

稲積 駿^{1,2}, 植田 暢大³, 吉野 幸一郎^{4,2,1}

1. 奈良先端科学技術大学院大学
2. 理化学研究所 ガーディアンロボットプロジェクト, 3. 京都大学, 4. 東京科学大学

実世界対話の意味理解

テキストの照応・述語項構造解析

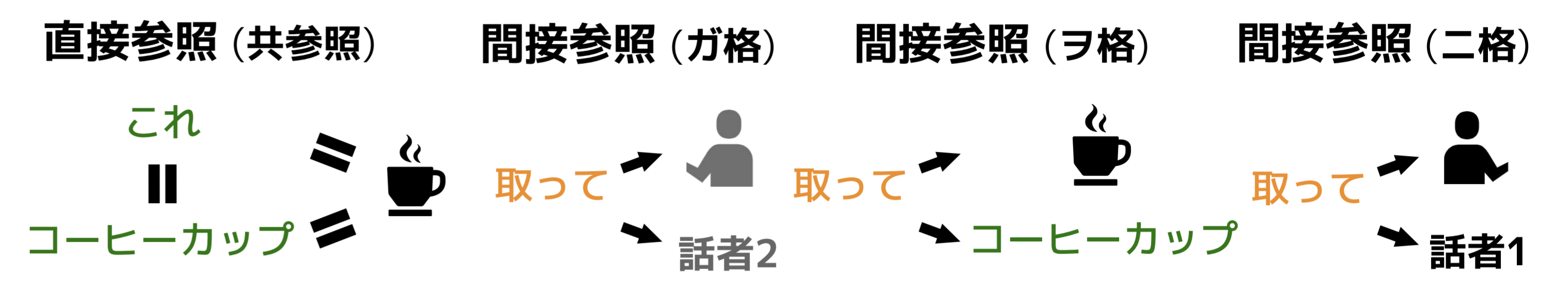
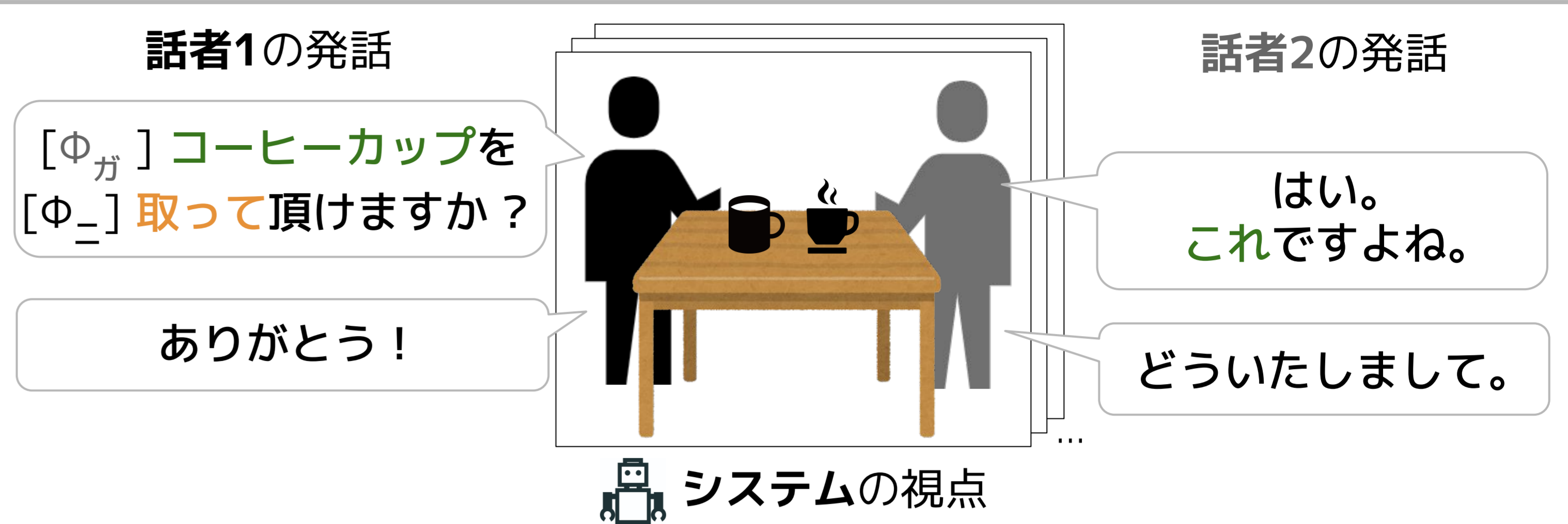
語句間の参照関係を特定するタスク

- 直接参照: 共参照と対応 コーヒーカップ = これ
- 間接参照: 述語項構造や橋渡し照応と対応 取って → コーヒーカップ

マルチモーダル参照解析 [Ueda+, 2024]

- 語句が指す物体の参照関係を特定するタスク 取って →
- フレーズグラウンディング: 語句から物体の直接参照のみを特定する場合 コーヒーカップ =
- テキストに含まれるイベント「誰が誰に何をどうする」を物体と紐付けて理解するシステムが実現可能

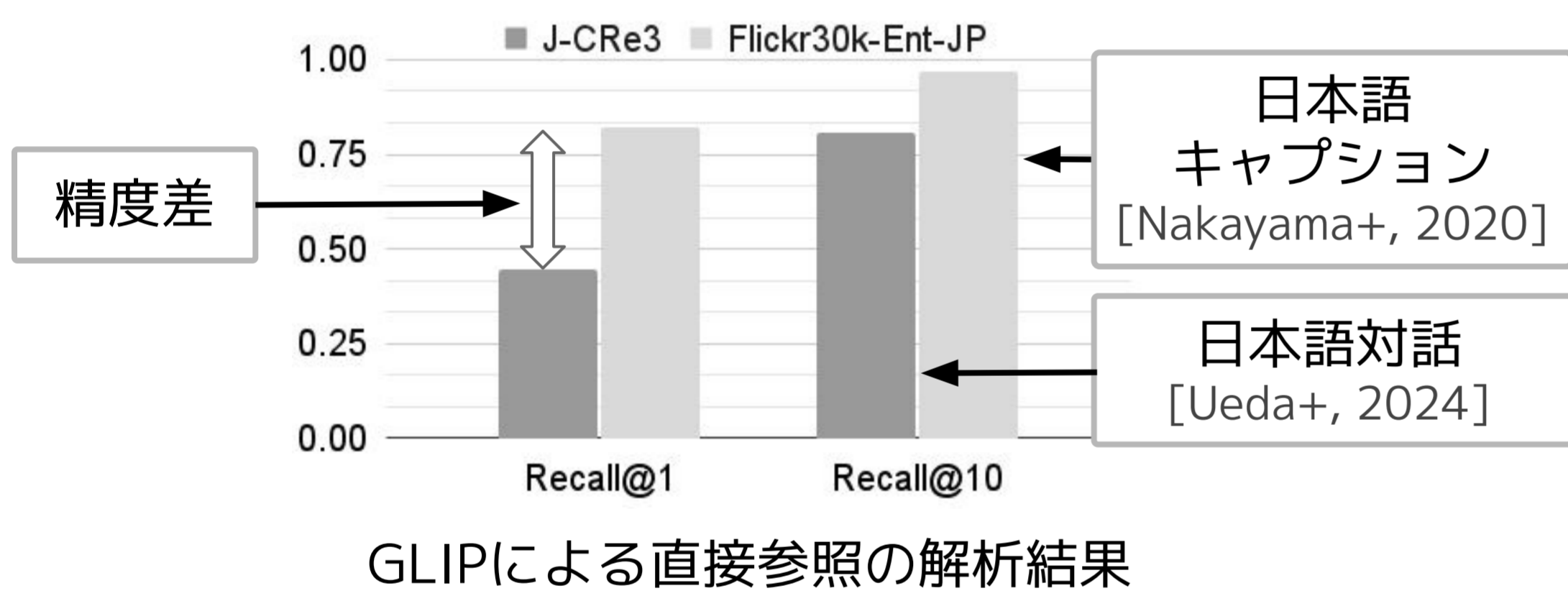
2者の実世界対話をシステムが解析する例



既存モデルの課題

- GLIP [Li+, CVPR2022]: 語句-物体間の直接参照を大量の画像とそのキャプション・クラスラベルで学習
- 対話の解析においてGLIPは:

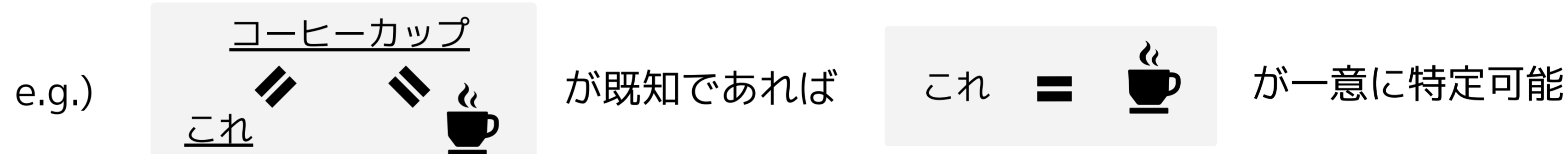
- 指示詞に対する直接参照の解析が困難 これ =
- 間接参照の解析を扱えず省略された語句の解析が困難 取って →



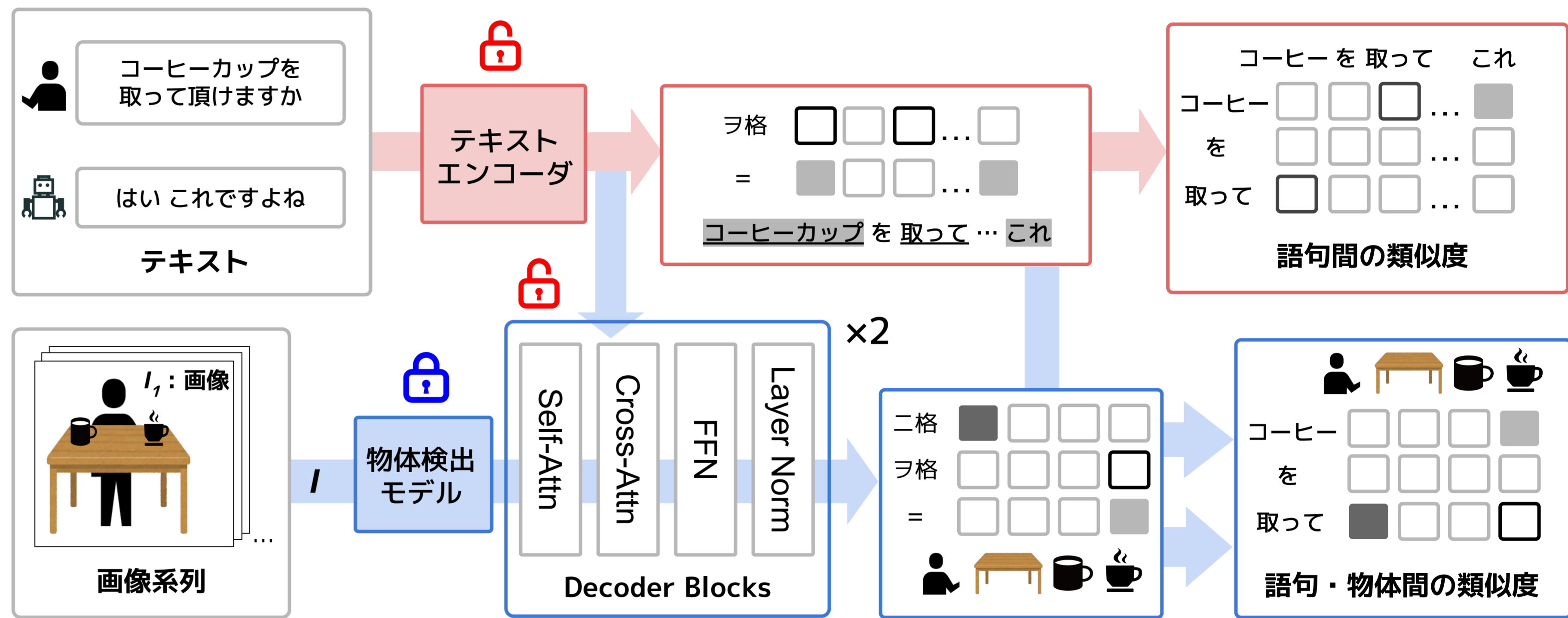
①と②の曖昧性解消を通して対話へのマルチモーダル参照解析の性能向上を目指す

提案手法

照応・述語項構造の知識を考慮して 語句-物体間の解析精度向上を図る



照応・述語項構造解析 と マルチモーダル参照解析 を統合的に扱う枠組みの提案



フレーズグラウンディングの結果

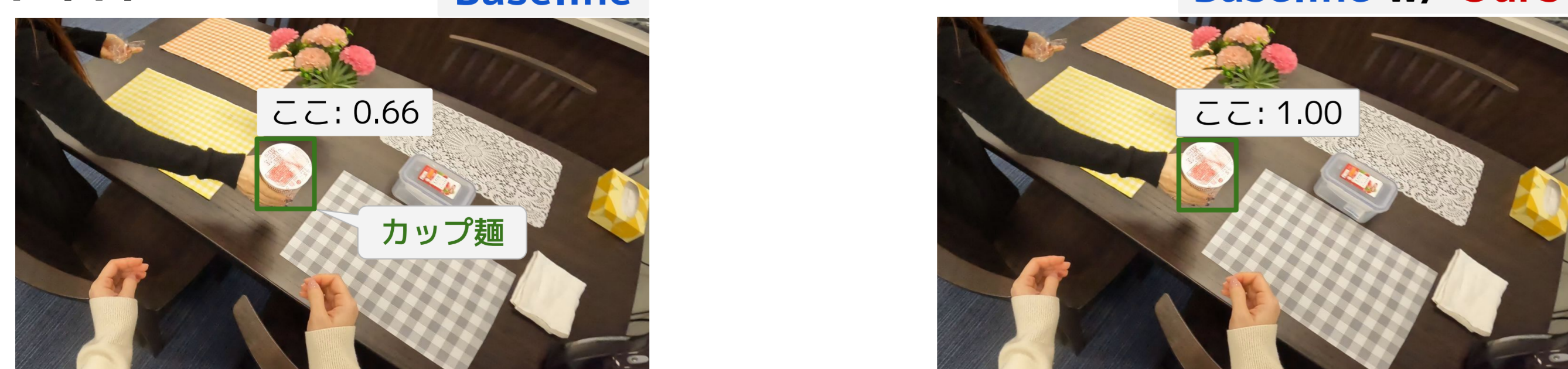
比較モデル

- Baseline
 - Baseline w/ Ours
 - Baseline w/ KWJA [Ueda+, 2023]
 - GLIP [Li+, CVPR2022]
1. 共参照関係でテキスト解析モデルを学習
2. 参照解析モデルを追加で学習
- 英語データ [Krishna+, 2017, Hudson+, 2019] による事前学習あり

定量評価



定性評価



お湯が沸いたら、**ここ**に入れてくれる?

共参照解析により:

- 指示詞に対する解析精度が向上
- 指示詞から物体への予測の確信度が強まる

マルチモーダル参照解析の結果

解析対象

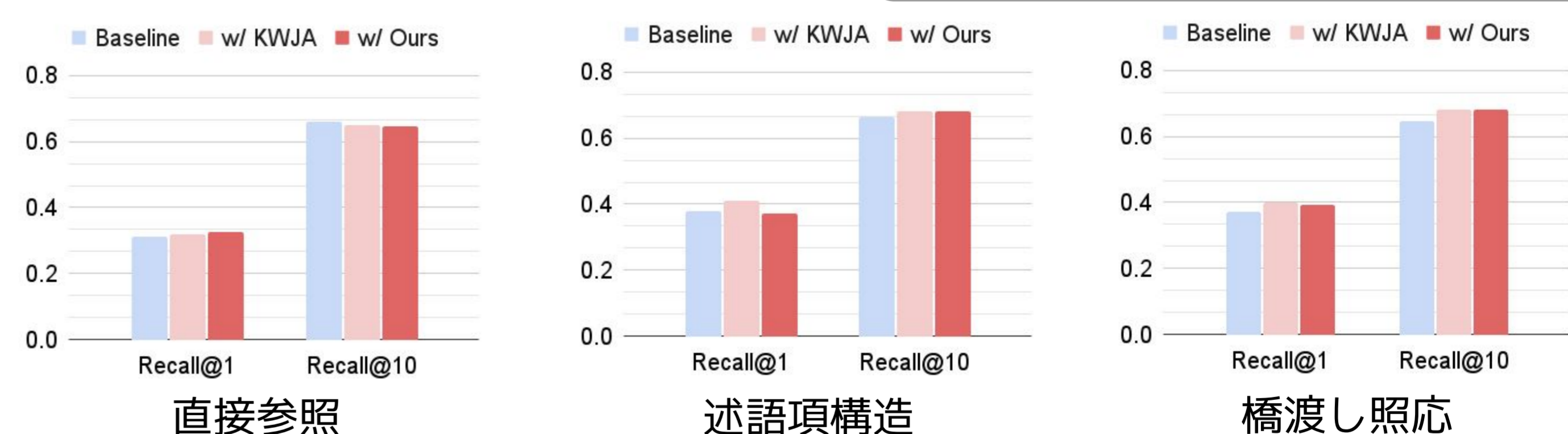
- 直接参照
- 間接参照
 - 述語項構造 (ガ格, ヲ格, ニ格, デ格)
 - 橋渡し照応 (ノ格)

比較モデル

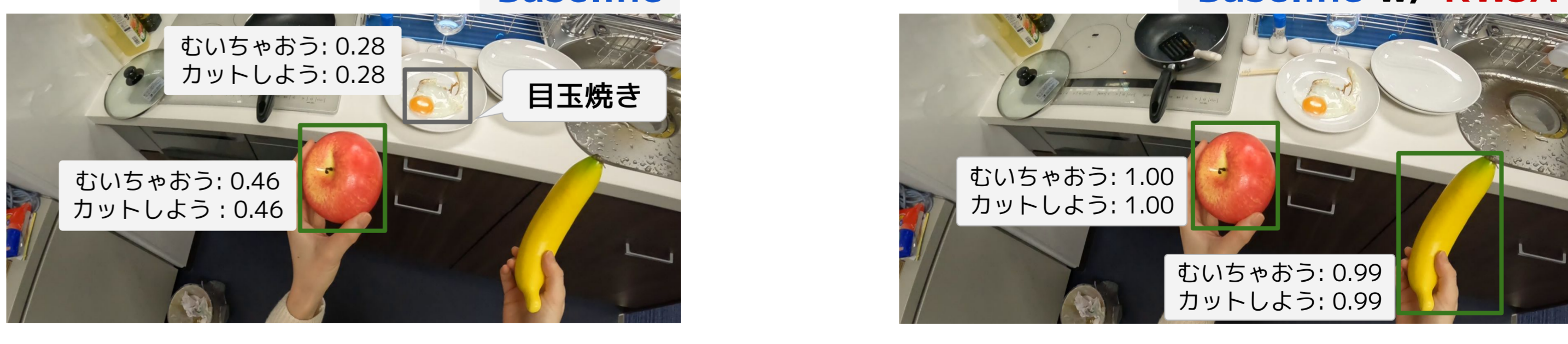
- Baseline
- Baseline w/ Ours
- Baseline w/ KWJA

照応・述語項構造解析を経た参照解析モデル

定量評価



定性評価



[バナナとリンゴを] 両方むいちゃおうか。で、3人分に**カット**しよう。

照応・述語項構造解析により:

- 語句-物体間の間接参照の解析精度が向上
- 述語から物体への予測の確信度が強まる