



Define, Evaluate, and Improve Task-Oriented Cognitive Capabilities for Instruction Generation Models



Lingjun Zhao



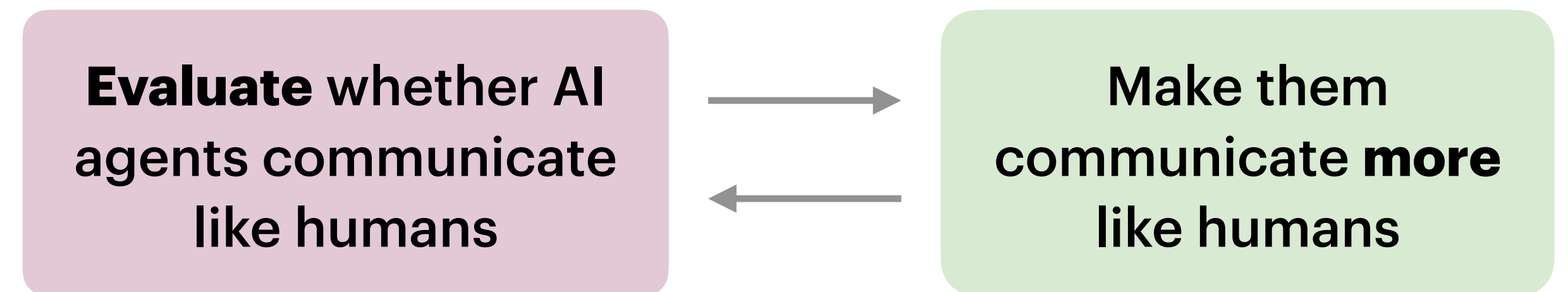
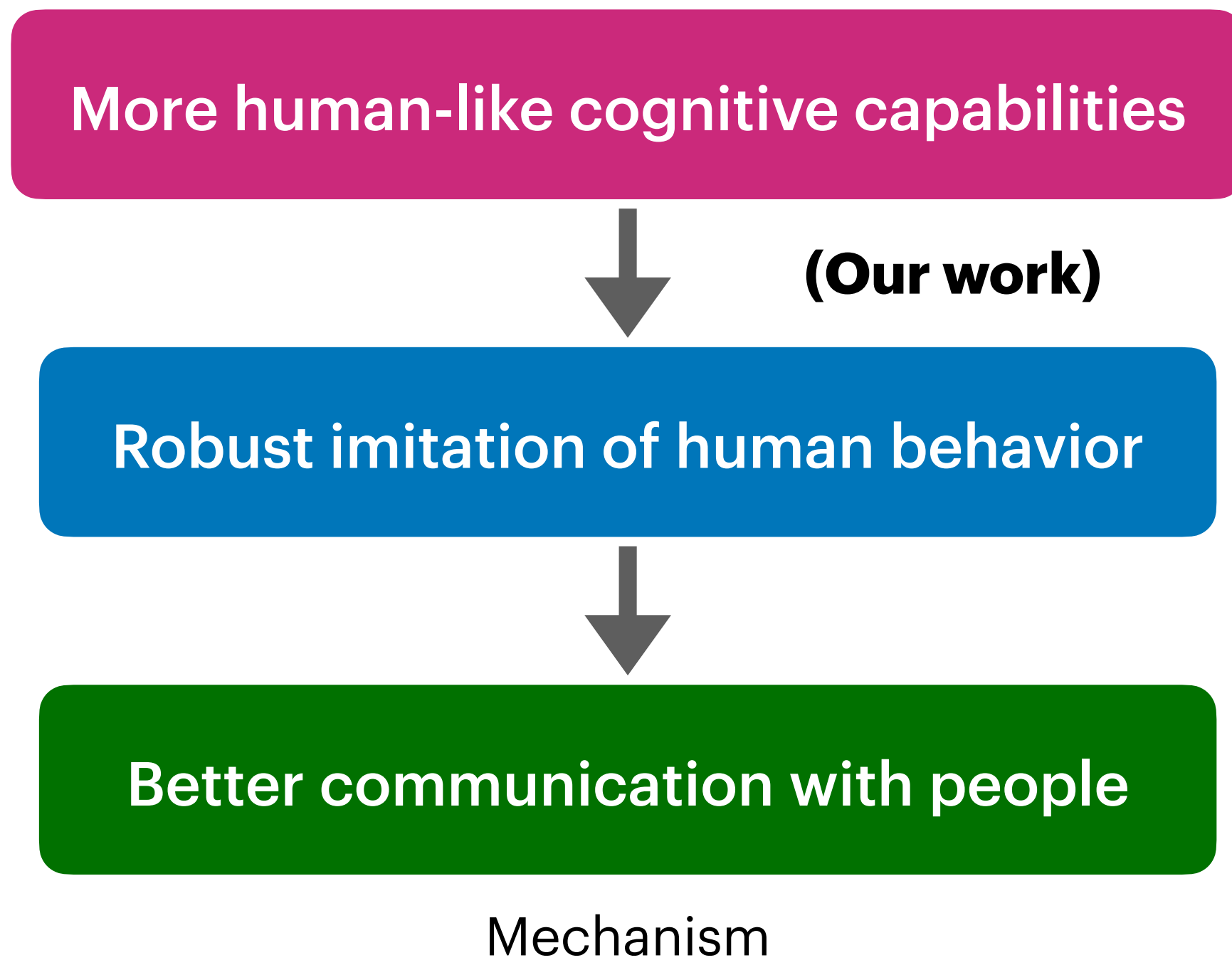
Khanh Nguyen



Hal Daumé III

Motivation: More human-like cognition leads to better communication

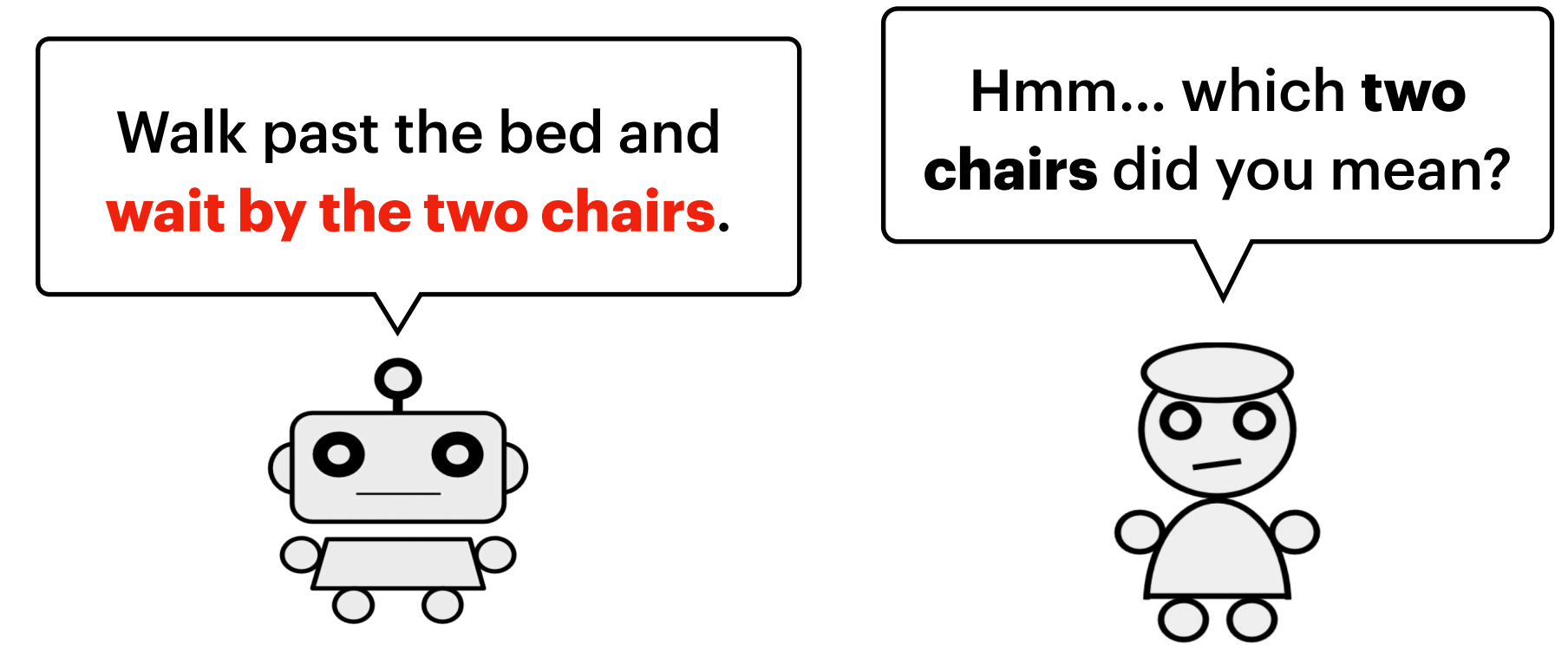
- By aligning AI agents with humans: how to perceive and describe the world



The evaluation step is difficult for black-box models

Problem: How to generate navigation instructions for people to follow

- Instructions generated by vanilla instruction generation (*speaker*) models fail to communicate well with humans
- How to generate better instructions by **reasoning pragmatically**?
- How to **evaluate cognitive capabilities** of speaker models?

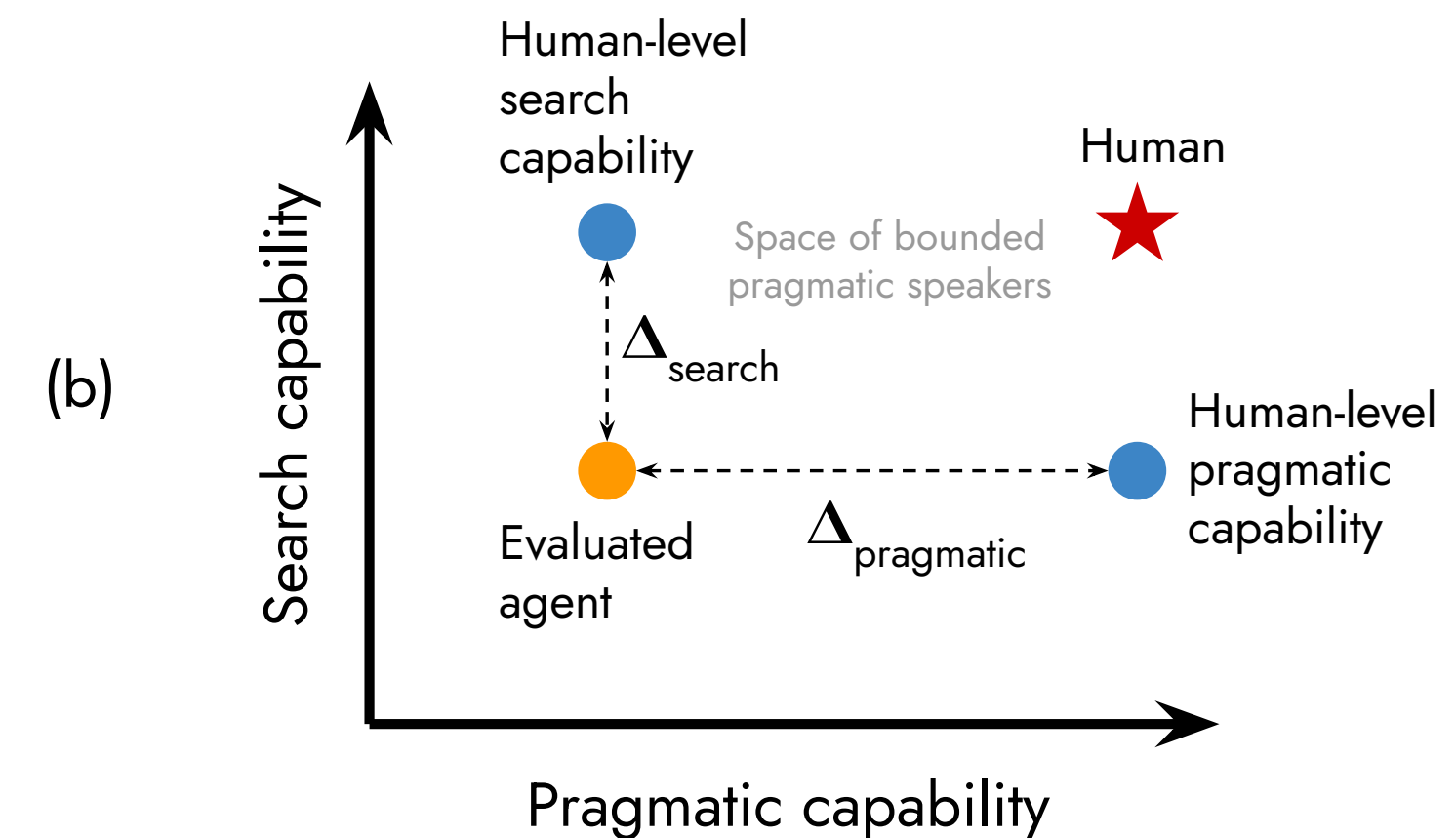
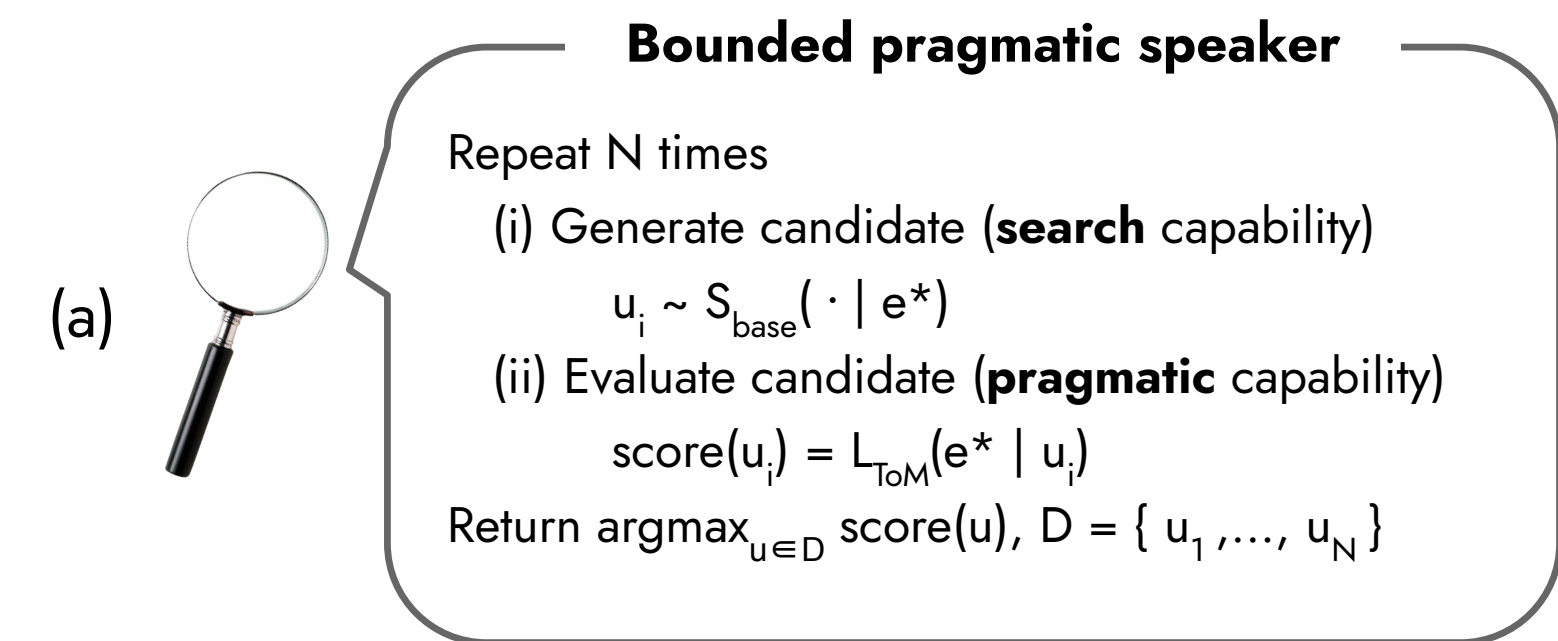


Contributions

- A new scheme for evaluating task-oriented cognitive capabilities in instruction generation models
- An 11% success rate improvement in guiding real humans in photorealistic environment, by equipping vanilla speakers with theory-of-mind capabilities
- A call to construct better theory-of-mind models for improving the instruction generation models

Distinguishing two capabilities: ToM and Search

- Humans are bounded pragmatic speakers (Sanborn and Chater 2016)
- Two **cognitive capabilities**:
 - * **Search**: evaluate whether can generate relevant instructions
 - * **Theory-of-Mind**: evaluate whether can simulate how human interprets the instructions

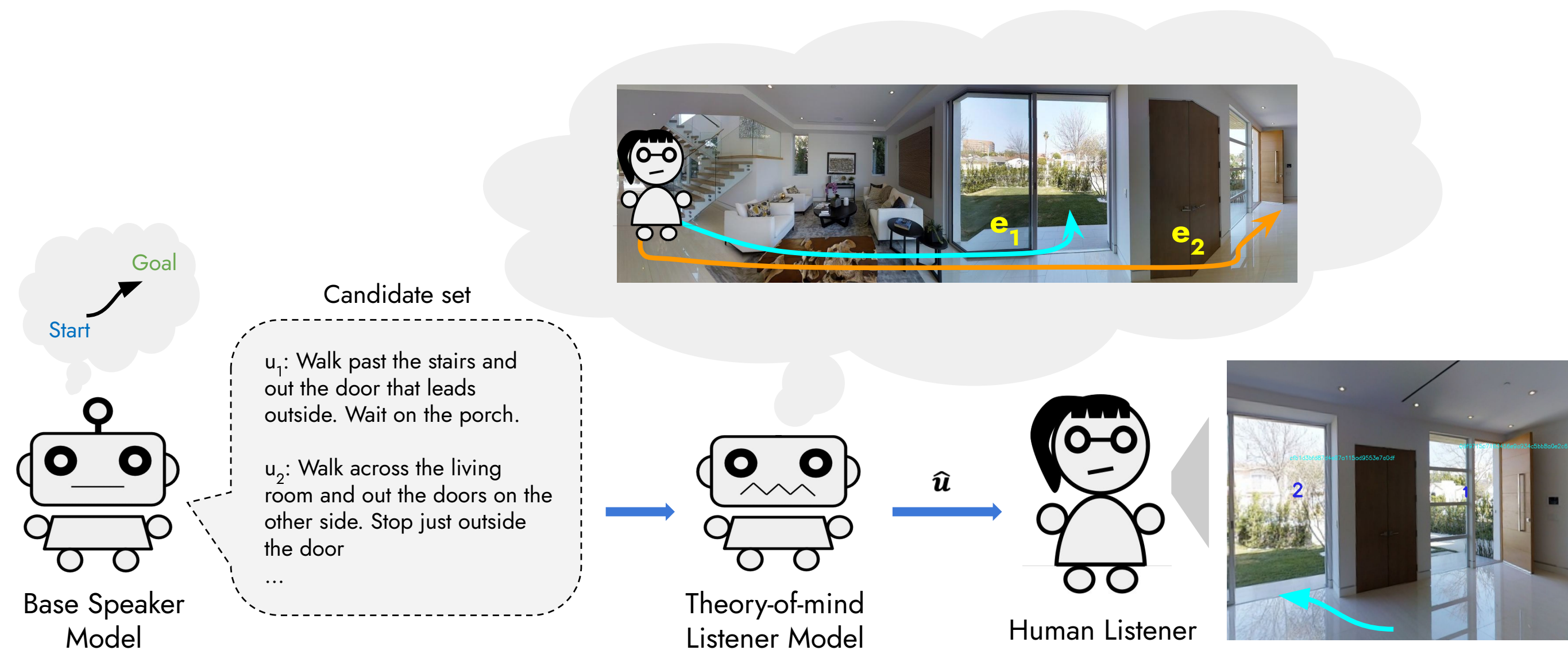


Recommendation:

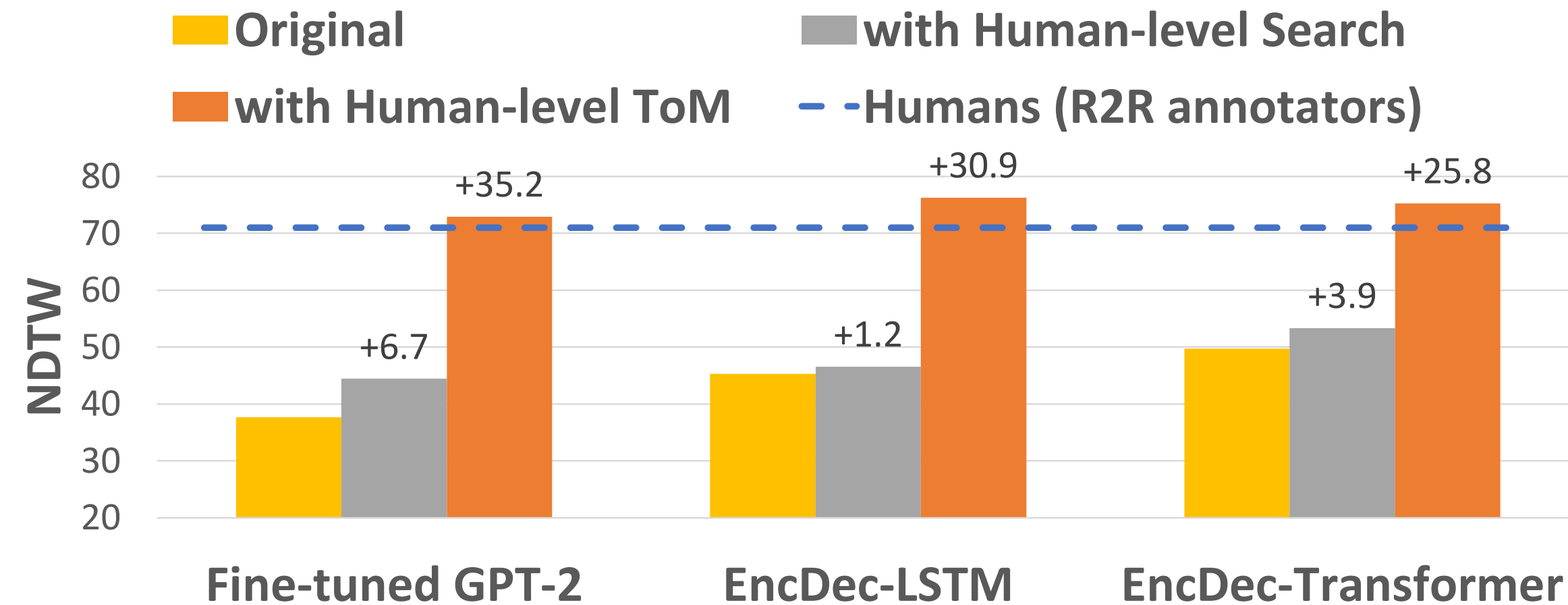
- (c)
- Large Δ_{search} , small $\Delta_{\text{pragmatic}}$ \Rightarrow improve inference algorithm
 - Large $\Delta_{\text{pragmatic}}$, small Δ_{search} \Rightarrow enhance model of listener

Bounded Pragmatic Speaker: Incorporate bounded Theory-of-mind into instruction generation

- **Base Speaker:** generates a set of relevant candidate instructions for a path
- **Theory-of-mind Listener:**
 - * RL agent(s) simulating how human would follow the instructions
 - * Select the instruction with simulated path most similar to the intended path
- **Human Listener:** follow the selected instruction in the environment



Pragmatic capability (theory-of-mind evaluation) is more deficient than Search capability (candidate generation)



Performance of the base speakers and their human-augmented versions

Experimental Settings

- **Speaker model dataset (reverse Matterport Room2Room dataset):**

Train:14k, Dev: 4k, Test: 1k

- **Evaluation: measure human's success in following generated instructions**

- * Give instructions to real humans
- * Measure similarity between human-generated and intended paths:

Normalized dynamic time warping (**NDTW** ↑)

- **Models:**

- * Finetuned GPT-2
- * EncDec-LSTM
- * EncDec-Transformer
- * Pragmatic Speakers

Using ensemble followers as theory-of-mind model can improve base speakers significantly to communicate with humans

ToM listener L_{ToM}	Base speaker S_{base}		
	Fine-tuned GPT-2	EncDec-LSTM	EncDec-Transformer
None	37.7 (▲ 0.0)	45.3 (▲ 0.0)	49.4 (▲ 0.0)
Single VLN-BERT (Majumdar et al., 2020)	38.9 (▲ 1.2)	39.8 (▼ 5.5)	46.2 (▼ 3.2)
Ensemble of 10 EnvDrop-CLIP (Shen et al., 2022)	37.8 (▲ 0.1)	53.1 [†] (▲ 7.8)	57.3 [†] (▲ 7.9)
Ensemble of 10 VLN \odot BERT (Hong et al., 2021)	43.4 (▲ 5.7)	56.4 [‡] (▲ 11.1)	54.2 (▲ 4.8)
Humans (skyline)	72.9 [‡] (▲ 35.2)	76.2 [‡] (▲ 30.9)	75.2 [‡] (▲ 25.8)

Performance of the speakers (NDTW) when equipped with different Theory-of-mind listener models

Shrink the gap with humans by 36%!

Takeaways

- Using ensemble followers as theory-of-mind model can **improve** base speakers trained with MLE objective
- Better task-oriented **theory-of-mind model** is needed to bridge the communication gap between AI and humans
- To develop safe and helpful AI requires quantifying the gaps between an AI agent and human