

# Flexible Intelligence By Integrating Language and Perception

## Humans use language in extremely diverse ways



- Recognition
- Retrieval
- Generation
- Question answering
- Disambiguation
- Language acquisition
- Follow commands
- Paraphrasing
- Common sense reasoning
- Planning

$$P(\text{sentence}, \text{video})$$

$$\text{argmax}_{v \in V} P(s, v)$$

$$\text{argmax}_{s \in L} P(s, v)$$

$$\text{argmax}_{s \in L} P(Q(s, s_q), v)$$

$$\text{argmax}_{i \in \text{parser}(s)} P(i, v)$$

$$\text{argmax}_{\theta} \prod_{s, v} P(s(\theta), v)$$

$$\text{argmax}_p \int_{v^+} P(C(s), v^+ | v) E(v^+, p, v)$$

$$\int_v |P(s, v) - P(s', v)|$$

$$\text{argmax}_{s \in L} \int_v P(s_q, v) P(Q(s, s_q), v)$$

$$\text{argmax}_{s \in L} \int_v P(s, v_0 : v : v_n)$$

Caption, answer questions, understand a description, explain it to someone, engage in a conversation, give agents commands, imagine something different, recognize its description in a story, rewrite that description in another language, understand if someone is missing the point, reproduce it, intervene, etc.

## New tasks like visual disambiguation and zero-shot SNLI

Narayanaswamy et al. 2014 Generated 250 sentences with 2 to 6 interpretations and filmed each. Over 2000 clips.

Barret et al. 2016

Danny approached the chair with a yellow bag.

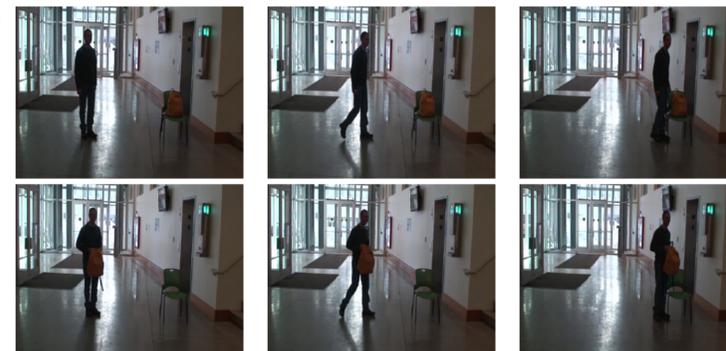
Narayanaswamy et al. 2014

Barbu et al. in prep.

Berzak et al. 2015

Yu et al. 2015

Paul et al. 2017



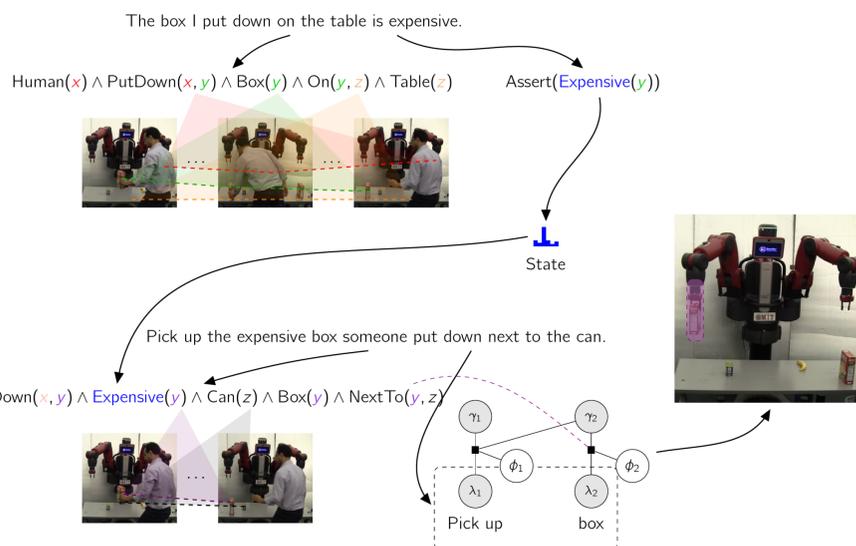
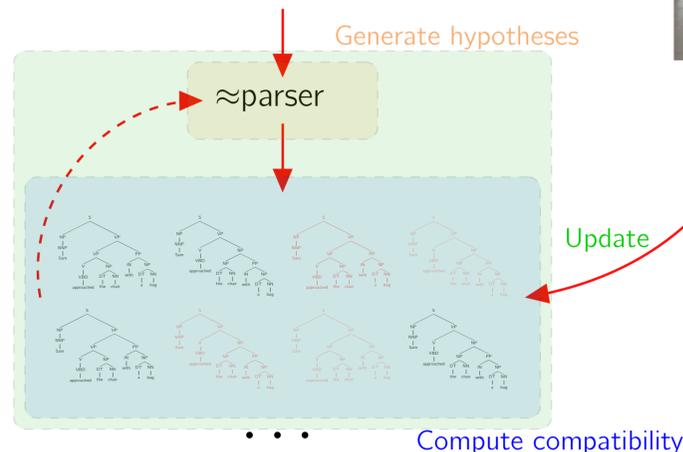
Comparing sentences vs Parikh et al. 2016 + ELMo

Does the sentence above imply the one below, and vice versa?

	Ground	Ours	Theirs
Alice carried the chair	↓ Y	Y	N
Alice held the chair	↑ N	N	Y
Alice carried the chair towards Ben	↓ Y	Y	N
Alice approached Ben	↑ N	N	Y?
Alice carried the chair towards Ben	↓ N	N	Y?
Alice left Ben	↑ N	N	Y?
Alice picked up the chair, and Ben put down the bag	↓ N	N	Y
Ben picked up the chair, and Alice put down the bag	↑ N	N	Y

## Models of language learning that are grounded in perception

Danny approached the chair with a bag.



2D and 3D velocities & distances are not invariant as the viewpoint changes.  
The contact, forces, and support are invariant for many actions.

