



Socially Intelligent Digital Humans

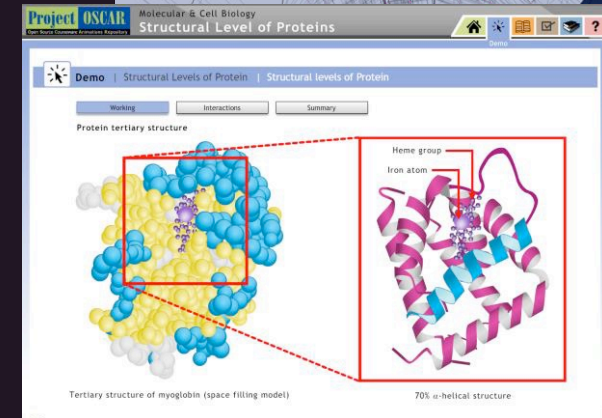
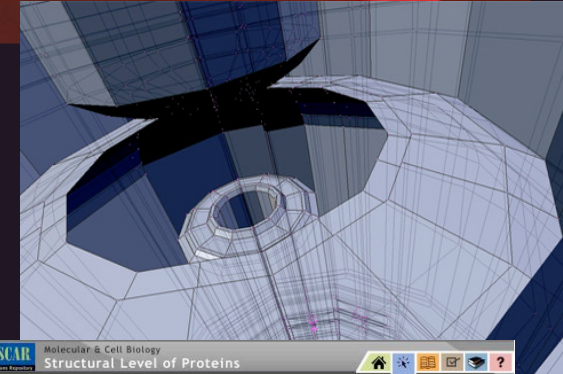
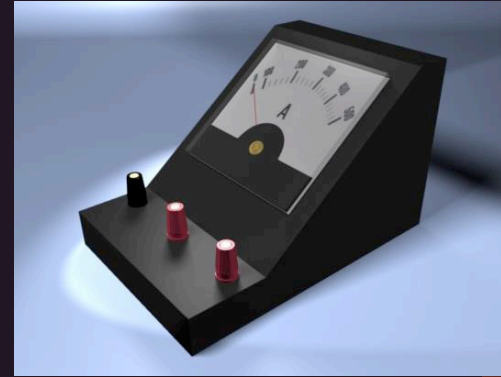
Evaluating and Generating
Multiparty Behavior

Chirag Raman
Asst. Professor, Tapri Lab

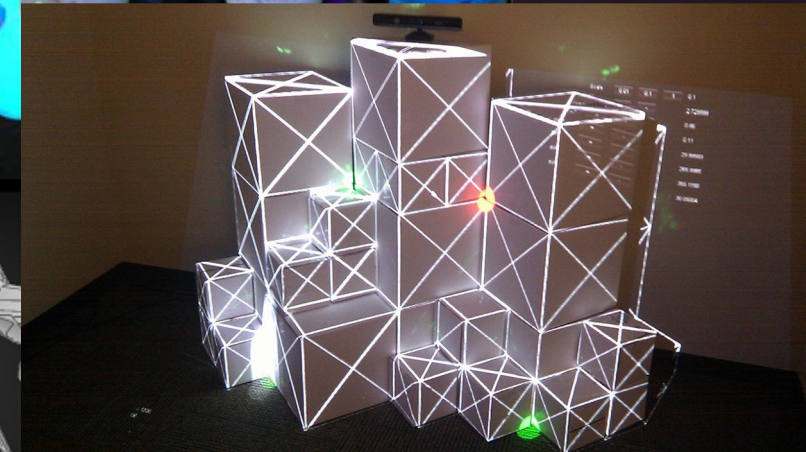
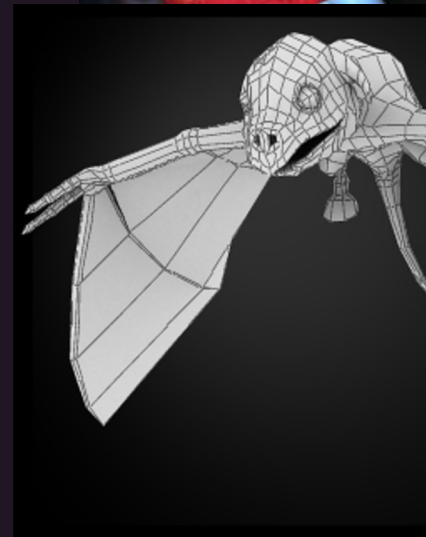
Spring School on Social XR, 2026



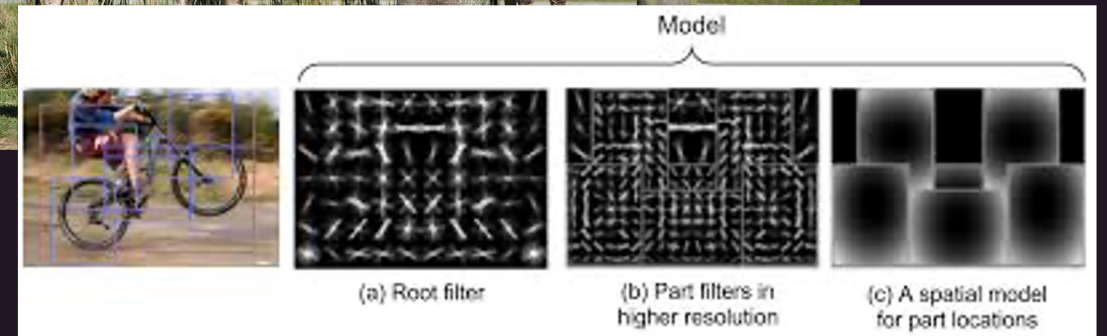
Introduction



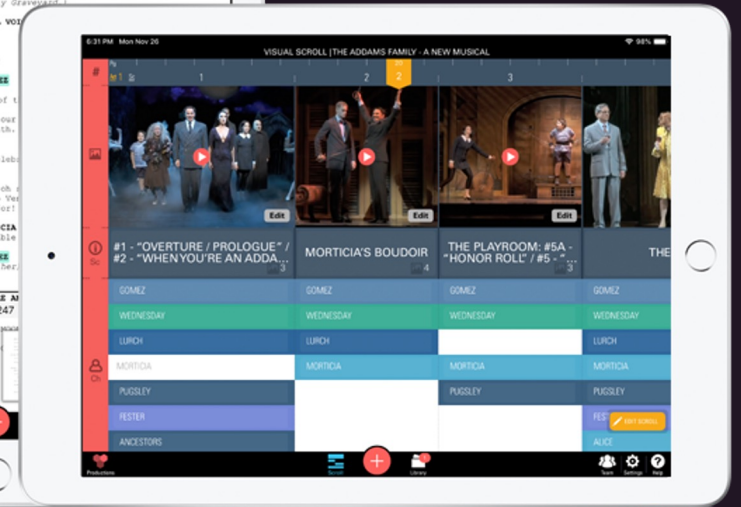
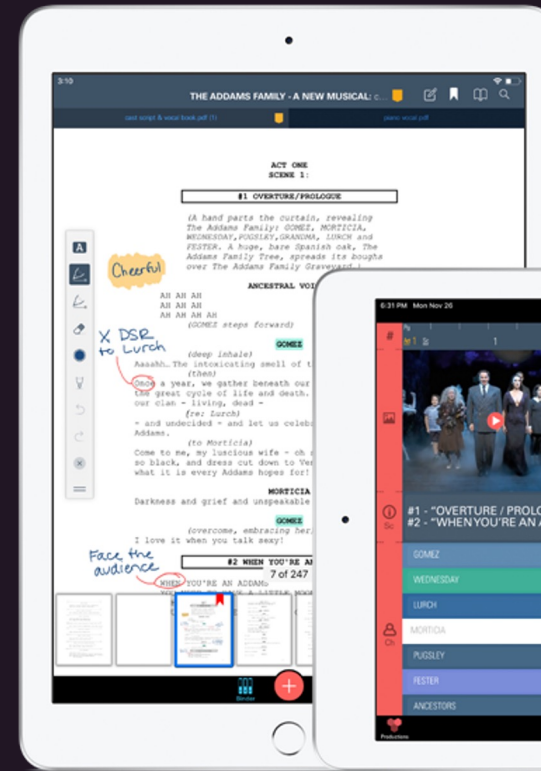
Introduction



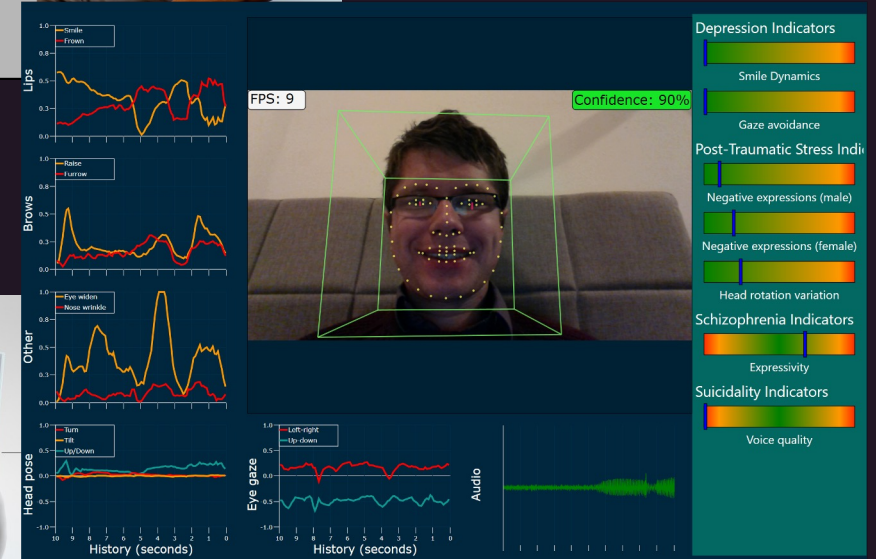
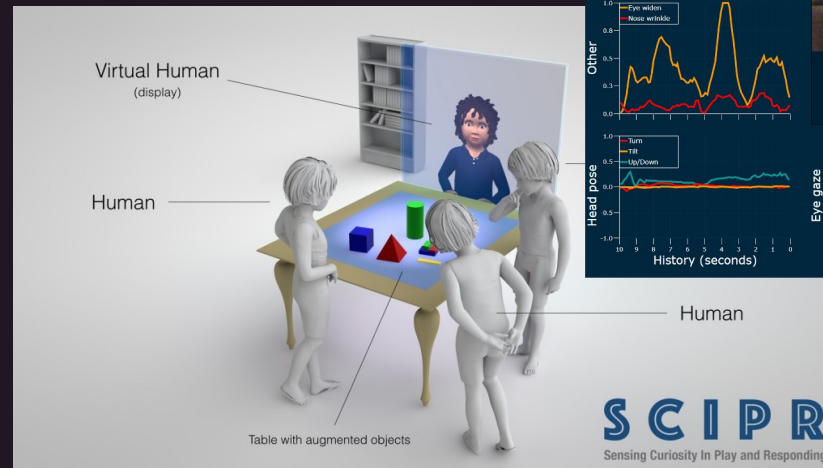
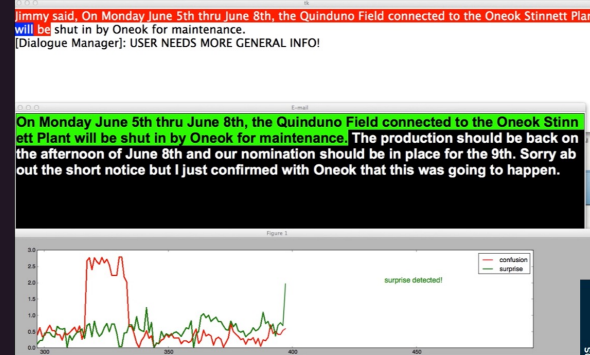
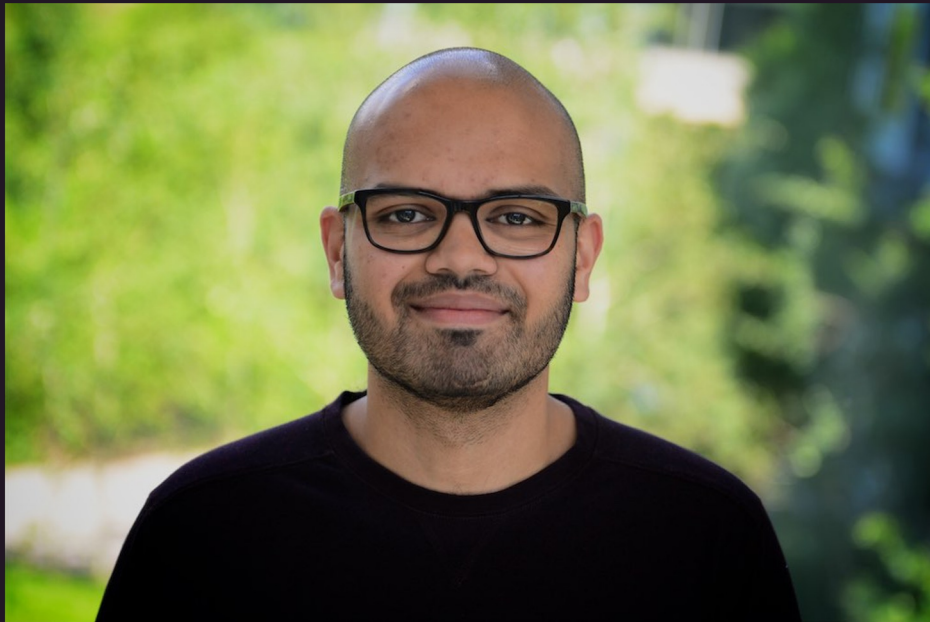
Introduction



Introduction



Introduction



Introduction



A photograph showing two hands wearing blue surgical gloves. The hand on the right is holding a pair of surgical forceps. The background is dark, and the lighting highlights the texture of the gloves and the metallic sheen of the forceps.

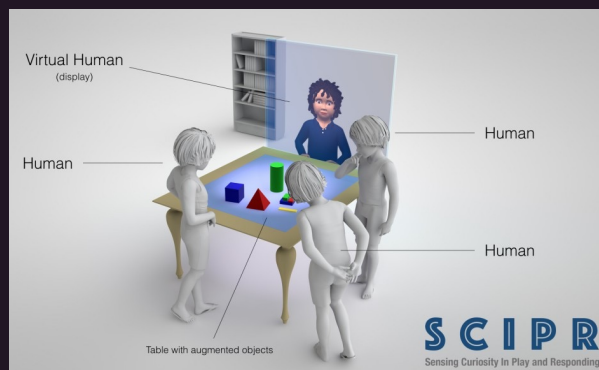
0.22 sec.

Social intelligence is inseparable from
task-based competence

The Future – Human-AI Coupled Systems



ESI Lab, TU Delft



SCIPR, CMU



Perfect World Games, NVIDIA ACE



m-Team, Michigan Med

Healthcare



Teacher Training, TU Delft

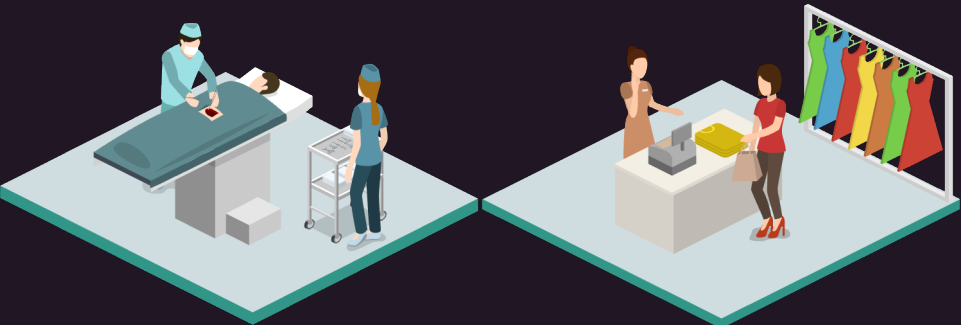
Education



Unreal Metahumans

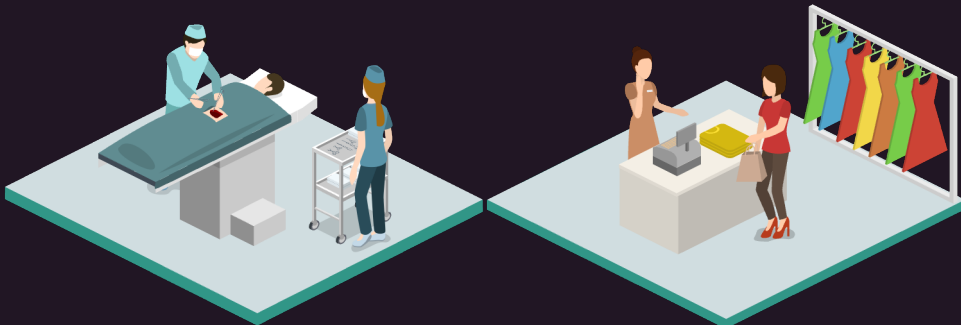
Entertainment

Humans generalize & evolve. AI reacts after retraining.



Humans infer hidden physical and social causes

Gerstenberg, T., & Tenenbaum, J. B. (2017). "Intuitive theories". *The Oxford Handbook of Causal Reasoning*.

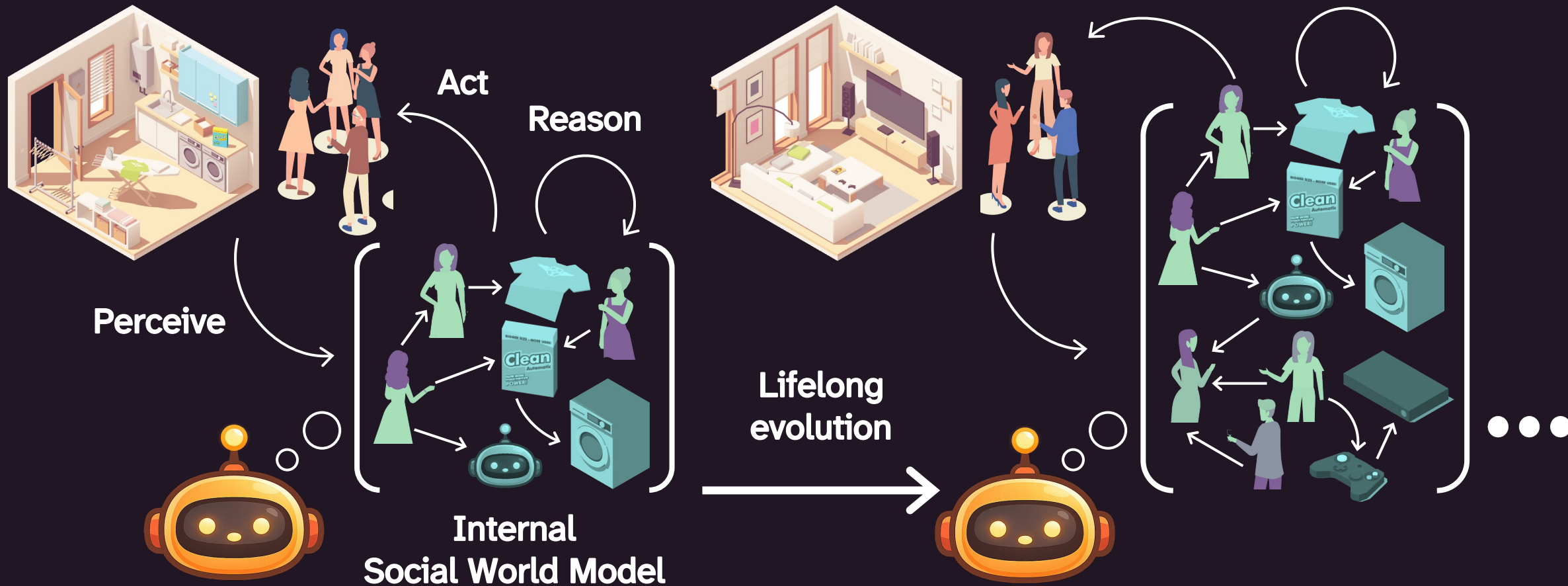


Need retraining per context
Fail to infer mental states

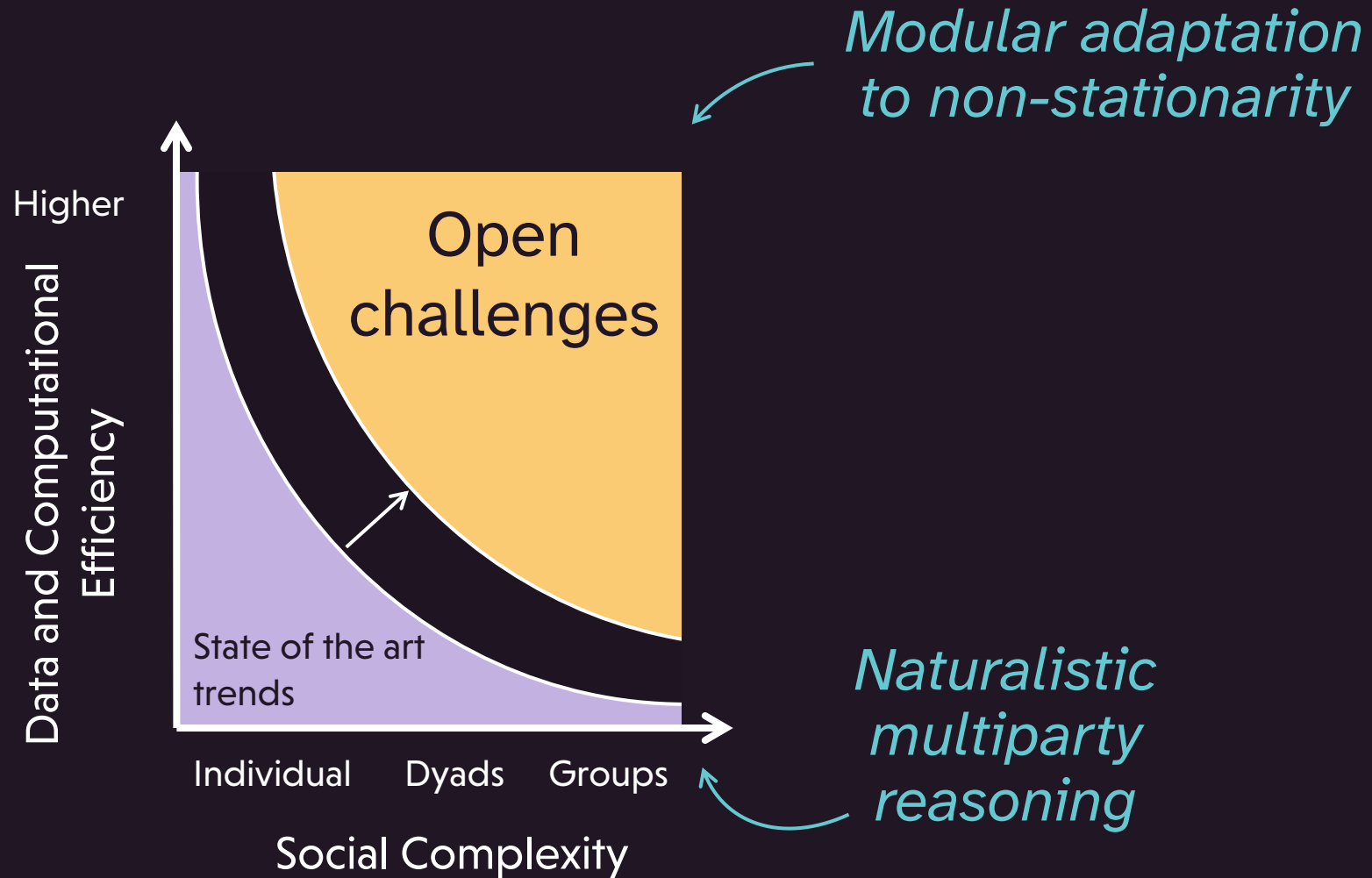
Duñez-Guzmán, E. A., et al. (2023). "A social path to human-like artificial intelligence". *Nature Machine Intelligence*.



**AI must simulate social worlds,
not just predict patterns from static data**

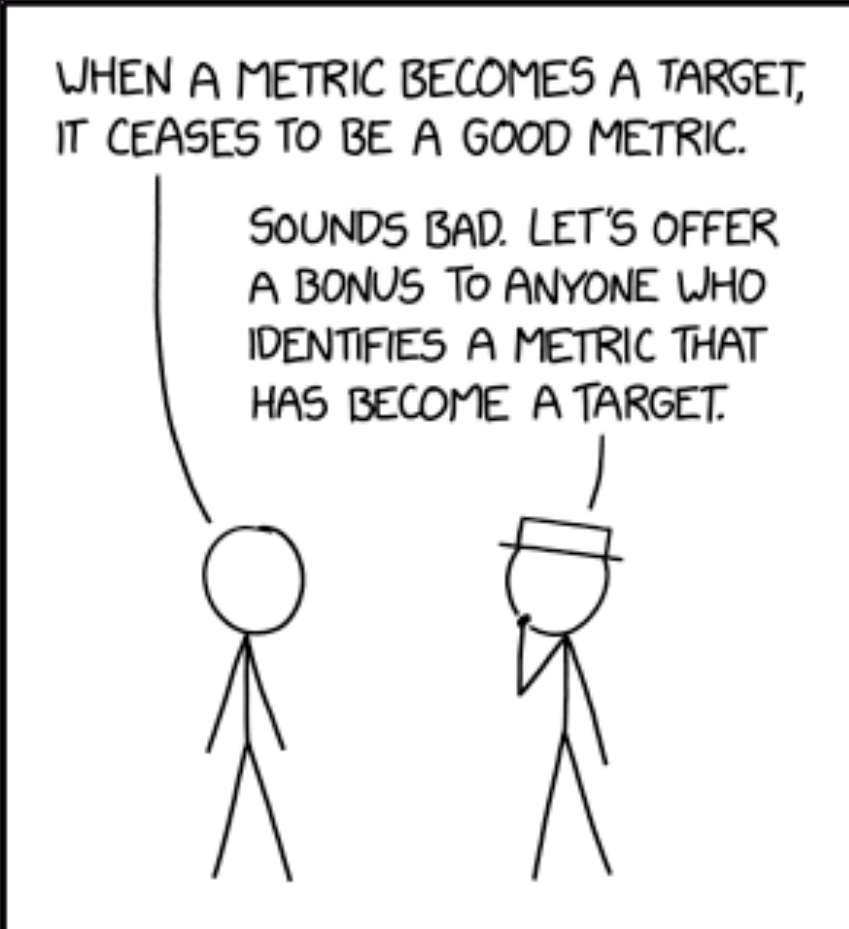


Toward Socially Intelligent, Efficient AI



Evaluation of Multiparty Social Behavior

Core challenges



[later] I'm pleased to report we're now identifying
and replacing hundreds of outdated metrics per hour

Building on evaluation quicksand

On the state of evaluation for language models.

NATHAN LAMBERT

OCT 16, 2024

Core challenges

Building on evaluation quicksand

On the state of evaluation for language models.

NATHAN LAMBERT

OCT 16, 2024

Behavior is stochastic: multiple valid responses for a stimulus

- No internal consistency between questionnaire items¹
- Objective metrics do not agree with subjective evaluations!²

¹ Wolfert, P., et al. (2022) "A Review of Evaluation Practices of Gesture Generation in Embodied Conversational Agents." *IEEE Transactions on Human-Machine Systems*.

² Kucherenko, T. et al. (2022) "The GENE Challenge 2022: A large evaluation of data-driven co-speech gesture generation." *ICMI*.

Research Questions

1. Which **objective measures** are appropriate for evaluation and training of models for multimodal behavior generation?
2. How do these objective measures relate to **subjective perception** of synthesized behavior?

Coordination to Perception

- Do you see the subtle change in Dua's expressions?
- The coordination of her head bobs change towards the end
- Both inter/intra-person coordination affects perception
e.g., freezing → perceived anxiety

Treffner, P., Peter, M. & Kleidon, M. Gestures and phases: The dynamics of speech-hand communication (2008).

Roelofs, K., Hagens, M. A. & Stins, J. Facing freeze: social threat induces bodily freeze in humans: Social threat induces bodily freeze in humans (2010).

Bernieri, F. J. & Rosenthal, R. Interpersonal coordination: Behavior matching and interactional synchrony. *Fundamentals of nonverbal behavior*. 511, 401-432 (1991).



Objective Metrics of Behavior Coordination

Cross Recurrence
Quantification
Analysis (CRQA)

Transient synchrony

Measures stability and
predictability

Beat Consistency

Empirical Mode
Decomposition (EMD)
across different
rhythmic frequencies

Handles non-
stationary signals
without assuming
fixed basis

Soft-Dynamic Time
Warping

Compares the shape of
motion or pitch
contours rather than
rigid clock time

Subjective Measures of Perception

Perceived Conversation
Quality

Inter-personal Relationships

Nature of Interaction

Equal Opportunity

Artificial Social Agent
Questionnaire

Human-likeness

Natural Appearance

Social Presence

Raman, Chirag, Navin Raj Prabhu, and Hayley Hung. "Perceived Conversation Quality in Spontaneous Interactions." *IEEE Transactions on Affective Computing*, 2023, 1-13. <https://doi.org/10.1109/taffc.2023.3233950>.

Fitrianie, Siska, Merijn Bruijnes, Fengxiang Li, Amal Abdulrahman, and Willem-Paul Brinkman. "The Artificial-Social-Agent Questionnaire: Establishing the Long and Short Questionnaire Versions." In *Proceedings of the 22nd ACM International Conference on Intelligent Virtual Agents*. New York, NY, USA: ACM, 2022. <https://doi.org/10.1145/3514197.3549612>.

How do we understand coordination?

Objective Metrics

CRQA

Beat Consistency

Soft-DTW

Subjective Measures

Perceived Conversation Quality

Artificial Social Agent
Questionnaire

How do we understand coordination?

Objective Metrics

CRQA

Beat Consistency

Soft-DTW



Subjective Measures

Perceived Conversation Quality

Artificial Social Agent
Questionnaire

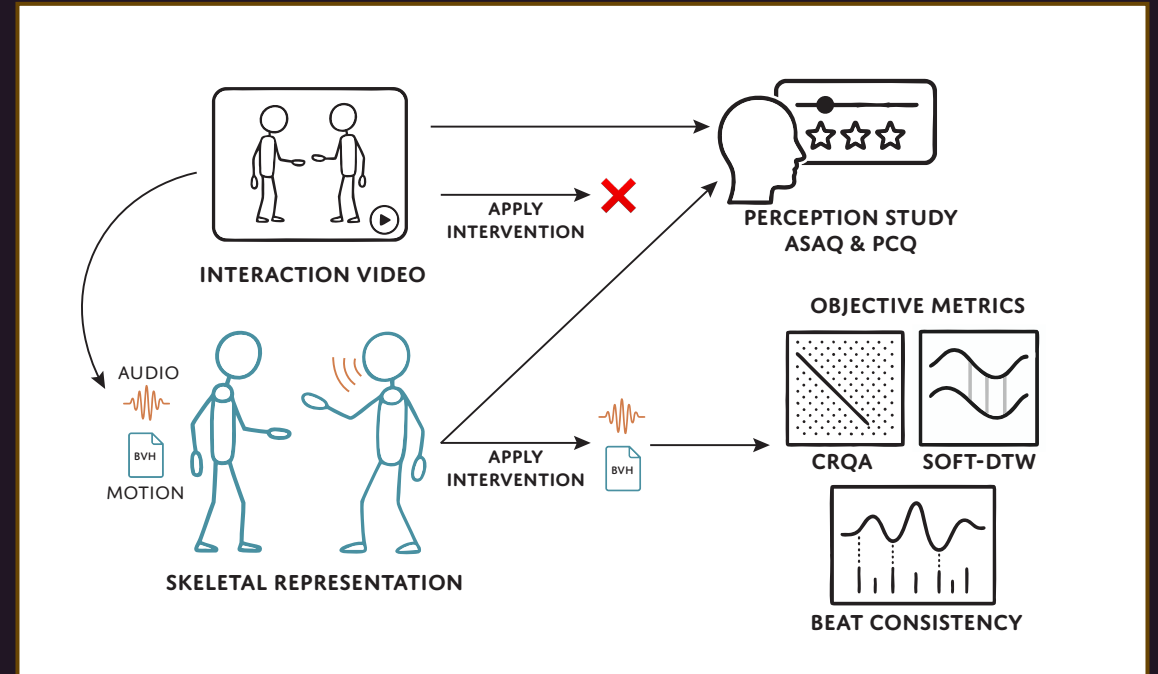
But First, Metric Sensitivity

We investigate how the metrics affect each other

Interventional Study!

- Audio-Gesture Desync
- Pitch Variance Limits
- Gesture Dampening

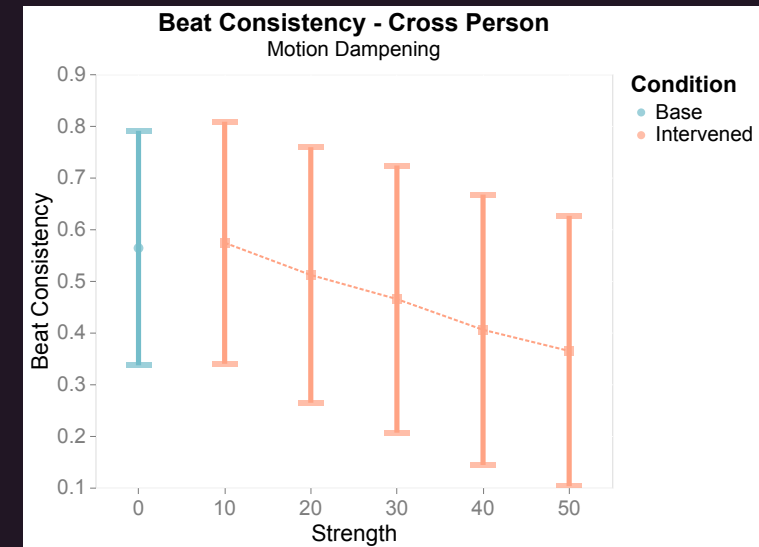
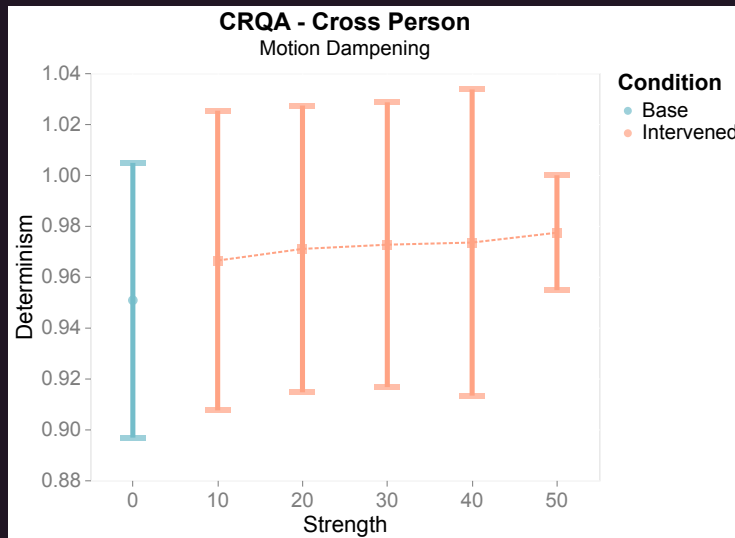
Applied at various strength levels



Metric Sensitivity: What did we find?

Kinematic Dampening (Smoothing hand movements):

- CRQA: determinism (%DET) increased
- Beat Consistency: decreased

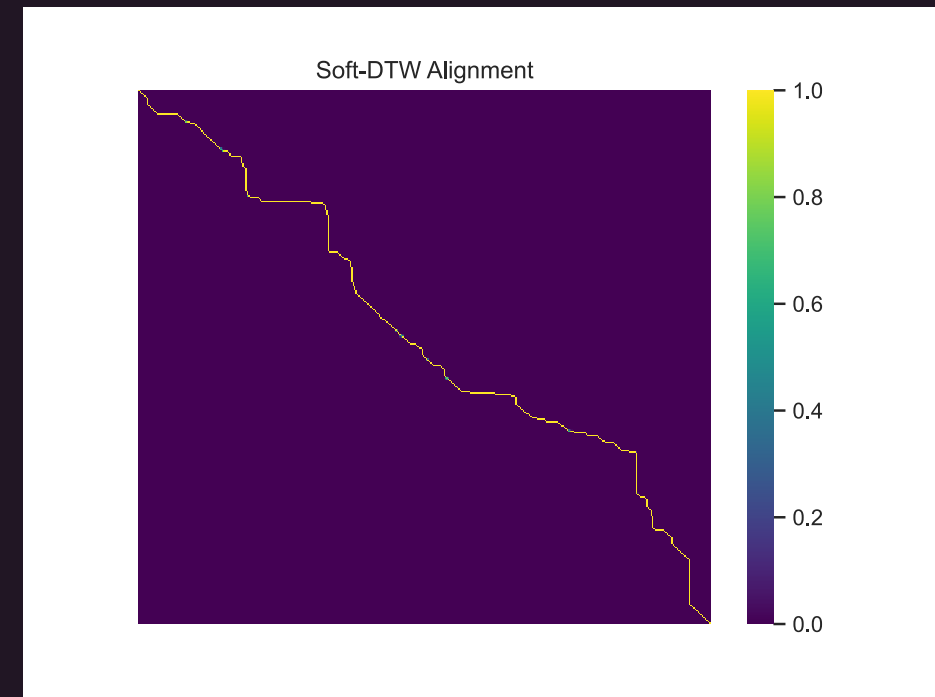
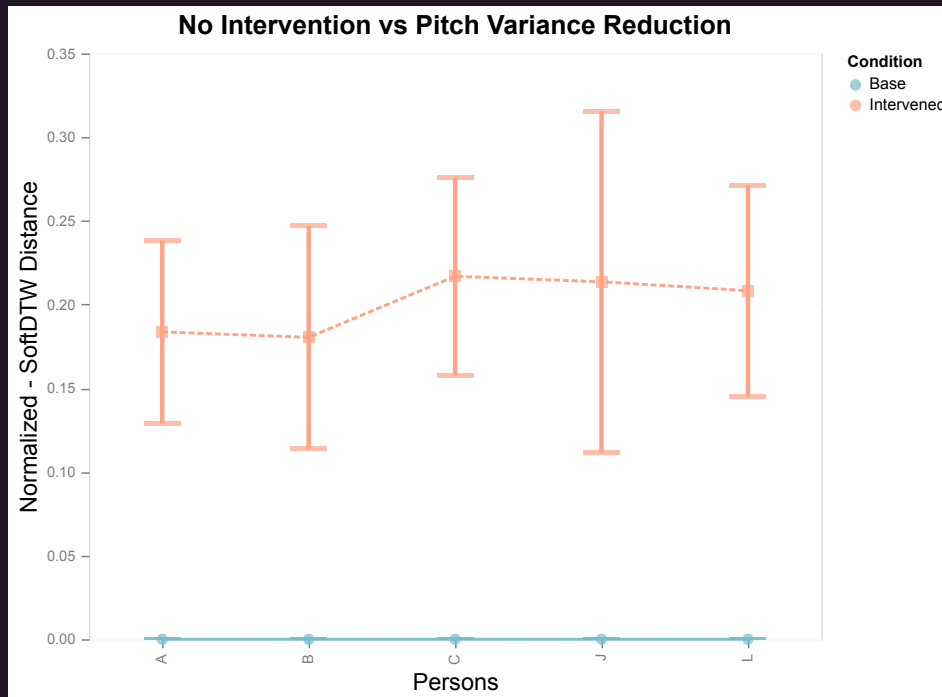


Shirekar, Ojas, Wim Pouw, Chenxu Hao, Vrushank Phadnis, Thabo Beeler, and Chirag Raman. "Multimodal Quantitative Measures for Multiparty Behavior Evaluation." In Proceedings of the 27th International Conference on Multimodal Interaction, 249–64. New York, NY, USA: ACM, 2025. <https://doi.org/10.1145/3716553.3750752>.

Metric Sensitivity: What did we find?

Pitch-Variance Reduction (flattening prosody):

- Soft-DTW (alignment): decreased



Shirekar, Ojas, Wim Pouw, Chenxu Hao, Vrushank Phadnis, Thabo Beeler, and Chirag Raman. "Multimodal Quantitative Measures for Multiparty Behavior Evaluation." In Proceedings of the 27th International Conference on Multimodal Interaction, 249–64. New York, NY, USA: ACM, 2025. <https://doi.org/10.1145/3716553.3750752>.

Metric Sensitivity: What did we find?

Speech-Gesture Delays:

- Beat Consistency: decreased, but effect is less significant



Shirekar, Ojas, Wim Pouw, Chenxu Hao, Vrushank Phadnis, Thabo Beeler, and Chirag Raman. "Multimodal Quantitative Measures for Multiparty Behavior Evaluation." In Proceedings of the 27th International Conference on Multimodal Interaction, 249–64. New York, NY, USA: ACM, 2025. <https://doi.org/10.1145/3716553.3750752>.

What's next?

Objective Metrics

CRQA

Beat Consistency

Soft-DTW



Subjective Measures

Perceived Conversation Quality

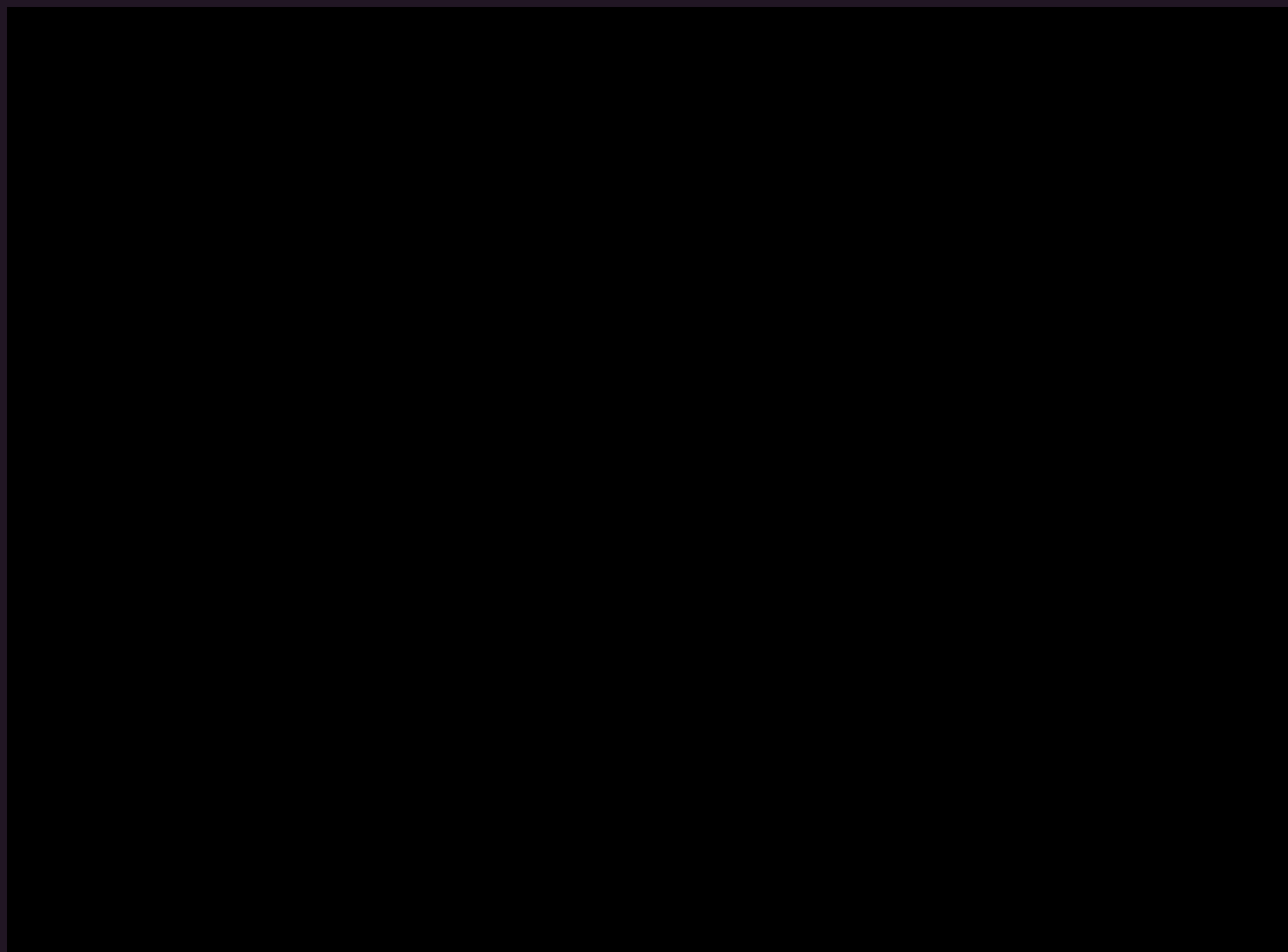
Artificial Social Agent
Questionnaire

Large-scale citizen-science
perception study

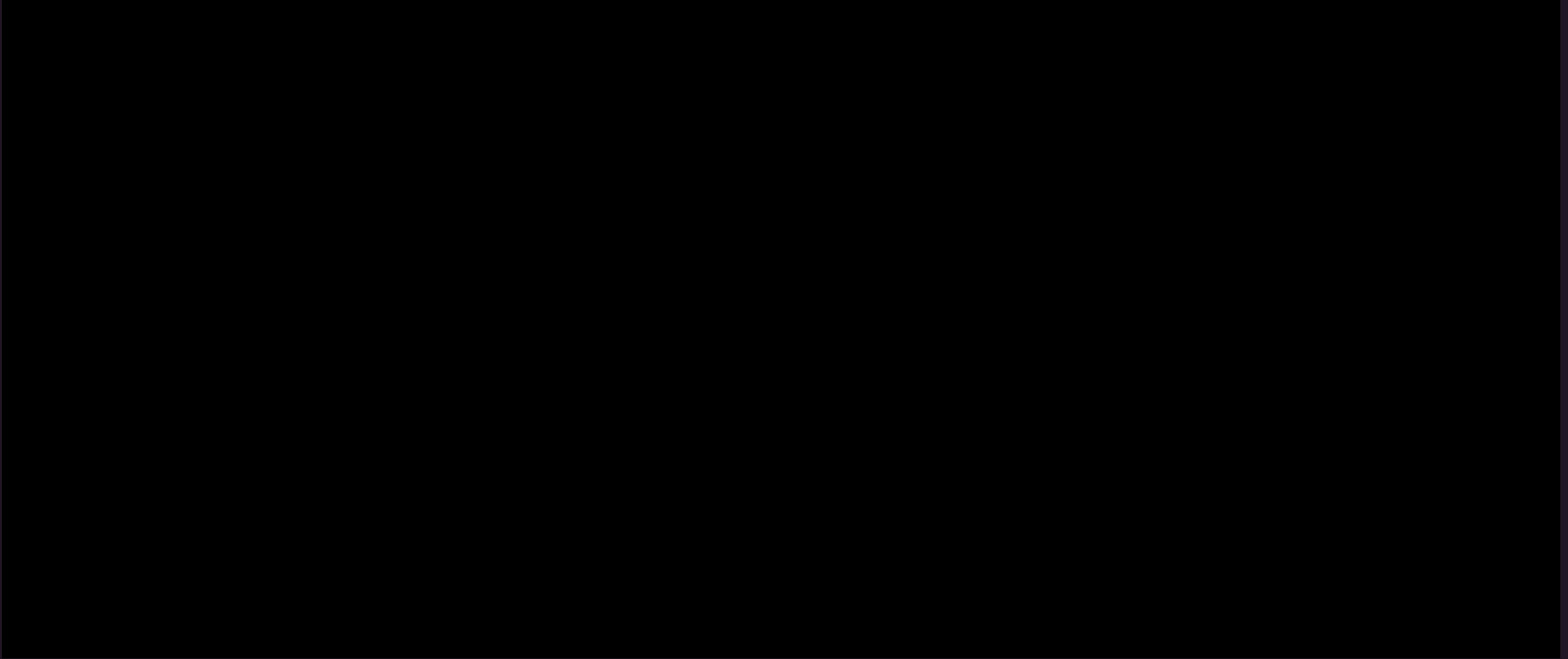


Dr. Wim Pouw
Cognitive Science
Tilburg University

Sample Stimulus

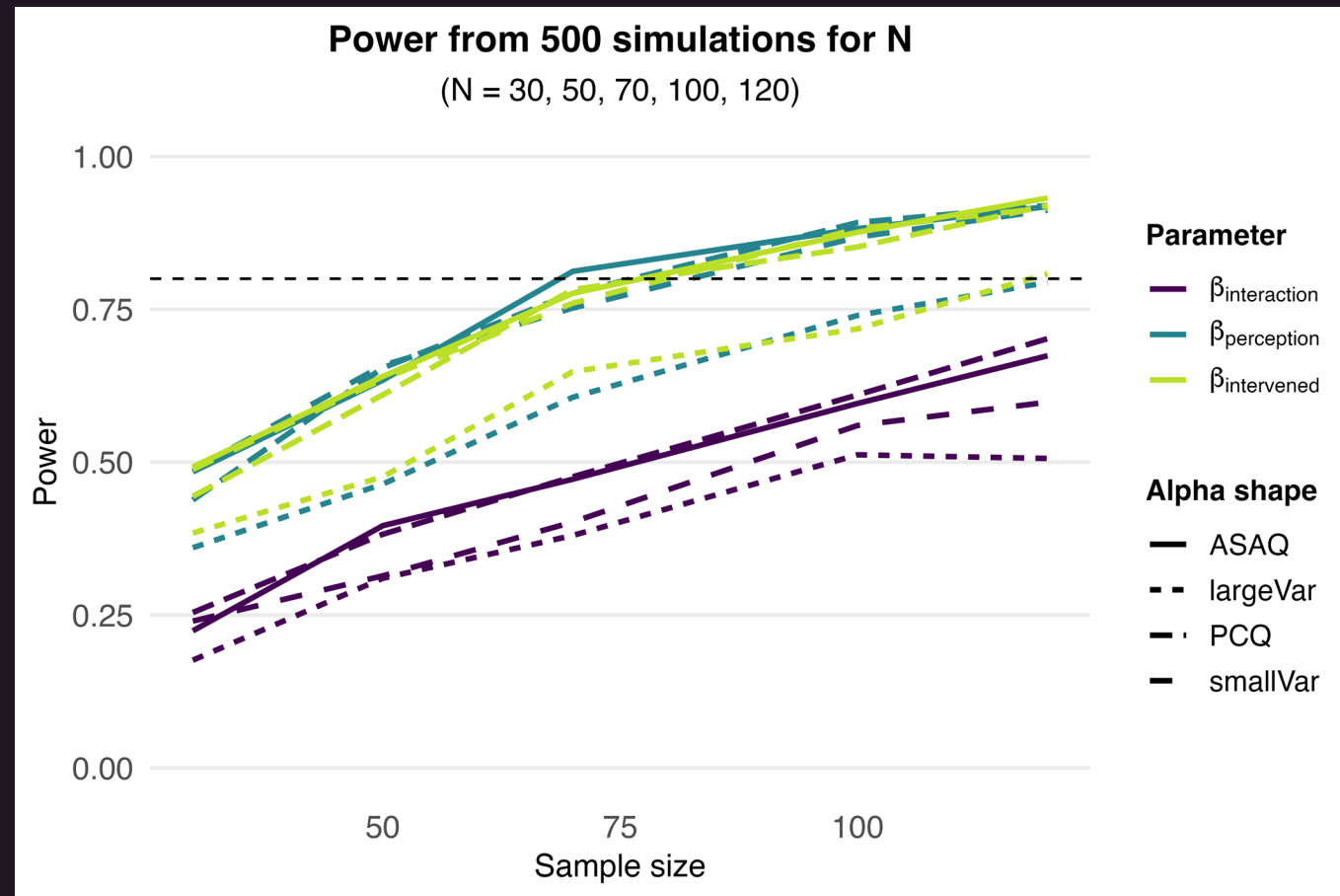


Sample Stimulus



Power Analysis

We need roughly 128 people!



Generating Multiparty Social Behavior

Modeling in-the-wild interactions

MatchNMingle Dataset



...



...



...



observed

future

Modeling in-the-wild interactions

MatchNMingle Dataset



...

...



...

observed

future

Modeling in-the-wild interactions

MatchNMingle Dataset



...

...



...



observed

future



...

→

group leaving : $y = 1$

200 instances over 90 mins
[Van Doorn 2018]

$X := \{ \text{cues over observed seq.} \}$

Modeling in-the-wild interactions

MatchNMingle Dataset



observed

future

group leaving : $y = 1$

200 instances over 90 mins
[Van Doorn 2018]



$X := \{ \text{cues over observed seq.} \}$



$Y := \{ \text{cues over future seq} \}$

Key Observation 1

The *social signal*¹ –the high-level attitudes and social meaning transferred in interactions–is embedded in the low-level cues²

1. Ambady et al. 2000
2. Vinciarelli et al. 2009

Challenges and Design Considerations

Multiple socially valid futures possible for an observed sequence

$$Y = f(X) \rightarrow \text{model } p(Y|X)$$

Participants' behaviors are interdependent - [Goffman 1963, Kendon 1990]

Capture uncertainty at *global* (group) level –
jointly forecast one future for all participants at a time

Key Observation 2

Social dynamics & behavior coordination are unique for every unique grouping of individuals

Moore 2013, Kendon 1990

Key Idea

View unique groups as meta-learning *tasks*
to generalize to unseen groups

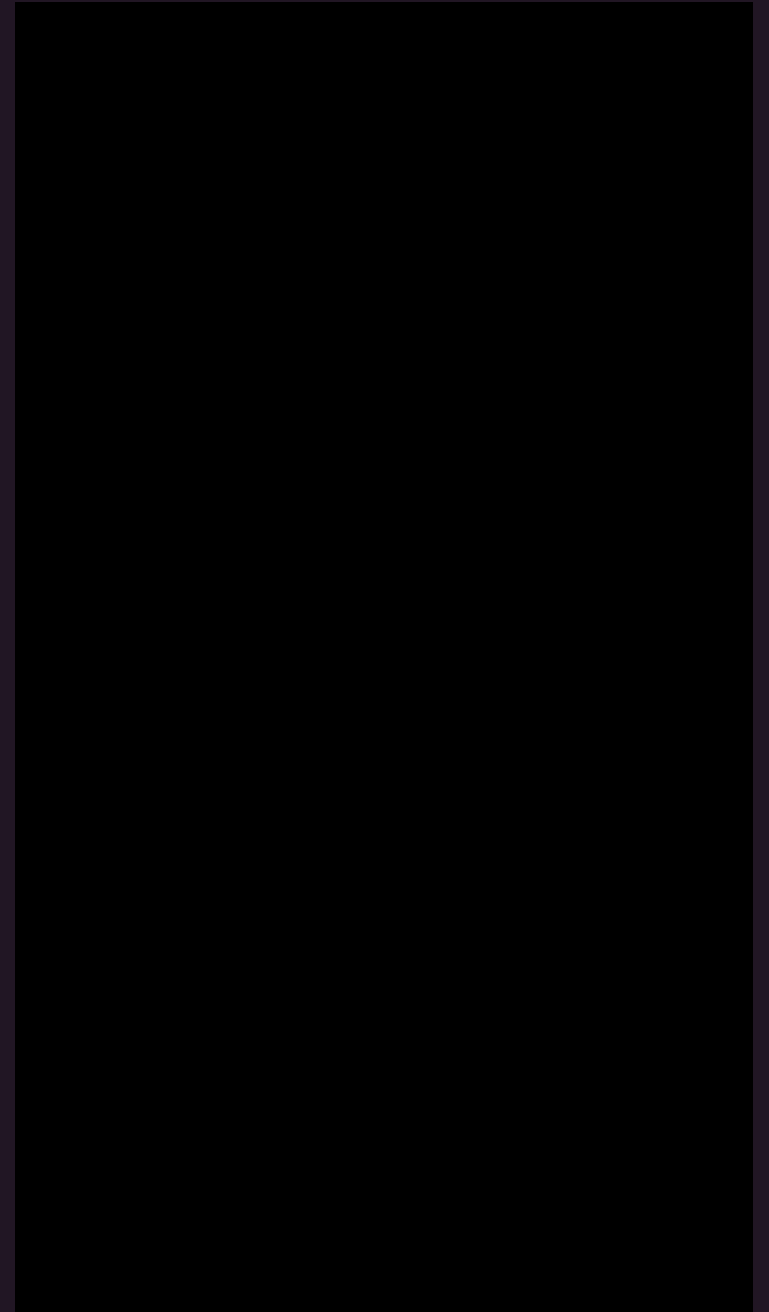
Model $p(\mathbf{Y}|\mathbf{X}, \mathcal{C})$,

where \mathcal{C} denotes a set of context
interaction sequences for a group

Behavior Interdependence

What do we see here?

- First case: group affects individual
- Second case: individual affects the group



Hierarchical Sequential Neural Processes

Behaviour of all participants

Y_t



z_t^i



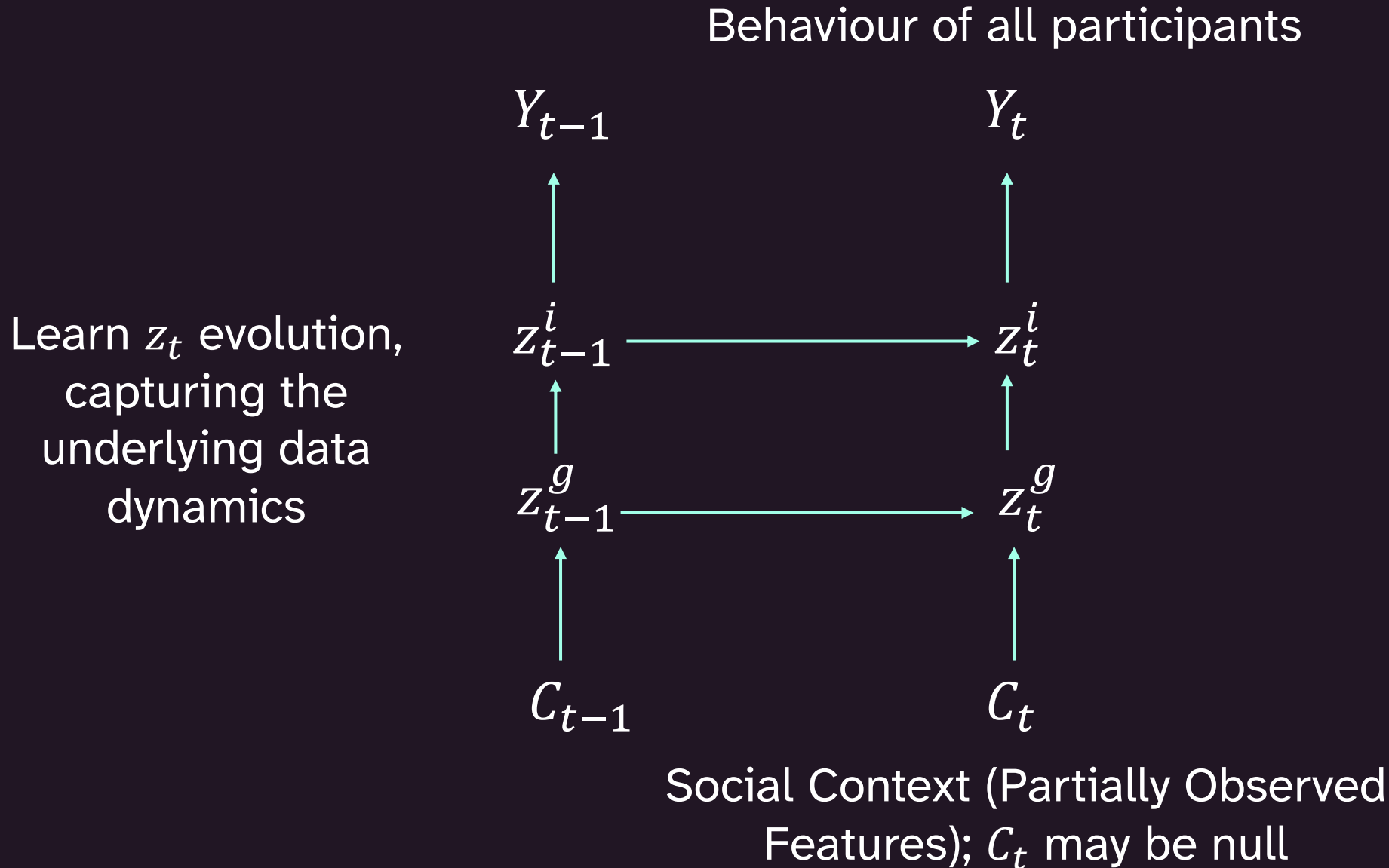
z_t^g



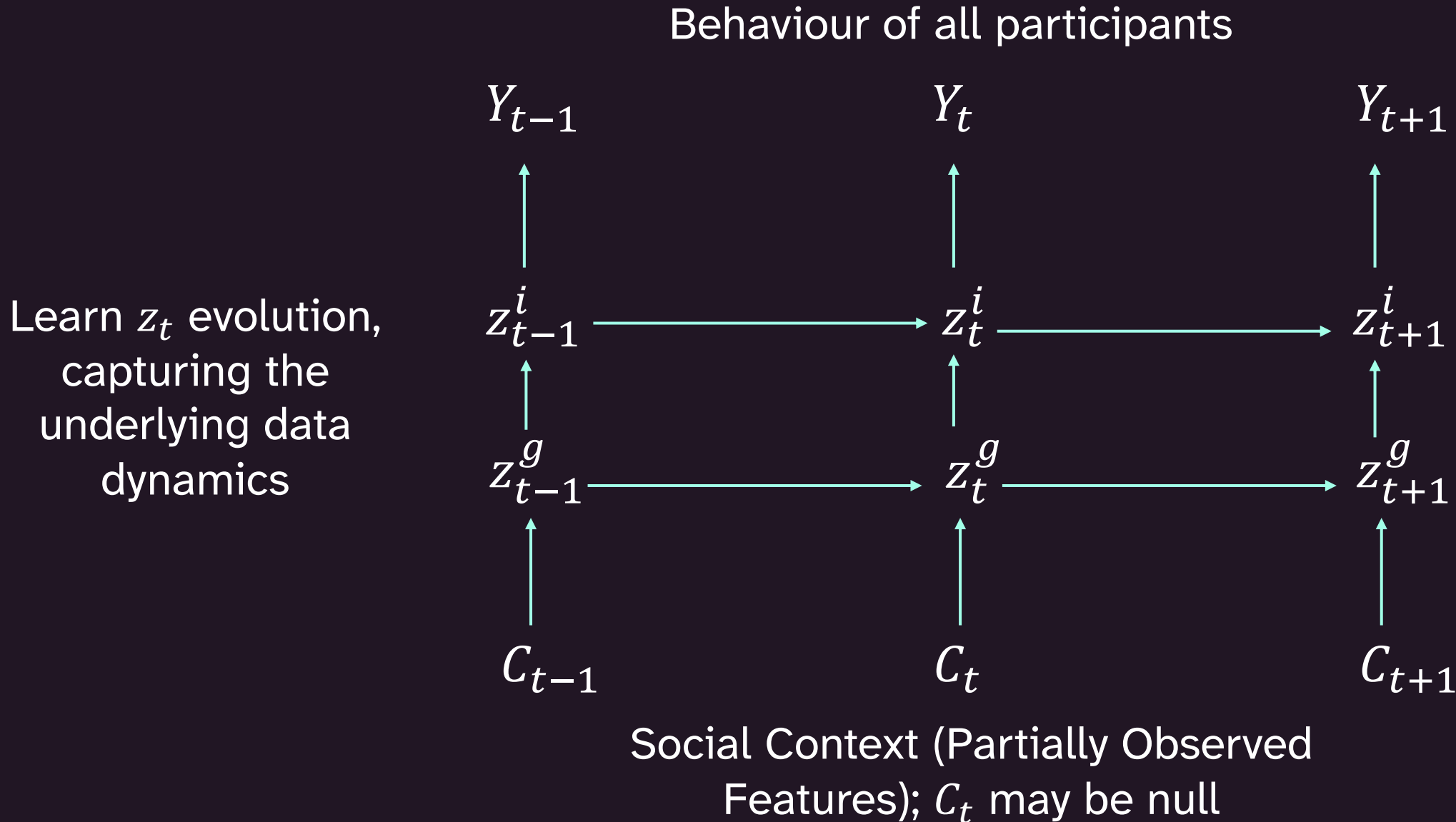
C_t

Social Context (Partially Observed Features); C_t may be null

Hierarchical Sequential Neural Processes

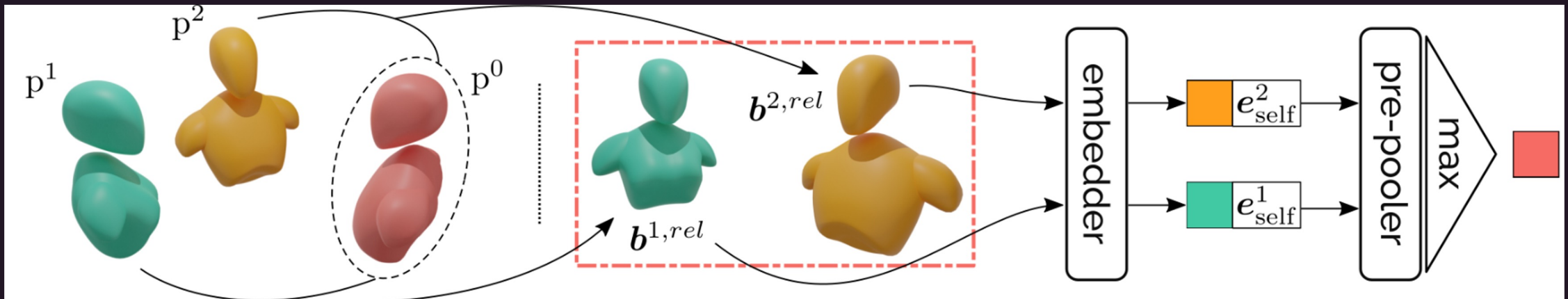


Hierarchical Sequential Neural Processes

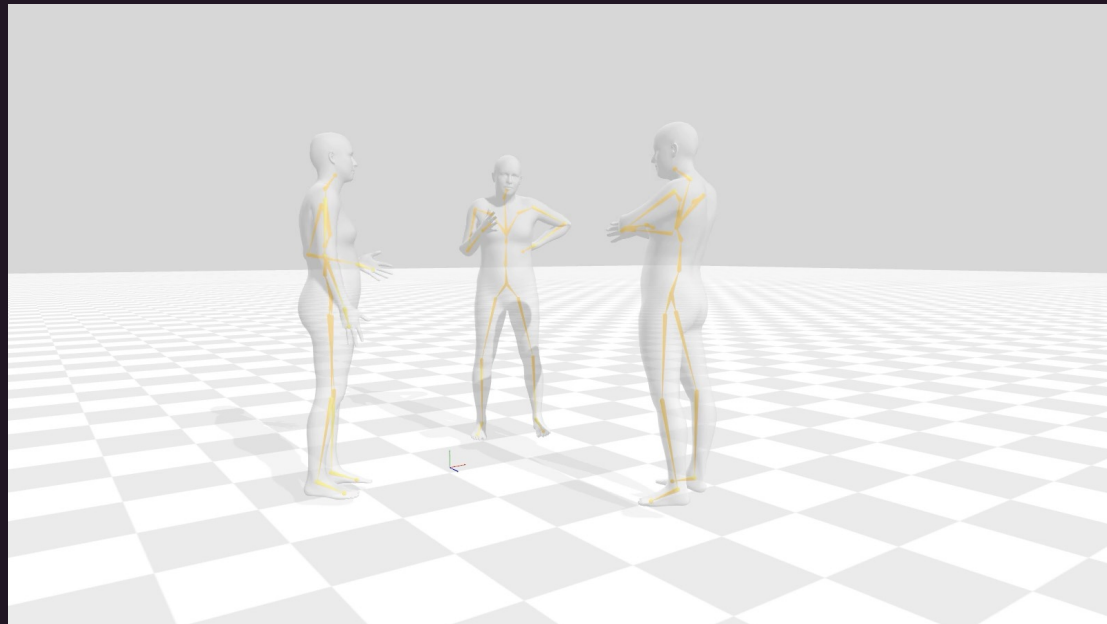


Incorporating partner behavior

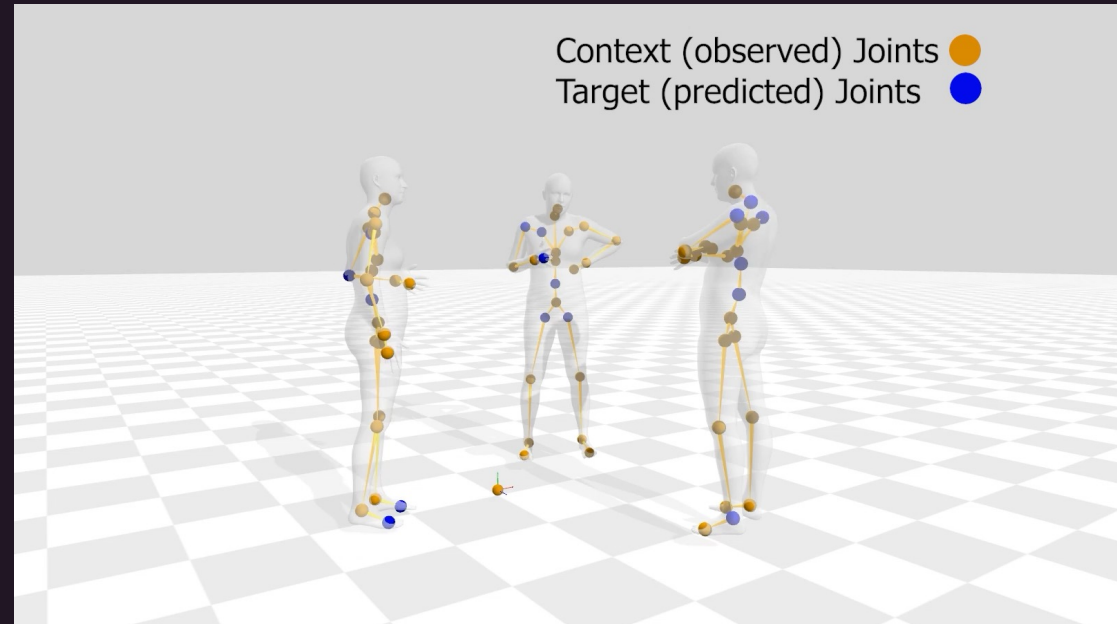
Encode the interaction from the perspective of each participant



WIP Predictions



Ground Truth



Generation

Recap

- Socially intelligent agents need to simulate social worlds and support lifelong evolution from experience
- Evaluation of simulated multiparty behavior remains challenging
 - Metric Sensitivity [CRQA, Beat Consistency, Soft-DTW]
 - Mapping low-level metrics to higher-order perception
- A hierarchical, sequential latent-variable model for capturing individual and group dynamics in generating behavior

Thanks! Questions?



Ojas Shirekar
Computer Science
TU Delft



Dr. Chenxu Hao
Cognitive Psychology
TU Delft



Dr. Vrushank Phadnis
UX & HCI
Google Research



Dr. Wim Pouw
Cognitive Science
Tilburg University