

Explicit World Models for Reliable Human-Robot Collaboration

**Kenneth Kwok¹, Basura Fernando¹, Qianli Xu²,
Vigneshwaran Subbaraju¹, Dongkyu Choi³, Boon Kiat Quek¹**

¹Institute of High Performance Computing, Agency for Science, Technology and Research, Singapore

²Institute for Infocomm Research (I²R), Agency for Science, Technology and Research, Singapore

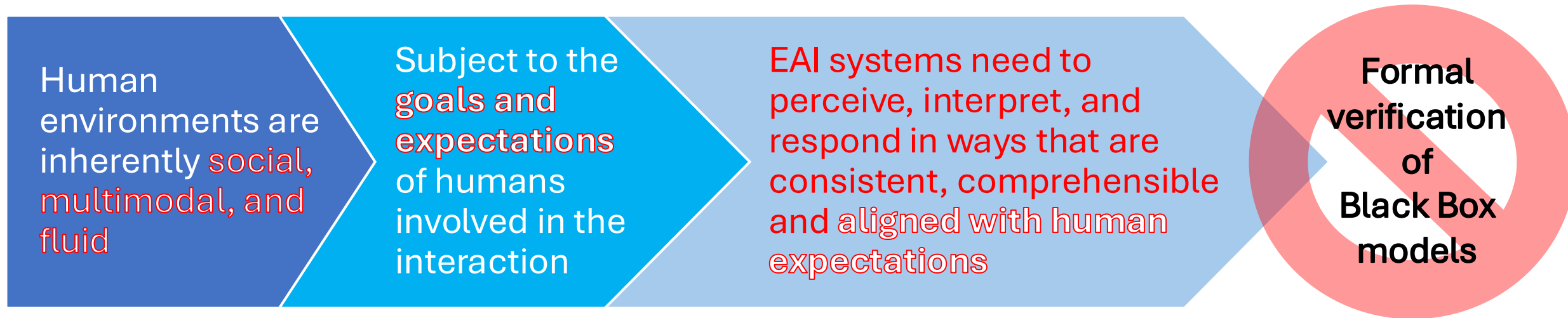
³A-FAB Technology Team, Mechatronics Research, Samsung Electronics, Korea

kenkwok@a-star.edu.sg, fernando_basura@a-star.edu.sg, xu_qianli@a-star.edu.sg,
vigneshwaran_subbaraju@a-star.edu.sg, edc.choi@samsung.com, quekbk@a-star.edu.sg



Reliable HRI with Embodied AI

Reliability in HRI is **contextually determined**:

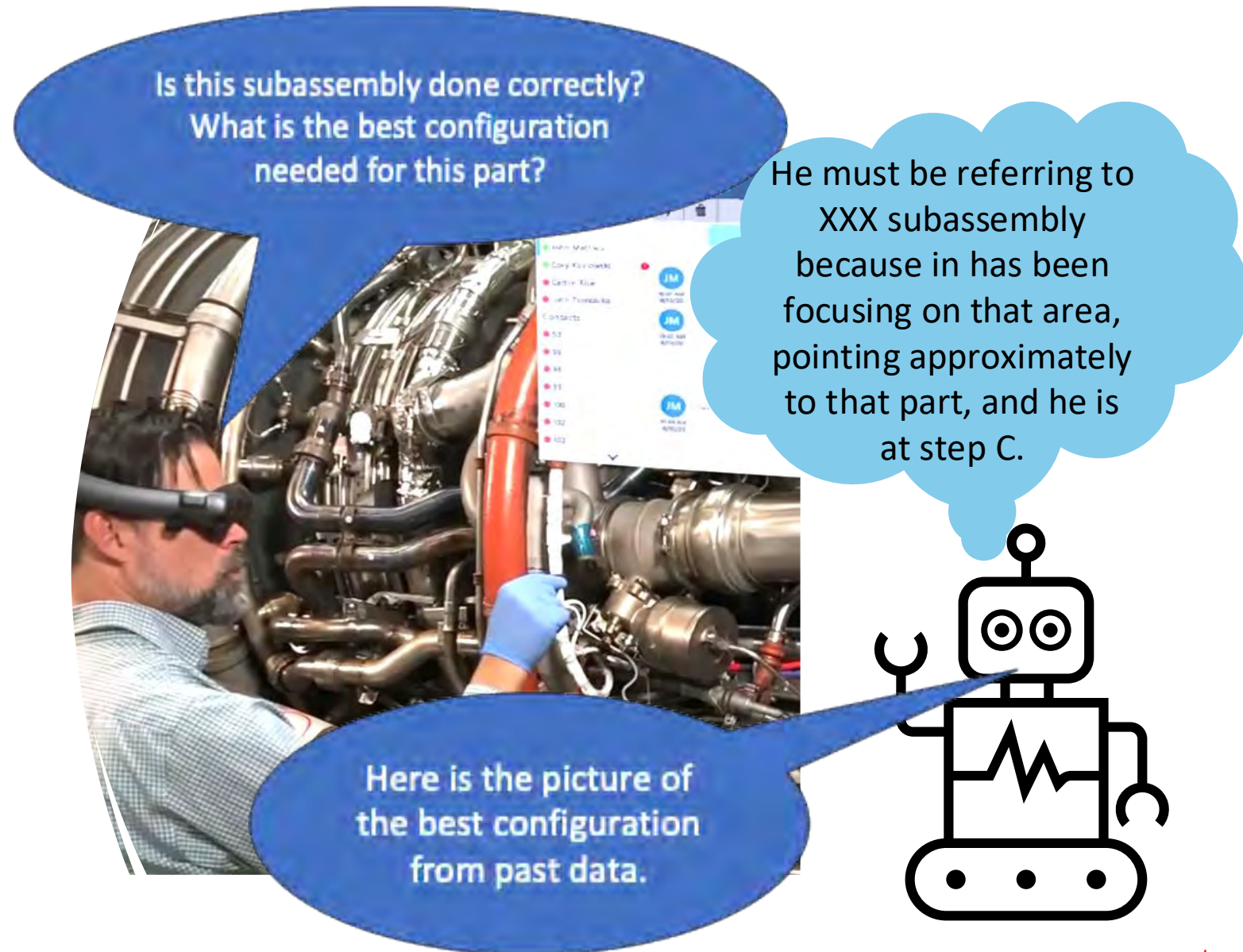


Proposed Approach:

HRI centred on building and updating an accessible **explicit world model**, representing the **common ground** between human and EAI to **align robot behaviours with human expectations**.

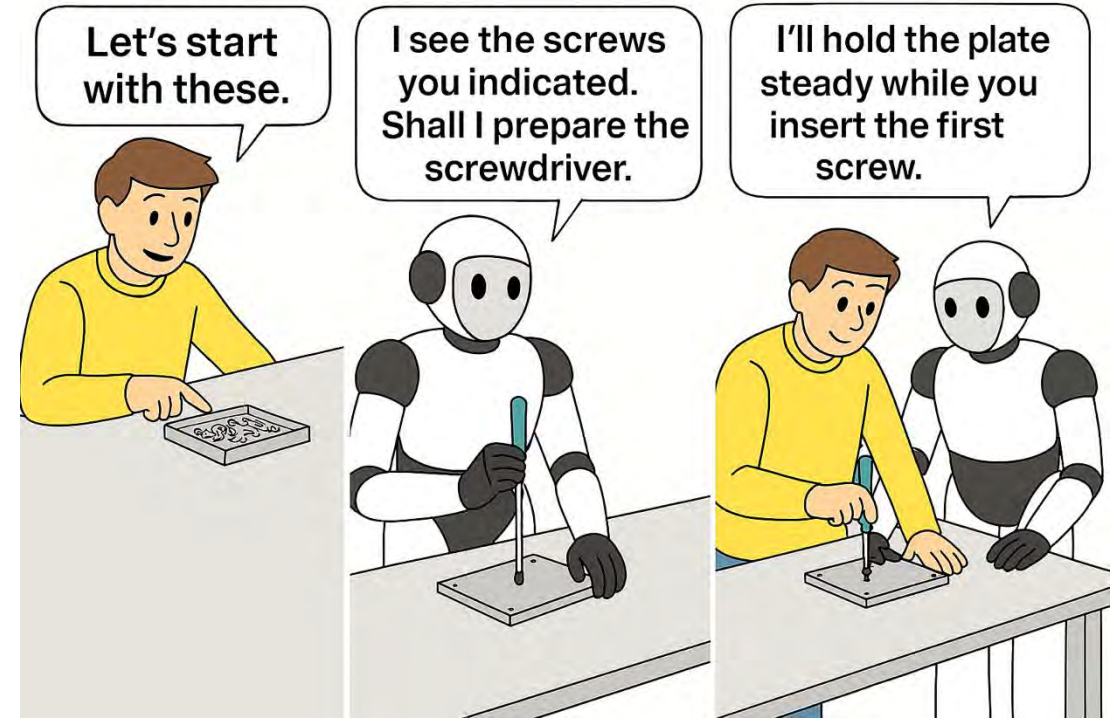
Common Ground in Human-Robot Collaboration

- Common ground = shared understanding of tasks, communications, and environments between agents (Dillenbourg and Traum, 2006)
- Aligns human and robots in terms of perception, cognition, and embodiment
- Forms the basis for cooperative action in Human-Robot Collaboration



Building Common Ground

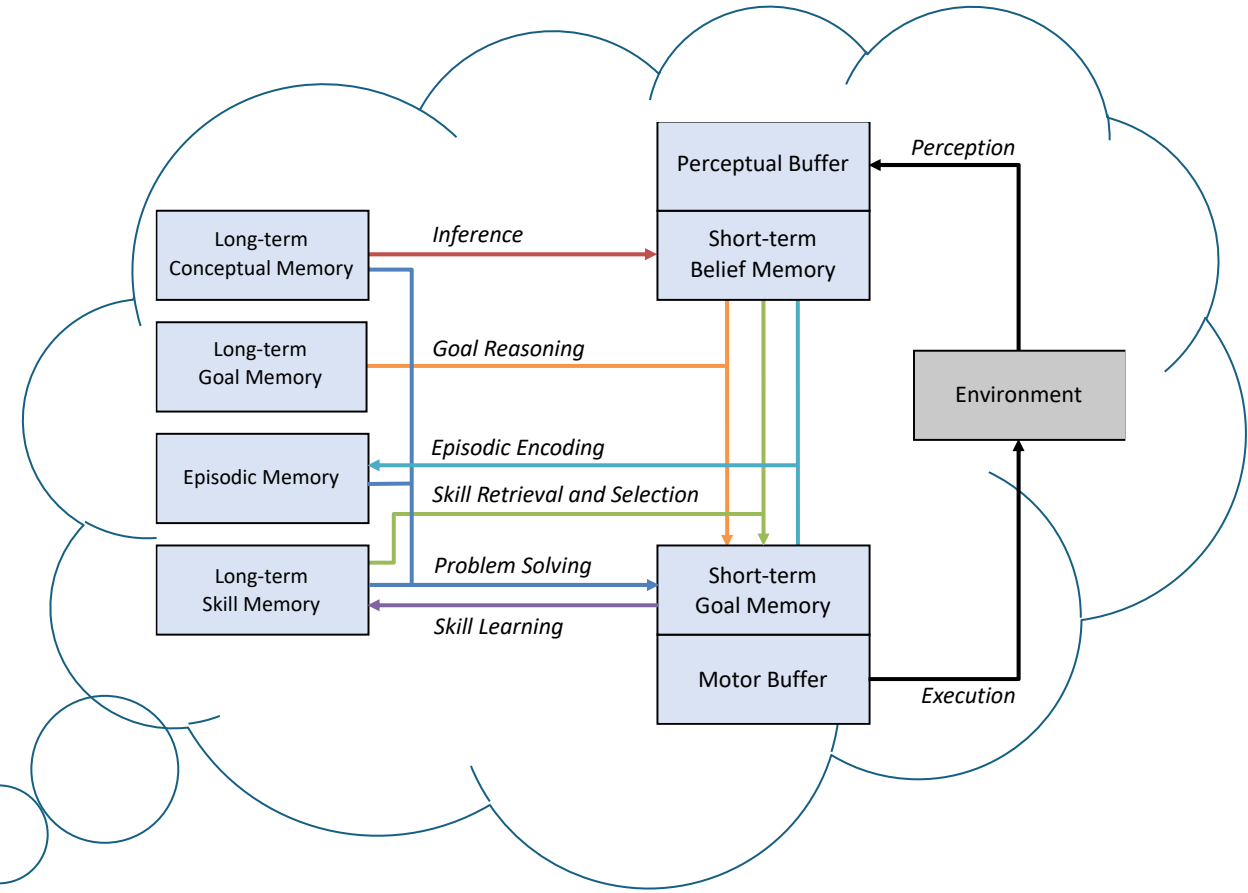
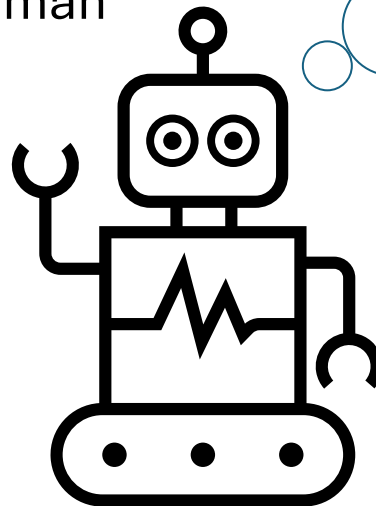
- Perceptual Grounding
 - Objects
 - Actors
 - Activity (Task), Situation (State)
- Joint Attention and Embodied Communication
 - Monitoring speech and behaviours (gaze, gestures etc.)
 - Continuous intention inference
 - Acting expressively to convey intentions



STEP	HUMAN ACTION	ROBOT ACTION	COMMON GROUND
1	Points to screws, says "Let's start with these."	Interprets gesture + speech, confirms understanding	Initial alignment on task goal
2	Nods, positions base plate	Retrieves screwdriver, avoids workspace obstruction	Dynamic coordination based on cues
3	Prepares to insert screw	Holds plate steady, verbal update	Maintains shared world model and predictability

Explicit World Model

- **Cognitive Architectures (CAs)**
 - Use explicit world models to represent the environment, relational concepts, and executable procedures
 - Raw sensory data is transformed into symbols and rules – a high level, interpretable and structured abstraction of the environment – for cognitive processing
 - Heavily dependent on human handcrafting

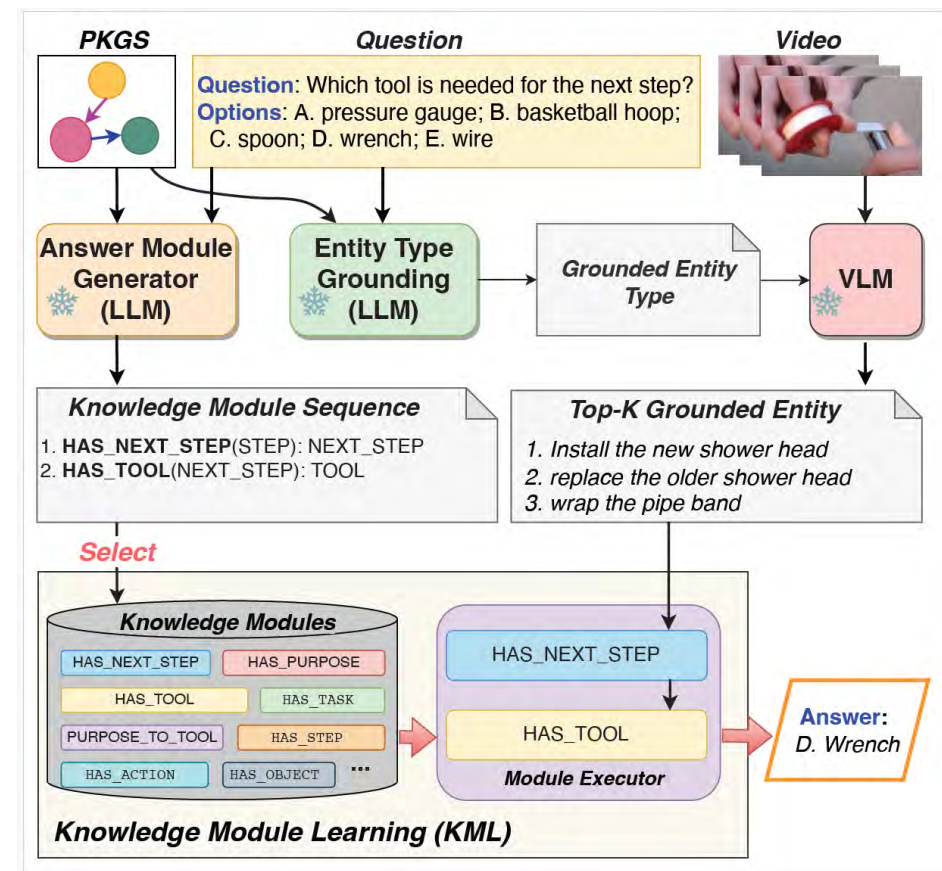
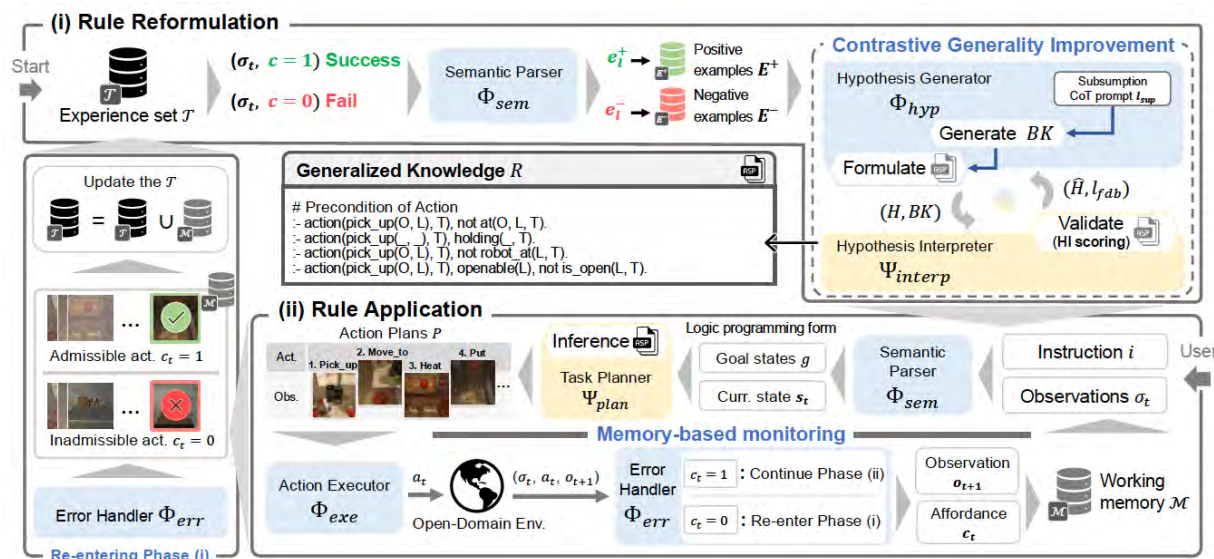


Langley, P. & Choi, D (2006). "A Unified Cognitive Architecture for Physical Agents." In *Proceedings of the 21st National Conference on Artificial Intelligence*, Vol 2, 1468-1474

Explicit World Model

• Neuro-symbolic Architectures

- Explicit, interpretable, and self-evolving world models
- Develop internal representations without human handcrafting



Knowledge Module Learning (Nguyen et al. 2025)

NeSyC (Choi et al., in press)

Choi, W.; Park, J.; Ahn, S.; Lee, D.; and Woo, H. In Press. "A Neuro-Symbolic Continual Learner for Complex Embodied Tasks in Open Domains." In 2025 13th International Conference on Learning Representations (ICLR).

Nguyen, T.-S.; Yang, H.; Neoh, T. Y.; Zhang, H.; Ee, Y. K.; and Fernando, B. 2025. "Neuro-Symbolic Knowledge Reasoning for Procedural Video Question Answering with Knowledge Module Learning (KML)." ArXiv, abs/2503.14957.



Spatial-Temporal World Models

1. World models as interface (Li Fei Fei's Marble)
 - Words and flat media → 3D assets humans can edit and share
2. World models as simulator (Deepmind's Genie)
 - Continuous, controllable video worlds that agents can interact with
 - SIMA 2-style agents built on top
3. World models as cognition (LeCun-style architectures)
 - Multimodal perception → latent variables and transition functions
 - Internal predictive state



Figure 3. Rich Supervision of 3D World Modeling for Physical Interactions, when conditioned on 3D robot point flows and partial observable RGB-D. The 3D world modeling objective enjoys dense pixel-level supervision while encoding a wide range of capabilities central to robotic manipulation. To predict full-scene evolution, the model needs to implicitly segment objects of interest, identify material property and/or articulation structure, perform implicit shape completion for contact reasoning, propagate robot-object interaction for object-object dynamics, and simultaneously considering the effects of gravity, encapsulated all in a single forward pass of the learned model.

Huang, W., Chao, Y. W., Mousavian, A., Liu, M. Y., Fox, D., Mo, K., & Fei-Fei, L. (2026). PointWorld: Scaling 3D World Models for In-The-Wild Robotic Manipulation. arXiv preprint arXiv:2601.03782.

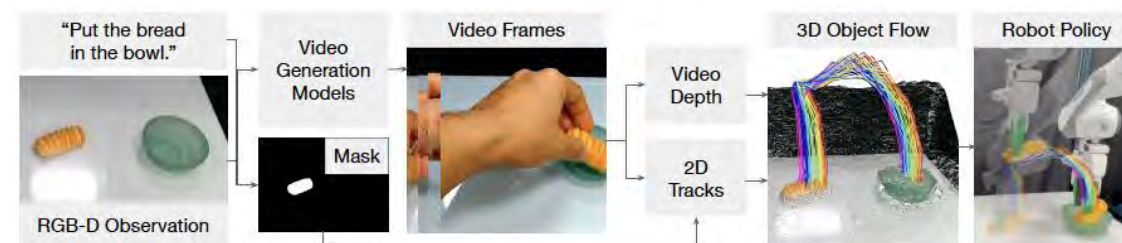
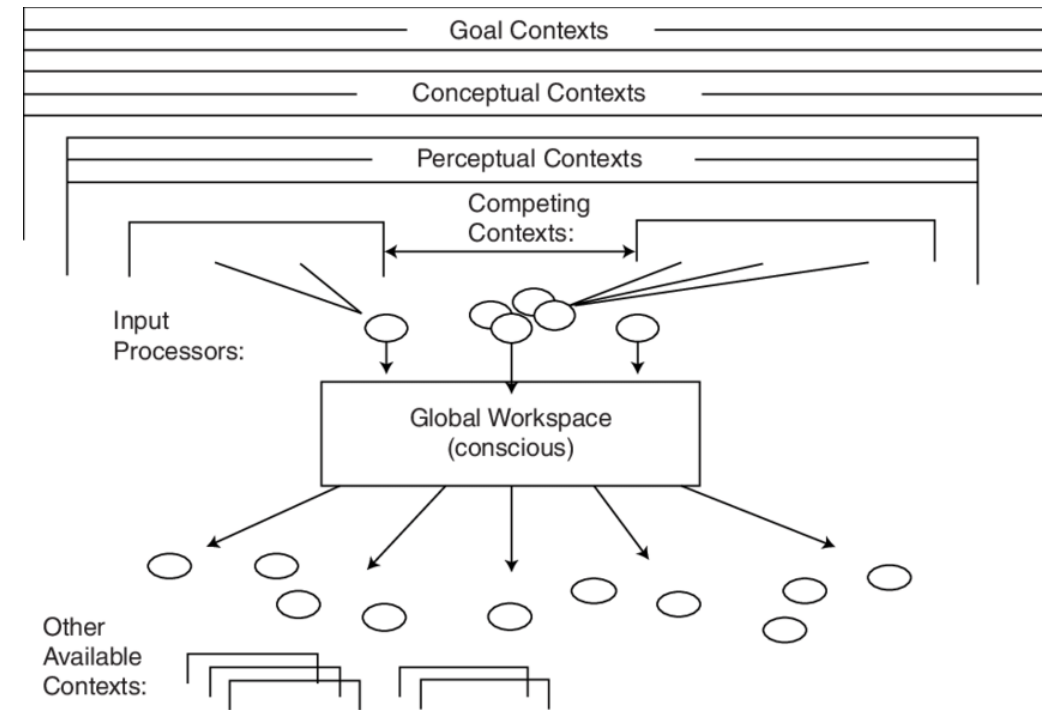


Fig. 2: An overview of Dream2Flow. Given a task instruction and an initial RGB-D observation, an image-to-video model synthesizes video frames conditioned on the instruction. We additionally obtain object masks, video depth, and point tracking from vision foundation models, which are used to reconstruct 3D object flow. Finally, a robot policy generates executable actions that track the 3D object flow using trajectory optimization or reinforcement learning.

Dharmarajan, K., Huang, W., Wu, J., Fei-Fei, L., & Zhang, R. (2025). Dream2Flow: Bridging Video Generation and Open-World Manipulation with 3D Object Flow. arXiv preprint arXiv:2512.24766

World Models as Global Workspaces

- **Global Workspace Theory** (Baars 2005) is a foundational cognitive architecture that describes *global-information sharing* mechanisms in the brain.
- World Models function as global workspaces to provide a unified, interpretable representation of the current state of perception, goals, and context that is broadcast across specialised cognitive modules — essentially acting as a *shared internal model* that enables coordination and decision-making.



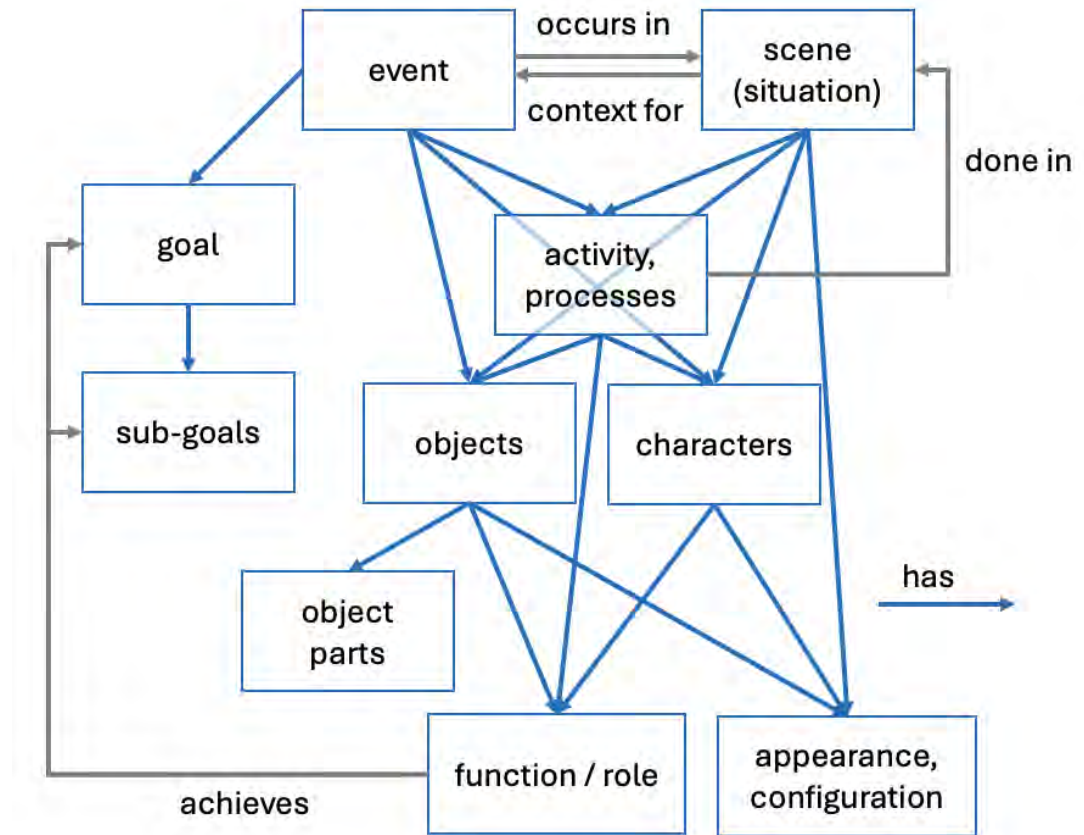
Baars Global Workspace Theory as depicted in Sun, et.al. (2007)

Baars, B. J. (2005). Global workspace theory of consciousness: Toward a cognitive neuroscience of human experience. In S. Laureys (Ed.), *Progress in brain research* (Vol. 150, pp. 45–53). Elsevier. [https://doi.org/10.1016/S0079-6123\(05\)50004-9](https://doi.org/10.1016/S0079-6123(05)50004-9)

Sun, Ron & Franklin, Stan. (2007). *Computational Models of Consciousness: A Taxonomy and Some Examples*. Cambridge Handbook of Consciousness. 10.1017/CBO9780511816789.008.

World Models as Global Workspaces

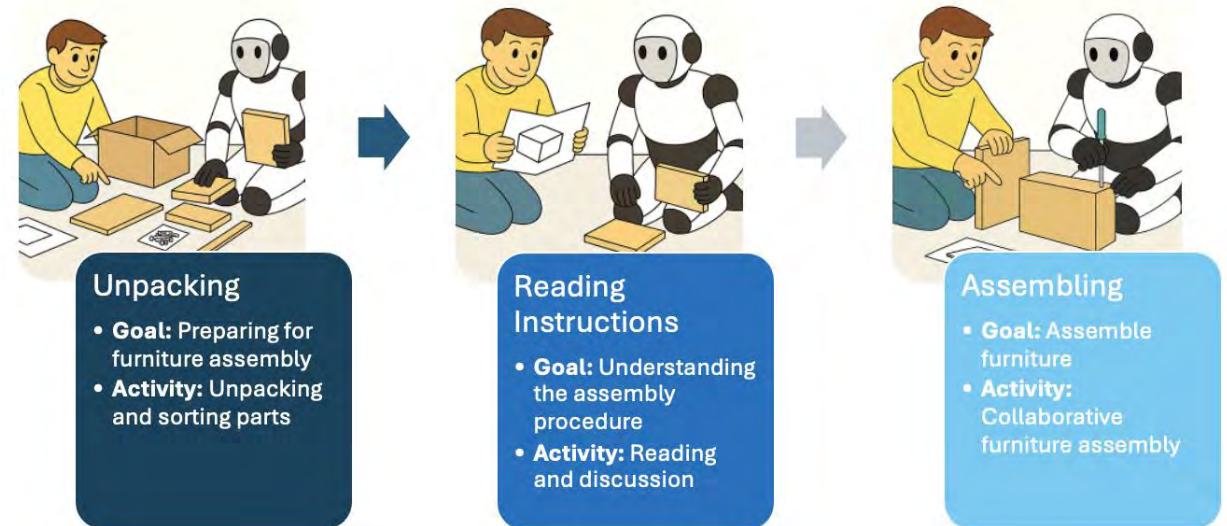
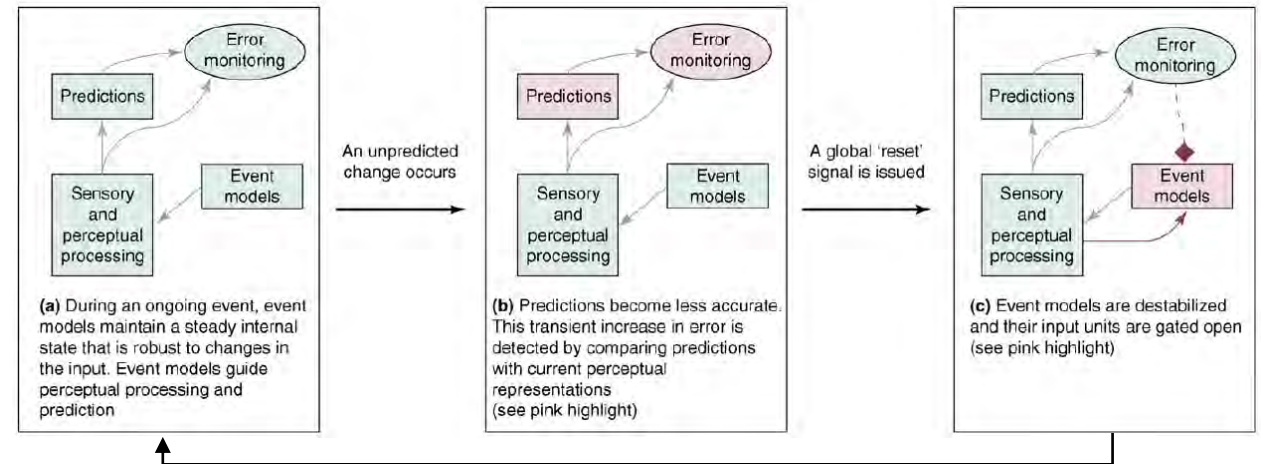
- Episodic (event-centric, relational) vs Spatial-Temporal (flow, physical)
- Event-centric Representation
 - Event, Scene (time, space, causes, characters and goals)
 - Hierarchical time scales
 - Extended (coarse grained)
 - Brief (fine grained)
- Constructing World Models
 - Event Segmentation
 - Goal Inference
 - Event Model



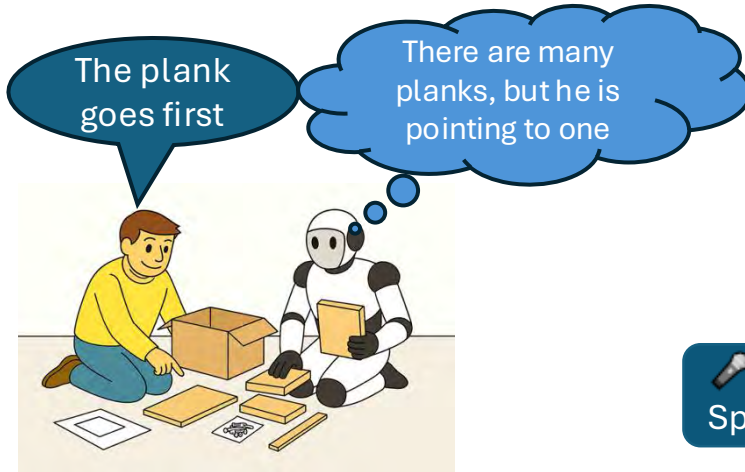
Event Perception and Segmentation

Event Segmentation Theory (EST) proposes that perceptual systems spontaneously segment an activity into events as a side effect of trying to anticipate upcoming information.

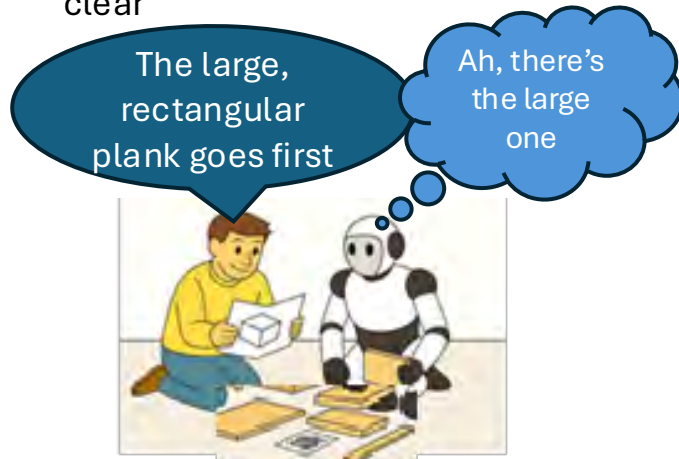
- Event boundaries are associated with changes in time, space, causes, characters and goals
- Actions are hierarchically organized by goals and subgoals
 - **Coarse-grained events** focus on objects, using more precise nouns and less precise verbs.
 - **Fine-grained events** focus on actions on those objects, using more precise verbs but specifying the objects less precisely.



Human State Sensing / Goal Inference

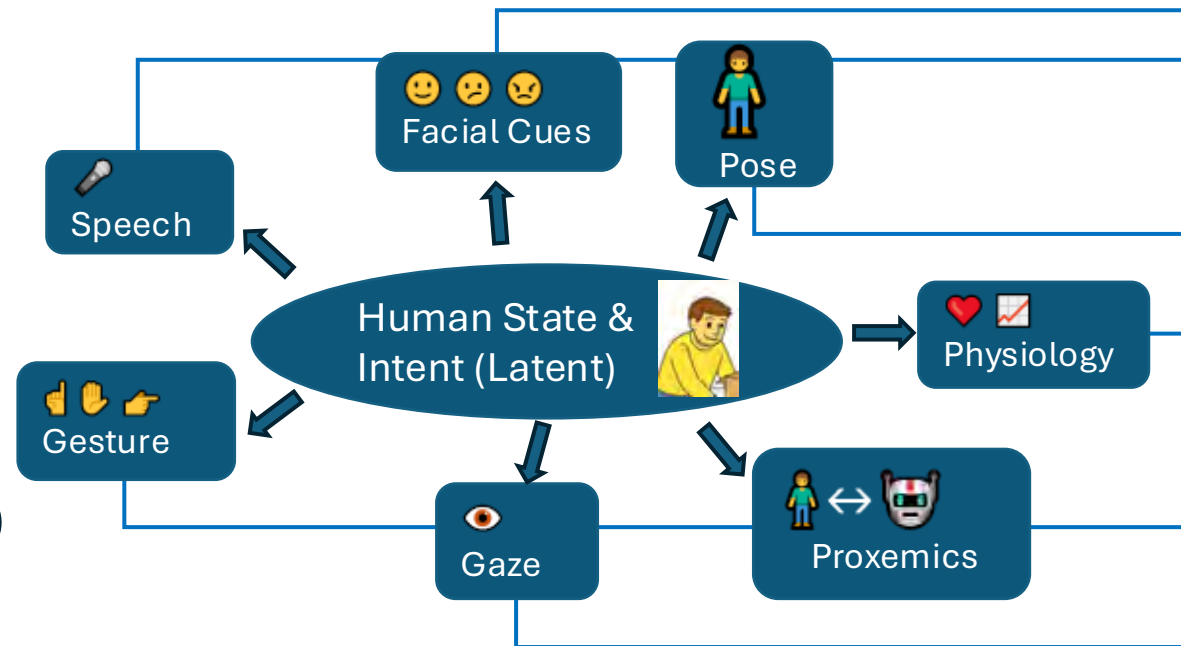


Ambiguous verbal instructions, but pointing/gaze makes it clear



Hands and eyes are busy reading but verbal instruction is specific

Human States / Intent expressed multimodally
Speech • Gesture • Gaze • Touch
Pose • Proxemics • Facial cues



Robot infers human state:

- Fatigue
- Cognitive load
- Anxiety
- Confusion
- Comfort
- Approval

Robot infers intent / goal:

- Object
- Action
- Task
- Situation

Monitoring human state (facial expressions, nodding etc.) can also help in predicting errors and misalignment.

➤ Enables flexible (non-task specific) error management framework

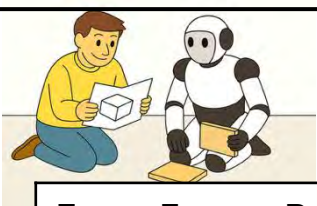
Stiber, M., Taylor, R. H., & Huang, C. M. (2023, March). On using social signals to enable flexible error-aware HRI. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot Interaction* (pp. 222-230).
 Bremers, A., Pabst, A., Parreira, M. T., & Ju, W. (2023). Using Social Cues to Recognize Task Failures for HRI: Overview, State-of-the-Art, and Future Directions. *arXiv preprint arXiv:2301.11972*.

Explicit World Model Representation

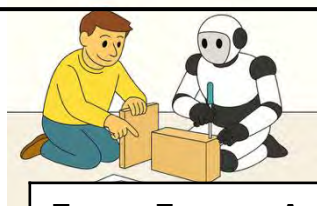
Event Frame: Assembling Furniture	
Time	Morning
Location	Living Room
Goal	Assemble Furniture
Characters	Human_1, Robot_1



Event Frame: Unpacking	
Time	{Time 1}
Space	{Configuration 1: characters, objects}
Sub-Goal	Identify and layout parts for assembly
States	{Pre: Parts in box; Post: Parts organised on the floor} <div><div>Human_1</div><div><div><div>Belief</div><div>PerceptionsStates</div></div><div><div>Desire</div><div>Goals</div></div><div><div>Intention</div><div>Plans</div></div></div></div>



Event Frame: Reading Instructions	
Time	{Time 2}
Space	{Configuration 2}
Sub-Goal	Understand assembly procedure
States	...



Event Frame: Assembling	
Time	{Time 3}
Space	{Configuration 3}
Sub-Goal	Assemble parts
States	...

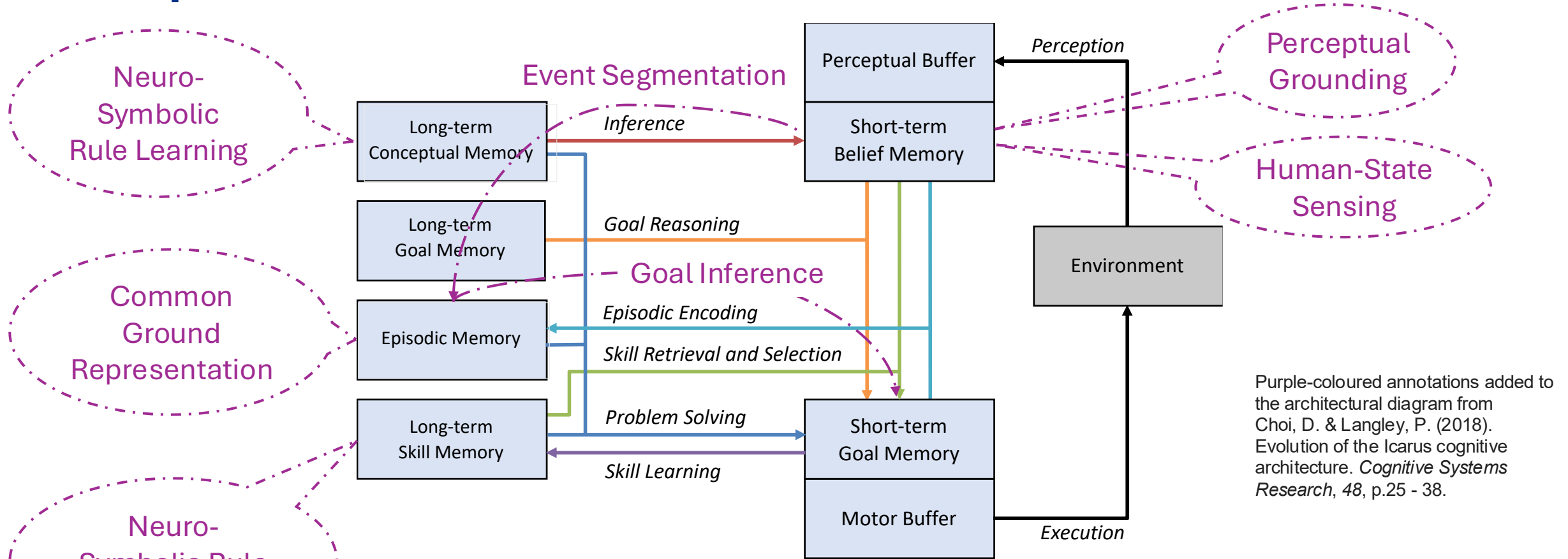
Event Frame: Step 1	
Time	{Time 3.1}
Space	{Configuration 3.1}
Sub-Goal	Assemble part 1
States	...

Event Frame: Step 2	
Time	{Time 3.2}
Space	{Configuration 3.2}
Sub-Goal	Assemble part 2
States	...

...

Event Frame: Step N	
Time	{Time 3.N}
Space	{Configuration 3.N}
Sub-Goal	Assemble parts together
States	...

Implications for Robotic Architectures



- With the help of neuro-symbolic rule learning of relations and procedures, robotic architectures get richer vocabulary of how to describe states and how to change certain states into another.
- These rules constitute an explicit world model that the robotic agents can use to infer abstract states, form a common ground with another agent, select relevant goals, and execute for the chosen goals.

Call To Action – Capability Gaps

- Robust Perceptual Grounding
 - Visual scene parsing
 - Situation recognition
 - Human state sensing/recognition
- Accessible Common Ground Representation
 - Expressive (capture rich social, multimodal, and fluid nature of HRI)
 - Efficient (lightweight enough for real-time updating)
- Implementation in Embodied Cognitive Architecture
 - Event perception and segmentation
 - Human goal inference (Theory of Mind)
 - Neuro-symbolic rule learning from continuous experience
 - Common ground influence on robot's decision making

Thank you!



AAAI 2026 Summer Symposium Series

Sponsored by the Association for the Advancement of Artificial Intelligence



June 22-June 24, 2026

Dongguk University | Seoul, South Korea

Architectures for Embodied Agents: A Synergy of Classic and Foundation Model Paradigms

Organised by:

Dongkyu Choi, Samsung Electronics

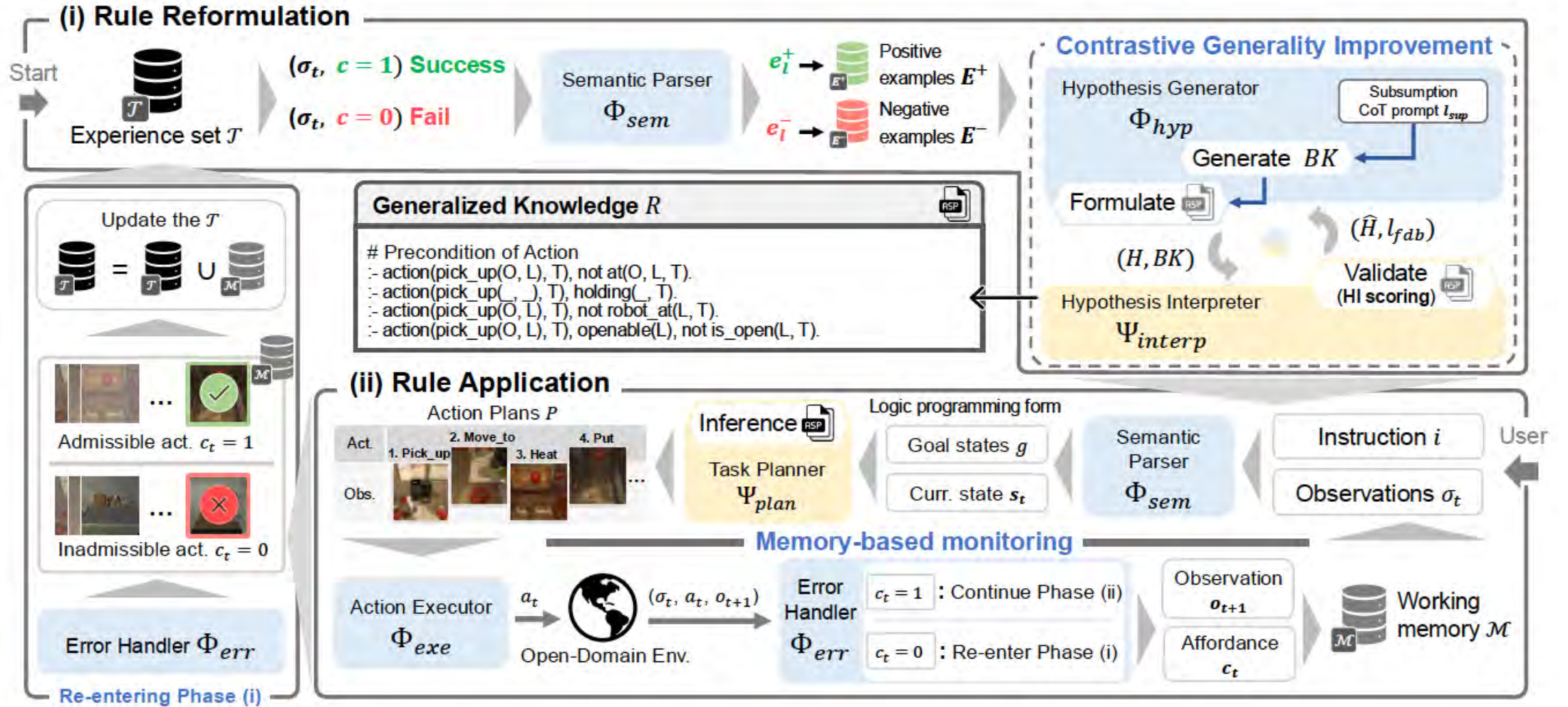
Pat Langley, Georgia Tech Research Institute

Jaeheung Park, Seoul National University

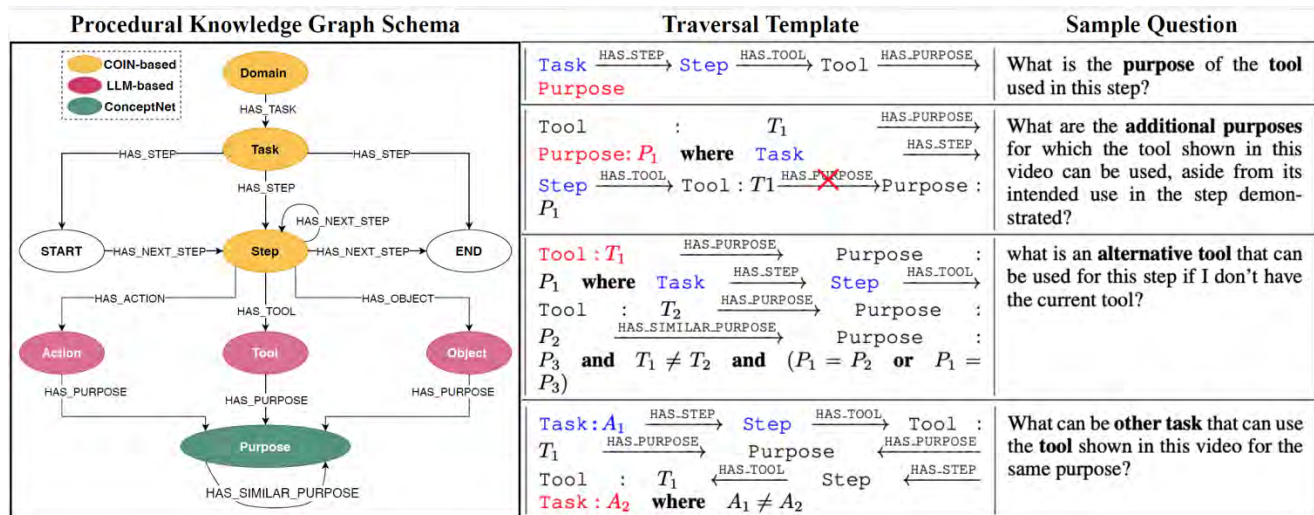
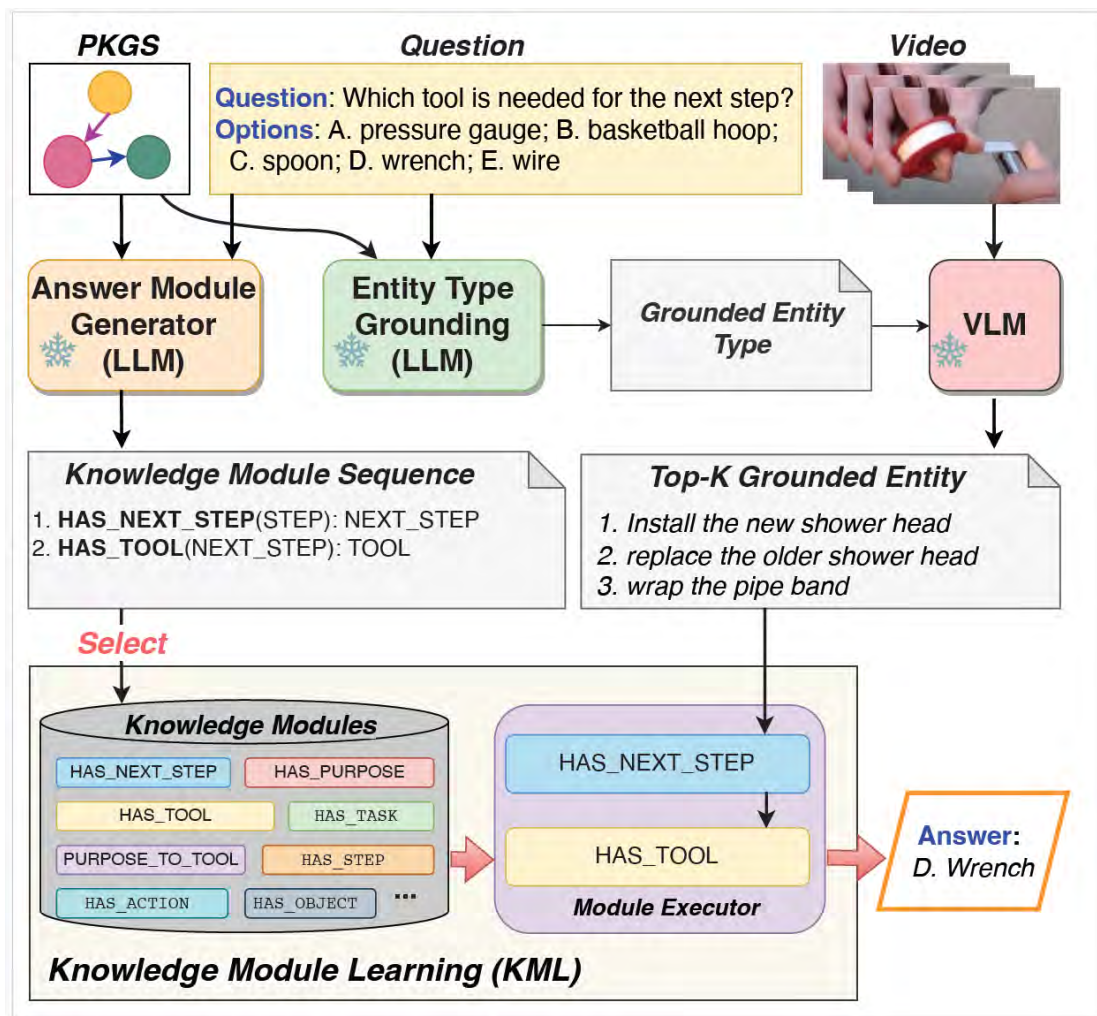
Kenneth Kwok, Agency for Science, Technology and Research

Sanjay Oruganti, Rensselaer Polytechnic Institute

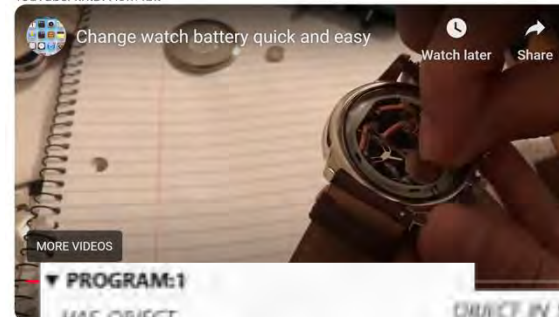
NESYC: A NEURO-SYMBOLIC CONTINUAL LEARNER FOR COMPLEX EMBODIED TASKS IN OPEN DOMAINS



Knowledge Module Learning



YouTube: krkDA49m4bk



Question

What is the other task that use the object in this video for the same purpose?

- ☐ Replace Battery On Key To Car
- ☐ Change Car Tire
- ☐ Clean Toilet
- ☐ Hang Wallpaper

MORE VIDEOS

PROGRAM:1

HAS_OBJECT

- back cover and waterproof ring: 0.6528
- back cover: 0.6800
- battery: 0.5318
- type of the back cover: 0.4789
- car key battery: 0.3872
- flat side of the new needle: 0.3847
- old memory chip: 0.3838
- screen connector: 0.3791
- fixed battery components and the back cover: 0.3624
- watch: 0.3388

DRIFT IN STEP

- install the back cover and waterproof ring: 0.4565
- open the back cover: 0.6129
- replace the battery: 0.5901
- put in the battery: 0.5347
- take out the battery: 0.5190
- install the new screen: 0.5128
- take out the car key battery: 0.4986
- unscrew the screws used to fix the screen: 0.4776
- check the type of the back cover: 0.4762
- put battery in: 0.4575

IN_TASK

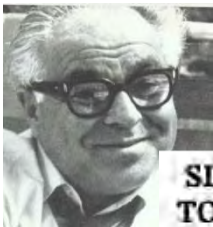
- ChangeBatteryOfWatch: 0.6564
- ReplaceLaptopScreen: 0.6241
- ReplaceBatteryOnKeyToCar: 0.4332
- ReplaceTyreValveStem: 0.4380
- ReplaceBatteryOnTVControl: 0.4298
- ReplaceFilterForAirPurifier: 0.4817
- ReplaceGraphicsCard: 0.3792
- ReplaceRefrigeratorWaterFilter: 0.3724
- RefillFountainPen: 0.3789
- InstallAirConditioner: 0.3658

Nguyen, T.-S.; Yang, H.; Neoh, T. Y.; Zhang, H.; Ee, Y. K.; and Fernando, B. 2025. "Neuro-Symbolic Knowledge Reasoning for Procedural Video Question Answering with Knowledge Module Learning (KML)." ArXiv, abs/2503.14957.



Event Model - Representation

- Scripts and Frames



SILVAN
TOMKIN

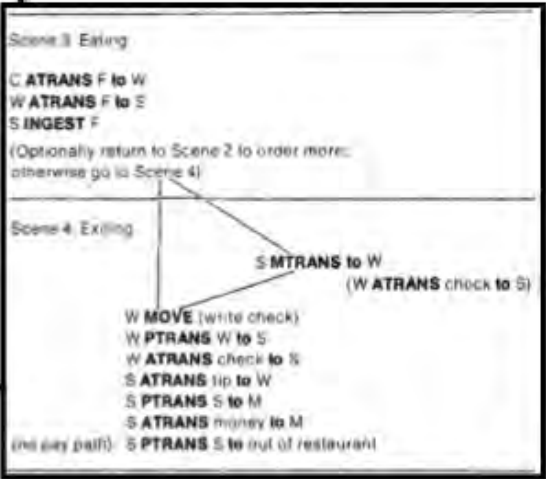
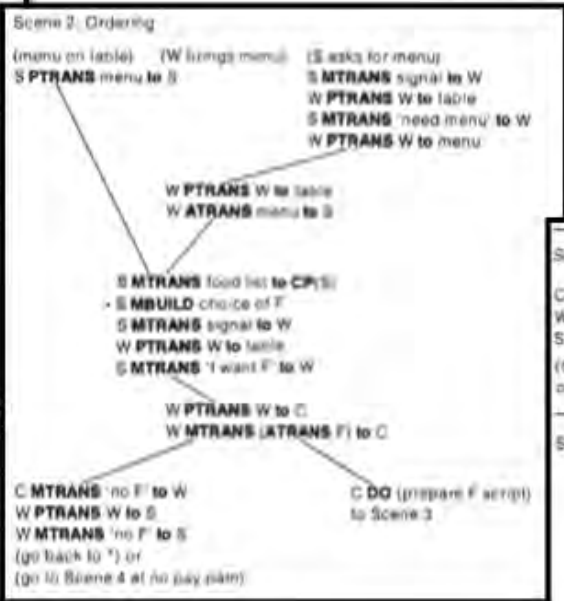


ROGER
SCHANK



ROBERT
ABELSON

Script: RESTAURANT	Rules: S-Customer
Track: Coffee Shop	W-Waiter
Props: Tables	C-Cook
Menu	M-Cashier
F-Food	O-Owner
Check	
Money	
Entry conditions: S is hungry.	Results: S has less money
S has money,	O has more money.
	S is not hungry
	S is pleased (optional)
Scene 1: Entering	
S PTRANS S into restaurant	
S ATTEND eyes to tables	
S MBUILD where to sit	
S PTRANS S to table	
S MOVE S to sitting position	



CHAIR frame	
A-kind-of:	furniture
number-of-legs:	an integer (default=4)
style-of-back:	straight, cushioned, ...
number-of-arms:	0,1,2

John's-chair frame	
a-kind-of:	furniture
number-of-legs:	4
style-of-back:	cushioned
number-of-arms:	0

Generic Restaurant Frame	
a-kind-of:	Business Establishment
Types:	Range: (Cafeteria, Seat-Yourself, Wait-to-be-seated, Fastfood) Default: IF plastic-orange-counter THEN fastfood IF stack-of-trays THEN cafeteria IF wait-for-waitress-sign OR reservation-made THEN wait-to-be-seated OTHERWISE seat_yourself
Location:	Range: an ADDRESS if-needed: (Look at the menu)
Name:	if needed: (Look at the menu)
Food-style:	Range: (Burgers, Chinese, American, Seafood, French) Default: Chinese if-added: (Update Alternative of Restaurant)
Time-of-Operation:	Range: a time-of-day Default: open evenings except Mondays
Payment form:	Range: (Cash, Credit Script)
Event-Sequence:	Default: Eat-at-Restau

Scripts Plans Goals and Understanding

An Inquiry into Human Knowledge Structures

Roger Schank
Robert Abelson

Psychology Press

A Framework for Representing Knowledge

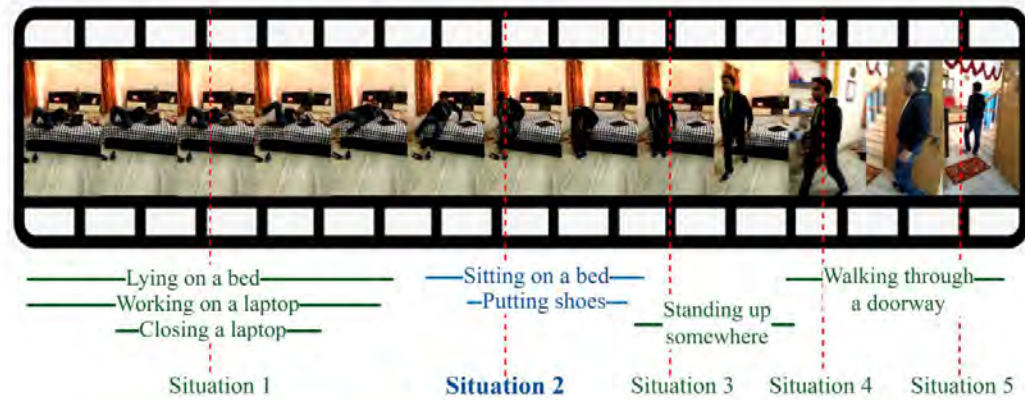
Marvin Minsky

MIT-AI Laboratory Memo 306, June, 1974.



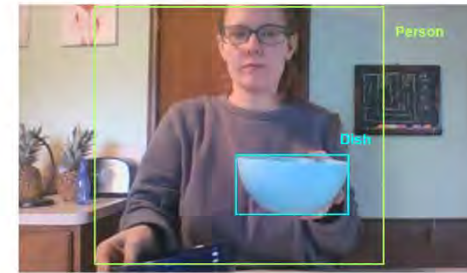
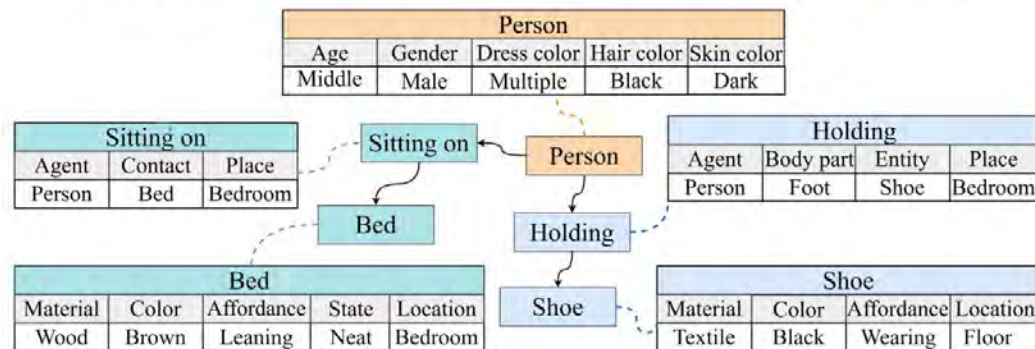
Event Model - Representation

- Situational Scene Graphs

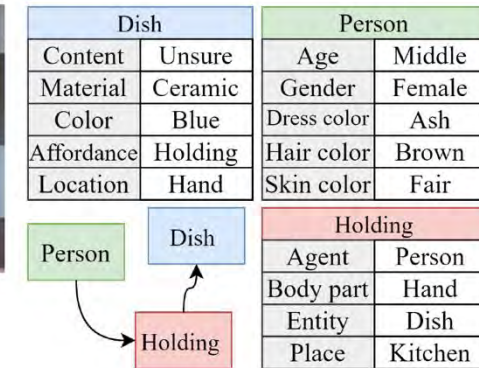


Action 1: Sitting on bed

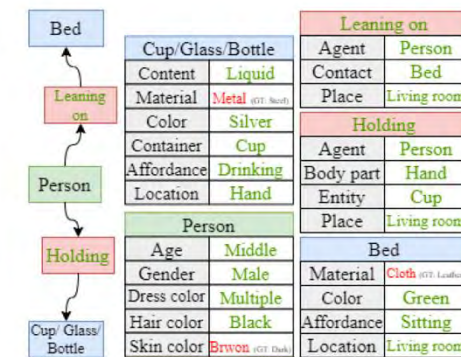
Action 2: Putting shoes



Action: Holding a dish



Actions: Lying on a bed, Holding a cup/glass/bottle.



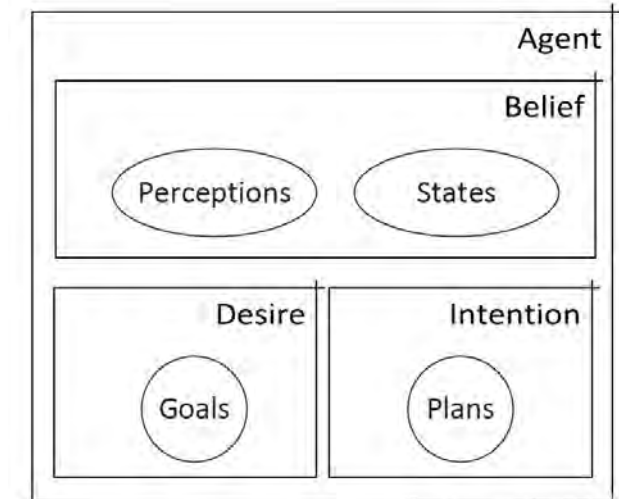
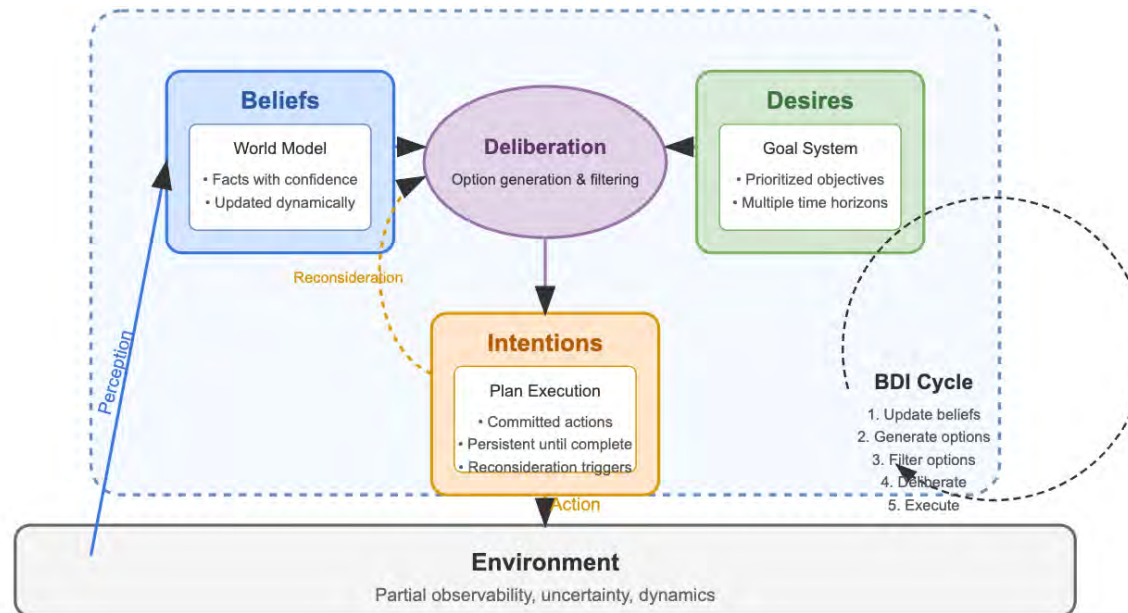
Event Model - Representation

Belief-Desire-Intention (BDI) Agent Model

- **Belief:** Agent's understanding of the world, including itself and other agents
- **Desire:** Agent's goals, preferences, and values
- **Intentions:** Agent's plans, strategies, and actions



Michael Bratman



Agent State expressed as Beliefs-Desires-Intentions (BDI)