# Encoding Gesture in Multimodal Dialogue: Creating a Corpus of Multimodal AMR

Kenneth Lai<sup>1</sup>, Richard Brutti<sup>1</sup>, Lucia Donatelli<sup>2</sup>, James Pustejovsky<sup>1</sup>

<sup>1</sup>Brandeis University, <sup>2</sup>Vrije Universiteit Amsterdam

LREC-COLING 2024

May 20-25, 2024





# Contributions

- We create a corpus of speech and gesture meaning in a task-based setting
  - Temporal alignment of meaning representations with speech signal and gesture morphology
  - Semantic alignment: cross-modal coreference and anaphora
- Our corpus allows for structured analysis of multimodal semantic phenomena, enabling downstream applications for multimodal dialogue understanding
- Corpus and annotation guidelines available at https://github.com/klai12/encoding-gesture-multimodal-dialogue





# Outline

- Introduction and Background
- Gesture AMR and Multi-sentence AMR
- Data and Annotation
- Results
- Discussion and Conclusion





#### Gesture

- Idea: Humans naturally communicate using multiple modalities beyond language
- Gesture: how people move their hands or other body parts when they speak or communicate information







#### Gesture

- Gestures can be classified according to:
  - Relation with speech: co-speech, post-speech, pro-speech
  - Function within discourse: referential, non-referential, interactive
- In this project, we focus on co-speech, referential gestures
  - Other kinds of gestures are future work





# Abstract Meaning Representation (AMR)

AMR is a graph-based meaning representation that expresses the meaning of a sentence in terms of its predicate-argument structure

- Relatively easy to annotate
- Readable by both humans and machines
- Existing community of researchers





# Abstract Meaning Representation (AMR)

- Push that block left.
- (p / push-01 :mode imperative :ARG0 (y / you) :ARG1 (b / block :mod (t / that)) :ARG2 (l / left))







#### **Gesture AMR**

Gesture AMR can represent the dialogue participants, and objects and actions being referred to, while abstracting away from physical descriptions

- (g / [gesture]-GA
  - :ARG0 [gesturer]
  - :ARG1 [content]
  - :ARG2 [addressee])

- Root node: gesture act type
  - ARGO: initiator of the gesture
  - $\circ$   $\hfill ARG1:$  semantic content of the gesture
  - ARG2: recipient of the gesture





#### **Gesture AMR**

(i / icon-GA
 :ARG0 (s / signaler)
 :ARG1 (p / push-01
 :direction (l / left))
 :ARG2 (a / actor))







#### **Gesture AMR**

- Basic gesture act relations
  - Iconic: describe concrete properties of objects or actions
  - Deictic: point to objects or locations
  - Emblematic: meanings set by convention
  - Metaphoric: describe abstract properties of concepts or ideas
    - No examples in our corpus
- Complex gestures
  - Gesture units: single gestures with multiple meaning components
  - Coordinated gestures: separate gestures that are simultaneous but independent





# Multi-sentence AMR

Multi-sentence AMR (MS-AMR) extends AMR with coreference, implicit roles, and bridging relations across sentences

Push that block left. 

#### (p / push-01)

- :mode imperative
- :ARG0 (y / you)
- :ARG1 (b / block

:ARG2 (1 / left))

- Gesture for "push left"
- (i / icon-GA
  - :ARG0 (s / signaler)
  - :ARG1 (p / push-01
    - :ARG0 (i / implicit-role: pusher)
- :mod (t / that)) :ARG1 (i2 / implicit-role: thing pushed)

:ARG2 (a / actor))



# **EGGNOG Corpus**

We annotate on top of the EGGNOG (Elicited Giant Gallery of Naturally Occurring Gestures) corpus

- Shared task of arranging block structures
  - Signaler gives instructions
  - Actor interprets and builds structures





## **EGGNOG Corpus**







## **EGGNOG Corpus**

- Time-stamped annotations of participant gestures, for both the gesturer's inferred intent (such as "here"), as well a physical description of the movement (such as "arms: apart, left; hands: claw, down;")
- We created speech transcripts for the videos using the Coqui speech-to-text toolkit and manually corrected the output





#### **Annotation Process**

- We dually annotated 21 videos of task-based dialogues with speech and gesture AMRs
- We calculated inter-annotator agreement (IAA) using Smatch and S<sup>2</sup>match

	The second secon					
File Edit Annotation Tier Type Search View	Options Window Help					
		Lexicon Comments Recognizers Metadata Controls				
		Volume:				
		100				
(P+0)	0 50 100					
		20160128 195448 00 Video.avi				
		Mute Solo 0 35 50 75 100				
		20160128 195441 00 Video avi				
		Mute Solo 0 15 50 75 100				
and the second sec		20160128 195448 00 Video way				
		10100110_1050+00_00_0000.wav				
00:00:41.735	Selection: 00:00:00.000 - 00:00:00.000 0					
H H 14 E4 H D DF DE D1 D DH	▶S 🖋 ⊨ ← → ↓ ↑ Selection Mode Loop Mode 🐗					
	ar a cha ann an an ann ann an ann an ann an an					
20160128 9 0/42.000 00:00:43.000 00:00:44.000	00:00:45.000 00:00:46.000 00:00:47.000 00:00:48.000 00:00:49.000	00.00.50.000 00.00.51.000 00.00.52.000 00.00.53.000 00.00.54.000 001				
<b>*</b>	·····					
		**************************************				
0.42.000 00.00.43.000 00.00.44.000	00:00:45.000 00:00:46.000 00:00:47.000 00:00:48.000 00:00:49.000	00:00:50.000 00:00:51.000 00:00:52.000 00:00:53.000 00:00:54.000 00:				
042.000 00:00:43.000 00:00:44.000 Speech down do	00:00-45:000 00:00-45:000 00:00-47:000 00:00-45:000 00:00-49:000	0000-50.000 0000-51.000 00:00-52.000 00:00-53.000 00:00-54.000 00+				
	00.00.45.000 00.00.45.000 00.00.47.000 00.00.45.000 00.00.40.000 00.00.45.000 000 00.000 00.000 000 00.000 0000	00:05:50:00 00:00:51:00 00:00:52:000 00:00:53:000 00:00:54:000 00: 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1				
	00:00:45:000 00:00:45:000 00:00:45:000 00:00:46:000 00:00:48:000 00:00:49:000 00:00:48:000 00:00:49:000 00:00:40:000 00:00:00:000 00:00:00:000 00:00:00	00055000 0005100 00052000 000553000 00054000 007				
0 42.000         0 00.01 + 3.000         0 00.01 + 4.000           Speech Ake         0 1 (see facts - 0.01 mole         nm           Speech Ake         0 (see facts - 0.01 mole         nm           Casture - Laby         Arrows, front, FRA, more, down;         ams. move           Casture - Laby         Arrows, front, FRA, more, down;         ams. move	00.00.45.000 00.00.47.000 00.00.45.000 00.00.45.000 00.00.00.00 00.00.00 00.00	00:00:50:00 00:00:51:00 00:00:52:00 00:00:51:00 00:00:54:00 00: 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1				





#### **Annotation Process**

- We then marked coreference and bridging relations between actions and objects across modalities using MS-AMR
- We calculated IAA using CoNLL-2012 F1 score







Step 1: 343 speech AMRs and 319 gesture AMRs

Gesture Act	Top-level	All
lcon	142	249
Deixis	63	129
Emblem	37	39
Gesture unit	50	-
Coordinated gestures	27	-





Step 1: Most common semantic contents of gestures

Ico	on	De	ixis	Emb	olem
put-01	40	location	99	yes	19
block	31	block	24	ok	6
slide-01	17	left	2	no	5





Step 1: Inter-annotator agreement (micro-averaged F1) scores

	Smatch	S <sup>2</sup> match
Speech	0.489	0.648
Gesture	0.575	0.715





Step 2: 436 relations: 388 coreference chains, 28 set-member relations, 20 part-whole relations

• ~5 mentions per coreference chain, covering ~50% of entities/events in the corpus

Inter-annotator agreement (CoNLL-2012 F1 score): 60.46





# Discussion

We observe a wide variety of communication styles

- Some participants spoke much more than they gestured; others, vice versa
- One participant's gestures were exclusively iconic, while another did not produce any purely iconic gestures at all
- ~1/3 of participant movements consist of beat or other non-referring gestures





# Discussion

Annotation is hard!

• Annotators were required to segment speech transcripts into utterances themselves, and could separate or combine gestures as they saw fit

Future work:

- Separate utterance/gesture segmentation from AMR annotation
- Automatic AMR parsing followed by manual correction





# Discussion

Our annotation scheme abstracts away from many spatiotemporal properties of gesture

- Iconic gestures vary in their speed, manner, spatial coordinates, etc.
- Deictic gestures vary in the perspective the speaker assumes, etc.

How detailed morphological differences in gesture map onto semantic meaning is often challenging to determine, so we choose to leave this for future work

• We can use roles such as :manner and :mode to more faithfully represent gesture morphology





# Conclusion

- We present an annotated corpus of gesture and speech AMR in a task-based setting, that enables analysis of how the semantics of gesture and natural language interact
- Although our current annotation scheme is rather coarse-grained, AMR proves to be a flexible representation scheme adaptable to this domain





# Acknowledgements

We would like to thank Dean Cahill, Gaby Dinh, Hayden McCormick, Ryan Partlan, Shiyi Shen, Christopher Tam, Alicia Tu, and Tali Tukachinsky for their assistance with this research.

This research is supported by the NSF National AI Institute for Student-AI Teaming (iSAT), under grant DRL 2019805.





# Thank you!





