

# Embodied Human-Computer Interactions through Situated Grounding

James Pustejovsky and Nikhil Krishnaswamy

IVA '20: ACM International Conference  
on Intelligent Virtual Agents

October 19–23, 2020

Glasgow, UK



# Situated Semantic Grounding and Embodiment

- Task-oriented dialogues are **embodied interactions** between agents, where language, gesture, gaze, and actions are situated within a common ground shared by all agents in the communication.
- Situated semantic grounding assumes **shared perception** of agents with **co-attention** over objects in a situated context, with **co-intention** towards a common goal.
- **VoxWorld** : a multimodal simulation framework for modeling **Embodied Human-Computer Interactions** and communication between agents engaged in a shared goal or task.
- **Embodied HCI** and robot control in action.

# Situated Meaning

Mother and son interacting in a shared task of icing cupcakes



## SITUATED MEANING IN A JOINT ACTIVITY

- SON: *Put it there (gesturing with co-attention)?*
- MOTHER: *Yes, go down for about two inches.*
- MOTHER: *OK, stop there. (co-attentional gaze)*
- SON: *Okay. (stops action)*
- MOTHER: *Now, start this one (pointing to another cupcake).*

# Situated Meaning

## Elements from the Common Ground

Agents	mother, son
Shared goals	baking, icing
Beliefs, desires, intentions	Mother knows how to ice, bake, etc. Mother is teaching son
Objects	Mother, son, cupcakes, plate, knives, pastry bag, icing, gloves
Shared perception	the objects on the table
Shared Space	kitchen

# Embodied Human-Computer Interaction

- Elements of Situated Meaning
  - Identifying the *actions and consequences* associated with objects in the environment.
  - Encoding a multimodal expression contextualized to the *dynamics of the discourse*
  - *Situated grounding*: Capturing how multimodal expressions are anchored, contextualized, and situated in context
- Modalities Deployed
  - gesture recognition and generation
  - language recognition and generation
  - affect, facial recognition, and gaze
  - action generation

# IVA in Embodied Environment

An encounter between two “people” with multimodal dialogue: language, gesture, gaze, action.



Figure: IVA Diana engaging in an embodied HCI with a human user.

▶ Link

# Affordance and Goal Recognition

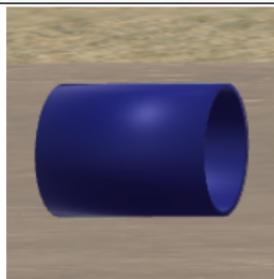
1. **Perceived purpose** is an integral component of how we interpret situations and reason about utterances in communicative contexts.
  - Events are purposeful and directed;
  - Places are functional;
  - Objects are usable and manipulable.
2. **Affordances** are latent action structures of how an agent interacts with objects in the environment, in different modalities:
  - language, gesture, vision, action;
3. **Qualia Structure** provides a link to such **latent actions structures** associated with objects in utterances and the context.

# Focus on Objects

- Context of objects is described by their properties.
- Object properties cannot be decoupled from the events they facilitate.
  - *Affordances* (Gibson, 1979)
  - *Qualia* (Pustejovsky, 1995)

“He **slid** the cup across the table. Liquid spilled out.”

“He **rolled** the cup across the table. Liquid spilled out.”



# Visual Object Concept Modeling Language (VoxML)

Pustejovsky and Krishnaswamy (2016)

- Encodes afforded behaviors for each object
  - **Gibsonian**: afforded by object structure (Gibson,1977,1979)
    - grasp, move, lift, etc.
  - **Telic**: goal-directed, purpose-driven (Pustejovsky, 1995, 2013)
    - drink from, read, etc.
- Voxeme
  - **Object Geometry**: Formal object characteristics in R3 space
  - **Habitat**: Conditioning environment affecting object **affordances** (behaviors attached due to object structure or purpose);
  - **Affordance Structure**:
    - What can one do to it
    - What can one do with it
    - What does it enable

# VoxML - cup

$$\left[ \begin{array}{l}
 \text{cup} \\
 \text{LEX} = \left[ \begin{array}{l} \text{PRED} = \text{cup} \\ \text{TYPE} = \text{physobj, artifact} \end{array} \right] \\
 \text{TYPE} = \left[ \begin{array}{l} \text{HEAD} = \text{cylindroid}[1] \\ \text{COMPONENTS} = \text{surface, interior} \\ \text{CONCAVITY} = \text{concave} \\ \text{ROTATSYM} = \{Y\} \\ \text{REFLECTSYM} = \{XY, YZ\} \end{array} \right] \\
 \text{HABITAT} = \left[ \begin{array}{l} \text{INTR} = [2] \left[ \begin{array}{l} \text{CONSTR} = \{Y > X, Y > Z\} \\ \text{UP} = \text{align}(Y, \mathcal{E}_Y) \\ \text{TOP} = \text{top}(+Y) \end{array} \right] \\ \text{EXTR} = [3] \left[ \text{UP} = \text{align}(Y, \mathcal{E}_{\perp Y}) \right] \end{array} \right] \\
 \text{AFFORD\_STR} = \left[ \begin{array}{l} \text{A}_1 = H_{[2]} \rightarrow [\text{put}(x, \text{on}([1]))] \text{support}([1], x) \\ \text{A}_2 = H_{[2]} \rightarrow [\text{put}(x, \text{in}([1]))] \text{contain}([1], x) \\ \text{A}_3 = H_{[2]} \rightarrow [\text{grasp}(x, [1])] \\ \text{A}_4 = H_{[3]} \rightarrow [\text{roll}(x, [1])] \end{array} \right] \\
 \text{EMBODIMENT} = \left[ \begin{array}{l} \text{SCALE} = \text{agent} \\ \text{MOVABLE} = \text{true} \end{array} \right]
 \end{array} \right]$$

## VoxML

## VoxML for Actions and Relations

$$\left[ \begin{array}{l}
 \mathbf{put} \\
 \text{LEX} = \left[ \begin{array}{l} \text{PRED} = \mathbf{put} \\ \text{TYPE} = \mathbf{transition\_event} \end{array} \right] \\
 \text{TYPE} = \left[ \begin{array}{l} \text{HEAD} = \mathbf{transition} \\ \text{ARGS} = \left[ \begin{array}{l} A_1 = \mathbf{x:agent} \\ A_2 = \mathbf{y:physobj} \\ A_3 = \mathbf{z:location} \end{array} \right] \\ \text{BODY} = \left[ \begin{array}{l} E_1 = \mathit{grasp}(x, y) \\ E_2 = \mathit{while}(\mathit{hold}(x, y), \mathit{move}(x, y)) \\ E_3 = \mathit{at}(y, z) \rightarrow \mathit{ungrasp}(x, y) \end{array} \right] \end{array} \right]
 \end{array} \right]$$
  

$$\left[ \begin{array}{l}
 \mathbf{on} \\
 \text{LEX} = \left[ \text{PRED} = \mathbf{on} \right] \\
 \text{TYPE} = \left[ \begin{array}{l} \text{CLASS} = \mathbf{config} \\ \text{VALUE} = \mathbf{EC} \\ \text{ARGS} = \left[ \begin{array}{l} A_1 = \mathbf{x:3D} \\ A_2 = \mathbf{y:3D} \end{array} \right] \\ \text{CONSTR} = \mathbf{y} \rightarrow \text{HABITAT} \rightarrow \text{INTR}[\mathit{align}] \end{array} \right]
 \end{array} \right]$$

## VoxML - grasp

$$\left[ \begin{array}{l} \mathbf{grasp} \\ \text{LEX} = \left[ \begin{array}{l} \text{PRED} = \mathbf{grasp} \\ \text{TYPE} = \mathbf{transition\_event} \end{array} \right] \\ \text{TYPE} = \left[ \begin{array}{l} \text{HEAD} = \mathbf{transition} \\ \text{ARGS} = \left[ \begin{array}{l} \text{A}_1 = \mathbf{x:agent} \\ \text{A}_2 = \mathbf{y:physobj} \end{array} \right] \\ \text{BODY} = \left[ \text{E}_1 = \mathit{grasp}(x, y) \right] \end{array} \right] \end{array} \right]$$

## VoxML - grasp cup

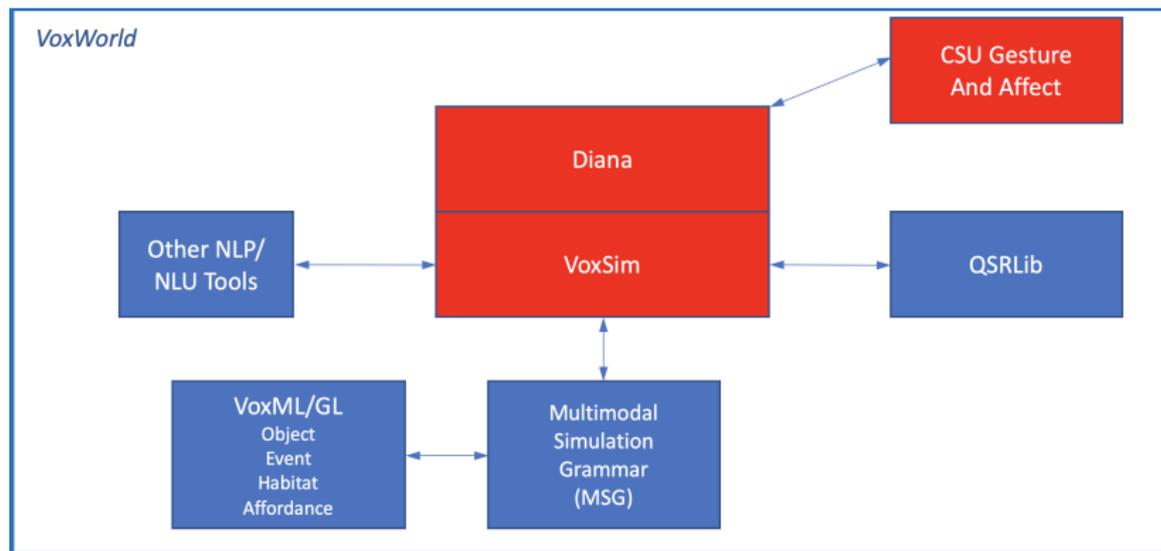
- Continuation-passing style semantics for composition
- Used within conventional sentence structures and between sentences in discourse in MSG

# Multimodal Simulations

- Human understanding depends on a wealth of **common-sense knowledge**; humans perform much reasoning **qualitatively**.
- To simulate events, every parameter must have a value
  - “Roll the ball.” How fast? In which direction?
  - “Roll the block.” Can this be done?
  - “Roll the cup.” Only possible in a certain orientation.
- VoxML: Formal semantic encoding of properties of objects, events, attributes, relations, functions.
- VoxSim: What can situated grounding do? (Krishnaswamy, 2017)
  - Exploit numerical information demanded by 3D visualization;
  - Perform qualitative reasoning about objects and events;
  - Capture semantic context often overlooked by unimodal language processing.

# VoxWorld: A Platform for Multimodal Simulations

## Interfacing Diana to CSU Gesture and Affect Systems



# Dynamic Discourse Interpretation

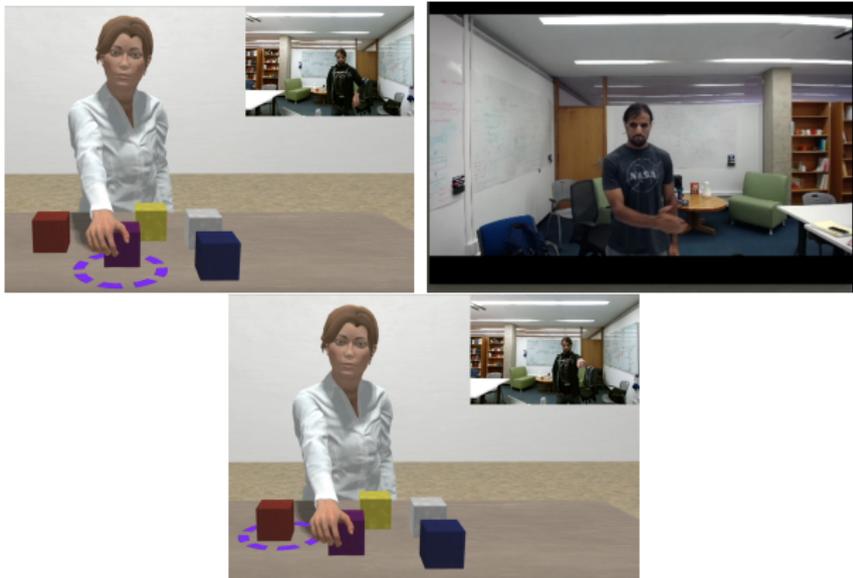
- Common Ground Structure
  - Co-belief
  - Co-perception
  - Co-situatedness
- Multimodal communication act:
  - language
  - gesture
  - action
- Dynamic tracking and updating of dialogue with:
  - Discourse Sequence Grammar
  - Gesture Grammar
  - Action Grammar

# Co-belief and Co-perception in the Common Ground

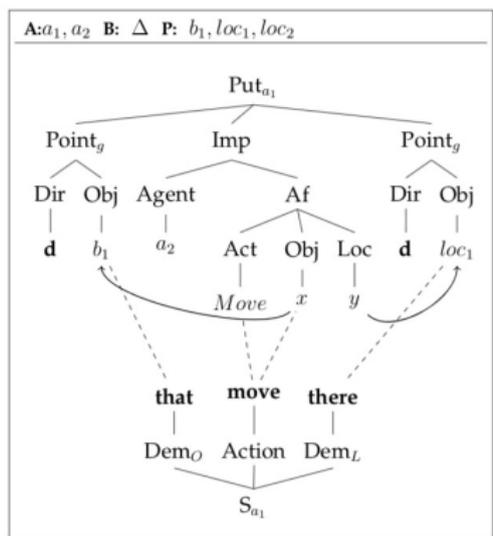
- *Public announcement logic (PAL)*
  - $[\alpha]\varphi$  denotes that an agent “ $\alpha$  knows  $\varphi$ ”.
  - Public Announcement:  $[\!\varphi_1]\varphi_2$
  - Any proposition,  $\varphi$ , in the common knowledge held by two agents,  $\alpha$  and  $\beta$ , is computed as:  $[(\alpha \cup \beta)^*]\varphi$ .
- *Public perception logic (PPL)*
  - $[\alpha]_\sigma\varphi$  denotes that agent “ $\alpha$  perceives that  $\varphi$ ”.
  - $[\alpha]_\sigma\hat{x}$  denotes that agent “ $\alpha$  perceives that there is an  $x$ .”
  - Public Display:  $[\!\varphi_1]_\sigma\varphi_2$
  - The co-perception by two agents,  $\alpha$  and  $\beta$  includes  $\varphi$  :  
 $[(\alpha \cup \beta)^*]_\sigma\varphi$

# Situated Meaning

Gesture and co-gestural speech imperative



$a_1$ : “That object  $b_1$  move  $b_1$  to there, location  $loc_1$ .”



$$\lambda k'_s \otimes k'_g. (\overline{\langle \text{that}, Point_1 \rangle \langle \text{move}, Move \rangle}) (\lambda r_s \otimes r_g. \langle \text{that}, Point_2 \rangle) \\ (\lambda k_s \otimes k_g. k'_s \otimes k'_g (k_s \otimes k_g r_s \otimes r_g))$$

# Transfer Learning of Object Affordances

- Gibsonian/Telic affordances are associated with abstract properties:
  - spheres **roll**, sphere-like entities probably do too;
  - small cups are **graspable**, small cylindroid-shaped objects probably are too.
- Similar objects have similar habitats/affordances:
- This informs the way you can talk about items in context:
  - **Q**: “What am I pointing at?”
  - **A**: “I don’t know, but it looks like {a ball/a container/etc.}”

# Transfer Learning of Object Affordances

Exploits the linkages between affordances and objects in VoxML

- Train over a sample of 17 different objects: blocks, KitchenWorld objects (apple, grape, banana, book, etc.)
- Trained 200 dimensional affordance and habitat embeddings using a Skip-Gram model, for 50,000 epochs with a window size of 3:
  - These embeddings serve as the inputs to the object prediction architectures
- Using the affordance embeddings in vector space, predict which object they belong to: using a 7-layer MLP; a 4-layer CNN with 1D convolutions

# Transfer Learning of Object Affordances

- The architectures:

MLP	CNN
Input	Input
<b>Dense (32 × tanh)</b>	<b>Conv1D (64 × ReLU)</b>
20% Dropout	ReLU
<b>Dense (196 × ReLU)</b>	20% Dropout
20% Dropout	<b>Conv1D (250 × ReLU)</b>
<b>Dense (92 × tanh)</b>	Global Max Pooling 1D
20% Dropout	20% Dropout
<b>Dense (196 × tanh)</b>	<b>Dense (196)</b>
<b>Dense (92 × ReLU)</b>	20% Dropout
<b>Dense (32 × tanh)</b>	ReLU
<b>Output (softmax)</b>	<b>Output (softmax)</b>
70,913 params	100,923 params

- Ground truth clusters generated by k-means clustering over human-annotated object similarity. Sample aggregate results:

Model	% predictions in correct cluster	% predictions always in correct cluster
MLP (Habitats)	78.82	27.06
MLP (Affordances)	<b>84.71</b>	38.82
CNN (Habitats)	78.82	27.06
CNN (Affordances)	81.18	<b>40.00</b>

- Object specific results (input: vectorized affordances for plate)

MLP (Habitats)	MLP (Affordances)	CNN (Habitats)	CNN (Affordances)
book, cup, bowl, bottle	cup, bottle, apple	book	cup, bottle

# Transfer Learning of Object Affordances



# Refactoring VoxWorld for Robot Navigation and Control

## Kirby's World

- Gesture and language communication with a Turtlebot-3:



- Fiducials represent registered proxies for object sorts in the environment:



Figure 2: Two fiducials.



## Conclusion - Embodied HCI

- VoxWorld facilitates experimentation with IVAs in embodied HCI contexts, using multiple modalities in diverse settings.
- An embodied HCI, such as that enabled by the simulation environment VoxWorld, provides a venue for the human and computer or robot to share an epistemic space,
- Any communicative modality that can be expressed within that space (e.g., linguistic, visual, gestural) enriches the ways in which a human and a computer or robot can communicate regarding objects, actions, and situation-based tasks.

# Thank You 😊

- **Brandeis LLC lab members:** Nikhil Krishnaswamy, Kyeongmin Rim, Mark Hutchens, Ken Lai, Katherine Krajovic, Daeja Showers, Eli Goldner, Kelley Lynch
- **CSU Vision lab members:** Ross Beveridge, Bruce Draper, Rahul Bangar, David White, Pradyumna Narayana, Dhruva Patil
- **University of Florida lab members:** Jaime Ruiz, Isaac Wang
- Funded by a grant from **DARPA** within the **CwC Program**

