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

Shun Inadumi [1,2], Nobuhiro Ueda [3,†] and Koichiro Yoshino [4,2,1]

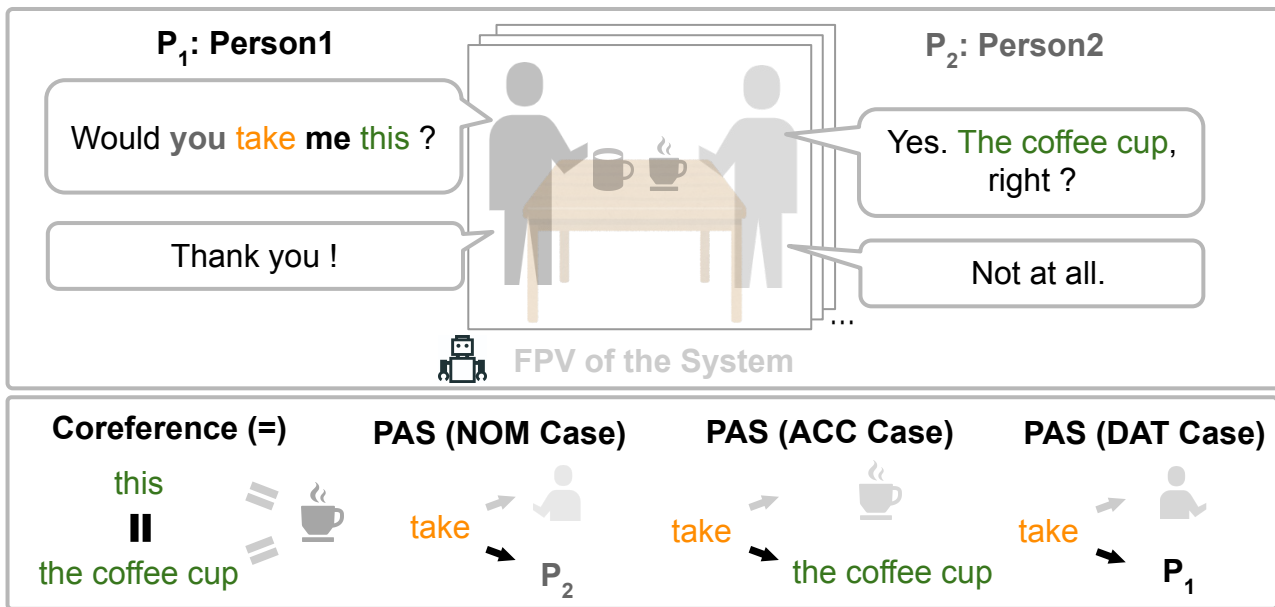1. NAIST, 2. RIKEN GRP, 3. Kyoto University, 4. Science Tokyo
[†]Currently at NEC Corporation.

# Understanding References in **Dialogue Text**

- **TRR: T**extual **R**eference **R**esolution:
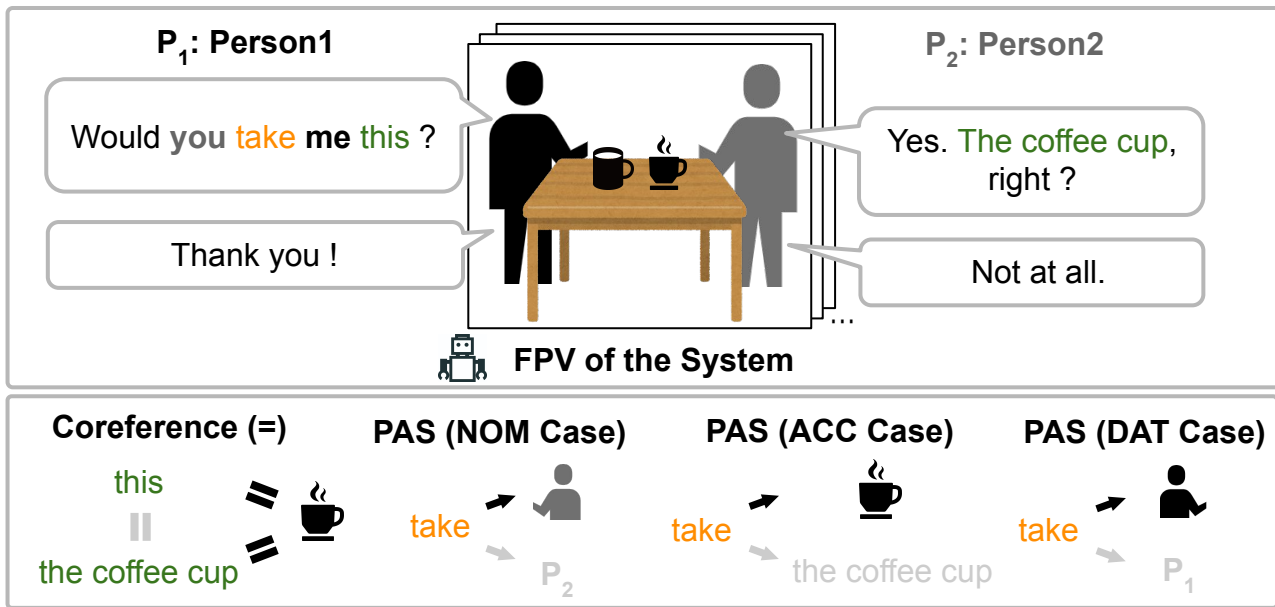  - Identify textual reference relations between phrases
  - TRR consists of coreference resolution, predicate-argument structure analysis and bridging anaphora resolution.



Example of the system analyzes a two-person dialogue

# Understanding References in **Visually Grounded Dialogues**

- **MRR: M**ultimodal **R**eference **R**esolution [Ueda+, 2024] :
  - Identify phrase-to-objects reference relations
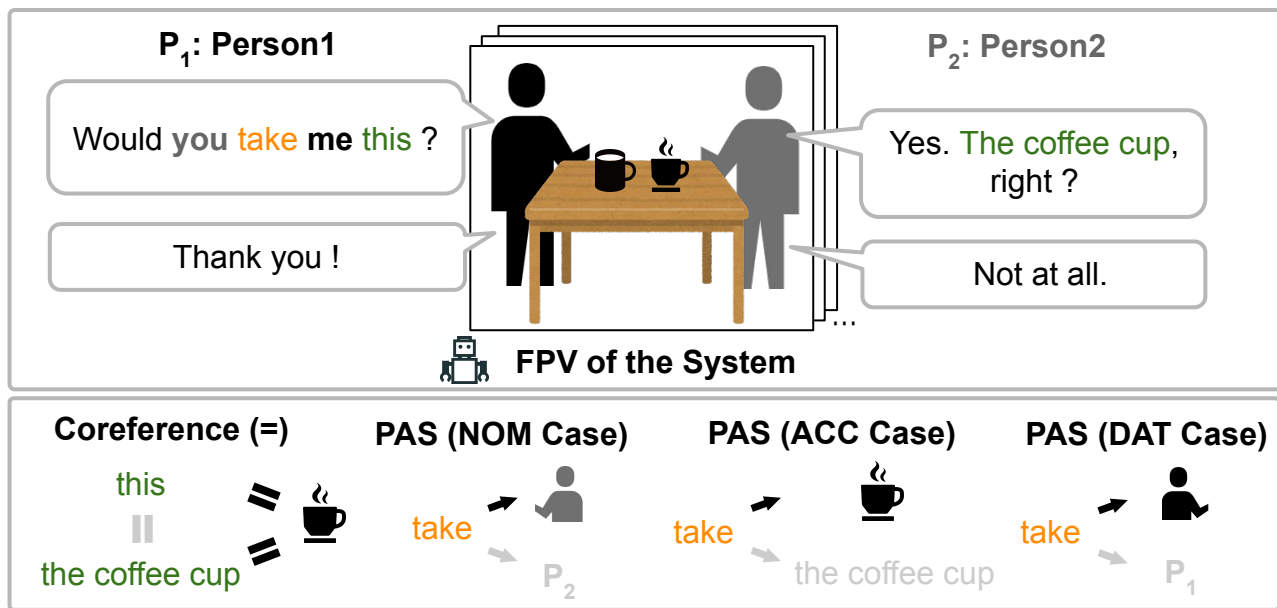  - Direct reference and Indirect reference



Example of the system analyzes a two-person dialogue

- **Multimodal Reference R**
  - Identify phrase-to-obj
  - Direct reference and m

> **Phrase grounding** [Plummer+, 2015] refers to the task of predicting only this type of reference.



**P$_1$: Person1**

Would **you** take **me** this ?

Thank you !

**P$_2$: Person2**

Yes. The coffee cup, right ?

Not at all.

**FPV of the System**

**Coreference (=)**

this

∥

the coffee cup

**PAS (NOM Case)**

take → P$_2$

**PAS (ACC Case)**

take → the coffee cup

**PAS (DAT Case)**

take → P$_1$

Example of the system analyzes a two-person dialogue
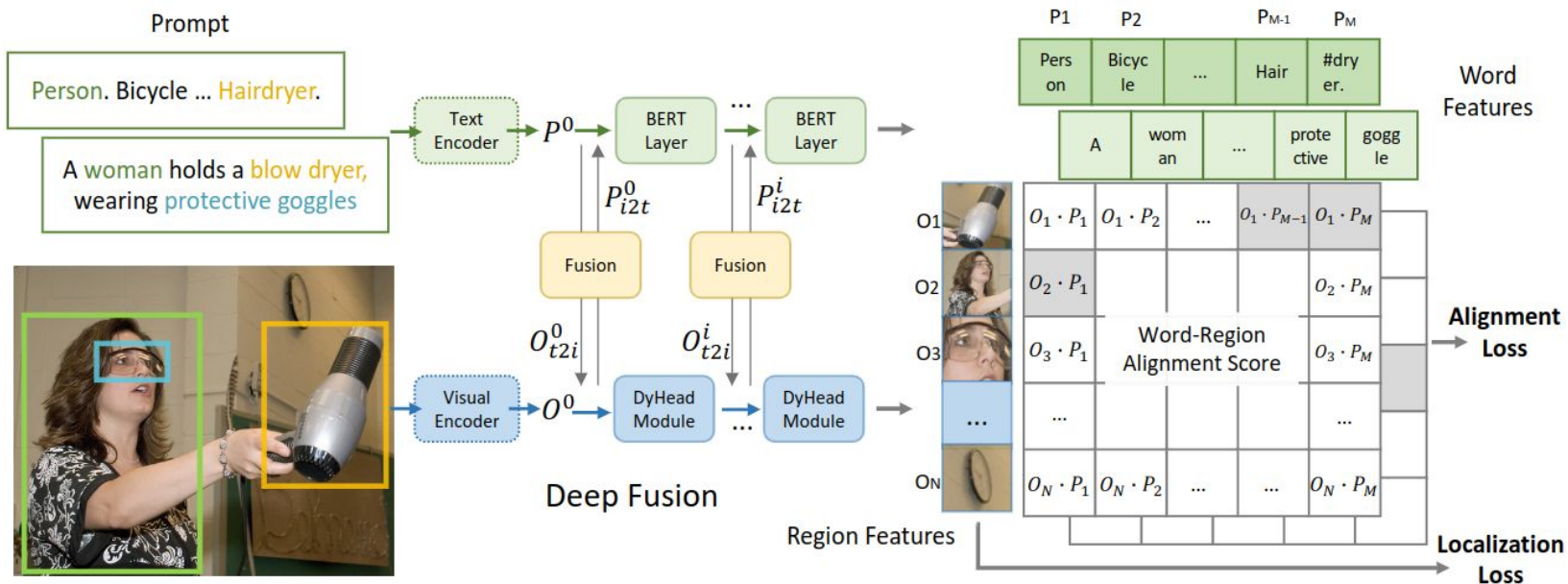
4

# Why is MRR important?

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



Example of the system analyzes a two-person dialogue
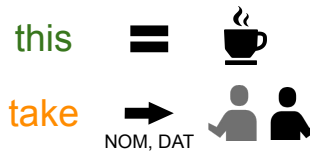
# Limitations of Existing Models

- GLIP (Grounded Language-Image Pretraining) [Li+, 2022] :
A phrase grounding model trained on large-scale **image-caption pairs** with
direct references.



Cited from [Li+, 2022].
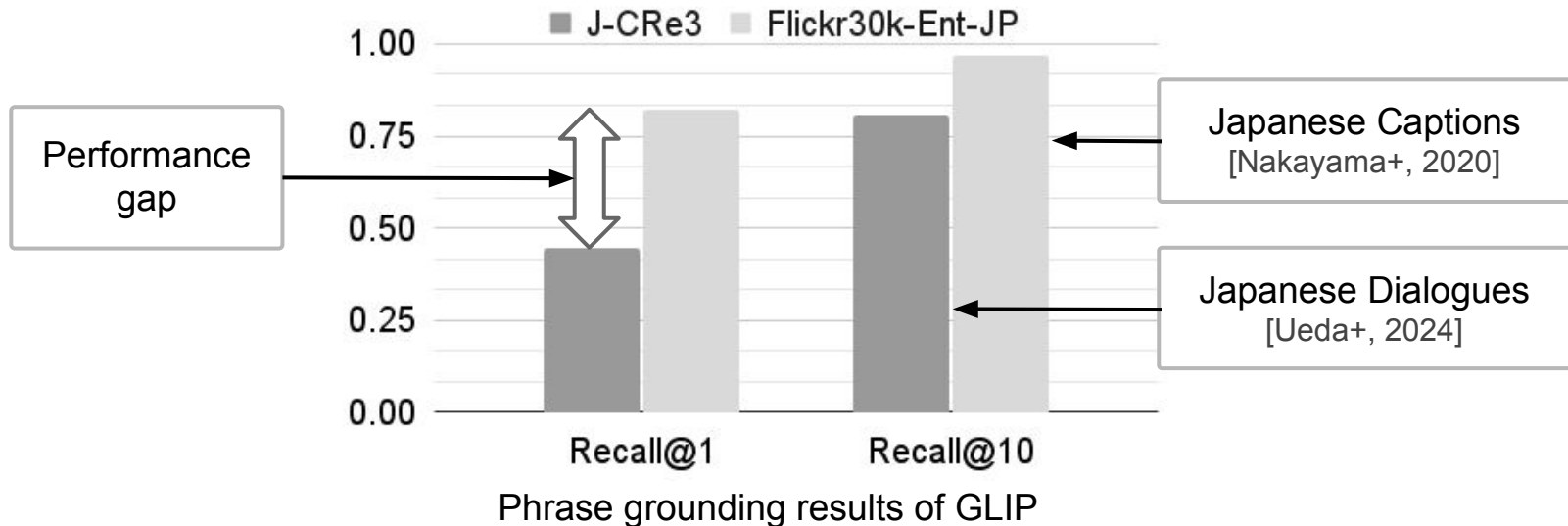
# Limitations of Existing Models

- In dialogue parsing, GLIP struggles with:
  - Resolving direct references made via **pronouns**
  - Parsing indirect references involving **ellipses**

this ═ ☕

take ➡ 🧑🧑
NOM, DAT

# Limitations of Existing Models

- In dialogue parsing, GLIP struggles with:
  - Resolving direct references made via **pronouns**
  - Parsing indirect references involving **ellipses**

this ⬛ ☕

take ➡ 🧍🧍
NOM, DAT



Performance gap

J-CRe3  Flickr30k-Ent-JP

1.00
0.75
0.50
0.25
0.00

Recall@1    Recall@10

Japanese Captions
[Nakayama+, 2020]

Japanese Dialogues
[Ueda+, 2024]

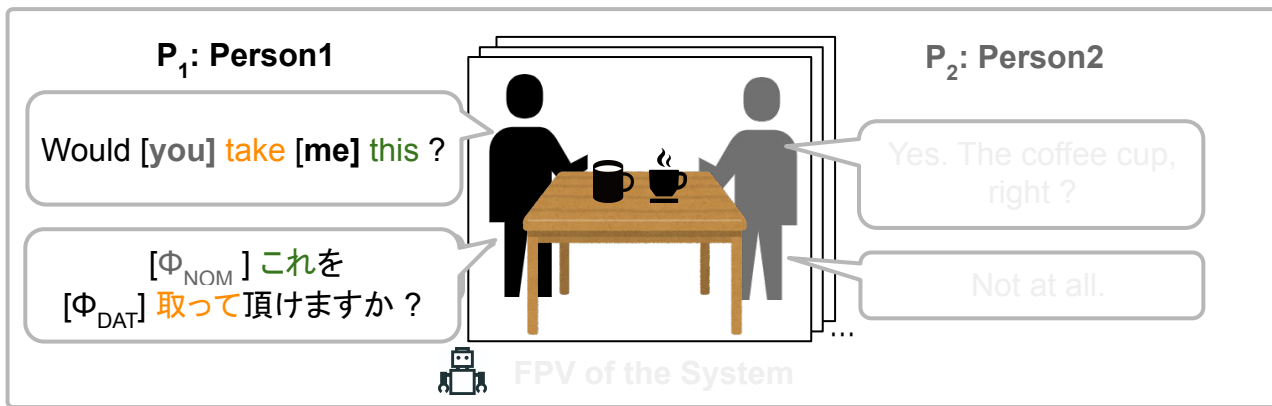Phrase grounding results of GLIP

# Limitations of Existing Models

- In dialogue parsing, GLIP struggles with:
  - Resolving direct references made via **pronouns**
  - Parsing indirect references involving **ellipses**

this  =  ☕

take  ➡  NOM, DAT

**P₁: Person1**

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

[Φ$_{NOM}$] これを
[Φ$_{DAT}$] 取って頂けますか ？

FPV of the System

**P₂: Person2**

Yes. The coffee cup, right ?

Not at all.

…

Pro-drop languages, such as Japanese, often omit **subjects** and **objects**.

9

# Limitations of Existing Models

- In dialogue parsing, GLIP struggles with:
    - Resolving direct references made via **pronouns**
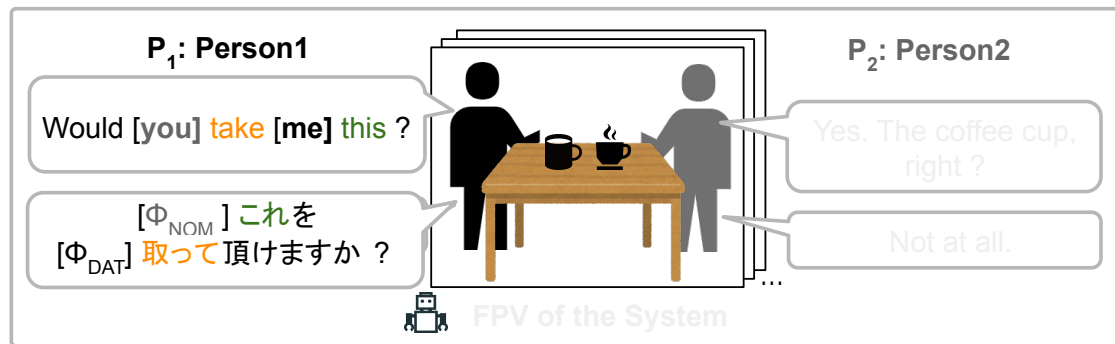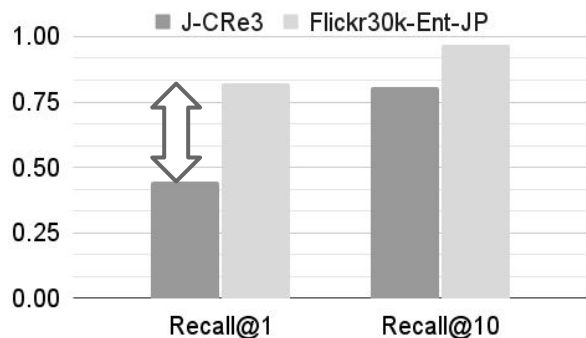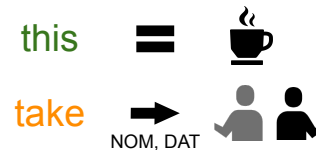    - Parsing indirect references involving **ellipses**



this $=$ ☕

take → 👤👤
NOM, DAT



**P₁: Person1**

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

[$\Phi_{NOM}$] これを
[$\Phi_{DAT}$] 取って 頂けますか ？

**P₂: Person2**

Yes. The coffee cup, right ?

Not at all.

FPV of the System

**By resolving these ambiguities**,
we aim to improve the understanding of real-world dialogues.

# Using Textual Reference Relations 💡

By incorporating <u>textual references</u>, we can improve MRR performance.

e.g.) If $\begin{array}{c}\text{the coffee cup}\\ \| \quad \| \,☕\\ \text{this}\end{array}$ is known, $\text{this} = ☕$ can be uniquely identified.
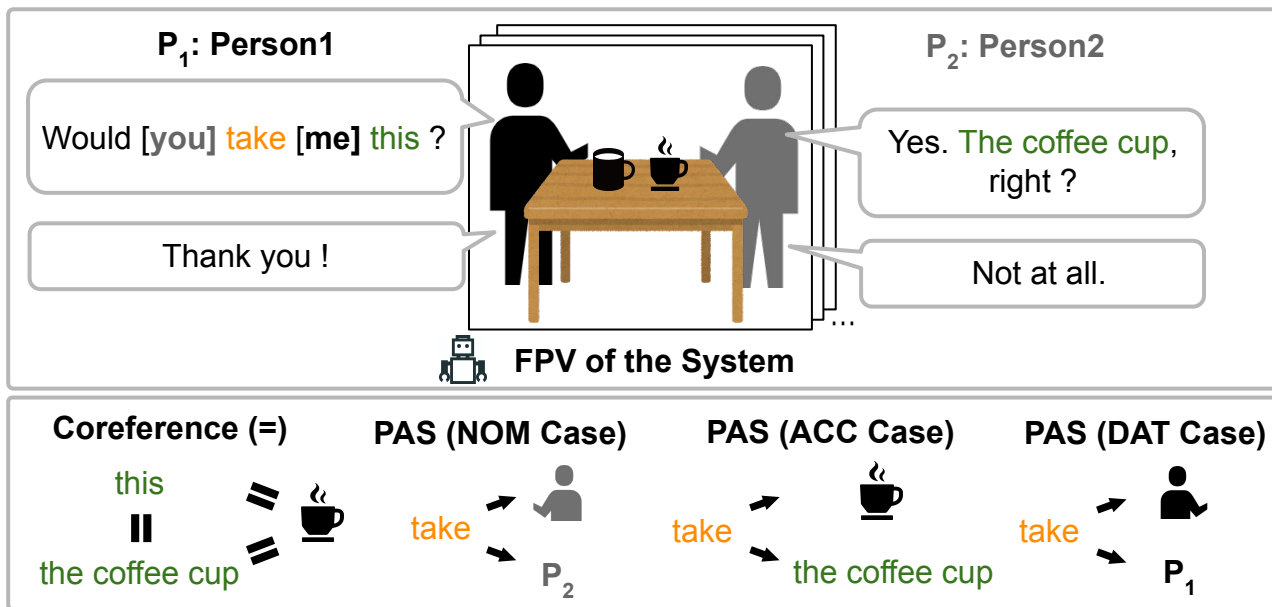
**P₁: Person1**

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

Thank you !

**P₂: Person2**

Yes. The coffee cup, right ?

Not at all.

...

🤖 **FPV of the System**

| Coreference (=) | PAS (NOM Case) | PAS (ACC Case) | PAS (DAT Case) |
|---|---|---|---|
| this | take → 🧑 | take → ☕ | take → 🧑 |
| ‖ | ↓ | ↓ | ↓ |
| ☕ | P₂ | the coffee cup | P₁ |
| the coffee cup | | | |

# Proposed Framework 💡

We propose a framework to jointly model TRR and MRR.

Learning to align phrase embeddings with object features



**Text**

Would you take me this ?

Yes. The coffee cup

**Seq. of Images**

Text Encoder

Object Detector

Self-Attn · Cross-Attn · FFN · Layer Norm

**Decoder Blocks**

DAT · ACC · = · take me this … coffee …

**Sim. matrix (phrase to phrase)**

take me this coffee

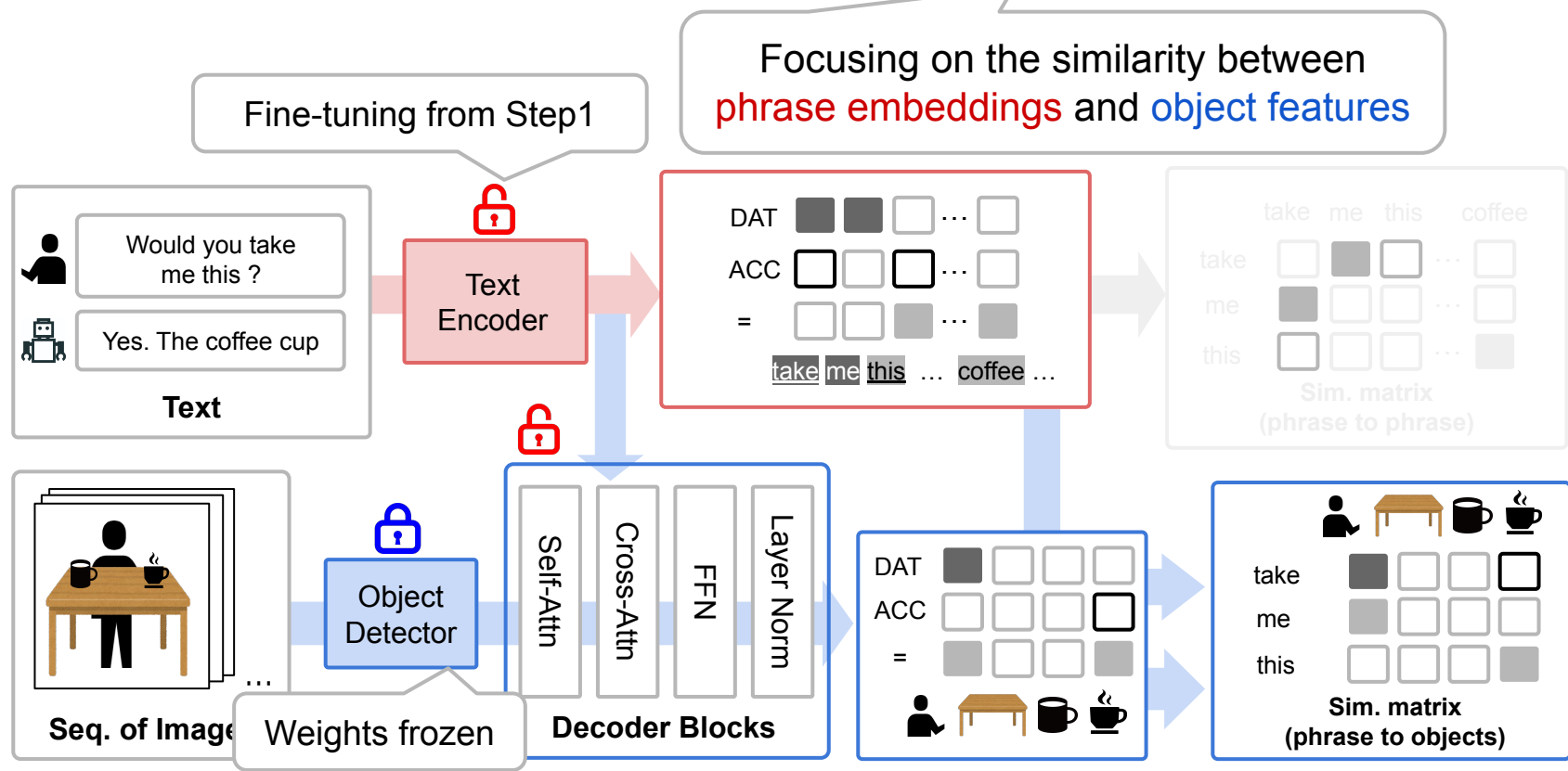**Sim. matrix (phrase to objects)**

take me this

# Step1: Textual Reference Resolution

We propose a framework to jointly model TRR and MRR.

Focusing on the similarity between phrase embeddings



Text

Would you take me this ?

Yes. The coffee cup

Text Encoder

DAT
ACC
=

take me this … coffee …

Sim. matrix (phrase to phrase)

take me this coffee

Seq. of Images

Object Detector

Self-Attn
Cross-Attn
FFN
Layer Norm

Decoder Blocks

DAT
ACC
=

Sim. matrix (phrase to objects)

# Step2: Multimodal Reference Resolution

We propose a framework to jointly model TRR and MRR.

Fine-tuning from Step1

Focusing on the similarity between phrase embeddings and object features

Would you take me this ?

Yes. The coffee cup

**Text**

Text Encoder

DAT
ACC
=

take me this … coffee …

Sim. matrix (phrase to phrase)

**Seq. of Images**

Object Detector

Weights frozen

Self-Attn | Cross-Attn | FFN | Layer Norm

**Decoder Blocks**

DAT
ACC
=

take
me
this
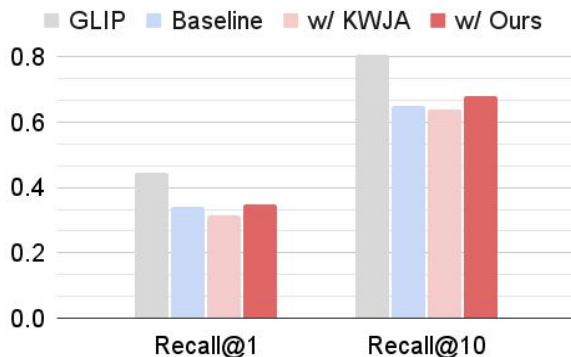
**Sim. matrix (phrase to objects)**

14

# Phrase Grounding Results
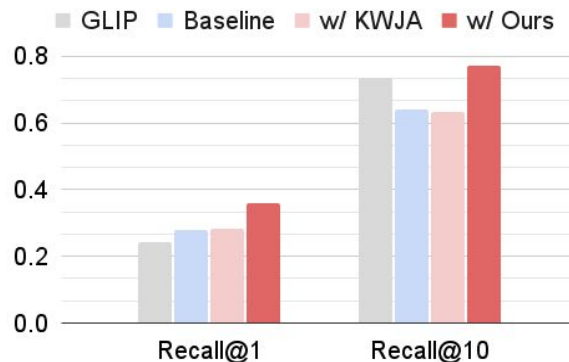
- Compared models:
  - Baseline
  - Baseline w/ Ours
  - Baseline w/ KWJA [Ueda+, 2023]
  - GLIP

Phrase grounding model with coreference resolution (fine-tuned on Japanese data [Nakayama+, 2020, Ueda+, 2024])

Pre-trained on English data [Krishna+, 2017, Hudson+, 2019]



Japanese Dialogue (Overall, 996)

Japanese Dialogue (**Pronouns**, 120/996)

Baseline w/ Ours achieved improved **pronoun phrase grounding** through **coreference resolution**.
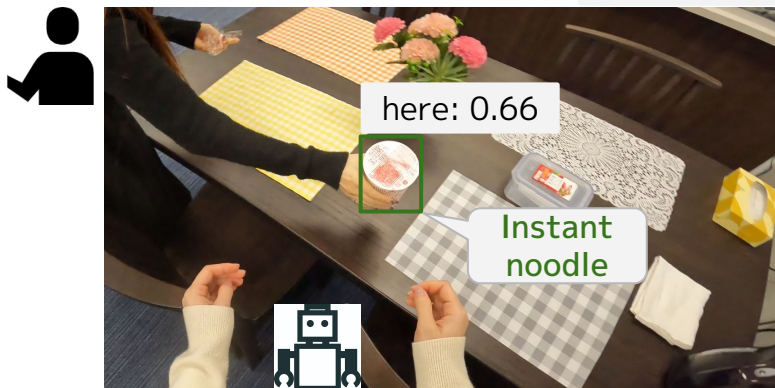
# Phrase Grounding Results

- Compared models:
  - Baseline
  - Baseline w/ Ours

Phrase grounding model with coreference resolution
(fine-tuned on Japanese data [Nakayama+, 2020, Ueda+, 2024])



**Baseline**

here: 0.66

Instant noodle

**Baseline w/ Ours**

here: 1.00

Can you put [the water] in here
since it comes up ?
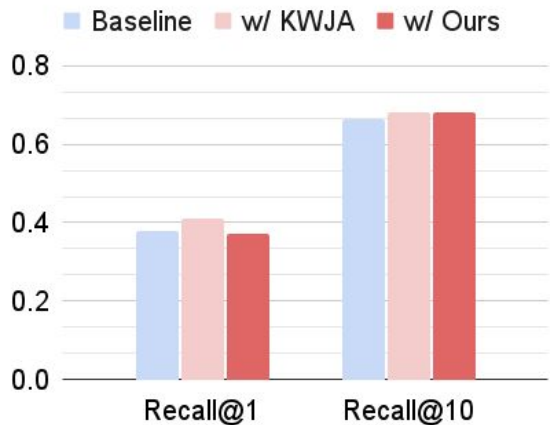
Japanese
omits [ ] phrases.

Coreference resolution strengthens confidence scores
in pronoun-to-object predictions.
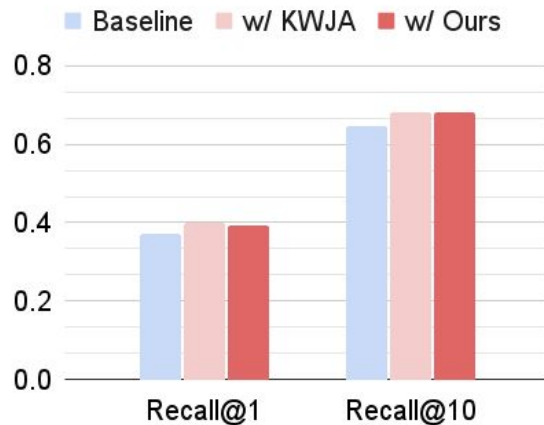
16

# Multimodal Reference Resolution Results

- Compared models:
  - Baseline
  - Baseline w/ Ours
  - Baseline w/ KWJA [Ueda+, 2023]

MRR model with TRR
(fine-tuned on Japanese data [Ueda+, 2024])



Indirect (Predicate-argument structure)
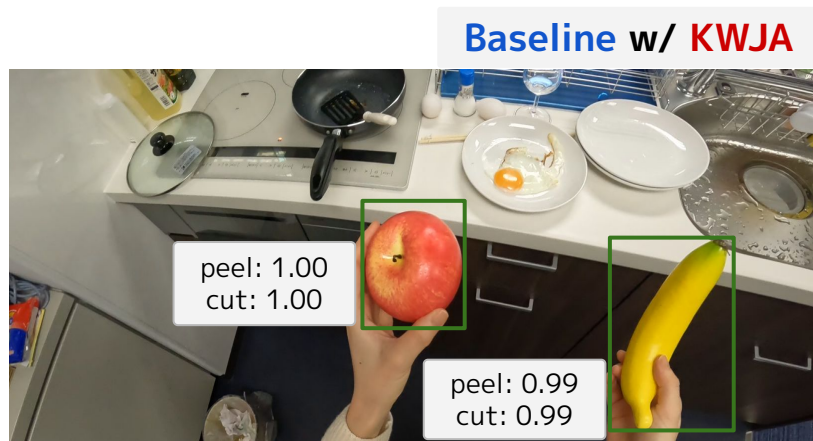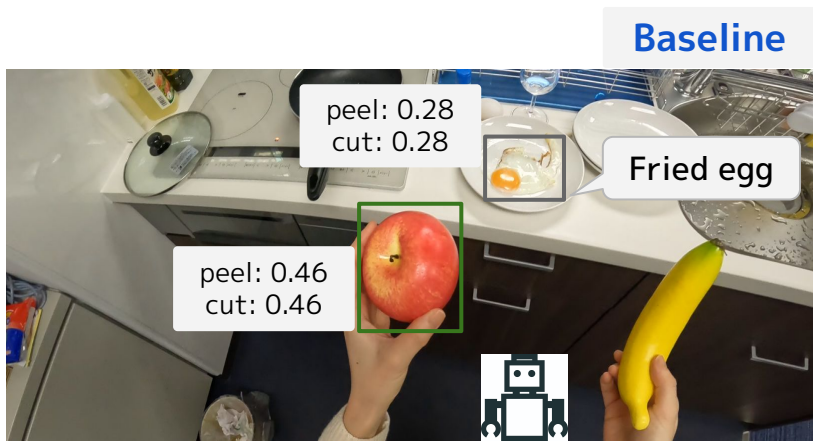


Indirect (Bridging anaphora)

Improved indirect reference performance
through textual reference resolution.

# Multimodal Reference Resolution Results

- Compared models:
  - Baseline
  - Baseline w/ KWJA [Ueda+, 2023]

MRR model with TRR
(fine-tuned on Japanese data [Ueda+, 2024])



**Baseline**

peel: 0.28
cut: 0.28

Fried egg

peel: 0.46
cut: 0.46

**Baseline w/ KWJA**

peel: 1.00
cut: 1.00

peel: 0.99
cut: 0.99

Shall we peel both [the apple and the banana]?
Then, let's cut [them] into portions for three people.

Japanese
omits [ ] phrases.

TRR strengthens confidence scores for predicates.
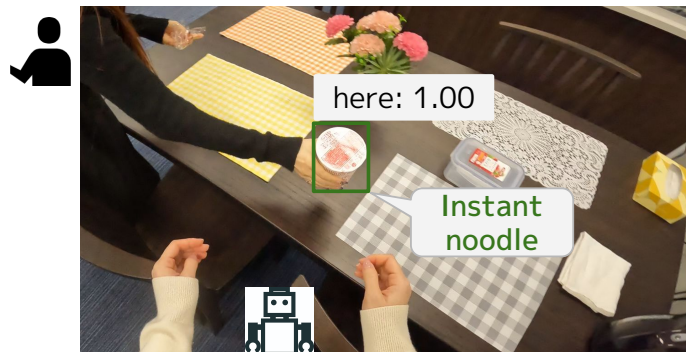
# Conclusion

**Purpose** :

- Through **resolving ambiguities in visually-grounded dialogues**,
  we aim to improve the understanding of real-world dialogues.

**Idea** 💡 :

- We propose a framework to jointly model textual
  reference resolution (TRR) and multimodal
  reference resolution.

**Main Results** :

- Improved pronoun phrase grounding through
  coreference resolution.
- Improved indirect reference performance through
  TRR.



here: 1.00

Instant noodle

Can you put [the water] in here since it comes up ?