

DOROTHIE: Spoken Dialogue for Handling Unexpected Situations in Interactive Autonomous Driving Agents

Ziqiao Ma, Ben VanDerPloeg*, Cristian-Paul Bara*, Yidong Huang*, Eui-In Kim, Felix Gervits, Matthew Marge, Joyce Chai

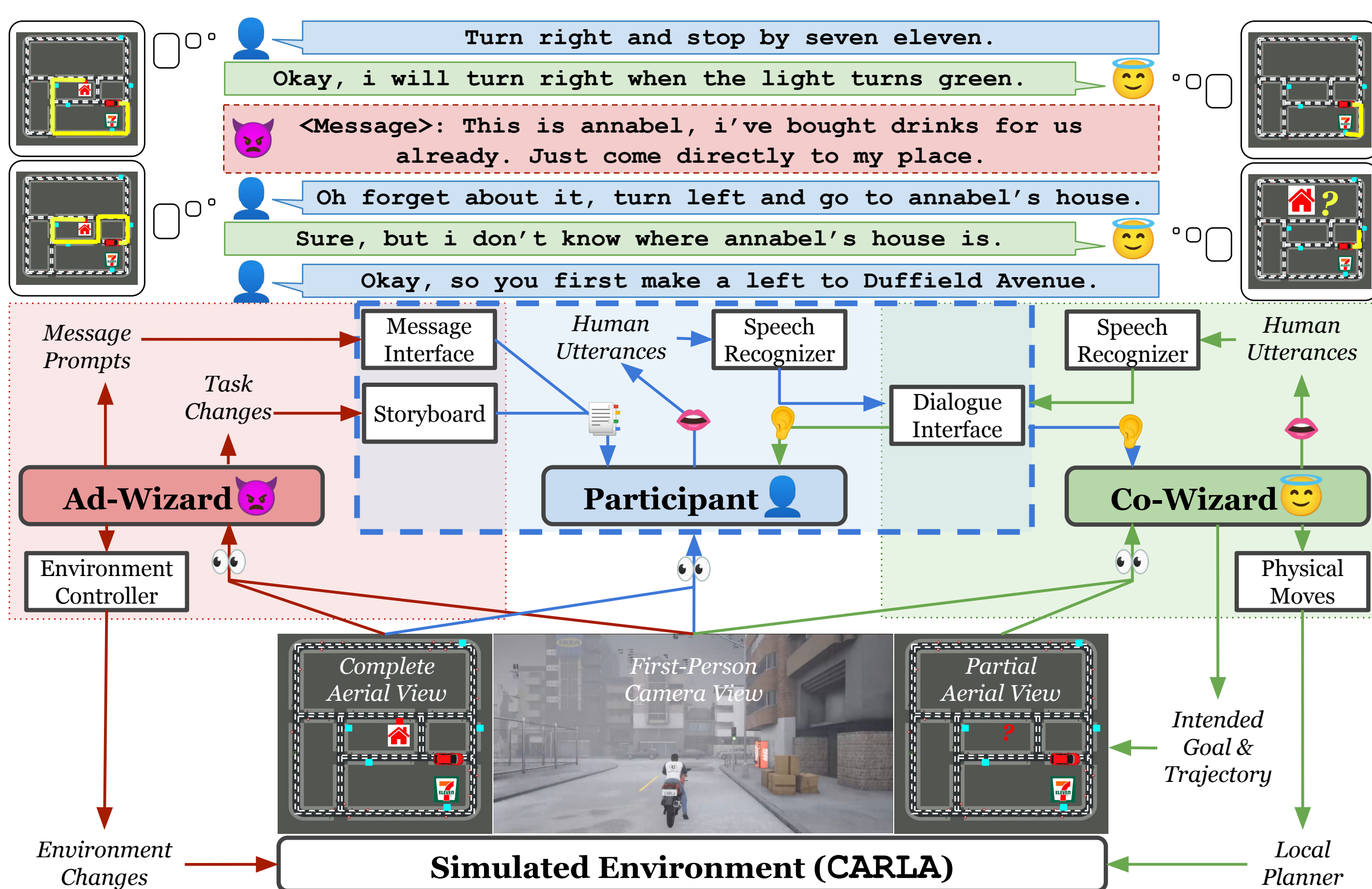


MOTIVATION

Current vision-and-language navigation (VLN) setups only reflect partial and simplified challenges compared to those faced by autonomous vehicles (AVs), which navigate in highly dynamic environment with continuous physical control. When unexpected situations arise, agents should collaborate with human operators in the form of spoken dialogue. It's thus important to empower AVs with the ability to harness human knowledge and expertise and to enable natural language communication and collaboration in tackling unexpected situations. In this work, we seek to enable AVs to navigate in **continuous** and **dynamic** environment, and communicate with human through **sensorimotor grounded dialogue**.

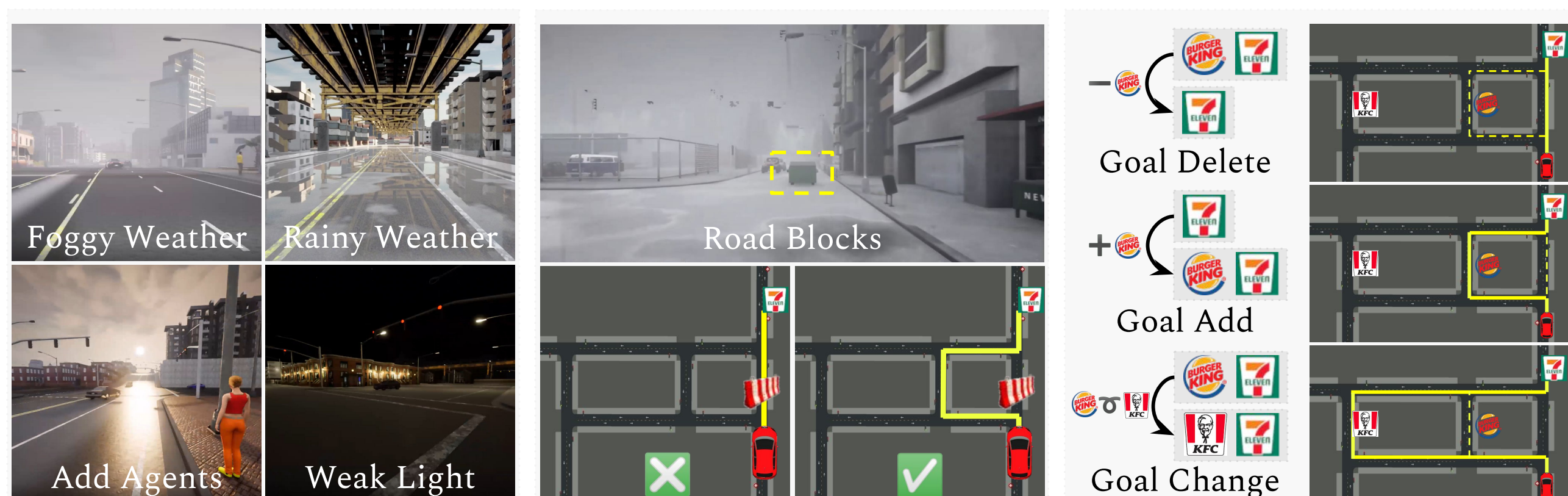
ENVIRONMENT AND HUMAN STUDIES

We developed a novel framework, **Dialogue On the ROad To Handle Irregular Events (DOROTHIE)** upon CARLA [1] to study situated human-vehicle communication based on the **Wizard-of-Oz (WoZ)** setting. We extend the traditional single-wizard to a duo-wizard setup approach by introducing a pair of Wizards.



- The **naïve participant** communicates with the vehicle to visit goal locations specified in a storyboard.
- The **Cooperative-Wizard** controls the agent's behaviors and carries language communication with the human participant to jointly achieve the goal.
- The **Adversarial-Wizard** controls the environment and task interface and introduces unexpected situation on-the-fly.

Physical Actions	Args	Descriptions
LaneFollow	-	Default behaviour, follow the current lane.
LaneSwitch	Angle (Rotation)	Switch to a neighboring lane.
JTurn	Angle (Rotation)	Turn to a connecting road at a junction.
UTurn	-	Make a U-turn to the opposite direction.
Stop	-	Brake the vehicle manually.
Start	-	Start the vehicle manually.
1-3 SpeedChange	Speed (± 5)	Change the desired cruise speed by 5 km/h.
LightChange	Light State (On/Off)	Change the front light state.
Mental Actions	Args	Descriptions
PlanUpdate	List[Junction ID]	Indicate intended trajectory towards a destination.
GoalUpdate	List[Landmark]	Indicate current goal as an intended landmark.
StatusUpdate	Tuple[Landmark, Status]	Indicate a change in task status.
KnowledgeUpdate	x, y	Guess the location of an unknown landmark.
Other	-	Other belief state updates.



DATASET

We recruit 40 human subjects and collect Situated Dialogue Navigation (SDN), a fine-grained navigation benchmark of 183 trials.

We annotate each dialogue session with 4 levels of linguistic units.

- Transaction Units (TUs)
- Exchange Units (EUs)
- Dialogue Moves
- Dialogue Slots

Metric	Value
Control Stream	18.7 h
Trimmed Audio	2.9 h
# Sessions	183
# Utterances	8415
# Words	50398
Vocabulary	1373
# Transactions	578
# Exchanges	4089
# Dialogue Moves	11623
# Slot Values	8618
# Physical Actions	9448

TASK DEFINITION

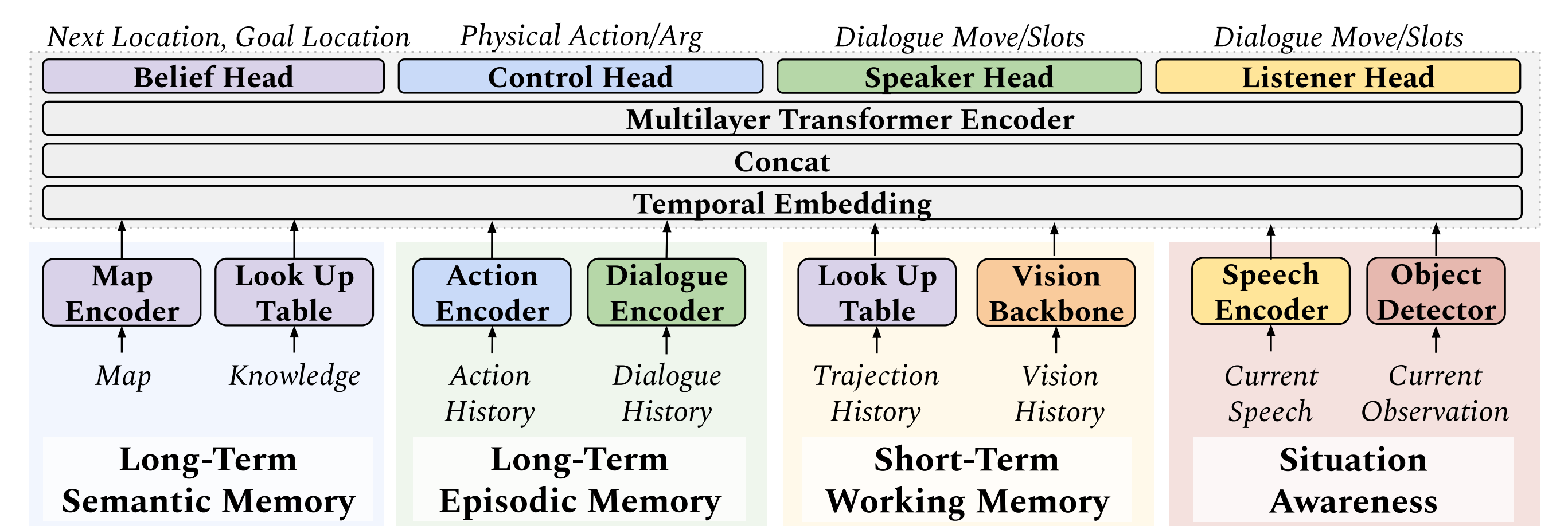
We evaluate the agent's ability in predicting dialogue moves from humans as well as generating its own dialogue moves and physical navigation actions.

- When: human speaks or agent selects a dialogue/navigation action
- Input: history of dialogue, RGB sensors, speech, and actions
- Output: human's current dialogue move/slot (UfD), agent's next dialogue move/slot (RfD), agent's next physical action (NfD)

Human Dialog Move/Slots $d_h, (s, v)$	Dialog History $\langle U^h_{human}, U^a_{agent} \rangle$	Agent Dialog Move/Slots $d_a, (s, v)$	Physical Action/Arg P_{move}, θ	Observation History O^a
QueryYN (A:UTurn, L:SE)	are you able to make a u-turn to seven eleven? no i don't think so.	ReplyN	Start	[0:38]
Acknowledge	okay.	Acknowledge	JTurn (+90)	[0:53]
Instruct (L:Ikea)	let's just go to ikea first then.	ReplyY	JTurn (0)	[1:17]
Instruct (A:UTurn)	turn left here actually.	Acknowledge	JTurn (0)	[1:27]
Instruct (A:UTurn)	make a u-turn some point.	ReplyY	Stop	[5:26]
QueryYN (A:UTurn)	is that possible?	Acknowledge	JTurn (-90)	[8:05]
Acknowledge	yes i do.	Check (L:Ikea, Stat:Ongoing)	Stop	
Instruct (A:UTurn, Str:Murfin)	then we are gonna turn right to murfin avenue.	ReplyN	JTurn (-90)	
ReplyN	but are we still going to ikea?	Explain (L:SE, Stat:Completed)	JTurn (-90)	
Clarify (Stat:Ongoing)	no.	QueryW (L:House)	Stop	
Instruct (A:Start, A:Stop)	take me to seven eleven.	ReplyW (A:House, Str:UpAnd)	Start	
QueryYN (L:House)	i think we are at seven eleven. where shall we go next?	Acknowledge	Stop	
ReplyW (A:House, Str:UpAnd)	um not yet.	Explain (L:SE, Stat:Completed)	Start	
ReplyW (A:UTurn, Str:UpAnd)	move forward a little bit and stop.	QueryW (L:House)	Stop	
Explain (A:UTurn, Str:UpAnd)	okay.	ReplyW (A:House, Str:UpAnd)	Start	
Instruct (A:House, Str:UpAnd)	can you take me to annabel's house?	ReplyW (A:UTurn, Str:UpAnd)	Start	
Instruct (A:UTurn, Str:Fuller)	where's annabel's house?	Acknowledge	Start	
Explain (Stat:Completed)	so let's just go straight on duffield ave, and turn right onto upland.	Acknowledge	Start	
	and once you turn right to upland, continue straight until fuller, and turn right onto fuller.	Acknowledge	Start	
	okay.	Acknowledge	Start	
	stop.	Acknowledge	Start	
	this is it.	Acknowledge	Start	
	where shall we go next?	QueryW (L:Queried)	Stop	

COMPUTATIONAL MODEL AND EXPERIMENTS

We present Temporally-Ordered Task-Oriented Transformer (TOTO), a Transformer-based baseline.



We summarize the experiment results in the table below.

Model	UfN (Seen)		RfN (Seen)		NfD (Seen)	
	Move Acc.	Slot F1	Move Acc.	Slot F1	Action Acc.	Act-Arg Joint Acc.
TOTO	40.9(± 3.9)	36.9(± 0.0)	29.2(± 0.7)	55.7(± 0.2)	41.2(± 2.5)	36.0(± 3.4)
TOTO (+ Belief Tracking)	39.5(± 2.2)	37.0(± 0.1)	28.8(± 0.9)	55.7(± 0.2)	40.7(± 3.6)	34.0(± 4.7)
TOTO (- Action History)	30.5(± 1.5)	36.9(± 0.0)	23.5(± 1.7)	55.7(± 0.0)	27.6(± 2.8)	24.6(± 4.0)
TOTO (- GT Transcript)	39.8(± 1.9)	36.9(± 0.1)	29.2(± 0.8)	55.6(± 0.1)	40.4(± 3.4)	31.6(± 4.3)
TOTO (- Object Detection)	42.5(± 2.8)	37.0(± 0.2)	30.4(± 0.7)	55.8(± 0.1)	39.2(± 3.5)	34.4(± 5.8)
TOTO (- Vision History)	41.9(± 1.3)	37.0(± 0.2)	29.1(± 0.5)	55.8(± 0.2)	42.0(± 3.1)	36.1(± 4.0)
TOTO (- Current Speech)	35.1(± 2.7)	36.7(± 0.5)	29.9(± 0.9)	55.9(± 0.2)	39.7(± 1.9)	33.7(± 3.0)
TOTO (- Map Knowledge)	42.6(± 1.2)	36.9(± 0.0)	29.3(± 0.9)	55.8(± 0.2)	44.6(± 3.3)	39.1(± 3.3)
Episodic Transformer	36.6(± 3.6)	37.0(± 0.2)	29.4(± 1.2)	55.9(± 0.2)	40.0(± 2.8)	32.2(± 4.0)
Fine-tuned BERT	66.8(± 2.0)	24.9(± 5.5)	52.7(± 1.0)	46.0(± 2.5)	32.4(± 2.1)	16.2(± 2.7)
Model	UfN (Unseen)		RfN (Unseen)		NfD (Unseen)	
	Move Acc.	Slot F1	Move Acc.	Slot F1	Action Acc.	Act-Arg Joint Acc.
TOTO	49.2(± 3.0)	26.2(± 0.0)	31.0(± 1.7)	54.0(± 0.7)	45.8(± 3.8)	41.1(± 2.8)
TOTO (+ Belief Tracking)	47.1(± 3.5)	26.2(± 0.0)	29.0(± 2.0)	53.7(± 0.7)	47.6(± 4.5)	38.8(± 3.1)
TOTO (- Action History)	35.5(± 3.2)	26.1(± 0.1)	28.2(± 3.9)	54.8(± 0.0)	36.8(± 0.8)	36.0(± 1.7)
TOTO (- GT Transcript)	46.7(± 2.4)	26.2(± 0.0)	31.6(± 2.6)	54.2(± 0.8)	46.2(± 5.9)	37.6(± 6.9)
TOTO (- Object Detection)	50.0(± 1.8)	26.2(± 0.1)	32.7(± 2.2)	53.8(± 1.2)	45.7(± 5.2)	40.3(± 5.4)
TOTO (- Vision History)	48.7(± 2.3)	26.2(± 0.1)	31.5(± 2.9)	54.3(± 0.7)	45.9(± 4.2)	42.3(± 3.5)
TOTO (- Current Speech)	42.8(± 2.5)	25.8(± 0.3)	33.8(± 1.4)	55.1(± 0.4)	46.5(± 4.9)	39.4(± 5.2)
TOTO (- Map Knowledge)	48.2(± 1.0)	26.2(± 0.1)	31.9(± 1.2)	54.9(± 0.8)	51.7(± 3.4)	46.0(± 4.0)
Episodic Transformer	45.1(± 3.8)	26.1(± 0.1)	33.4(± 2.2)	54.7(± 0.8)	46.6(± 3.3)	37.0(± 5.9)
Fine-tuned BERT	67.2(± 1.5)	16.2(± 3.5)	57.0(± 0.9)	46.9(± 2.2)	37.1(± 1.5)	19.6(± 3.6)

Overall, our preliminary experiment has shown that the tasks are challenging.

- TOTO is able to handle all tasks uniformly on both the seen and unseen splits of the test set, and outperform the majority of the unimodal baselines;
- The finetuned language model only masters dialogue move tasks and the Episodic Transformer underperforms language tasks.
- Action history is crucial in the understanding and prediction of navigation/dialogue actions;

REFERENCES

[1] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. CARLA: An open urban driving simulator. In Proceedings of the 1st Annual Conference on Robot Learning, pages 1–16, 2011

LINKS

