



The VoxWorld Platform for Multimodal Embodied Agents

Nikhil Krishnaswamy, William Pickard, Brittany Cates, Nathaniel Blanchard, and James Pustejovsky • nkrishna@colostate.edu

VoxWorld

A platform for multimodal agent behaviors, presented as a resource to the AI/NLP community

What Makes an Agent?

- Perceives through sensors and acts through actuators
- Epistemic point of view from which it observes the world
- Virtual world is mode of presentation, allows observer to see what agent does
- Embodied agents add new dimensions to human/agent interactions
- Must recognize and interpret inputs in multiple modalities (e.g., gesture, speech, gaze action)
- Solving these problems has driven development of VoxWorld

Theory

- **VoxML** modeling language and **VoxSim** event simulator
- Events composed of subevent semantics that decompose into minimal primitive set
- Objects provide minimal encoding of properties, e.g., **habitats** and **affordances**
- Relations sample from distributions under constraints
- Event, relations, and objects composed at runtime

```

put
BODY = [
  E1 = grasp(x, y)
  E2 = while(hold(x, y) ^ ~at(y, z))
    -> move(x, y, z, PO, (loc(y), z, y))
  E3 = if(at(y, z) -> ungrasp(x, y)
]

in front
CLASS = config
VALUE = RCC8.EC
TYPE = [
  ARGS = [
    A1 = x:physobj
    A2 = y:physobj
  ]
  CONSTR = Z(x) > Z(y)
]

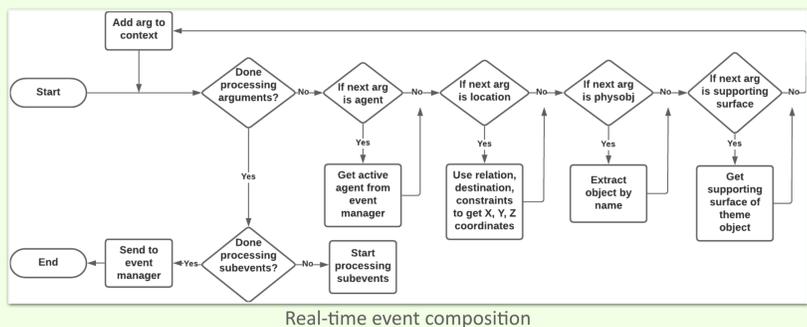
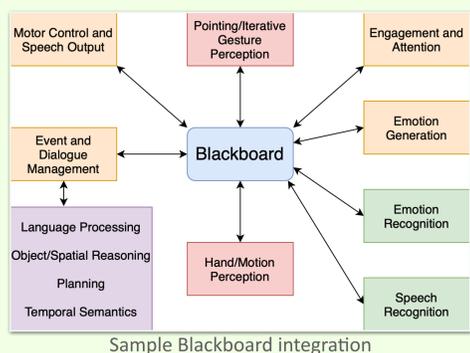
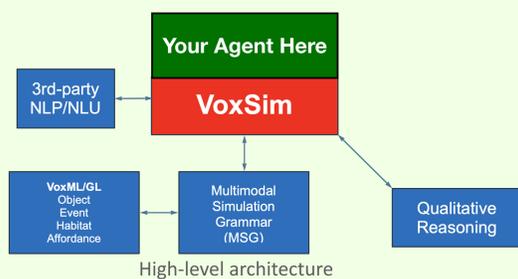
block
AFFORD_STR = [
  A1 = H2j -> [put(x, y, on([1]))]
  support([1], y)
  A2 = H2j -> [grasp(x, [1])]
  hold(x, [1])
  A3 = H2j -> [lift(x, [1])]
  hold(x, [1])
  A4 = H2j -> [ungrasp(x, [1])]
  release(x, [1])
]

```

Sample voxemes

Implementation

- Built on Unity game engine
- Accommodates qualitative calculi, machine learning inputs
- Interaction management via blackboard and pushdown automata
- Integrated with functional programming semantics
- Support arbitrary inputs and web deployment



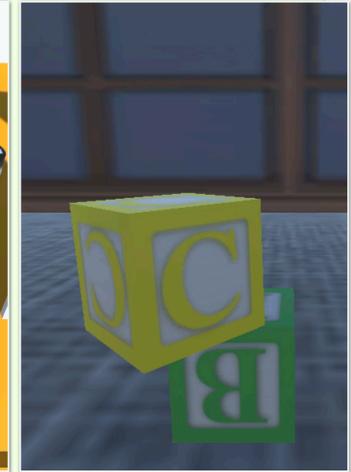
Agent Implementations



Diana: Multimodal Interaction



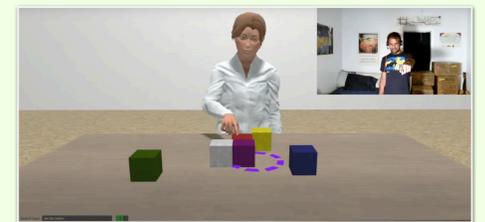
Kirby: Mobile Robotics



BabyBAW: Exploration with RL

Diana: interactive multimodal agent for peer-to-peer human-computer communication

- Diana interprets asynchronous speech and gesture
- Received mean 74.3 System Usability Score in User Study



Kirby: Navigating robot

- Same multimodal interface as Diana
- Integrates Robot Operating System, LIDAR, live camera feed
- Fiducial and object detection

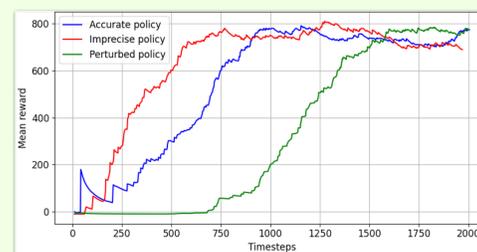


GoPiGo3 with LIDAR

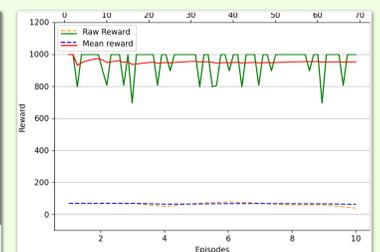


BabyBAW: Environmental exploration with RL

- Self-guided learner intended to approximate aspects of infant/toddler learning
- Integrates Unity ML-Agents, OpenAI Gym
- Ground actions and objects in world to learned labels



BabyBAW learns to stack!



Resources

- Asset package and bleeding-edge source code: <https://github.com/VoxML/VoxSim>
- Sample project: <https://github.com/VoxML/VoxWorld-QS>
- Documentation: <https://www.voxicon.net/api/>