  - Baseline w/ **KWJA**
- MRR models with **TRR**

### Main Results:



### Analysis:



Improvements from TRR:

- Improved indirect reference performance
- Increased confidence scores for **predicates**