

AI-based analysis of action and attention of child-parent interaction

Iba Baig ^{1,2}, Yanbin Xu ², Kevin Li ³, Seiji Cattelain ^{3,4}, Hayato Ono ⁴, Sho Tsuji ^{3,4}, Mingbo Cai ^{2,4}.

¹ Northeastern University, ² University of Miami, ³ Ecole normale supérieure, ⁴ University of Tokyo

Background

- Parent-child interactions--toy handling, spatial proximity, pose, gaze, language--are key in development.
- Manual annotation of observational videos is labor-intensive, time-consuming, and difficult to scale. Automation aids behavioral study (Weng et al. 2025).
- Vision Language Models (VLMs) offer opportunity to automate high-resolution video understanding and behavioral coding.

Research Questions

- Can a multimodal Video-LM reliably annotate parent-child behavior at a fine-temporal scale?
- What behavioral patterns, transitions, and temporal structures can be extracted automatically?

Methods

Variables:

Controlled vocab & prompts



Our Architecture:

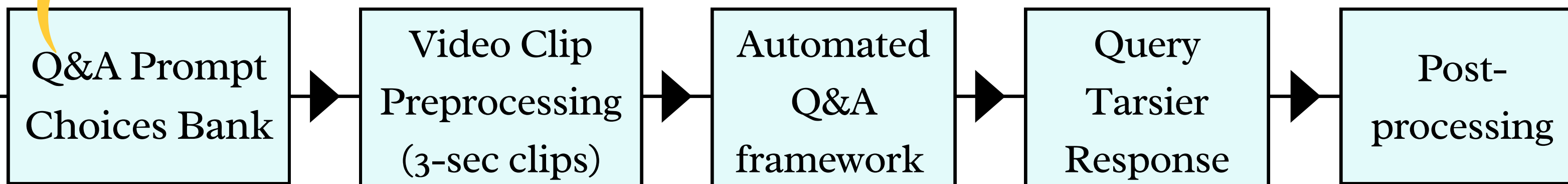
- Automated multi-turn conversational protocol to transform Tarsier outputs into a behavioral annotator

Tool: *Open-Source Pre-trained Video-Language Model Tarsier2-7B*

*Tarsier Model weights frozen

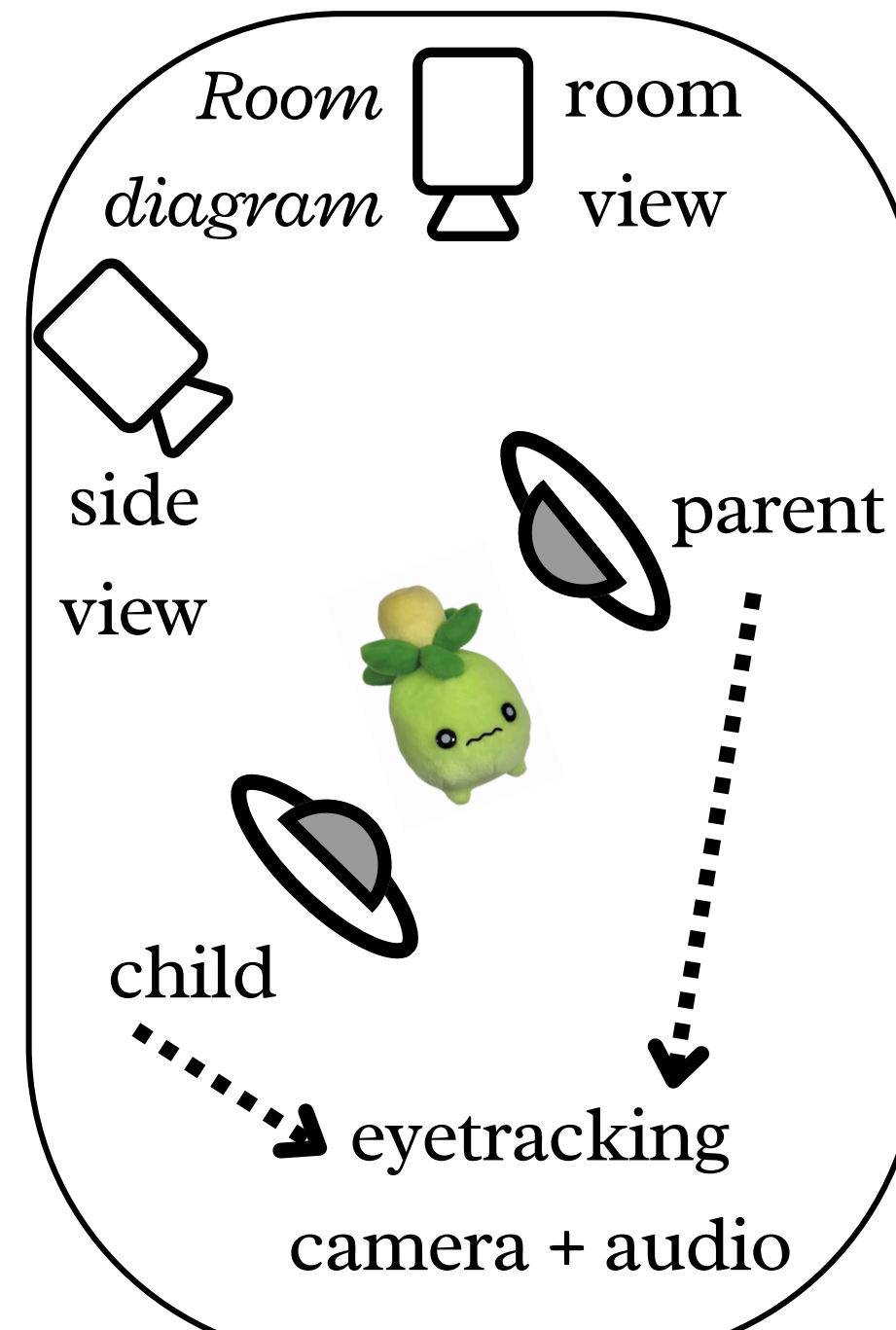


ByteDance (Yuan et al. 2025)

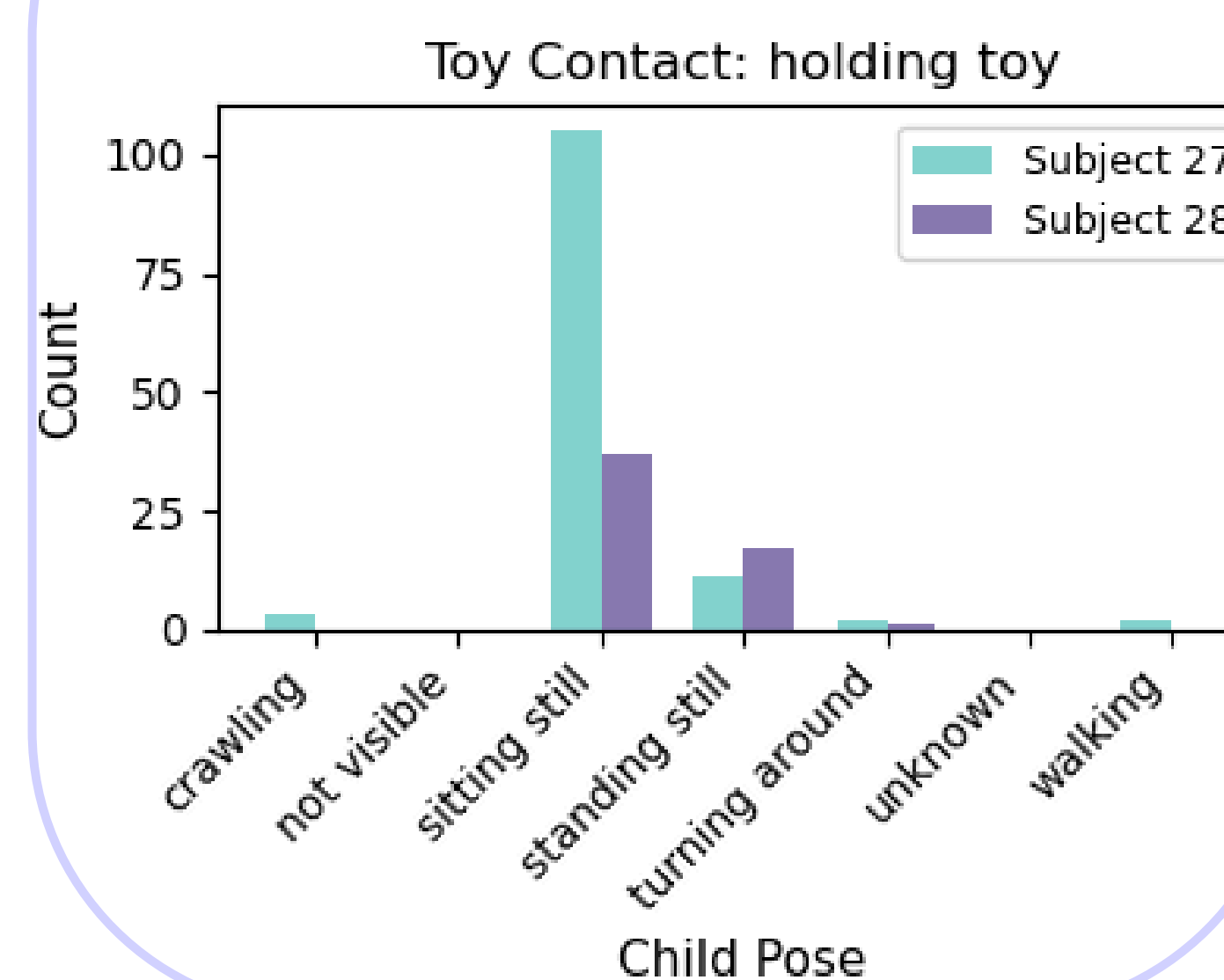


Dataset

- Parent-child pairs recorded with four-camera videos
- Sample: 31 caregiver-child pairs (mean age=3.4 yrs)
- Parents instructed to teach child toy's name, hobby, food preference, and personality per toy

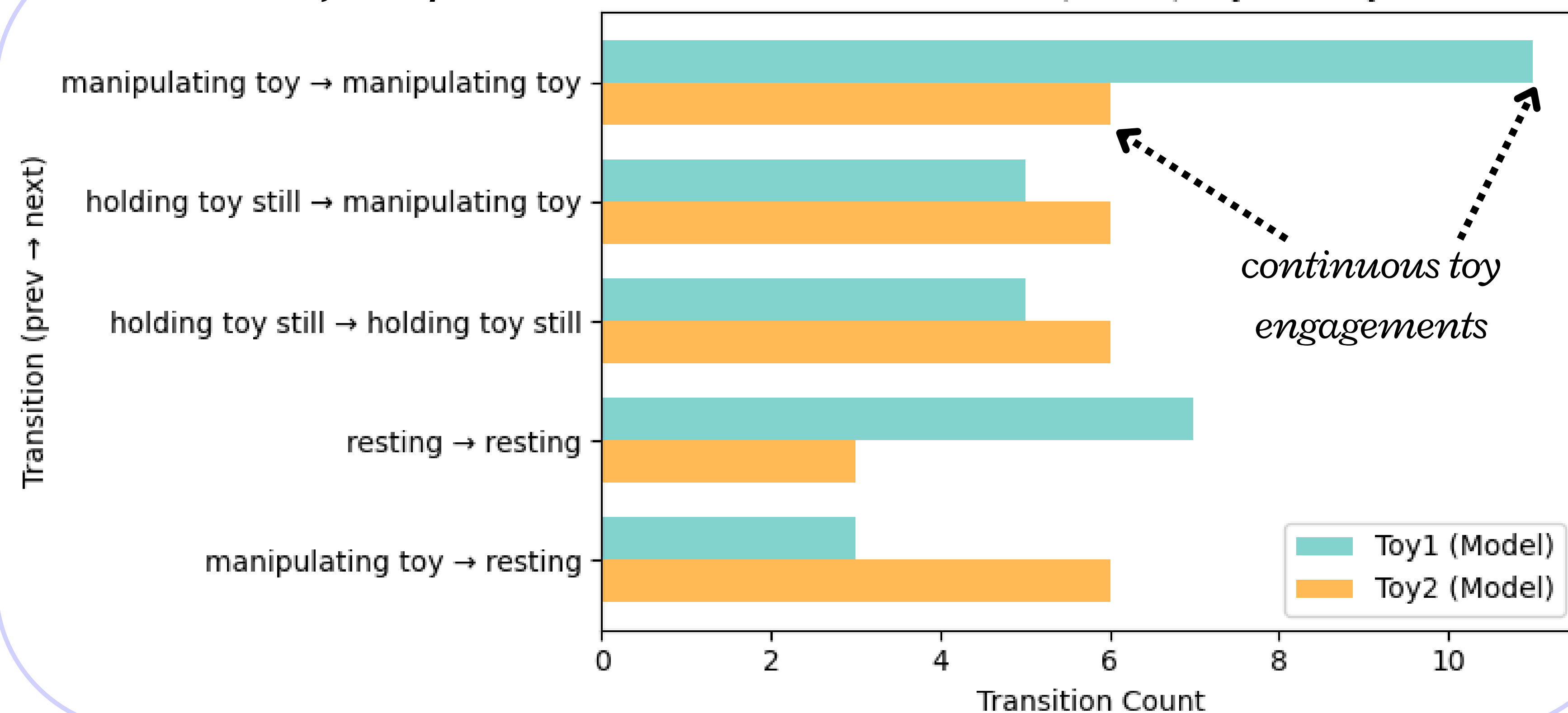


Subject 27 v.s. 28 Toy Contact + Movement



Subject 27

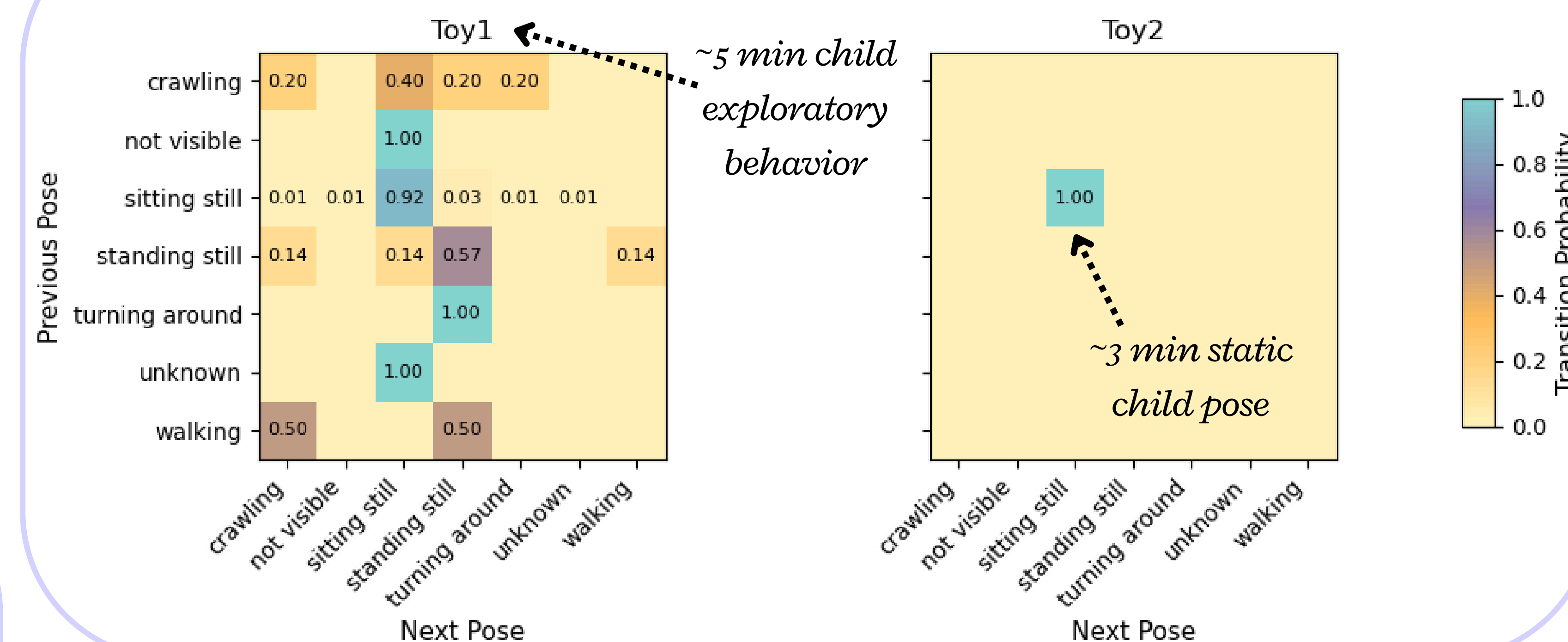
Child Hand Action Transitions (Model) Toy1 vs Toy2



Stimuli

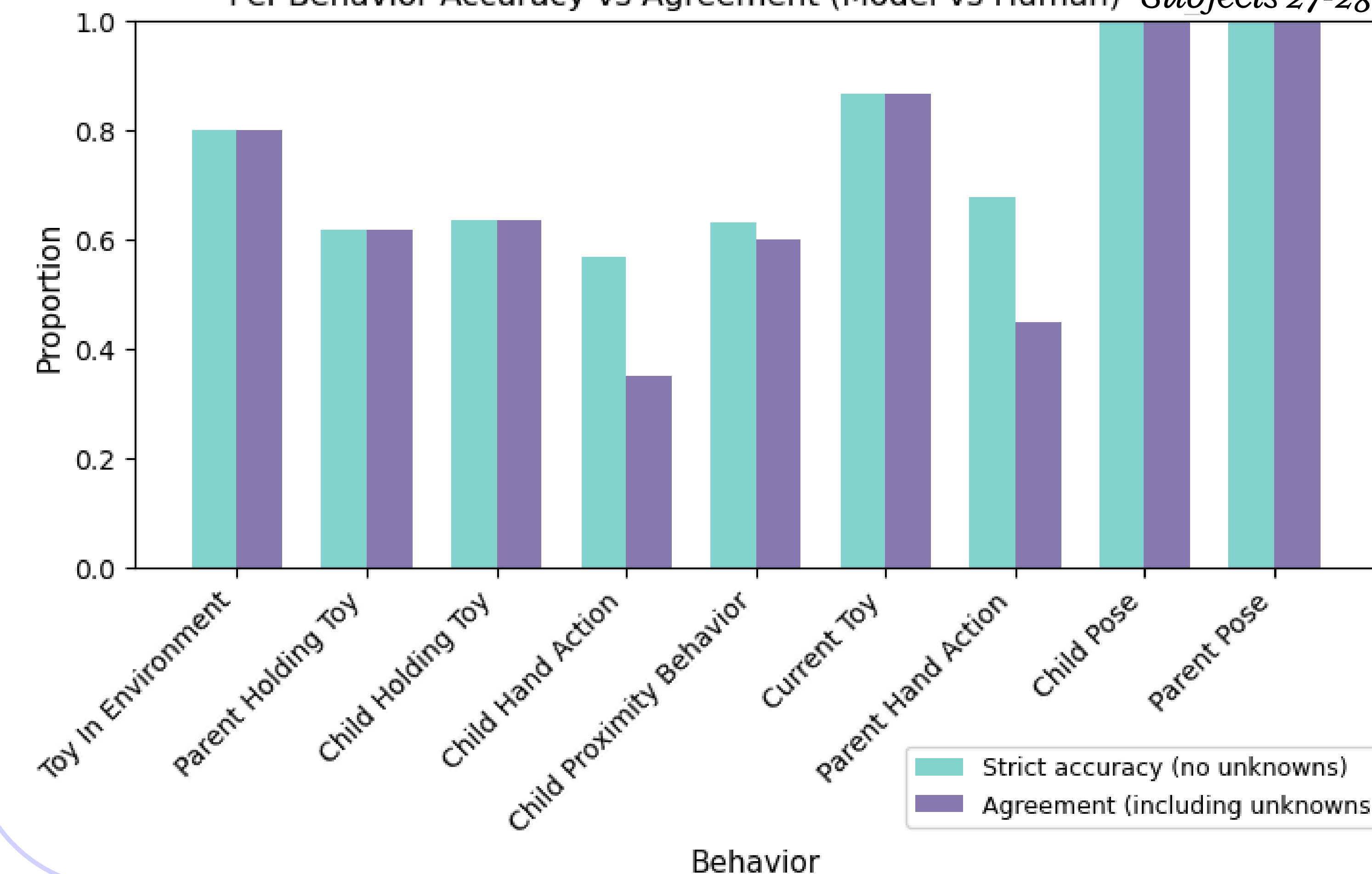


Child Pose Transition Probabilities (Model) Toy1 vs Toy2 [Subject 27]



Results

Per-Behavior Accuracy vs Agreement (Model vs Human) Subjects 27-28



Conclusion

- Multimodal video-LLMs can automatically annotate parent-child interaction videos with reasonable agreement to human annotations
- Best performance for pose estimation
- Lower accuracy in hand actions/spatial proximity
- Pre-trained video-LM performance

Future Direction

- Open source annotation tool!
- Integrate multimodal analysis with collaborators (side/top + egocentric + eye-tracking + audio-visual)
- Predictive modeling for individualized dynamics

SCAN ME

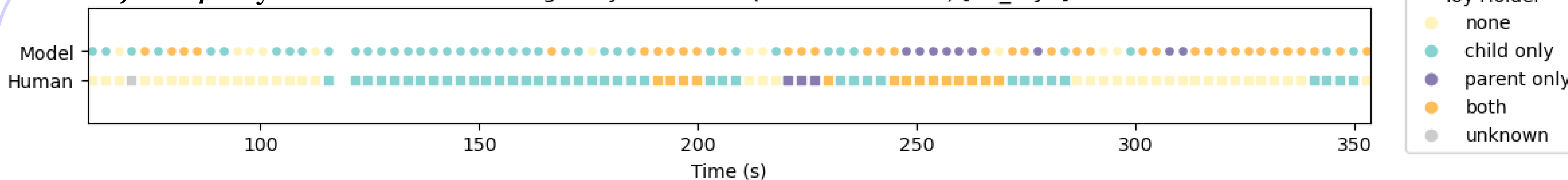


References

- Weng et al. (2025). Artificial intelligence-powered 3D analysis of video-based caregiver-child interactions. *Sci Adv*
- Yuan et al. (2025). Tarsier 2: Advancing Large Vision-Language Models from Detailed Video Description to Comprehensive Video Understanding. *arXiv:2501.07888*

Subject 27 Toy 1 & 2

Who Is Holding a Toy Over Time (Model vs Human) [27_toy1]



Who Is Holding a Toy Over Time (Model vs Human) [27_toy2]

