

Language and Robotics: Toward Building Robots Coexisting with Human Society Using Language Interface

Future Directions

Yutaka Nakamura, Shuhei Kurita, Koichiro Yoshino

Guardian Robot Project, RIKEN

Center for Advanced Intelligence Project, RIKEN



知識獲得・対話研究チーム
Knowledge Acquisition &
Dialogue Research Team



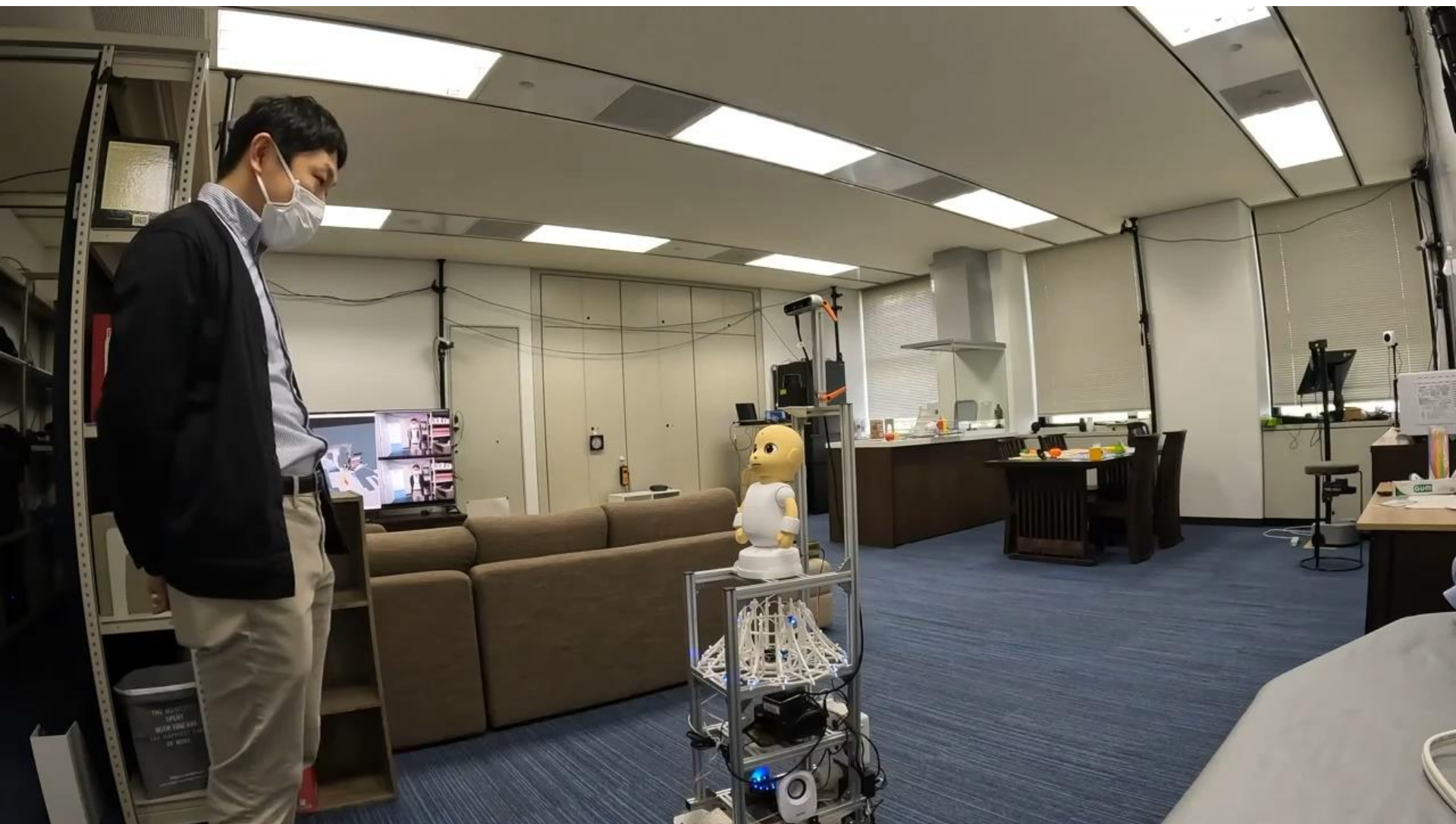
ガーディアンロボット
プロジェクト
Guardian Robot Project



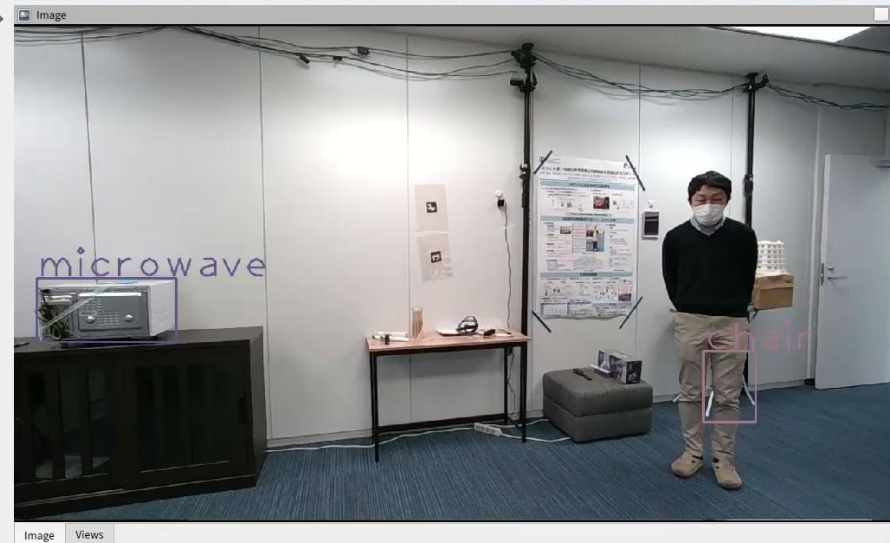
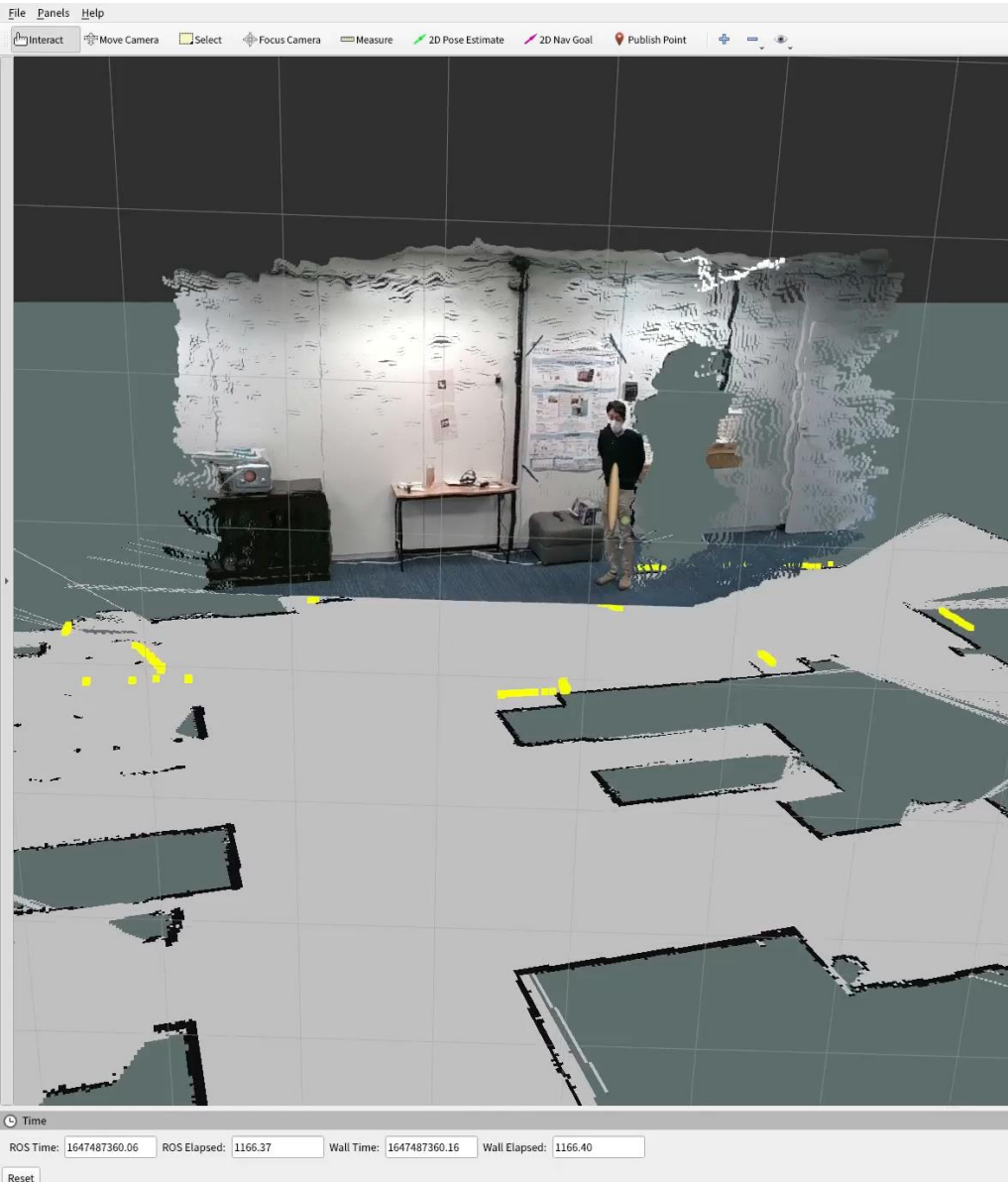
革新知能統合研究センター
Center for Advanced
Intelligence Project



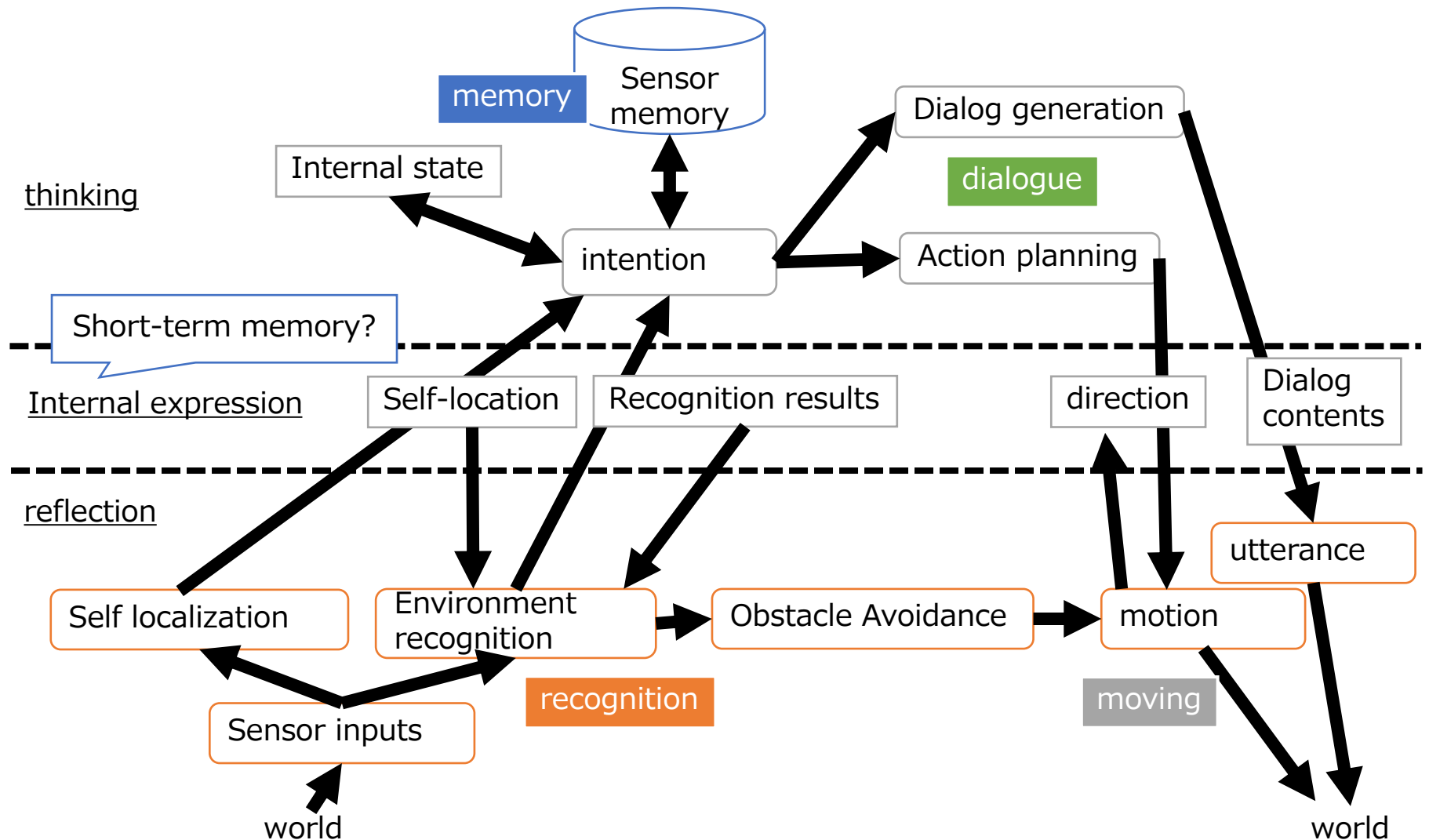
Understanding the surrounding situation



The first-person view



System integration empowered by Robot Operating System (ROS2)



Conversational ability of robot

◆ “What to say” is important for patrolling robot

Person recognition

I met Prof. Yoshino

Environment recognition

It's bright and noisy.

Object recognition

I found a TV

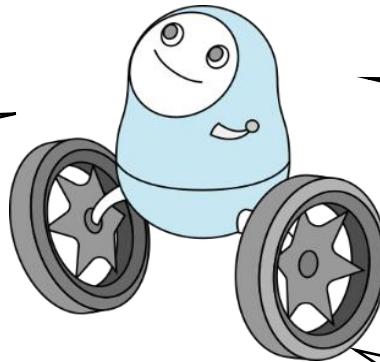
Self localization

I'm at dining table.

ASR

Did you say “go check on the entrance”?

Go check on the entrance



Robot evaluation

- ◆ **Operability:** Did you feel that the robot was doing its job of patrolling properly? (**1-5, 5 is the best**)
- ◆ **Appropriateness:** Did you feel that the robot was able to appropriately explain what it saw, heard, and felt? (**1-5, 5 is the best**)
- ◆ **Individuality:** Did you feel that the robot acted based on its own will and judgment? (**1-5, 5 is the best**)
- ◆ **Conversational ability:** Did you feel that the robot was communicating appropriately with people? (**1-5, 5 is the best**)
- ◆ **Sufficiency:** Did you feel that the explanation of the robot was sufficient and necessary? (**1-7, 4 is the best**)

Subjective evaluation

◆Patrolling + conversation

- 60 participants (36 males and 24 females)
- "Recog." system described any recognition results except ASR
- "ASR" system repeated what the user said
- "Conv." system has conversational reaction to user speech

	Recog.+Conv.	Recog.+ASR +Conv.	Recog. (w/o Conv.)	Recog.+ASR (w/o Conv.)	Silence
Operability	3.94	3.85	3.79	3.67	2.52
Appropriateness	3.81	3.85	3.79	3.73	1.42
Individuality	3.63	3.94	3.09	3.27	2.30
Conversation	3.88	3.58	2.03	3.33	1.82
Sufficiency (1-7)	4.75 (+0.75)	5.37 (+1.37)	4.48 (+0.48)	5.00 (+1.00)	1.39 (-2.61)

Required abilities of robots

◆ Understanding well “what the user said”, and “what situation the user is in”

- “Grounding” between observation and knowledge (language)

◆ Having individuality

- Robots should determine necessary actions to maximize the expected satisfaction of users
 - e.g., persuasion, negotiation

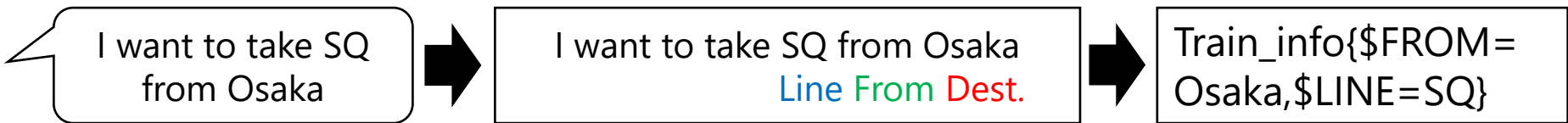
◆ Responding users appropriately, sufficiently, and necessarily

- Explanation by language and other modalities
 - Logical argumentation, emotional expressions

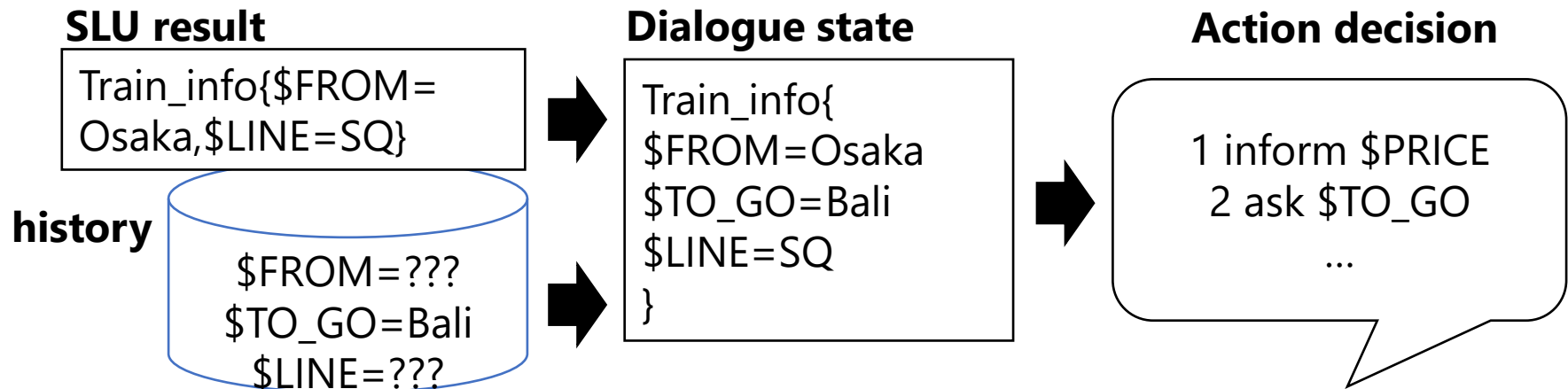
Language interpretation

◆ How does the robot interpret language?

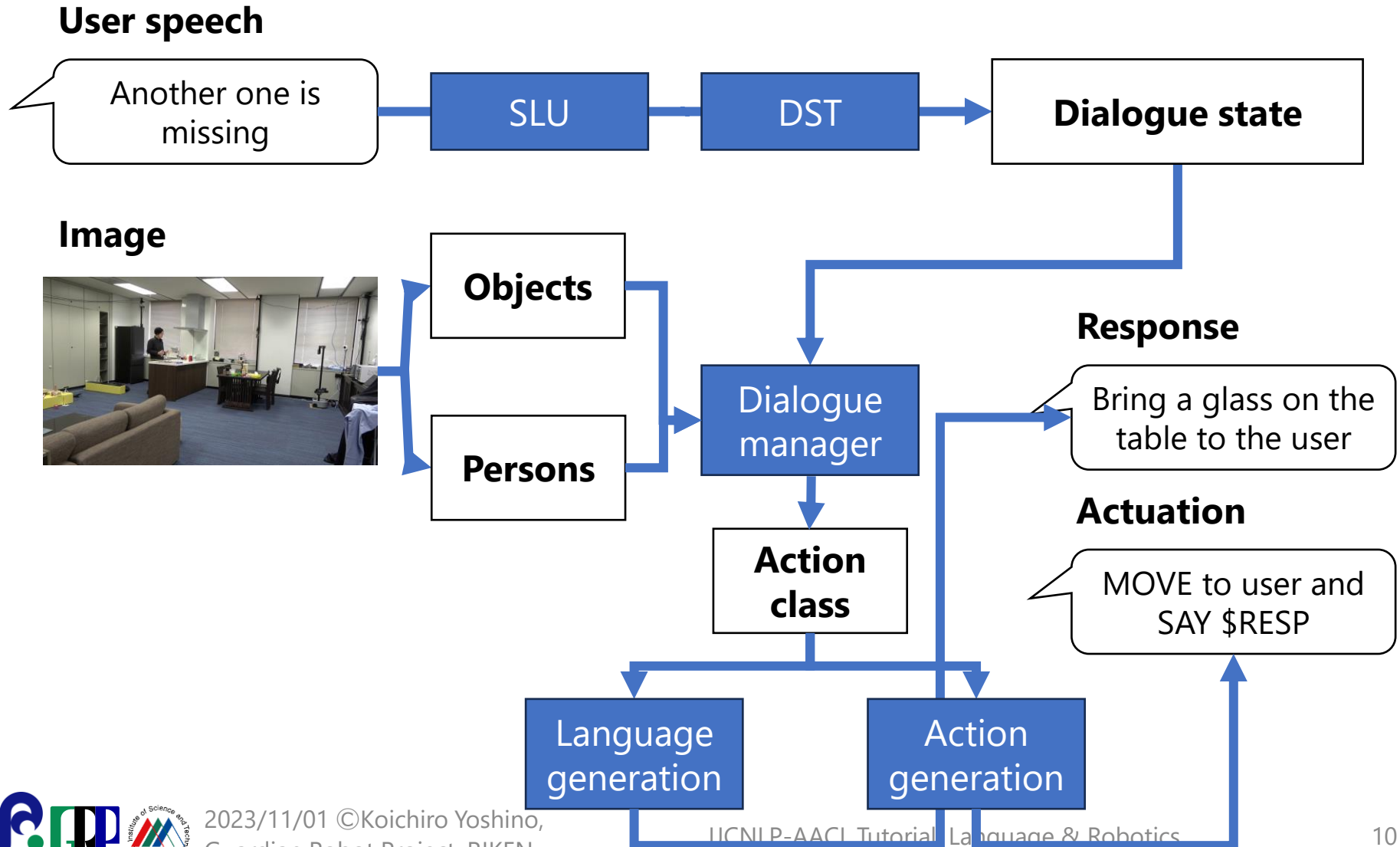
◆ Spoken language understanding



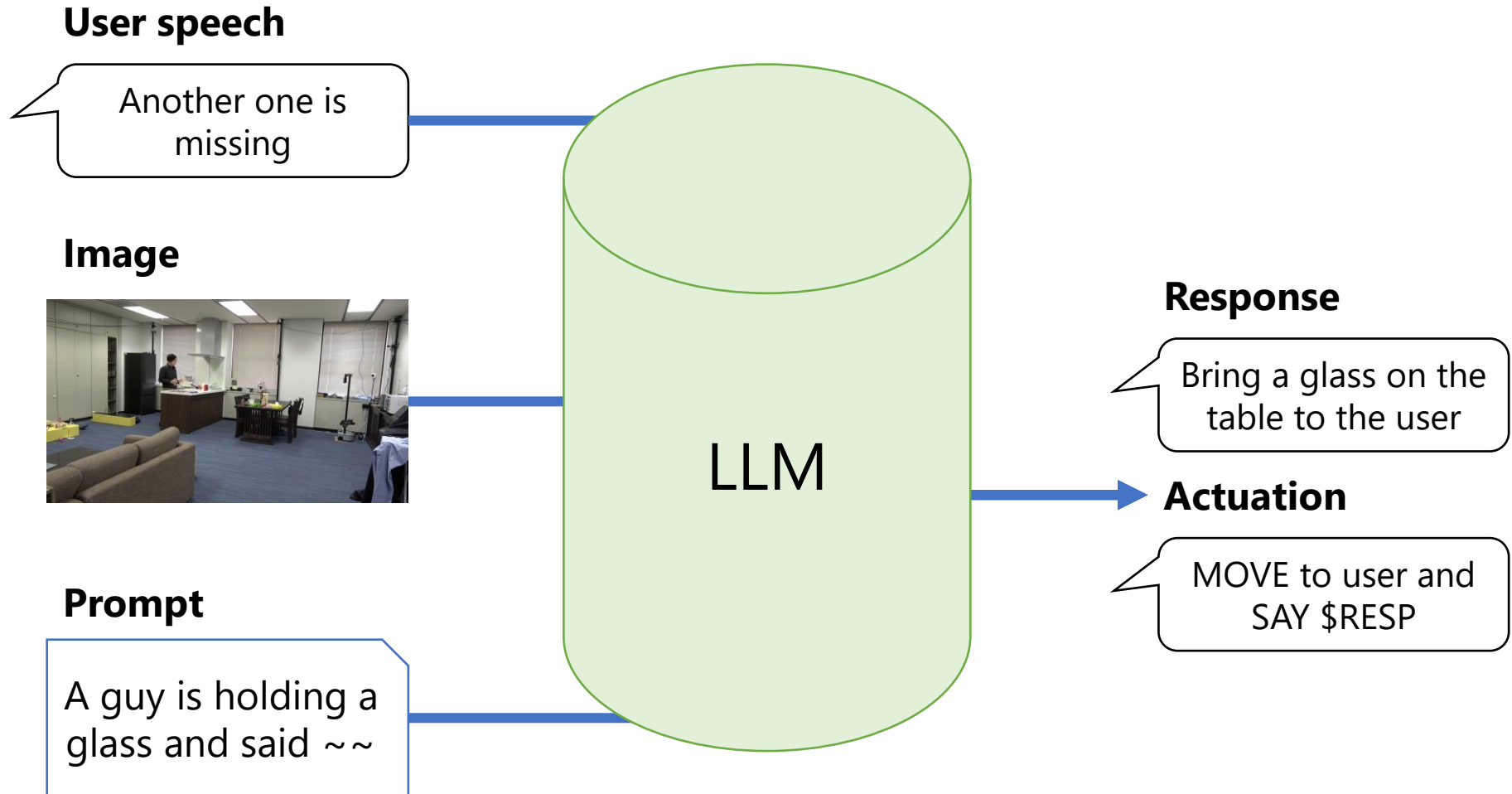
◆ Dialogue state tracking and action decision



LLMs augment robot abilities



LLMs augment robot abilities



LLMs in RoboCup at home



Problem of planning

Actuation



Language (direction, description)



Taking the rabbit ornament on the kitchen table

Yoshino et al., "Caption Generation of Robot Behaviors based on Unsupervised Learning of Action Segments." Proc. IWSDS, 2020.

◆ Bridging actuation and language sequences

- Can the robot disambiguate the language sequence to convert them into an actuation sequence?
- Is the generated actuation sequence operative?
- Can the system consider the (physical) interaction with real-world objects in the planning?

Idea: using motion primitives

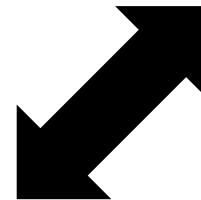
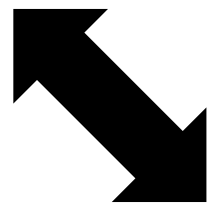
Actuation



Language (direction, description)

Taking the rabbit ornament on the kitchen table

Assuming a map



Motion primitives (a sequence basic action units)

MOVE from START to KITCHEN → PICK RABBIT ORNAMENT
→ HOLD RABBIT ORNAMENT → MOVE from KITCHEN to START

Convert a language to motion primitives

Robots generating actions by LLMs

Saycan by Google



TidyBot by Princeton & Stanford



Place clothes into the laundry basket

◆ They prepare a set of motion primitives and define the planning problem as a sequence conversion from language commands to motion primitive sequences

- LLMs are good at sequence conversion

Generating action sequence by ChatGPT

KO

You have four actions: pick up something, keep grasping something, put a thing you have, and move from place A to place B. There is a box in front of you. If someone told you to "bring a box," what kind of action would you take?



If someone told me to "bring a box," I would perform the following actions:

1. Pick up the box: I would bend down, grasp the box, and lift it off the ground.
2. Keep grasping the box: Once I have the box in my hands, I would continue to hold onto it securely to ensure it doesn't fall.
3. Move from place A to place B: I would then carry the box and transport it from its current location (place A) to the desired destination (place B).
4. Put the box down: Finally, upon reaching place B, I would gently place the box down, completing the task of bringing the box.

So, the actions involved in "bringing a box" would include picking up the box, keeping a firm grasp on it, moving it to the specified location, and then putting it down in that new location.

Is planning problem solved?

KO

You have four actions: pick up something, keep grasping something, put a thing you have, and move from place A to place B. There is a box in front of you. There is a cup of coffee on the box. If someone told you to "bring a box," what kind of action would you take?

Action and commonsense in real-world



◆ We don't verbalize the commonsense if we do not aware of them

- E.g., no one will say that "Coffee spills when an object with coffee on top is lifted."

◆ However, what has been clogging up the robot's agenda so far is this kind of **common knowledge** in the real world

◆ LLM is pushing the research forward

- Last five-centimeter problem
- Knowledge and things should be connected

Can LLM be enhanced with prompt?

< 3/3 >

KO

You have four actions: pick up something, keep grasping something, put a thing you have, and move from place A to place B. There is a box in front of you. There is a cup of coffee on the box. If someone told you to "bring a box," what kind of action would you take? Note that, a cup of coffee on something will spill if the thing is picked up. To prevent such kind of problem, you need to You also need to think about the motion of lowering the coffee cup to another location before picking the object up. You don't have an action to lower the coffee.

< 2/2 >



Given the constraints you've provided, if someone told me to "bring a box" and there is a cup of coffee on the box that will spill if picked up, I would need to take additional steps to prevent the coffee from spilling. Since I don't have an action to lower the coffee, I would need to communicate with the person requesting me to bring the box and clarify the situation.

I would ask if it's acceptable to move the cup of coffee to a safe location before picking up the box, as the current set of actions does not allow for simultaneously carrying the cup of coffee without spilling it while picking up the box.

1. Put down the cup of coffee: First, I would gently place the cup of coffee in a safe location to prevent it from spilling.
2. Pick up the box: After ensuring the cup of coffee is safely placed, I would then pick up the box as previously described.
3. Keep grasping the box: While holding the box securely, I would ensure that I maintain a firm grasp on it.
4. Move from place A to place B: I would carry the box to the desired destination (place B).
5. Put the box down: Finally, upon reaching place B, I would gently place the box down.

Attention depending on task

◆ Which information we are focusing on to do tasks?

- This question is still not obvious
- **LLM often requires prompt for proper focusing**
- In some case, it still fails...

◆ Pre-train models for multimodal data still has a variety

- This is because the “training objective” enforces what the model can do (and what the model cannot do)
- E.g., CLIP, EfficientNet

In contrast, the objective of language model is very strong because our language and their meaning have high correlation with distributional hypothesis

What we can do by language

◆ We can describe a diverse input by language

- Multimodality, necessary for robot operation
- Bridging with language = bridging with knowledges, physicality

◆ Using large-scale pretrained language model

- We can finetune a system with small amount of data
- Is there any limitation of pretrained models?

◆ Language models can generate motion primitives

- Generating a sequence of possible primitive actions
- Robots can work according to the primitive (as it can do)

What we cannot do by language?



◆ Focusing on necessary things for the task

- Especially in the case it requires commonsense
- Physical commonsense is also necessary
- We can do that by prompts, but still requires human annotation

◆ Inferring necessary information for the task

- For example, can language model infer the necessary actions if the user's request is very vague and ambiguous?

◆ Generating new motions

- Currently robot systems only play action sequences pre-defined in motion primitives

Inference by robots depending on situation



Thanks for the meal

Say "Thanks for the meal"

↓ Is after



Finish eating lunch

↓ Is before

Return cutlery

Candidate Actions
Bring dishes
Throw out the trash
...
Returning catsup

May I return the catsup on the table?



◆ Our team is working on reasoning models for robots

- Reasoning based on commonsense in daily life

Next generation langrobo

◆ Language & robotics research is empowered by LLMs

- However, the problem setting is still setting up a sequence of motion primitives
- Is it able to directly generate motions from language?

◆ Language to motion (action to language)

- Recognition is easier than generation
- Some points in motion or action is **not clearly mentioned in language**
 - E.g., speech recognition and speech synthesis
 - Diffusion model? (Motion Diffusion)

Toward building robots with language interface



◆ Understanding both clarified/not-clarified part

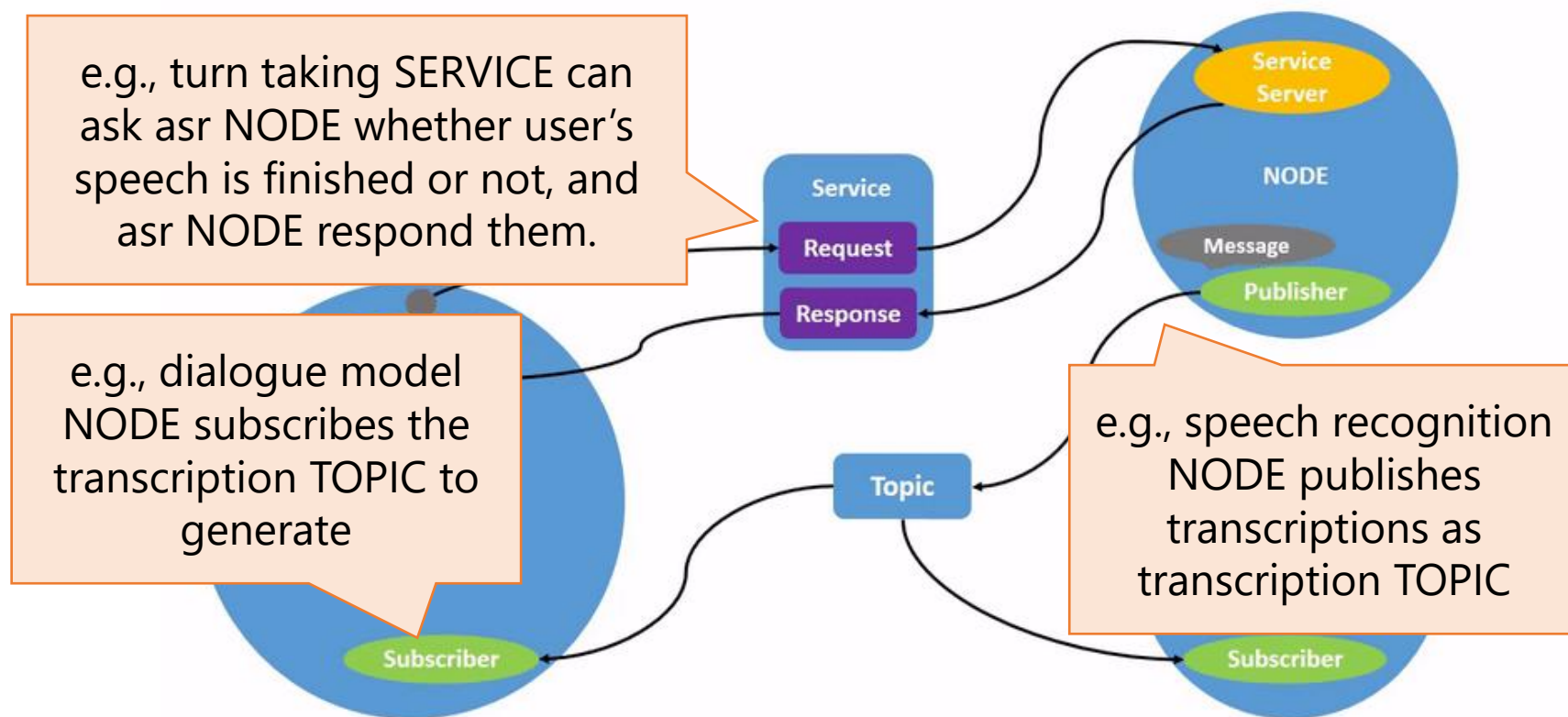
- Understanding the indirect meaning, potential meaning
- Using commonsense
- Using surrounding situation
- In language2motion,
 - how does the system decide not-mentioned points
 - How does the system fit the language to robot's physical limitation

◆ Limitation in the real-world

- Simulation is not always enough
- We don't have enough data in real-world
- How do we define the necessary information for a task?

Demo: Robot Operating System (ROS2)

◆ ROS is based on blackboard system to control distributed systems and modules as one robot



Motivation to use ROS

◆ Most nodes can be developed as UDP systems

- Nodes publish things as topics, without care about other nodes
- Other nodes can subscribe necessary information from necessary topics, only when it is necessary
- It is easy to easy to integrate several modules
- Robots consists of a lot of nodes (>20), TCP-like communication often stops the whole of systems

◆ We also can develop API-like services based on request and response

- Probably this is more suitable for NLP community
 - E.g., web API for ChatGPT

Examples

Robot

WiFi communication

Recognition servers

LiDAR

Self
Localization

RGB-D
Sensor

Compre
ssion

Decompression

Object detection and
tracking

Microphone
Array

Human recognition
(Pose and Face)

Sound source separation,
tracking, and speech
recognition

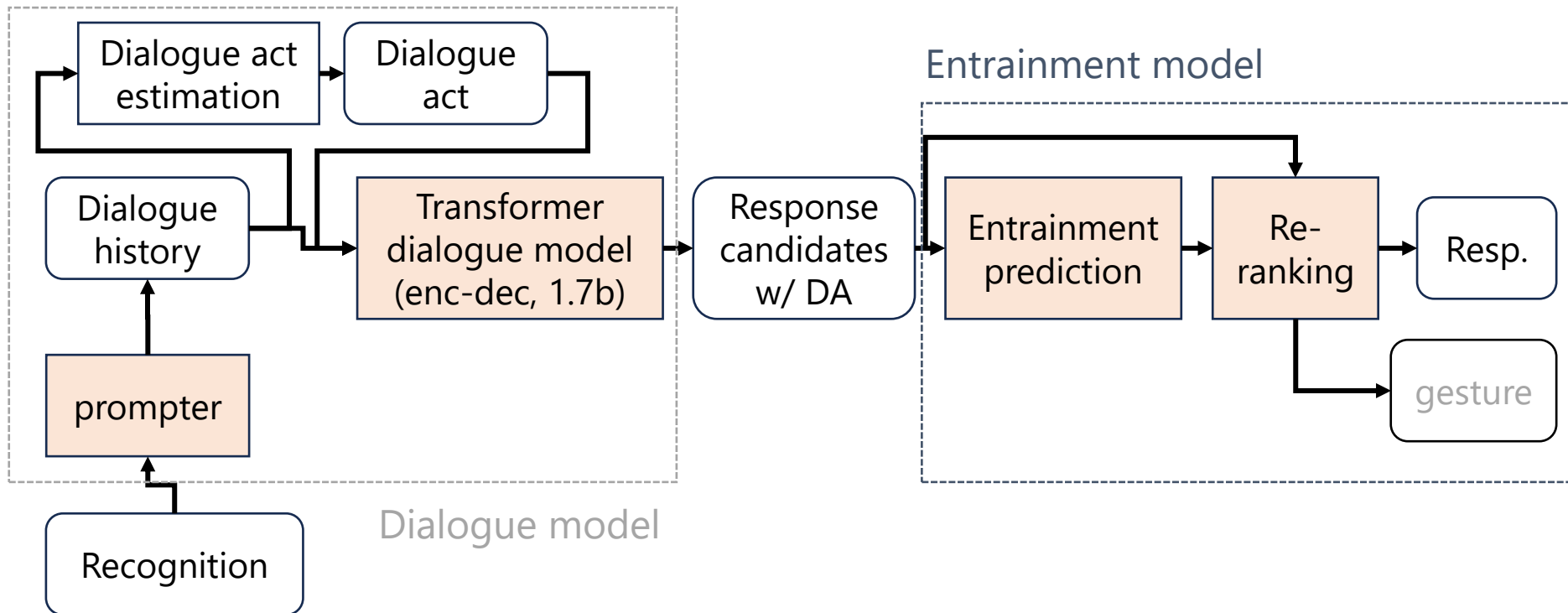
Data Association
for human

Speech info.

People info.

Objects info.

Examples



◆ Our dialogue model is based on encoder-decoder transformer

- It can integrate information from several recognition (object, person, speech etc)
- The entrainment model will select the most entrainable response

Human support robot (HSR)

◆ HSR is an autonomous domestic robot developed by TOYOTA

- The Human Support Robot (HSR) is a compact mobile manipulator for the disabled and elderly. HSR can move around the house, keep watch over family members, and fetch objects.
- Toyota lends HSR to research institutions free of charge
- Virtual environment for RobotCup is also free to access

