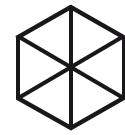




**MAX PLANCK INSTITUTE**  
FOR EVOLUTIONARY ANTHROPOLOGY



**LEUPHANA**  
UNIVERSITY LÜNEBURG

# Rethinking Variation in Social Cognition: Gaze Following across Individuals, Ages, and Communities

Von der Fakultät Nachhaltigkeit  
der Leuphana Universität Lüneburg zur Erlangung des Grades  
Doktorin der Psychologie  
– Dr. rer. nat. –

genehmigte Dissertation von

Julia Christin Prein  
geboren am 26.02.1995 in Berlin-Kreuzberg, Deutschland

Eingereicht am: 1. Oktober 2024

Erstbetreuer: Prof. Dr. Manuel Bohn, *Leuphana Universität Lüneburg*

Zweitbetreuer: Prof. Dr. Sebastian Wallot, *Leuphana Universität Lüneburg*

Erstgutachter: Prof. Dr. Manuel Bohn, *Leuphana Universität Lüneburg*

Zweitgutachter: Prof. Dr. Sebastian Wallot, *Leuphana Universität Lüneburg*

Drittgutachter: Prof. Dr. Daniel Haun, *Universität Leipzig & Max-Planck-Institut für Evolutionäre Anthropologie*

# Acknowledgments

This dissertation rests on the shoulders of many inspiring, thought-provoking, motivating, and caring individuals. The last four years have been many different things – but definitely neither boring nor stagnant. I want to thank everyone who has accompanied me on this journey.

My motivation to work on this dissertation is greatly due to my wonderful supervisors, Manuel and Daniel. Manuel, I cannot thank you enough for supporting and believing in me. You have taught me to critically challenge the status quo in research and to focus my energy on improving it – pragmatically, step by step. Even though I do know that you believe there *are* stupid questions, I felt very comfortable to ask you all along the way and I have benefited a lot from your patience with me. Thanks for roughly 200 meetings that have always left me more motivated and focused than I was before. I cannot imagine how science would be without you! Daniel, I am so happy and proud to have been a member of the CCP team. Thank you for offering me so many unique opportunities and unforgettable memories to travel to conferences and Uganda. I am deeply impressed by how quickly you get a grasp of a new project and know immediately which questions to ask. Thank you for being such an approachable director and having an open door to hear about my personal challenges. To my external supervisor, Robert: your teaching in Research Methods has turned my hate for statistics into a deep passion. Thanks for looking into open PhD positions with and for me, and accompanying me along the way. Without you, I would have never even considered doing a PhD.

Thanks to the new people that I met at the Leuphana. Sebastian, thank you for making it possible for me to pursue this doctoral title. Patricia, Franziska, Nele, and (again) Manuel, thanks for the welcoming and kind atmosphere in the LüneLütten lab. I am excited about what is to come next!

I am grateful to our wonderful cohort of PhDs at the MPI. You are a remarkable group of people, and I am very happy that we could share our daily struggles with each other. I will miss seeing you every day! Noemi, I am so glad that we got to share our office for a (way too short) period of time and successfully solved the Christmas murder mystery with Marie. Wilson, thanks for enduring our laughter through the office walls, for countless lunch breaks together, and for sharing random stories about your fieldwork (and Iceland)! Without my work wife Marie, my time during this dissertation would not have been even half as great as it has been. I could (and have and will) cry thinking about not sharing my daily life with you anymore – even though I completely trust in our relationship to live until we die. You have made me laugh so much, even if all I wanted to do was cry. Whatever it is that I am up to, I know you will always have my back and help me get back on my feet if I fail. I love your honesty, creativity, immense courage, and charming awkwardness, and feel so honored that you are my friend. Thank you for being you!

## *Acknowledgments*

My dear, dear family, I am so deeply grateful to know you by my side. Whenever I talk about individual differences, I think about you, Caro, and how we two sisters have the same parents and yet, are so different (in the very best way!!). I believe it is rare that siblings do not feel the need to compete with each other but can just love each other as they are. Thank you. To my mom, who told me right in the beginning that this journey would be a marathon and not a sprint. In some of my darker moments, you told me that even if it felt like I was about to fall, I could never — because there are so many caring people who would catch me beforehand. Thank you for always having an open home and lots of support (in the form of love, food, vacations, phone calls, and advice) for me. To my dad, thanks for showing up to Leipzig so often and establishing a little concert routine together! I believe I have learned a lot of my problem-solving skills and patience from you. Even if you do not remember saying “Mach Dinge richtig” in Caro’s and my childhood, I think my dissertation has greatly benefited from this sentence. Thank you, Louis, for helping with my move to Leipzig and spontaneously catching up over coffee in Leipzig or Berlin. To my grandpa, who told me he never worried that I would find my own way. And to my grandma, for all your curiosity, play, and love for children. To the whole Gerken-Prein clan (including our furry best friends Resi and Milou, reminding me to cuddle and catch some fresh air), I could not wish for a better family. Finally, we are all in the North again, reunited.

Thank you to all my amazing friends, who constantly remind me to live my life outside the office. To O74, OLAV, and Shubhangi, Schlotti, and Maló, for the cozy home we built together and all the little dance breaks, kitchen karaoke, reading clubs, and your patience with my mood swings and stealing lunch boxes. With you, even COVID isolation was fun. Thanks to Laura, for drying my tears, always having an open ear and heart, and helping with every single move. To Marvin, for all the long hours of deep talk and dancing. Thanks to Arne, for nerding about statistics and providing way too much green mouth refreshener. To the Osna Coxis, especially Merle, Leah, and Alissa, for all the New Year’s Eves, regular video calls, big hugs and words of comfort. To Niki, for constantly reaching out, joking around, and asking me how I’m doing. To Ani, Nica, and Schlotti for a lifelong friendship, wholehearted laughter, and a traveling diary. I am so grateful to have you all in my life.

And finally, Steven. You have taken over so many different roles in my life over the last four years. In the beginning, I was simply happy to get some technological support and coding enthusiasm from you. Thank you for pushing me beyond my levels of comfort and believing in my programming skills (... and for being my IT partner ;)). Quickly, however, coding became only one of the thousand topics that we would regularly discuss. Every morning, I looked forward to you making creepy, funny faces on our office door and having long, controversial, and nerdy conversations with you. I deeply miss you in my daily work life. And yet, I do appreciate all the changes that our relationship has gone through in recent years — through all the highs and lows. I am so sincerely grateful and glad to share my life with you. It is so easy being myself around you. Thank you for all the commitment, stability, honesty, relaxation, music, joy, and silly jokes that you bring into my life. Your love and support means the world to me. I am excited about our future in Hamburg!

# Copyright Notice

The research articles included in this cumulative dissertation have been or will be published in international peer-reviewed journals. Copyright of the text and illustrations lies with the author or authors of the respective chapter(s). The publishers own the exclusive rights to publish or use the text and illustrations for their own purposes. Reprinting any part of this dissertation thesis requires the permission of the copyright holder(s). The individual contributions of this cumulative thesis have been or will be published (in chronological order) as follows:

Prein, J. C., Kalinke, S., Haun, D. B. M.\* & Bohn, M.\* (2023). TANGO: A reliable, open-source, browser-based task to assess individual differences in gaze understanding in 3 to 5-year-old children and adults. *Behavior Research Methods*, 56(3), 2469–2485. <https://doi.org/10.3758/s13428-023-02159-5>

Prein, J. C., Maurits, L., Werwach, A., Haun, D. B. M.\* & Bohn, M.\* (2024). *Variation in gaze understanding across the life span: A process-level perspective*. PsyArXiv. <https://doi.org/10.31234/osf.io/dy73a>

Bohn, M.\*, Prein, J. C.\* Ayikoru, A., Bednarski, F. M., Dzabatou, A., Frank, M. C., Henderson, A. M. E., Isabella, J., Kalbitz, J., Kanngiesser, P., Keşşafoglu, D., Koymen, B., Manrique-Hernandez, M., Magazi, S., Mújica-Manrique, L., Ohlendorf, J., Olaoba, D., Pieters, W., Pope-Caldwell, S., ... Haun, D. (2024). *A universal of human social cognition: Children from 17 communities process gaze in similar ways* [Manuscript submitted for publication]. PsyArXiv. <https://doi.org/10.31234/osf.io/z3ahv>

Prein, J.C., Bednarski, F. M., Dzabatou, A., Frank, M. C., Henderson, A. M. E., Kalbitz, J., Kanngiesser, P., Keşşafoglu, D., Köymen, B., Manrique-Hernandez, M. V., Magazi, S., Mújica-Manrique, L., Ohlendorf, J., Olaoba, D., Pieters, W. R., Pope-Caldwell, S., Sen, U., Slocombe, K., Sparks, R. Z., Stengelin, R., Sunderarajan, J., Sutherland, K., Tusiime, F., Vieira, W., Zhang, Z., Zong, Y., Haun, D. B. M.\* Bohn, M.\* (2024). *Measuring variation in gaze following across communities, ages, and individuals – a showcase of the TANGO-CC* [Manuscript submitted for publication]. PsyArxiv. <https://doi.org/10.31234/osf.io/fcq2g>

Further articles that were written in the context of this dissertation can be found in the Appendix and have been or will be published as follows (in chronological order):

Schuwerk, T., Kampis, D., Baillargeon, R., Biro, S., Bohn, M., Byers-Heinlein, K., Dörrenberg, S., Fisher, C., Franchin, L., Fulcher, T., Garbisch, I., Geraci, A., Grosse Wiesmann, C., Hamlin, K., Haun, D. B. M., Hepach, R., Hunnius, S., Hyde, D. C., Karman, P., ..., Prein, J., ... Rakoczy, H. (2021). *Action anticipation based on an agent's epistemic state in toddlers and adults*. *Child Development* [In-Principle Acceptance of Registered Report Stage 1: Study Design]. PsyArXiv. <https://doi.org/10.31234/osf.io/x4jbm>

Bohn, M., Prein, J. C., Engicht, J., Haun, D., Gagarina, N., & Koch, T. (2023). PREVIC: An adaptive parent report measure of expressive vocabulary in children between 3 and 8 years of age [Manuscript submitted for publication]. PsyArXiv. <https://doi.org/10.31234/osf.io/hvncp>

Bohn, M.\*, Prein, J.\*, Koch, T., Bee, R. M., Delikaya, B., Haun, D., & Gagarina, N. (2024). oREV: An item response theory-based open receptive vocabulary task for 3- to 8-year-old children. *Behavior Research Methods*, 56(3), 2595–2605. <https://doi.org/10.3758/s13428-023-02169-3>

Steffan, A., Zimmer, L., Arias-Trejo, N., Bohn, M., Dal Ben, R., Flores-Coronado, M. A., Franchin, L., Garbisch, I., Grosse Wiesmann, C., Hamlin, J. K., Havron, N., Hay, J. F., Hermansen, T. K., Jakobsen, K. V., Kalinke, S., Ko, E.-S., Kulke, L., Mayor, J., Meristo, M., ..., Prein, J., ..., Schuwerk, T. (2024). Validation of an open source, remote web-based eye-tracking method (WebGazer) for research in early childhood. *Infancy*, 29(1), 31–55. <https://doi.org/10.1111/infa.12564>

\* denotes shared first or last authorship.

## Disclaimer

In this dissertation, the plural pronoun “we” highlights the collaborative efforts within the presented scientific projects. The contributions of all co-authors of the included manuscripts are transparently and accurately attributed using the Contributor Roles Taxonomy (CRediT; <https://credit.niso.org>) and can be found for each respective manuscript in [Appendix A]. In situations where the primary responsible person is not explicitly defined, “we” refers to the efforts of this dissertation’s author alone, as the plural pronoun “we” reflects standard practice in psychological scientific writings. This use does not diminish the author’s significance in conceptualizing, executing, documenting, and interpreting the presented research. The singular pronoun “I” explicitly denotes personal opinions, reflections, and decisions regarding this dissertation’s theoretical framework paper and its structure.

# **Abstract**

Prein et al. (2023)

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Luctus venenatis lectus magna fringilla urna. Laoreet non curabitur gravida arcu ac tortor. Faucibus vitae aliquet nec ullamcorper sit amet risus nullam eget. Cras semper auctor neque vitae tempus quam pellentesque. Imperdiet massa tincidunt nunc pulvinar sapien et ligula ullamcorper. Et tortor at risus viverra adipiscing at. Congue nisi vitae suscipit tellus mauris. Habitant morbi tristique senectus et netus et malesuada fames ac. Eget mauris pharetra et ultrices neque. Aenean et tortor at risus viverra. Tempor orci dapibus ultrices in iaculis nunc sed augue. Euismod lacinia at quis risus sed vulputate odio ut enim. Id eu nisl nunc mi ipsum faucibus. Est lorem ipsum dolor sit amet. Eget velit aliquet sagittis id consectetur purus. Faucibus ornare suspendisse sed nisi. Pellentesque diam volutpat commodo sed egestas egestas fringilla.



# Zusammenfassung

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Luctus venenatis lectus magna fringilla urna. Laoreet non curabitur gravida arcu ac tortor. Faucibus vitae aliquet nec ullamcorper sit amet risus nullam eget. Cras semper auctor neque vitae tempus quam pellentesque. Imperdiet massa tincidunt nunc pulvinar sapien et ligula ullamcorper. Et tortor at risus viverra adipiscing at. Congue nisi vitae suscipit tellus mauris. Habitant morbi tristique senectus et netus et malesuada fames ac. Eget mauris pharetra et ultrices neque. Aenean et tortor at risus viverra. Tempor orci dapibus ultrices in iaculis nunc sed augue. Euismod lacinia at quis risus sed vulputate odio ut enim. Id eu nisl nunc mi ipsum faucibus. Est lorem ipsum dolor sit amet. Eget velit aliquet sagittis id consectetur purus. Faucibus ornare suspendisse sed nisi. Pellentesque diam volutpat commodo sed egestas fringilla.



# Contents

<b>Acknowledgments</b>	<b>iii</b>
<b>Copyright Notice</b>	<b>v</b>
<b>Disclaimer</b>	<b>vi</b>
<b>Abstract</b>	<b>vii</b>
<b>Zusammenfassung</b>	<b>ix</b>
<b>List of Abbreviations</b>	<b>xiii</b>
<b>Structure of this Dissertation</b>	<b>xv</b>
<b>1 General Introduction</b>	<b>1</b>
1.1 Social Cognition . . . . .	2
1.1.1 Terminology . . . . .	2
1.1.2 Developmental Milestones . . . . .	3
1.1.3 Individual Differences . . . . .	3
1.1.4 Cross-cultural Research . . . . .	4
1.2 Gaze Following . . . . .	4
1.2.1 Terminology . . . . .	4
1.2.2 Developmental Milestones . . . . .	4
1.2.3 Individual Differences . . . . .	4
1.2.4 Cross-cultural Research . . . . .	4
1.2.5 Research Gaps . . . . .	4
1.3 Methodological Considerations . . . . .	4
1.3.1 Terminology . . . . .	4
1.3.2 Existing Social Cognition Measures . . . . .	5
1.3.3 The Need for a Paradigm Shift . . . . .	6
1.3.4 A Potential Way Forward . . . . .	7
<b>2 General Approach</b>	<b>9</b>
2.1 Aims of this Dissertation . . . . .	9
2.1.1 Aims of Study I . . . . .	9
2.1.2 Aims of Study II . . . . .	10
2.1.3 Aims of Study III . . . . .	10
2.1.4 Aims of Study IV . . . . .	11
2.2 Research Practices . . . . .	11
2.2.1 Accessibility of Study Materials . . . . .	11
2.2.2 Presentation Mode . . . . .	12
2.2.3 Programming Framework . . . . .	13
2.2.4 Data Processing . . . . .	14

<b>3 General Discussion</b>	<b>15</b>
3.1 Summary of Results . . . . .	15
3.1.1 Results of Study I . . . . .	15
3.1.2 Results of Study II . . . . .	15
3.1.3 Results of Study III . . . . .	15
3.1.4 Results of Study IV . . . . .	15
3.2 Research Contributions . . . . .	15
3.3 Limitations . . . . .	15
3.3.1 The task itself . . . . .	15
3.3.2 Modeling . . . . .	16
3.3.3 Cross-cultural work . . . . .	16
3.4 Outlook . . . . .	16
3.4.1 Possible spin-offs of the TANGO . . . . .	16
3.4.2 Open research questions . . . . .	16
3.5 Conclusion . . . . .	17
<b>References</b>	<b>19</b>
<b>Appendix A — Main Publications</b>	<b>23</b>
Study I . . . . .	24
Study II . . . . .	41
Study III . . . . .	101
Study IV . . . . .	131
<b>Appendix B — Further Publications</b>	<b>163</b>
Action anticipation based on an agent's epistemic state in toddlers and adults . . . . .	163
PREVIC: An adaptive parent report measure of expressive vocabulary in children between 3 and 8 years of age . . . . .	164
oREV: An item response theory-based open receptive vocabulary task for 3- to 8-year-old children . . . . .	165
Validation of an open source, remote web-based eye-tracking method (WebGazer) for re- search in early childhood . . . . .	166
<b>Appendix C — Social Cognition Survey</b>	<b>167</b>
Documentation of the Expert Online Survey “Defining Social Cognition” . . . . .	167
<b>Selbstständigkeitserklärung</b>	<b>177</b>

# List of Abbreviations

**TANGO** Task for Assessing iNdividual differences in Gaze understanding - Open

**TANGO—CC** TANGO—Cross-Cultural

**ToM** Theory of Mind



# **Structure of this Dissertation**

This dissertation is embedded in an overarching, still ongoing project led by Prof. Dr. Manuel Bohn (Leuphana University Lüneburg) and Prof. Dr. Daniel Haun (Max Planck Institute for Evolutionary Anthropology Leipzig) that will create a new task battery for social cognition. This task battery aims to capture individual differences in children from diverse communities. The present thesis describes the steps that have already been undertaken in this endeavor: during my time as a doctoral candidate, we have designed, implemented, and validated the first task of the battery. We have decided to focus the first task on one of the most fundamental social-cognitive abilities: gaze following.

The present dissertation seeks to approach the development of the new gaze following paradigm from two sides: from a contentual and from a methodological perspective. The first part will focus on the psychological construct(s) underlying gaze following: I will describe the importance of human's social cognition and its milestones during child development. I will continue to elaborate on the construct of gaze following, its developmental milestones, and the current state of research. I will end both these sections by considering studies that focus on individual differences and, subsequently, studies that were conducted cross-culturally.

The second part will focus on the methodological considerations for constructing a new task battery. I will examine existing social cognition measures, identify common challenges and shortcomings, and establish the need for an individual differences- and cross-cultural research perspective.

The remaining chapters of this dissertation will focus on our approach to develop a new gaze following measure that tackles these recognized methodological challenges. I will summarize our study findings, discuss their limitations and provide an outlook for future studies in social cognition, focusing on individual differences and cross-cultural research.



# 1 General Introduction

Humans share roughly 90 percent of their DNA with that of cats (Pontius et al., 2007), and even close to 99 percent match our closest living relatives, the chimpanzees (Waterson et al., 2005). In many ways, we are comparable to other mammals: for example, we have eyes, ears, and noses that sensorily perceive the surroundings, and we have a neocortex that processes this information and performs higher-order brain functions like spatial reasoning and motor commands. In other domains, the animal kingdom far exceeds humans in their abilities and characteristics. Cheetahs can run up to 120 km/h, and Peregrine falcons reach over 380 km/h when diving for prey. Bats and dolphins use echolocation to navigate in complete darkness. Ocean Sunfish have extraordinary reproductive abilities and can lay up to 300 million eggs at once, and the Greenland Shark can live as long as 400 years. So many other species are faster, produce more offspring, and live longer than humans. And yet, the majority of the earth's surface gets populated, exploited, and domesticated by mainly *one* mammal: the *Homo sapiens*. Which abilities brought humans to their ecological success? What makes us uniquely human?

In 1977, NASA launched the Voyager 1 and 2 spacecraft with a golden time capsule as a bold interstellar greeting. This Golden Record aimed to present the essence of Earth and humanity to potential extraterrestrial receivers. It included natural sounds, a selection of music from different cultures, greetings in 55 languages, and a photographic gallery. Among several pictures of mathematical formulas, scientific diagrams and planetary constellations, a large proportion of the included photos focused on the most fundamental human experiences. These pictures showed humans hunting and playing music together, eating and arguing at a full dinner table, collaboratively rowing a boat, competitively running against each other in marathons, nursing their children, and teaching them how to write. Strikingly, humans were rarely presented alone. These photographs, attempting to capture a comprehensive snapshot of Earth's environment, focused on humans interacting with each other in social and cultural groups.

Our ability to engage in social interactions and communicate with others has also been highlighted in prominent evolutionary and psychological theories (e.g., Tomasello, 2019). The number of publications in different psychological domains highlights that many psychologists agree on the importance of social cognition. A quick search on the APA PsycNet (on June 25, 2024, at 17:00 CET) revealed 28,594 results for the search term "social cognition", while only 3,095 results appeared for "spatial cognition" and 95 appeared for "physical cognition". Humans have been described as "ultra-social primates" (Grossmann, 2022, p. 1). We form social bonds (e.g., Bowlby, 1958), live (mostly) cooperatively (e.g., Tomasello, 2020), exchange information and learn through observing others (e.g., Csibra & Gergely, 2009). By relying on these strategies, knowledge and skills could be passed on from generation to generation — which could not have possibly been acquired by one single individual (Henrich, 2015). Sociality — and the cognitive abilities enabling it — might very well be human's superpower.

A seminal paper by Herrmann et al. (2007) tested the hypothesis that humans have unparalleled social-cognitive abilities. The researchers administered a cognitive test battery with 16 different tasks to chimpanzees, orangutans, and 2.5-year-old human infants. Indeed, the results showed that all three groups showed similar cognitive abilities within the physical world (space, quantities, causality), while human infants outperformed the chimpanzees and orangutans when it came

to the social world (social learning, communication, Theory of Mind). This study clearly underlines humans' remarkable abilities in the social realm.

## 1.1 Social Cognition

### 1.1.1 Terminology

So far, I have argued for the importance of human's social-cognitive abilities. Before going deeper into the developmental literature on social cognition, let us clarify the term: What is social cognition? Interestingly, there is not one agreed-upon definition accepted by most psychologists. Rather, cynical tongues might say that there are as many definitions of social cognition as there are research articles on it. As Ostrom put it in the Foreword of the Handbook of Social Cognition: "I regard single-sentence definitions of social cognition to be slightly offensive [...]. It would be more accurate to say that my preferred definition is the entire Handbook of Social Cognition, both this and the first edition combined." (Ostrom, 1994, p. vii). Indeed, researchers have argued that the non-specific and heterogeneous vocabulary hinders empirical and theoretical advancement: sometimes, a single concept is described by several different labels (convergence of meaning), while other concepts share the same label but refer to different constructs (divergence of meaning) (Quesque et al., 2024; Quesque & Rossetti, 2020).

To get a clearer grasp on what developmental psychologists understand under the term social cognition, we conducted a short expert survey at the beginning of my Ph.D. (in Autumn/Winter 2020). We distributed an online questionnaire via personal contacts and emailing lists (e.g. cogdevsoci mailing list) to researchers broadly working in cognitive, comparative, cross-cultural and/or developmental psychology. Out of the 100 experts that completed the survey, more than half held a professorship, and the majority focused on basic research with children. First, we asked experts to define social cognition as a psychological construct. Exactly one third of the sample (33%, n = 33) chose a definition by Glynn & Watkiss (2016): "Social cognition is the process by which actors, at individual or collective levels, decode and encode their social world, using mental models, knowledge structures and cultural understandings to process information, extract meaning and determine appropriate action." (p. 1). The next two most frequently chosen definitions focused similarly on the information processing aspects in environments with other agents ("[...] encompasses all the information-processing mechanisms that underlie how people capture, process, store, and apply information about others to navigate social situations" (Decety, 2020, p. ix); "[...] is concerned with the study of the thought processes, both implicit and explicit, through which humans attain understanding of self, others, and their environment" (Moskowitz, 2013, p. 1)). In the next step, we asked experts to name which dimensions they mentioned most often when discussing or writing about social cognition. The dimensions that were chosen by more than half of the experts were beliefs, knowledge, perspective-taking and intentions. Finally, we compiled a list of 21 social-cognitive abilities and asked experts to rate the likelihood that children of the same age would vary to perform these abilities. The experts assumed that children varied mostly in whether they manipulate emotions, deceive others, simulate others' reasoning processes, take another's perspective, and understand the subjectivity of knowledge states. Abilities that were assumed to vary the least between children of the same age were recognizing agents and following gaze. It appeared that experts assumed variation in many social-cognitive abilities and judged the more room for individual differences, the more complex a given social-cognitive ability. Please note that more information on the expert survey can be found in [Appendix C].

- distinction theory of mind vs social cognition needed?
- ToM Abkürzung einführen

A related term to social cognition is Theory of Mind (ToM). While these terms are sometimes used synonymously, social cognition functions as the umbrella term that encompasses ToM. ToM itself describes a system organized around beliefs, desires, and actions that we use in order to make sense of other agents. We try to understand actions by attributing mental states to ourselves and others – summarized by Wellman as follows: “in our everyday thinking we construe people as engaging in acts they *think* will get them what they *want*” (Wellman, 2013, p. 69).

### 1.1.2 Developmental Milestones

Importantly, human infants do not come into the world already possessing all these social-cognitive abilities (Tomasello, 2020). Much of the developmental psychology research today focuses on how children acquire the social-cognitive abilities needed to become full members of our society.

- developmental milestones in social cognition anreißen

### 1.1.3 Individual Differences

Many decades of research in social cognition have taught us much about developmental milestones and at which average children possess certain social-cognitive abilities (for a review, see Rakoczy, 2022). For example, we have learned that children are around four to five years of age when they understand that other agents can have false beliefs that differ from the child’s own belief and from reality (Wellman et al., 2001).

Anyone who has ever interacted with children of a similar age must have noticed that they can still differ in their abilities: some children might be faster in learning how to walk, while others might be quicker in learning how to talk.

*Individual differences* describe the features in which agents vary (“traits or other characteristics by which individuals may be distinguished from one another”, (American Psychological Association, n.d.-a)). *Inter-individual* differences focus on the variability *between* children, while *intra-individual* differences focus on the variability *within* the same child. Asendorpf (1992) have argued that inter-individual differences are often interpreted as trait-like characteristics (e.g. personality) that are more or less stable over time, while intra-individual differences might capture state-like characteristics (e.g. motivation) or development within a given construct. The present thesis and the associated publications focus on inter-individual differences, often abbreviated to individual differences.

A prominent framework that could explain differences between individuals is Vygotsky’s sociocultural theory of cognitive development (Vygotsky, 1962, 1978, 1987). This theory emphasizes the role of social interaction in children’s cognitive development and as such, learning as a social process. In other words, it is assumed that social interaction shapes social cognition.

In order to empirically test if social interaction shapes social cognition, we need to study children’s daily environments and cognitive abilities on the individual level.

Studies focussing on individual differences aim to measure variation between children in a given social-cognitive ability and, potentially, study the co-development of certain cognitive abilities or

link this variation to internal and external factors that influence its development. As such, individual differences studies act as a powerful tool to build and test developmental theories and the structure of social cognition (Happé et al., 2017). Additionally, individual differences studies are needed to design (clinical or educational) intervention programs. Capturing change before and after an intervention helps to evaluate its efficacy.

- also for modeling, understanding underlying processes?

Examples include studies conducted with twins (Hughes et al., 2005; Ronald et al., 2006).

- example from comparative work: some researchers focus on what chimps in average do, others focus on the very end of the distribution: what are they capable of?
- cross-cultural variation: not only average from small samples

False-belief understanding can predict children's pretense use, abilities for social interaction, peer popularity, which clearly underlines the relevance of individual differences in social-cognitive abilities for real-life outcomes (Wellman, 2013).

#### **1.1.4 Cross-cultural Research**

Wellman (2013) has argued that children's social-cognitive development is a particularly informative domain to assess in which ways people all over the world resemble and differ from one another.

### **1.2 Gaze Following**

#### **1.2.1 Terminology**

#### **1.2.2 Developmental Milestones**

#### **1.2.3 Individual Differences**

#### **1.2.4 Cross-cultural Research**

#### **1.2.5 Research Gaps**

### **1.3 Methodological Considerations**

#### **1.3.1 Terminology**

In the context of this work, the terms *variability*, *variation* and *variance* need further explanation. According to the APA Dictionary of Psychology, the definitions of variability and variation greatly resemble each other (variability: “1. the quality of being subject to change or variation in behavior or emotion, 2. the degree to which members of a group or population differ from each other, as measured by statistics such as the range, standard deviation, and variance” (American Psychological Association, n.d.-d); variation: “1. the existence of qualitative differences in form, structure,

behavior, and physiology among the individuals of a population, whether due to heredity or to environment. Both artificial selection and natural selection operate on variations among organisms, but only genetic variation is transmitted to offspring. 2. in statistics, the degree of variance or dispersion of values that is obtained for a specific variable” (American Psychological Association, n.d.-f)). Both definitions focus on differences between individuals that can be applied to the behavioral level. As in most of the developmental literature, the terms variability and variation will be used interchangeably in this work.

As the definitions of variability and variation indicate, the term *variance* denotes a statistical measure for the degree to which group members differ from one another. It describes the dispersion of a distribution by calculating how far the values of a sample spread out from their average value. Consequently, a smaller variance means less differences between individuals in a given sample (a measure of the spread, or dispersion, of scores within a sample or population, whereby a small variance indicates highly similar scores, all close to the sample mean, and a large variance indicates more scores at a greater distance from the mean and possibly spread over a larger range”, (American Psychological Association, n.d.-e)).

*Validity* concerns the truthfulness of a measure or the degree to which a measure captures what it is designed to measure conceptually (“1. the characteristic of being founded on truth, accuracy, fact, or law. 2. the degree to which empirical evidence and theoretical rationales support the adequacy and appropriateness of conclusions drawn from some form of assessment [...]”, (American Psychological Association, n.d.-c)). There are several different types of validity. In the context of this work, we focus on *convergent validity* (*i.e.*, the degree to which two measures that should be theoretically related are indeed related), *predictive validity* (*i.e.*, the degree to which a measure relates to other measures separated by a determined period (= predicting scores in the other measure)), and *external validity* (*i.e.*, the degree to which the results of a measure can be generalized beyond the original sample).

*Reliability* refers to the degree to which a measure produces consistent, stable results when repeated again (“the trustworthiness or consistency of a measure, that is, the degree to which a test or other measurement instrument is free of random error, yielding the same results across multiple applications to the same sample”, (American Psychological Association, n.d.-b); “). Two main types of reliability are often assessed when evaluating a given measure. *Internal consistency* focuses on the relationship between items in the same measure to ensure they measure the same construct. (*Test-Retest reliability* assesses a measure’s stability over time and estimates the relationship of the measurement scores on two different administration occasions (e.g., test day one followed by test day two three weeks later).

### 1.3.2 Existing Social Cognition Measures

Unfortunately, there are not too many measures that were designed to capture individual differences in social cognition in children and according to Beaudoin et al. (2020), social cognition studies rarely report psychometrical information.

Reasons for such an insufficiency in individual differences measures might be the inherent challenges in methodological advancements. Developing a reliable and valid measure is time-consuming, and (due to some features of the publishing system) not overly rewarding.

An example is a longitudinal study by Sodian et al. (2016) who have investigated how goals, beliefs and desires relate to moral ToM. They reported that “Contrary to expectations, neither goal encoding nor implicit false belief understanding was correlated with the Mo[ral] ToM false belief question” (p. 1227).

Two main interpretations come to mind when reading about these results. On one hand, the hypothesis that goal encoding and implicit FB understanding should be related to moral ToM could be simply incorrect. On the other, more optimistic hand, the non-existent relationship could be explained by methodological issues: potentially, the applied measures could not capture variation between children. Interestingly, Sodian et al. (2016) stated two sentences later: “[...] verbal IQ was positively correlated with almost all the assessed variables” (p. 1227). Why should we care about this relationship? The measures for verbal IQ were - by definition - designed to capture individual differences in linguistic intelligence. Possibly, the measures that we use for correlational analyses need to leave room for individual differences in order to find them.

What does it mean for a measure to be designed for individual differences? To drive the point home, let us imagine the opposite first. Many developmental study paradigms present children with two objects. For example, the Sally-Anne-task lets children witness how a ball gets moved from one box to a basket while another agent does or does not witness the change of location. Children then get asked where the agent believes the ball to be. Possible responses are (1) the box or (2) the basket. Similarly, in the traditional gaze following paradigm, an agent moves a toy into one of two boxes. Children themselves do not see where the toy was put but only see whether the agent looks to the (1) left- or (2) right-hand side. Common outcome measures focus on children’s looking or reaching behavior toward the two boxes. What becomes apparent is that these common social cognition tasks rely on dichotomous measures: either the child passes or fails the task. The child’s ability in question cannot be depicted on a continuous scale. A measure designed to capture individual differences would benefit from a continuous outcome measure, or, at least, more than one trial per child.

### 1.3.3 The Need for a Paradigm Shift

Classical experimental paradigms often give little weight to individual differences. When research questions focus on condition differences or the average age in which children pass a given social-cognitive task, variation between children is often handled by excluding outliers and calculating mean scores. As such, individual differences are regarded as noise or measurement error (Kidd et al., 2018).

I will argue that we need to question the interpretation that variation between individuals equals noise and should be thrown away. Indeed, some research agendas rely on the assumption that meaningful individual differences exist. For example, correlational studies aim to quantify the degree of a relationship between two variables. They assess whether children that perform well in one task also perform well in another task. In order to assess such a relationship, an individual’s ability must be measured accurately on both scales. Indeed, the correlation between two measures cannot be greater than the reliability of the individual measures. For example, let us assume we apply measure A with a known reliability estimate of .5 and measure B with a reliability estimate of .6 ...

Many research questions in developmental psychology are indeed individual differences questions.

Research focusing on the (dis)continuity of development must also rely on measures that can accurately track children's abilities across time. A textbook example is the question whether children develop some abilities similar to a tree that grows continuously over time or rather like a butterfly that develops in distinct stages (from caterpillar, to cocoon, to butterfly).

- An example of a gradual process developmental theory comes from Vygotsky. Theories within this realm assume that change occurs slowly and gradually.
- An example of a stage-like developmental theory comes from Piaget who assumes that children ... Children are assumed to go through these qualitatively different stages in a set universal order.
- First step is describing change, second step is explaining change! What changes over time and why? Process models.
- Universality in development? Maybe absolute differences in test scores but similar patterns overall?
- Interestingly, sometimes the same study paradigms are used for correlational research that were previously designed for experimental research focusing on group averages. It becomes apparent that this approach is

### 1.3.4 A Potential Way Forward

That this endeavor is nevertheless worth it can be seen from a historical perspective on science. Often, only after new measurement techniques, new knowledge could be gained and theories could be advanced. In other words: newly designed methods led to new opportunities and research questions. [todo example, maybe Teleskop]

- correlations only as large as the least reliable measure; we correlate lots but rarely know the psychometrics of the tasks
- potential issues:
  - group-level averages
  - small sample sizes
  - low trial numbers
  - dichotomous measures
  - missing or unsatisfactory psychometric properties
- Poor measurement on an individual level might conceal relationships between different aspects of cognition and may obscure developmental change
- Wishlist:
  - Objective, standardized measures (compared to anecdotal evidence, diaries)
  - Parametric measures (= continuous instead of dichotomous. Avoid floor & ceiling effects) (Schaafsma, Pfaff, Spunt, and Adolphs, 2015)
  - Satisfactory reliability estimates (Beaudoin et al., 2020; Hughes & Devine, 2015)
  - Induce variation across individuals and age groups (Repacholi, 2003)

## *1 General Introduction*

- Behavior measured should represent behavior in real world (Repacholi, 2003)
- many seemingly unrelated issues: replicability crisis, no correlations when theoretically expected, overreliance on Global North samples
- how can we address these issues? maybe for many similar solution, namely robust methods
- what makes a method robust? vali & reli
- but also: to move science forward, collaborative endeveaours are needed. Open science, share the task, so that others can reproduce results and test theories & generalizability
- we need new measures to capture individual differences in social cognition
- this thesis does so for a fundamental social-cognitive ability: gaze following
- paradigm shift: here example of one construct how this might be possible
- capture ind diff reliably. check process behind it and related constructs. see whether this is universally applicable

## 2 General Approach

### 2.1 Aims of this Dissertation

The overarching aim of this dissertation is to deepen our understanding of variability in a core social-cognitive ability, namely gaze following. In the Introduction, I have highlighted calls for better (*i.e.*, valid, reliable, generalizable) social cognition measures. The projects included in this dissertation follow this call — and try to demonstrate that, even if demanding, methodological advancement is feasible. By focusing on an individual differences measure, this dissertation aims to question the conception that variation between participants equals noise and should be disregarded. By applying the newly designed gaze following measure, I hope to show-case how individual differences research on social cognition helps us to answer fundamental questions about the development, cognitive processes and universality of human's extraordinary social-cognitive abilities. In the following, I will outline the role that each of the four included studies play in reaching these overall objectives.

#### 2.1.1 Aims of Study I

Much of the existing social cognition research relies on study designs that were created for group-level analyses. Yet, many research questions focus on individual differences. For example, do siblings have a positive influence on children's social-cognitive development (Devine & Hughes, 2018; Perner et al., 1994; Peterson, 2000; Stewart, 1983)? How are gaze following and language abilities related (Brooks & Meltzoff, 2005, 2015; Macdonald & Tatler, 2013; Okumura et al., 2017)? Can a child be trained to improve their perspective-taking ability (Lecce et al., 2014; Ruffman et al., 2018; Zhang et al., 2021)?

In order to relate cognitive abilities and external factors to each other, social cognition tasks need to be able to measure differences between children (or within children over a certain period of time). Tasks that are designed to do so remain rare. This study aimed at approaching this issue and described the development and validation of our new gaze following task, the Task for Assessing iNdividual differences in Gaze understanding - Open (TANGO). We have designed a playful study in which participants are asked to locate a balloon. The twist is that the participants cannot see the balloon itself. However, they see how another agent gazes out of a window, into the direction of the balloon. Participants have to follow the agent's gaze to locate the balloon. Their precision in doing so is measured continuously, namely in the distance between the actual target's location and their assumed location.

Study I consisted of three main parts. In the first part, we tested the task's usability in a sample of 3- to 5-year-old children and adults. We assessed whether the TANGO captured individual differences and whether data collection mode (*i.e.*, remotely vs. in-person) influenced children's performance. In the second part, we aimed at assessing whether these individual differences were reliable. We repeated data collection two weeks after the initial session and estimated retest-reliability and internal consistency. In the third part, we attempted to gauge the TANGO's validity. For its construct validity, we explored the relationship between the TANGO and variables of children's daily social environment. For its predictive validity, we examined the TANGO's relationship with a language

## *2 General Approach*

measure six months later. Together, these studies aimed at thoroughly assessing the TANGO as a new task to measure individual differences in gaze following.

### **2.1.2 Aims of Study II**

Study II built upon the results of Study I and adopted the TANGO to examine three research questions related to individual differences. We aimed at assessing (1) how gaze following develops across the lifespan, (2) how people process gaze cues, and (3) how gaze following relates to other (social-) cognitive abilities.

In attempting to answer the first question, we remotely collected data from 3- to 80-year-olds and examined changes in their ability to follow gaze. While previous gaze following research often focused on the youngest age in which children can follow gaze (D'Entremont et al., 1997; e.g., D'Entremont, 2000), we focused on the development (in precision) after this initial milestone.

Second, we developed a computational cognitive model to explain the underlying processes of how people process gaze. We hypothesized that participants observe the pupil location within the eye and calculate gaze vectors that point toward the attentional focus. As this process was likely to be noisy, individual differences were assumed in the participant's uncertainty surrounding this gaze vector. We judged the evidence for this gaze model by comparing it to two alternative explanations of the data: a center bias and a random guessing model. Due to its underlying geometrical features, the gaze model predicted that uncertainties around the gaze vector increased the further an agent looks to the side. We planned to check for this signature pattern in the data by assessing the participants' imprecision levels in each target bin.

Third, we intended to test whether an individual's precision in gaze following related to other (social-) cognitive abilities. Consistent with our gaze model, we assessed whether gaze following was related to non-social vector following. To test this claim, we designed a new vector following task that matched the TANGO as closely as possible but within a non-social context. The vector following task relied on the concept of magnetism and asked children to locate a magnet based on the location of a gearwheel. Furthermore, we applied four tasks out of a ToM task battery (Wellman et al., 2001) and two perspective-taking tasks (Flavell, Everett, et al., 1981; Flavell, Flavell, et al., 1981). These comparisons aimed at disentangling the components that make up gaze following. In sum, Study II aimed at providing a comprehensive picture of the lifelong development, cognitive processes, and components of gaze following.

### **2.1.3 Aims of Study III**

Studies I and II presented above relied on child samples from the Global North (Leipzig, Germany). Developmental theories often assume that central features of social cognition hold universally true across human cultures (Wellman, 2013). With Study III, we aimed at directly evaluating this claim. To broaden our perspective, we collected data of children's gaze following abilities from 17 diverse communities spanning five different continents. To our knowledge, this attempt constituted the largest international sample on children's gaze following abilities to this date.

We aimed at analyzing the data in three different ways. First, we examined the developmental trajectory of gaze following across communities. Second, we explored whether absolute differences

in gaze following imprecision related to variables concerning the data collection mode and/or children's daily environments. Third, we checked whether the signature pattern predicted by our gaze model in Study II could be found in all communities studied. With this approach, we aimed to investigate whether the same underlying model could describe the ways in which children worldwide process gaze. All in one, this study intended to provide first evidence for the universality or variability of gaze following across cultures.

### **2.1.4 Aims of Study IV**

Study IV intended to strengthen the basis for using the newly developed gaze following task for cross-cultural research. The paper had three main aims. First, we described the process of designing the cross-cultural gaze following task, the TANGO—Cross-Cultural (TANGO—CC). We aimed to provide a roadmap for other researchers how to pragmatically develop a new social cognition task that can be used for studying diverse communities.

Second, we provided a tutorial how other researchers can use and customize the TANGO—CC according to their needs. This section aimed to highlight the components of the task that are constant and those that are variable and can be adapted to diverse languages and communities.

Third (and most importantly), we aimed at examining the TANGO—CC's reliability across diverse communities. While the psychometric properties of the TANGO were previously tested in one setting (Leipzig, Germany; Study I), we cannot simply generalize these findings and assume the task's validity and reliability based on a mono-cultural sample. For this purpose, we aimed at further analyzing the data set from Study III. We assessed performance across trial types, within- and between-community variation, and internal consistency measures in all of the 17 communities studied. Together, this study aimed at providing a direction of how researchers can study variation in social cognition at an individual- and community-level.

## **2.2 Research Practices**

In the sections above, I have highlighted the theoretical approaches and aims of the publications included in this dissertation. All of these studies build upon the same underlying testing infrastructure. As this dissertation has a strong focus on methodological advances, I want to outline our main technological decisions within the process of designing this new testing infrastructure. In-depth information regarding each specific study can be found in the corresponding manuscripts in Appendix A.

### **2.2.1 Accessibility of Study Materials**

Elson et al. (2023) figuratively stated that “psychological measures aren’t toothbrushes” (p.1). The authors argue that psychological tasks are often used only a handful of times – namely by the same researcher(s) who designed the measure. However, this reluctance to share own measures and reuse measures of other researchers hinders cumulative science. Elson and colleagues, therefore, call for more transparency and standardization in the psychological tasks we implement and apply. We took this call seriously and strongly adhered to Open Science principles in the construction of our new standardized gaze following task.

The projects included in this dissertation pre-registered studies with their planned sample sizes, study procedures and analyses for hypothesis testing before data collection, as well as shared the study materials, analysis scripts and data sets after data collection. This dissertation and its included publications have published all related information in online repositories on GitHub (<https://github.com> and the Open Science Framework (<https://osf.io>; links leading to specific projects can be found in the respective manuscripts in the Appendix A).

Articles were written in R Markdown which allowed us to combine continuous text and results of statistical analyses into one coherent document. The aim was to reduce the chances to accidentally copy-paste incorrect or outdated findings into the manuscript. We tracked all analysis scripts and the manuscript texts in a version-controlled manner. This enables other researchers to transparently retrace decisions and changes along the way. The newly implemented task(s) were made openly available to other researchers to use and share as they wish. The source code of the website(s) have also been published on GitHub. By cloning the corresponding repository, motivated researchers can modify and further develop the task to their needs.

As the final project, Study IV incorporated refinements that have further optimized the usability of the gaze following task. While adapting the source code for cross-cultural stimulus presentation, I have refactored the code base; that is, I have improved the internal structure, readability, and maintainability of the software without altering its behavior or functionality). I have added comments to the code base using JSDoc (<https://jsdoc.app>), a markup language developed to document and annotate JavaScript source code files. Furthermore, I have collected frequently asked questions and created a manual for other researchers using the task (<https://ccp-odc.eva.mpg.de/tango-cc/manual.html>). Compared to the first version of the TANGO, these changes served to improve the task's sustainability over time.

## **2.2.2 Presentation Mode**

Before designing the TANGO itself, we had to choose a presentation mode for the experimental stimuli. Many traditional developmental studies have designed interactive tasks in which children interact with the experimenter or their caregiver. This makes standardization difficult as experimenters (or caregivers) have individual communication styles, might accidentally alter the instruction script due to memory gaps, or need replacement at some point (e.g., when temporary work contracts of research assistants end). In addition, this approach makes scaling data collection difficult as training experimenters requires substantial temporal and financial resources. Since studies investigating individual differences rely on large sample sizes, we decided to take another route. We designed an online testing platform that hosts studies in all common web browsers. It can be accessed on all devices (e.g., laptop, tablet) and does not require prior installation. Even though the task was programmed as a website, it does not require a working WiFi connection (for more information, see Programming Framework further down).

The design as an online study has several advantages. First, it allows for maximally standardized procedures as audio instructions are pre-recorded and automated. Second, data does not need to be entered manually or coded from video. Instead, the website stores response times and clicking behavior of participants automatically. This greatly reduces potential coding biases and rater errors. Third, the presentation as a website allows for in-person and remote study participation. Families can choose whether they want to participate by coming to the research lab or kindergartens, or whether they want to participate within the comfort of their own home. In this case, no appointments or rooms in the lab need to be scheduled. The families receive a personalized invitation link

and can access the online study whenever and wherever it suits them best. This increases the inclusiveness for families with busy schedules, longer traveling distances to the lab or children who might be shy to interact with strangers. While studies conducted in a research lab might be prone to over-study families with certain characteristics (e.g., families living closer by, having a higher SES or more free time), online studies might make study participation more accessible for a wider range of families – which, in turn, would increase the representativeness of the sample.

URL workings todo

We have consciously chosen to implement an active behavioral measure as we wanted to circumvent issues in interpretability of implicit measures. For example, it is unclear whether looking time measures capture children's level of surprise, attention, memory formation or other cognitive processes, and whether they show convergent validity (Aslin, 2007; Dörrenberg et al., 2018). While our outcome response does not provide a direct measure of underlying cognitive processes, it focuses on children's active behavior which might translate more directly into their actions in daily life.

The task itself is presented as an animated, interactive picture book. Audio instructions guide children through the task and provide a description of the presented events. Trial types build up on each other and familiarize children with the response format.

### 2.2.3 Programming Framework

The online study was programmed in JavaScript (functionality of the website), HTML (content of the website), CSS (styling of the website) and PHP (server-side interaction for down-/ uploading data). We decided to refrain from using web development frameworks like Next.js or React: as web development is a very fast moving domain with frequently appearing technology advances, projects relying on new frameworks often require regular monitoring. Our goal was to create a durable code base that avoids breaking changes through large updates. Additionally, we aimed at keeping the code base as easily accessible as possible so that researchers with less programming experience can also read through and adjust the code. We reasoned that this could be better achieved by relying on the so-called "Vanilla JavaScript" without additional frameworks. An exception to this was the one animation library we used. The GreenSock Animation Platform (GSAP; <https://gsap.com>) is a JavaScript library for industry-standard, high-performance web animations. However, GSAP is widely spread, well documented, and easy to read so including this library rather increased user friendliness.

We used one further tool, namely a module bundler for JavaScript. Parcel (<https://parceljs.org>), and, in later projects, webpack (<https://webpack.js.org>) tracked modules and dependencies, optimized website performance and enhanced the developer experience (e.g., by enabling a development server and hot reloading). However, please note that these bundlers only alter the source code on the sever-side – the functions defined in Vanilla JavaScript remained in their human-readable format. During website development, we made use of the development server made available by the module bundler. This means that we could test the website on our local computer and only uploaded content to the web server once the website was fully implemented. The bundler builds a so-called "dist" folder that contains all website content in a compressed and optimized format. This folder is essential for both online and offline data collection. For online data collection, the dist folder gets uploaded to the web server and the study can be accessed under its corresponding URL. For offline data collection, a live server can be installed which then hosts the dist folder locally. This is especially practical in settings with unreliable or non-existent WiFi connection: data collection in communities living in remote

## *2 General Approach*

villages becomes viable. Additionally, this approach is helpful for data collection with children in (German) kindergartens as these institutions often do not have public internet connection or stable telephone reception.

### **2.2.4 Data Processing**

After a child completed the task, the website automatically saves all responses and response times. The website stores all click data and, if desired, some additional device information (e.g., touch screen yes/no, what type and version of web browser) in a text format. During initial development and the first published article (see Study I), we used JSON (JavaScript Object Notation) files which are commonly used for transmitting data in web applications. In later stages of the project (see Study III), we instead used CSV (Comma-Separated Values) files as these are more commonly used among psychology researchers. In these text files, each trial is represented by one row of data. Columns depict the collected variables; for example, subject id, trial type, target location, and response. If requested, a webcam recording can be captured that films the study participation via the front camera. We implemented this feature especially for remote data collection. If caregivers consented to the webcam use, we could use the webcam recordings to ensure that indeed the children themselves carried out the task and did not get any external help. Depending on the choice of the experimenter and the data collection setup, data could be uploaded to secure servers located within the Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany, or could be downloaded to the testing device (e.g., tablet or computer).

Data pre-processing was greatly simplified as the study data does not require any further response coding. Individual subject files were merged into one final data set using the programming language R. Data visualization and statistical analyses were conducted within the same program.

# **3 General Discussion**

## **3.1 Summary of Results**

### **3.1.1 Results of Study I**

The task design

### **3.1.2 Results of Study II**

### **3.1.3 Results of Study III**

### **3.1.4 Results of Study IV**

## **3.2 Research Contributions**

- also include orev & previc as examples how the platform idea is useful tool (gleiche Plattform, Prinzipien der Entwicklung)

## **3.3 Limitations**

### **3.3.1 The task itself**

- social enough? still on screen, only observing
- idea for first social interaction, where participant guides where agent looks. more like a second-person perspective
- change in framing: gaze understanding to gaze following. why?

Developmental studies often report high number of drop-out rates, for example because of children's fuzziness (todo: example). Administering the TANGO has proven very helpful in this regard: across communities, children tended to complete the task with ease and even asked for a repetition. Drop-out rates could thus be greatly reduced. As the data collection procedures varied slightly from community to community, we cannot report exact numbers of drop-outs but rely on the experience of experimenters.

### 3.3.2 Modeling

- why is gafo still social if vector?
- why model needed? cc signature pattern, other interpretation than raw scores, in future incorporate more complex mental states
- why new? compared to others, only reliance on eyes not head. plus ind diff

### 3.3.3 Cross-cultural work

- do not want to signal cc is easy or we propose parachute science
- technology effect: yes but still ind diff in communities where touchscreen = 1
- more helpful to check relative imprecision, not absolute values within or between communities
- idea for touchscreen familiarization: bubble wrap game
- (Wellman, 2013) how can social cognition be universal and yet vary? framework understanding
- importantly, focusing on signature patterns helps us to study underlying cognitive processes, which raw absolutes score could not have revealed
- don't just compare urban global north to rural global south, which is common pitfall in CC work (Barrett 2020)

## 3.4 Outlook

Applying the tool: Chapter on how Gafo can be extended to answer new research questions.

### 3.4.1 Possible spin-offs of the TANGO

- gafo alien to answer why we need to integrate information from 2 eyes. probably important for side locations in TANGO?
- circle like clock to circumvent precision with distance diffusion
- different starting directions to avoid center bias

### 3.4.2 Open research questions

- but also: all methods sanity checks so that now theory testing can start
- do social interactions shape social cognition abilities? we need good measures on both ends, machine learning as promising avenue
- back to uniquely human: what about great apes? how would they perform in gafo?
- now that we can study vector following, we can also go down this path: is vector estimation needed in other social cognitive abilities? what about for example action prediction, reaching at sth?

### **3.5 Conclusion**



# References

- American Psychological Association. (n.d.-a). Individual Differences. In *APA Dictionary of Psychology*. Retrieved June 26, 2024.
- American Psychological Association. (n.d.-b). Reliability. In *APA Dictionary of Psychology*. Retrieved June 26, 2024.
- American Psychological Association. (n.d.-c). Validity. In *APA Dictionary of Psychology*. Retrieved June 26, 2024.
- American Psychological Association. (n.d.-d). Variability. In *APA Dictionary of Psychology*. Retrieved June 26, 2024.
- American Psychological Association. (n.d.-e). Variance. In *APA Dictionary of Psychology*. Retrieved June 26, 2024.
- American Psychological Association. (n.d.-f). Variation. In *APA Dictionary of Psychology*. Retrieved June 26, 2024.
- Asendorpf, J. B. (1992). Beyond stability: Predicting inter-individual differences in intra-individual change. *European Journal of Personality*, 6(2), 103–117. <https://doi.org/10.1002/per.2410060204>
- Aslin, R. N. (2007). What's in a look? *Developmental Science*, 10(1), 48–53. <https://doi.org/10.1111/j.1467-7687.2007.00563.x>
- Beaudoin, C., Leblanc, É., Gagner, C., & Beauchamp, M. H. (2020). Systematic Review and Inventory of Theory of Mind Measures for Young Children. *Frontiers in Psychology*, 10, 2905. <https://doi.org/10.3389/fpsyg.2019.02905>
- Bowlby, J. (1958). The nature of the child's tie to his mother. *The International Journal of Psychoanalysis*, 39, 350–373.
- Brooks, R., & Meltzoff, A. N. (2005). The development of gaze following and its relation to language. *Developmental Science*, 8(6), 535–543. <https://doi.org/10.1111/j.1467-7687.2005.00445.x>
- Brooks, R., & Meltzoff, A. N. (2015). Connecting the dots from infancy to childhood: A longitudinal study connecting gaze following, language, and explicit theory of mind. *Journal of Experimental Child Psychology*, 130, 67–78. <https://doi.org/10.1016/j.jecp.2014.09.010>
- Csibra, G., & Gergely, G. (2009). Natural pedagogy. *Trends in Cognitive Sciences*, 13(4), 148–153. <https://doi.org/10.1016/j.tics.2009.01.005>
- D'Entremont, B. (2000). A perceptual–attentional explanation of gaze following in 3- and 6-month-olds. *Developmental Science*, 3(3), 302–311. <https://doi.org/10.1111/1467-7687.00124>
- D'Entremont, B., Hains, S. M. J., & Muir, D. W. (1997). A demonstration of gaze following in 3- to 6-month-olds. *Infant Behavior and Development*, 20(4), 569–572. [https://doi.org/10.1016/S0163-6383\(97\)90048-5](https://doi.org/10.1016/S0163-6383(97)90048-5)
- Decety, J. (2020). Preface. In J. Decety (Ed.), *The Social Brain: A Developmental Perspective* (pp. ix–xii). The MIT Press.
- Devine, R. T., & Hughes, C. (2018). Family Correlates of False Belief Understanding in Early Childhood: A Meta-Analysis. *Child Development*, 89(3), 971–987. <https://doi.org/10.1111/cdev.12682>
- Dörrenberg, S., Rakoczy, H., & Liszkowski, U. (2018). How (not) to measure infant Theory of Mind: Testing the replicability and validity of four non-verbal measures. *Cognitive Development*, 46, 12–30. <https://doi.org/10.1016/j.cogdev.2018.01.001>
- Elson, M., Hussey, I., Alsalti, T., & Arslan, R. C. (2023). Psychological measures aren't toothbrushes. *Communications Psychology*, 1(1), 1–4. <https://doi.org/10.1038/s44271-023-00026-9>
- Flavell, J. H., Everett, B. A., Croft, K., & Flavell, E. R. (1981). Young children's knowledge about visual

## References

- perception: Further evidence for the Level 1–Level 2 distinction. *Developmental Psychology*, 17(1), 99–103. <https://doi.org/10.1037/0012-1649.17.1.99>
- Flavell, J. H., Flavell, E. R., Green, F. L., & Wilcox, S. A. (1981). The Development of Three Spatial Perspective-Taking Rules. *Child Development*, 52(1), 356–358. <https://doi.org/10.2307/1129250>
- Glynn, M. A., & Watkiss, L. (2016). Social Cognition. In M. Augier & D. J. Teece (Eds.), *The Palgrave Encyclopedia of Strategic Management* (pp. 1–4). Palgrave Macmillan UK. [https://doi.org/10.1057/978-1-349-94848-2\\_614-1](https://doi.org/10.1057/978-1-349-94848-2_614-1)
- Grossmann, T. (2022). Becoming uniquely human? Comparing chimpanzee to human infancy. *Developmental Science*, 25(1), e13142. <https://doi.org/10.1111/desc.13142>
- Happé, F., Cook, J. L., & Bird, G. (2017). The Structure of Social Cognition: In(ter)dependence of Sociocognitive Processes. *Annual Review of Psychology*, 68(1), 243–267. <https://doi.org/10.1146/annurev-psych-010416-044046>
- Henrich, J. (2015). *The Secret of Our Success: How Culture Is Driving Human Evolution, Domesticating Our Species, and Making Us Smarter*. Princeton University Press. <https://doi.org/10.1515/9781400873296>
- Herrmann, E., Call, J., Hernández-Lloreda, M. V., Hare, B., & Tomasello, M. (2007). Humans have evolved specialized skills of social cognition: The cultural intelligence hypothesis. *Science (New York, N.Y.)*, 317(5843), 1360–1366. <https://doi.org/10.1126/science.1146282>
- Hughes, C., Jaffee, S. R., Happé, F., Taylor, A., Caspi, A., & Moffitt, T. E. (2005). Origins of individual differences in theory of mind: From nature to nurture? *Child Development*, 76(2), 356–370. <https://doi.org/10.1111/j.1467-8624.2005.00850.x>
- Kidd, E., Donnelly, S., & Christiansen, M. H. (2018). Individual Differences in Language Acquisition and Processing. *Trends in Cognitive Sciences*, 22(2), 154–169. <https://doi.org/10.1016/j.tics.2017.11.006>
- Lecce, S., Bianco, F., Devine, R. T., Hughes, C., & Banerjee, R. (2014). Promoting theory of mind during middle childhood: A training program. *Journal of Experimental Child Psychology*, 126, 52–67. <https://doi.org/10.1016/j.jecp.2014.03.002>
- Macdonald, R. G., & Tatler, B. W. (2013). Do as eye say: Gaze cueing and language in a real-world social interaction. *Journal of Vision*, 13(4), 6–6. <https://doi.org/10.1167/13.4.6>
- Moskowitz, G. B. (2013). Social Cognition. *Obo in Psychology*. <https://doi.org/10.1093/obo/9780199828340-0099>
- Okumura, Y., Kanakogi, Y., Kobayashi, T., & Itakura, S. (2017). Individual differences in object-processing explain the relationship between early gaze-following and later language development. *Cognition*, 166, 418–424. <https://doi.org/10.1016/j.cognition.2017.06.005>
- Ostrom, T. M. (1994). Foreword. In S. R. Wyer & T. K. Srull (Eds.), *Handbook of Social Cognition: Applications*. (2nd ed., Vol. 2, pp. vii–xii). Lawrence Erlbaum Associates.
- Perner, J., Ruffman, T., & Leekam, S. R. (1994). Theory of Mind Is Contagious: You Catch It from Your Sibs. *Child Development*, 65(4), 1228–1238. <https://doi.org/10.2307/1131316>
- Peterson. (2000). Kindred spirits: Influences of siblings' perspectives on theory of mind. *Cognitive Development*, 15(4), 435–455. [https://doi.org/10.1016/S0885-2014\(01\)00040-5](https://doi.org/10.1016/S0885-2014(01)00040-5)
- Pontius, J. U., Mullikin, J. C., Smith, D. R., Team, A. S., Lindblad-Toh, K., Gnerre, S., Clamp, M., Chang, J., Stephens, R., Neelam, B., Volfovsky, N., Schäffer, A. A., Agarwala, R., Narfström, K., Murphy, W. J., Giger, U., Roca, A. L., Antunes, A., Menotti-Raymond, M., ... O'Brien, S. J. (2007). Initial sequence and comparative analysis of the cat genome. *Genome Research*, 17(11), 1675–1689. <https://doi.org/10.1101/gr.6380007>
- Prein, J. C., Kalinke, S., Haun, D. B. M., & Bohn, M. (2023). TANGO: A reliable, open-source, browser-based task to assess individual differences in gaze understanding in 3 to 5-year-old children and adults. *Behavior Research Methods*, 56(3), 2469–2485. <https://doi.org/10.3758/s13428-023-02159-5>

- Quesque, F., Apperly, I., Baillargeon, R., Baron-Cohen, S., Becchio, C., Bekkering, H., Bernstein, D., Bertoux, M., Bird, G., Bukowski, H., Burgmer, P., Carruthers, P., Catmur, C., Dziobek, I., Epley, N., Erle, T. M., Frith, C., Frith, U., Galang, C. M., ... Brass, M. (2024). Defining key concepts for mental state attribution. *Communications Psychology*, 2(1), 1–5. <https://doi.org/10.1038/s44271-024-00077-6>
- Quesque, F., & Rossetti, Y. (2020). What Do Theory-of-Mind Tasks Actually Measure? Theory and Practice. *Perspectives on Psychological Science*, 15(2), 384–396. <https://doi.org/10.1177/1745691619896607>
- Rakoczy, H. (2022). Foundations of theory of mind and its development in early childhood. *Nature Reviews Psychology*, 1(4), 223–235. <https://doi.org/10.1038/s44159-022-00037-z>
- Ronald, A., Viding, E., Happé, F., & Plomin, R. (2006). Individual differences in theory of mind ability in middle childhood and links with verbal ability and autistic traits: A twin study. *Social Neuroscience*, 1(3-4), 412–425. <https://doi.org/10.1080/17470910601068088>
- Ruffman, T., Puri, A., Galloway, O., Su, J., & Taumoepeau, M. (2018). Variety in parental use of “want” relates to subsequent growth in children’s theory of mind. *Developmental Psychology*, 54(4), 677. <https://doi.org/10.1037/dev0000459>
- Sodian, B., Licata, M., Kristen-Antonow, S., Paulus, M., Killen, M., & Woodward, A. (2016). Understanding of Goals, Beliefs, and Desires Predicts Morally Relevant Theory of Mind: A Longitudinal Investigation. *Child Development*, 87(4), 1221–1232. <https://doi.org/10.1111/cdev.12533>
- Stewart, R. B. (1983). Sibling Interaction: The Role of the Older Child as Teacher for the Younger. *Merrill-Palmer Quarterly*, 29(1), 47–68. <https://www.jstor.org/stable/23086191>
- Tomasello, M. (2019). *Becoming Human: A Theory of Ontogeny*. Harvard University Press. <https://doi.org/10.4159/9780674988651>
- Tomasello, M. (2020). The adaptive origins of uniquely human sociality. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 375(1803), 20190493. <https://doi.org/10.1098/rstb.2019.0493>
- Vygotsky, L. S. (1962). *Thought and language* (E. Hanfmann & G. Vakar, Eds.; pp. xxi, 168). MIT Press. <https://doi.org/10.1037/11193-000>
- Vygotsky, L. S. (1978). *Mind in Society: Development of Higher Psychological Processes* (M. Cole, V. Jolm-Steiner, S. Scribner, & E. Souberman, Eds.). Harvard University Press. <https://doi.org/10.2307/j.ctvjf9vz4>
- Vygotsky, L. S. (1987). Thinking and speech. In R. W. Rieber & A. S. Carton (Eds.), *The collected works of L.S. Vygotsky, Volume 1: Problems of general psychology* (Vol. 11, pp. 39–285). Plenum Press. (Original work published 1934).
- Waterson, R. H., Lander, E. S., Wilson, R. K., & The Chimpanzee Sequencing and Analysis Consortium. (2005). Initial sequence of the chimpanzee genome and comparison with the human genome. *Nature*, 437(7055), 69–87. <https://doi.org/10.1038/nature04072>
- Wellman, H. M. (2013). Universal Social Cognition: Childhood Theory of Mind. In M. R. Banaji & S. A. Gelman (Eds.), *Navigating the Social World: What Infants, Children, and Other Species Can Teach Us* (p. o). Oxford University Press. <https://doi.org/10.1093/acprof:oso/9780199890712.003.0014>
- Wellman, H. M., Cross, D., & Watson, J. (2001). Meta-analysis of theory-of-mind development: The truth about false belief. *Child Development*, 72, 655–684. <https://doi.org/10.1111/1467-8624.00304>
- Zhang, Z., Yu, H., Long, M., & Li, H. (2021). Worse Theory of Mind in Only-Children Compared to Children With Siblings and Its Intervention. *Frontiers in Psychology*, 12, 5073. <https://doi.org/10.3389/fpsyg.2021.754168>



# Appendix A — Main Publications

This dissertation includes four main publications that were either published (Study I) or under review (Study II, Study III, Study IV) at the time of the dissertation submission. The full texts of these publications are provided below. For the accepted manuscript, the published version is provided. For manuscripts under review, the submitted versions are provided which are published online as pre-prints.

**Study I:** Prein, J. C., Kalinke, S., Haun, D. B. M.\* & Bohn, M.\* (2023). TANGO: A reliable, open-source, browser-based task to assess individual differences in gaze understanding in 3 to 5-year-old children and adults. *Behavior Research Methods*, 56(3), 2469–2485. <https://doi.org/10.3758/s13428-023-02159-5>

**Study II:** Prein, J. C., Maurits, L., Werwach, A., Haun, D. B. M.\* & Bohn, M.\* (2024). *Variation in gaze understanding across the life span: A process-level perspective*. PsyArXiv. <https://doi.org/10.31234/osf.io/dy73a>

**Study III:** Bohn, M.\*, Prein, J. C.\*, Ayikoru, A., Bednarski, F. M., Dzabatou, A., Frank, M. C., Henderson, A. M. E., Isabella, J., Kalbitz, J., Kanngiesser, P., Keşsafoğlu, D., Koymen, B., Manrique-Hernandez, M., Magazi, S., Mújica-Manrique, L., Ohlendorf, J., Olaoba, D., Pieters, W., Pope-Caldwell, S., ... Haun, D. (2024). *A universal of human social cognition: Children from 17 communities process gaze in similar ways*. PsyArXiv. <https://doi.org/10.31234/osf.io/z3ahv>

**Study IV:** Prein, J.C., Bednarski, F. M., Dzabatou, A., Frank, M. C., Henderson, A. M. E., Kalbitz, J., Kanngiesser, P., Keşsafoğlu, D., Köymen, B., Manrique-Hernandez, M. V., Magazi, S., Mújica-Manrique, L., Ohlendorf, J., Olaoba, D., Pieters, W. R., Pope-Caldwell, S., Sen, U., Slocombe, K., Sparks, R. Z., Stengelin, R., Sunderarajan, J., Sutherland, K., Tusiime, F., Vieira, W., Zhang, Z., Zong, Y., Haun, D. B. M.\* Bohn, M.\* (2024). *Measuring variation in gaze following across communities, ages, and individuals – a showcase of the TANGO-CC*. PsyArxiv. <https://doi.org/10.31234/osf.io/fcq2g>

## Study I

Behavior Research Methods  
<https://doi.org/10.3758/s13428-023-02159-5>



### TANGO: A reliable, open-source, browser-based task to assess individual differences in gaze understanding in 3 to 5-year-old children and adults

Julia Christin Prein<sup>1</sup> · Steven Kalinke<sup>1</sup> · Daniel B. M. Haun<sup>1</sup> · Manuel Bohn<sup>1,2</sup>

Accepted: 2 June 2023  
© The Author(s) 2023

#### Abstract

Traditional measures of social cognition used in developmental research often lack satisfactory psychometric properties and are not designed to capture variation *between* individuals. Here, we present the TANGO (Task for Assessing Individual differences in Gaze understanding-Open); a brief (approx. 5–10min), reliable, open-source task to quantify individual differences in the understanding of gaze cues. Localizing the attentional focus of an agent is crucial in inferring their mental states, building common ground, and thus, supporting cooperation. Our interactive browser-based task works across devices and enables in-person and remote testing. The implemented spatial layout allows for discrete and continuous measures of participants' click imprecision and is easily adaptable to different study requirements. Our task measures inter-individual differences in a child ( $N = 387$ ) and an adult ( $N = 236$ ) sample. Our two study versions and data collection modes yield comparable results that show substantial developmental gains: the older children are, the more accurately they locate the target. High internal consistency and test-retest reliability estimates underline that the captured variation is systematic. Associations with social-environmental factors and language skills speak to the validity of the task. This work shows a promising way forward in studying individual differences in social cognition and will help us explore the structure and development of our core social-cognitive processes in greater detail.

**Keywords** Social cognition · Individual differences · Gaze cues · Cognitive development

#### Introduction

Social cognition—representing and reasoning about an agent's perspectives, knowledge states, intentions, beliefs, and preferences to explain and predict their behavior—is among the most-studied phenomena in developmental research. In recent decades, much progress has been made in determining the average age at which a specific social-cognitive ability emerges in development (Gopnik & Slaughter, 1991; Peterson et al., 2012; Rakoczy, 2022; Wellman et al.,

Haun, Daniel B. M. and Bohn, Manuel contributed equally to this work.

Julia Christin Prein  
julia\_prein@eva.mpg.de

<sup>1</sup> Department of Comparative Cultural Psychology, Max Planck Institute for Evolutionary Anthropology, Deutscher Platz 6, 04103 Leipzig, Germany

<sup>2</sup> Institute of Psychology, Leuphana University Lüneburg, Lüneburg, Germany

2001; Wellman & Liu, 2004). Yet, there are always individual differences. Identifying variability in social-cognitive abilities and factors influencing their development is vital in theory building (e.g., to test causal predictions) and designing interventions (Happé et al., 2017; Kidd et al., 2018; Lecce et al., 2014; Mundy et al., 2007; Underwood, 1975).

Numerous studies have already examined individual differences in social cognition (for an overview, see Hughes & Devine, 2015; Slaughter, 2015). The most common, recurring research questions are concerned with the developmental sequence of social-cognitive abilities (e.g., Wellman & Liu, 2004), and which factors drive the development of social cognition (Devine & Hughes, 2018; Gola, 2012). For example, Okumura and colleagues asked how early gaze-following and object processing relate to later language development (Okumura et al., 2017). In general, individual differences studies often focus on the relationship between social-cognitive abilities and: (1) family influences, (2) other cognitive constructs, and (3) social behavioral outcomes (for an overview, see Slaughter and Repacholi, 2003). Studies

on social-cognitive abilities and family influences include the effect of parenting practices (for a review, see Pavarini et al., 2013), attachment quality (e.g., Astor et al., 2020), mental state talk (Gola, 2012; Hughes et al., 2011; Lecce et al., 2014), and family background as parental education, occupation, sibling interaction and childcare (Bulgarelli & Molina, 2016; Cutting & Dunn, 1999; Dunn et al., 1991). Another group of individual differences studies focuses on the interplay of social and physical cognition (Herrmann et al., 2010), executive functions (Benson et al., 2013; Buttelmann et al., 2022; Carlson & Moses, 2001; Carlson et al., 2004; Hughes & Ensor, 2007), and language abilities (McEwen et al., 2007; Milligan et al., 2007; Okumura et al., 2017). Studies on social behavioral outcomes measured the interplay of social cognition and prosociality (for a review, see Imuta et al., 2016; Walker, 2005), stereotypes, resource allocations (Rizzo & Killen, 2018), and moral intentions (Sodian et al., 2016).

However, developmental psychologists are frequently surprised to find minor or no association between measures of social cognition that are thought to be theoretically related – cross-sectionally and/or longitudinally (e.g., Poulin-Dubois et al., 2023; Sodian, 2023; Sodian et al., 2016). This might be because traditional measures of social cognition are not designed to capture variation *between* children: they often rely on low trial numbers, small sample sizes, and dichotomous measures. A recent review showed that many studies on social cognition measures failed to report relevant psychometric properties at all (Beaudoin et al., 2020) or – when they did – showed mixed results on test–retest reliability (Hughes et al., 2000; Mayes et al., 1996).

To give an example: the most commonly applied prototypical measure for social cognition is the change-of-location false belief task (Baron-Cohen et al., 1985; Wimmer & Perner, 1983). Here, children watch a short sequence of events (often acted out or narrated by the experimenters). A doll called Sally puts her marble into a basket. After Sally leaves the scene, a second doll named Anne takes the marble and moves it into a box. Participants then get asked where Sally will look for her marble once she returns. The outcome measures false belief understanding in a dichotomous way: children pass the task if they take the protagonist's epistemic state into account and answer that she will look into the basket. Many years of research utilizing these verbal change-of-location tasks suggest that children develop belief-representing abilities at four to five years of age (for a review, see Wellman et al., 2001). Several cross-cultural studies supported this evidence (Barrett et al., 2013; Callaghan et al., 2005; cf. Mayer & Träuble, 2015).

However, from this age onwards, the change-of-location task shows ceiling effects and has very limited diagnostic value (Repacholi, 2003). Thus, this task seems well suited

to track a particular group-level developmental transition, yet it fails to capture individual differences (cf. “reliability paradox,” Hedge et al., 2018). As Wellman (2012) put it, “it’s really only passing/failing one sort of understanding averaged across age” (p. 317). This has profound implications for what studies on individual differences using this task (or others) can show. Poor measurement of social cognition on an individual level is likely to conceal relations between different aspects of cognition and may obscure developmental change. For example, Sodian et al. (2016) neither found a correlation between two moral Theory of Mind False Belief and Intention tasks at 60 months, nor a relationship between these two factors and implicit False Belief understanding at 18 months.

The “Sandbox task” is one of the few tasks that attempt to overcome these methodological challenges (Begeer et al., 2012; Bernstein et al., 2011; Coburn et al., 2015; Mahy et al., 2017; Sommerville et al., 2013). This continuous FB task measures the degree to which the estimate of another’s belief is biased by one’s own knowledge. Recent work questions the interpretation of this measure (Samuel et al., 2018a, b): it is unclear whether a smaller egocentric bias can be directly translated into a better mental state reasoning ability. Another evaluation criterion should, therefore, be whether a task captures meaningful variability in performance; that is, differences in test scores should correspond to differences in the social-cognitive ability in question.

Thus, developmental psychology faces a dilemma: many research questions rely on measuring individuals’ development, yet, there is a lack of tasks to measure these individual differences reliably. To capture the emergence of social-cognitive abilities and their relation to social factors in greater precision and detail, we must consequently address the methodological limitations of existing study designs (Hughes et al., 2011; Hughes & Leekam, 2004).

Schaafsma et al., (2015) compiled a “wish list” for new social-cognitive paradigms. They advocated for parametric – instead of dichotomous – measures covering proficiency as a range, avoiding floor and ceiling effects, and showing satisfactory test–retest reliability estimates (see also Beaudoin et al., 2020; Hughes & Devine, 2015). New tasks should capture variation across age groups, including older children and adults (Repacholi and Slaughter, 2003). Another goal in creating new tasks should be to focus on the “face value”: measures should probe the underlying social-cognitive ability as straight-forward and directly as possible. Keeping task demands minimal is also beneficial for using the paradigm in a variety of different cultural, clinical, and demographic contexts (Molleman et al., 2019). The task should serve as a proxy for behavior as it appears in the real world and should be validated in relation to real-world experiences (Repacholi and Slaughter, 2003).

### A new measure of gaze understanding

Our goal was to design a new measure of social cognition that captures individual differences across age groups in a systematic, reliable, and valid way. We focused on a fundamental ability implicated in many social-cognitive reasoning processes: gaze understanding – the ability to locate and use the attentional focus of an agent. The first component of this ability is often termed gaze following – turning one’s eyes in the same direction as the gaze of another agent – and has been studied intensively (Astor et al., 2021; Byers-Heinlein et al., 2021; Coelho et al., 2006; Del Bianco et al., 2019; Frischen et al., 2007; Hernik & Broesch, 2019; Itakura & Tanaka, 1998; Lee et al., 1998; Moore, 2008; Shepherd, 2010; Tomasello et al., 2007). In our definition, gaze understanding goes one step further by including the *acting on the gaze-cued location* – therefore, using the available social information to guide one’s behavior as needed in real-life conditions.

Following an agent’s gaze provides insights into their intentions, thoughts, and feelings by acting as a “front end ability” (Brooks & Meltzoff, 2005, p. 535). Gaze is integral for many more sophisticated social-cognitive abilities, for example, inferences about knowledge states. As such, the eyes have been regarded as a “window into the mind” (Shepherd, 2010). Monitoring another’s attention also supports building a common ground, which is important for action coordination and cooperative social interactions (Bohn & Köymen, 2018; Tomasello et al., 2007). In addition, gaze and language development seem to be related (Brooks & Meltzoff, 2005). Gaze facilitates word learning by helping to identify the referent of a new word and has been regarded as a crucial signal of nonverbal communication (Hernik & Broesch, 2019; Macdonald & Tatler, 2013).

While the emergence of gaze following has been well established, less is known about the developmental trajectory throughout childhood and adolescence. One possibility is that our social-cognitive ability in question is fully developed once emerged in infancy. However, many cognitive abilities continue to develop beyond early childhood (e.g., Gathercole et al., 2004 for working memory; Raviv & Arnon, 2018 for visual statistical learning). Therefore, children could potentially improve in understanding gaze, fine-tuning the performance of the already existing skill. Consequently, we aimed to assess the differentiation of the ability to understand gaze. Our goal was *not* to establish the youngest age at which children understand gaze cues. Rather, we wanted to examine how that ability changes with age. To accurately measure developmental change, we were interested in capturing individual variability.

To address the psychometric shortcoming of earlier work, we implemented the following design features: First,

we used a continuous measure which allowed us to capture fine-grained individual differences at different ages. Second, we designed short trials that facilitate more than a dozen replicates per subject. The result is more precise individual-level estimates. Third, we systematically investigated the psychometric properties of the new task.

Designing this task required a new testing infrastructure. We designed the task as an interactive web application. Previous research has successfully used online study implementations that compare well to in-person data collection (Bohn et al., 2021a, b; Frank et al., 2016). This greatly increased the flexibility with which we could modify the stimuli on a trial-by-trial basis. Furthermore, because the task is largely self-contained, it is much more controlled and standardized. Most importantly, it makes the task portable: testing is possible in-person using tablets but also remotely via the internet (no installation needed). As such, it provides a solid basis to study individual differences in gaze understanding across ages at scale. We make the task and its source code openly accessible for other researchers to use and modify.

### Task design

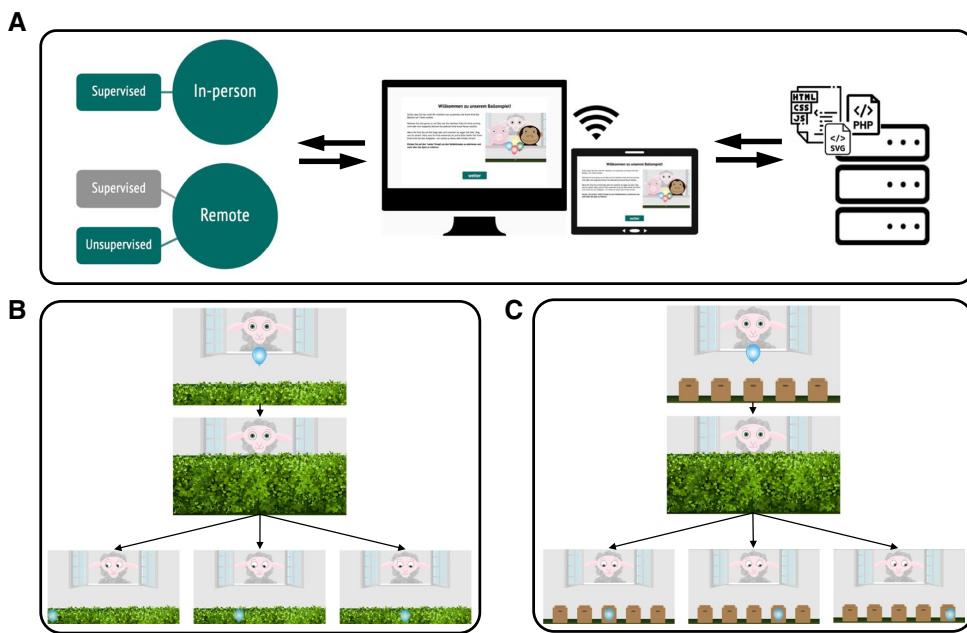
#### Implementation

The code is open-source (<https://github.com/ccp-eva/tango-demo>), and a live demo version can be found under: <https://ccp-odc.eva.mpg.de/tango-demo/>.

The web app was developed using JavaScript, HTML5, CSS, and PHP. For stimulus presentation, a scalable vector graphic (SVG) composition was parsed. This way, the composition scales according to the user’s viewport without loss of quality while keeping the aspect ratio and relative object positions constant. Furthermore, SVGs allow us to define all composite parts of the scene (e.g., pupil of the agent) individually. This is needed for precisely calculating the exact pupil and target locations and sizes. Additionally, it makes it easy to adjust the stimuli and, for example, add another agent to the scene. The web app generates two file types: (1) a text file (.json) containing metadata, trial specifications, and participants’ click responses, and (2) a video file (.webm) of the participant’s webcam recording. These files can either be sent to a server or downloaded to the local device. Personalized links can be created by passing on URL parameters.

#### Stimuli

Our newly implemented task asks children and adults to search for a balloon. The events proceed as follows (see Fig. 1B and C). An animated agent (a sheep, monkey, or



**Fig. 1** Study setup. **A** Infrastructure for online testing. (i) Subjects aged 3 to 99+ can participate. Data collection can take place anywhere: online, in kindergartens, or in research labs. (ii) The task is presented as a website that works across devices. (iii) The scripts for the website and the recorded data are stored on secure in-house servers. **B** Hedge version (continuous) of the TANGO. (i) The agent stands in a window with the target in front of them. (ii) A hedge grows and covers the target. (iii) The target falls to a random loca-

tion on the ground. The agent's eyes track the movement of the target. Three exemplary target locations are shown to depict how indicative the agent's gaze cues are in determining the target's location. The transparent target is only shown for an illustrative purpose (not visible during the test). **C** Box version (discrete) of the TANGO. Number of boxes (min. 1; max. 8) as potential hiding locations can be set according to the researcher's need

pig) looks out of a window of a house. A balloon (i.e., target; blue, green, yellow, or red) is located in front of them. The target then falls to the ground. At all times, the agent's gaze tracks the movement of the target: the pupils and iris move so that their center aligns with the center of the target. While the distance of the target's flight depends on the final location, the target moves at a constant speed. Participants are then asked to locate the target: they respond by touching or clicking on the screen. Visual access to the target's true location is manipulated by a hedge. Participants either have full, partial, or no visual access to the true target location. When partial or no information about the target location is accessible, participants are expected to use the agent's gaze as a cue.

To keep participants engaged and interested, the presentation of events is accompanied by cartoon-like effects. Each trial starts with an attention-getter: an eye-blinking sound plays while the pupils and iris of the agent enlarge (increase to 130%) and change in opacity (decrease to 75%) for 0.3 s. The landing of the target is accompanied by a tapping sound. Once the target landed, the instructor's voice asked "Where is the balloon?". To confirm the participant's click, a short

plop sound plays, and a small orange circle appears at the location of choice. Participants do not receive differential feedback so that learning effects are reduced, and trials stay comparable across the sample. If no response is registered within 5 s after the target landed, an audio prompt reminds the participant to respond.

### Trials

Trials differ in the amount of visual access that participants have to the final target position. Before the test trials start, participants complete four training trials during which they familiarize themselves with touching the screen. In the first training trial, participants have full visual access to the target flight and the target's end location and are simply asked to click on the visible balloon. In the second and third training trials, participants have partial access: they witness the target flight but cannot see the target's end location. They are then asked to click on the hidden balloon, i.e., the location where they saw the target land. In test trials, participants have no visual access to the target flight or the end location. Participants are expected to use the agent's gaze as a cue

to locate the target. The first trial of each type comprises a voice-over description of the presented events. The audio descriptions explicitly state that the agent is always looking at the target (see Supplements for audio script). After the four training trials, participants receive 15 test trials. The complete sequence of four training trials and 15 test trials can be easily completed within 5–10 min.

### Study versions

We designed two study versions that differ in the target's final hiding place and, consequently, in the outcome measure: a *hedge version* (continuous) and a *box version* (discrete). Both versions use the same first training trial and then differ in the consecutive training and test trials. In the hedge version, participants have to indicate their estimated target location directly on a hedge. Here, the dependent variable is imprecision, which is defined as the absolute difference between the target center and the *x* coordinate of the participant's click. In the box version, the target lands in a box, and participants are asked to click on the box that hides the target. Researchers can choose how many boxes are shown: one up to eight boxes can be displayed as potential hiding locations. Here, we use a categorical outcome (i.e., which box was clicked) to calculate the proportion of correct responses. Note that in the test trials of both versions, the target flight is covered by a hedge. In the hedge version, the hedge then shrinks to a minimum height required to cover the target's end location. In the box version, the hedge shrinks completely. The boxes then hide the target's final destination (see Fig. 1B and C).

### Randomization

All agents and target colors appear equally often and are not repeated in more than two consecutive trials. The randomization of the target end location depends on the study version. In the hedge version, the full width of the screen is divided into ten bins. Exact coordinates within each bin are then randomly generated. In the box version, the target randomly lands in one of the boxes. As with agent and color choice, each bin/box occurs equally often and can only occur twice in a row.

### Individual differences

Our first aim was to assess whether the TANGO captures inter-individual variation in a child and adult sample. Furthermore, we were interested in whether and how the data collection mode (in-person vs. remote) influences responses. Since we expected a greater difference in responses between

the two data collection modes for children, the analysis of data collection mode was restricted to a child sample.

Task design, data collection, and sample sizes were pre-registered: <https://osf.io/snju6> (child sample) and <https://osf.io/r3bhn> (adult sample). The analyses reported here were not pre-registered but followed the structure of the ones specified in the above pre-registrations (see Footnotes for deviations). The additional analyses mentioned in the pre-registrations (e.g., computational model) address separate research questions (e.g., process-level account of gaze understanding) and will be reported elsewhere. In this paper, we focus on the methodological and psychometric aspects of our task.

The study design and procedure obtained ethical clearance by the MPG Ethics commission Munich, Germany, falling under a packaged ethics application (Appl. No. 2021\_45), and was approved by an internal ethics committee at the Max Planck Institute for Evolutionary Anthropology. The research adheres to the legal requirements of psychological research with children in Germany.

Participants were equally distributed across the two study versions. Data were collected between May and October 2021.

### Participants

We collected data from an in-person child sample, a remote child sample, and a remote adult sample. In-person testing with children took place in kindergartens in Leipzig, Germany. The in-person child sample consisted of 120 children, including 40 3-year-olds (mean = 41.45 months, SD = 3.85, range = 36–47, 22 girls), 40 4-year-olds (mean = 54.60 months, SD = 3.10, range = 48–59, 19 girls), and 40 5-year-olds (mean = 66.95 months, SD = 3.39, range = 60–71, 22 girls).

We pre-registered the replacement for participants that finished fewer than four test trials. This was not the case for any participant. One child stopped participation after 12 test trials but was included in the sample due to the pre-registered replacement rule. Two additional participants were recruited but not included in the study because the participant did not feel comfortable interacting with the tablet alone ( $n = 1$ ), or due to an originally miscalculated age of the child ( $n = 1$ ).

For our remote child sample, we recruited families via an internal database of children living in Leipzig, Germany, whose parents volunteered to participate in child development studies and who indicated an interest in online studies. Families received an email with a short study description and a personalized link. If they had not participated in the study within 2 weeks, they received a reminder via e-mail. The response rate to invitations after the reminder was ~50%.

The remote child sample included 147 children, including 45 3-year-olds (mean = 42.62 months, SD = 3.35, range = 36–47, 14 girls), 47 4-year-olds (mean = 52.64 months, SD = 3.40, range = 48–59, 25 girls), and 55 5-year-olds (mean = 65.11 months, SD = 3.77, range = 60–71, 27 girls). Of these, three families participated twice. In these cases, we only kept the data sets from the first participation.

Four additional participants were recruited but not included in the study because they were already part of the in-person kindergarten sample ( $n = 3$ ), or because of unknown age ( $n = 1$ ).

Please note that we did not collect participant-specific demographics. In the following, we aim to provide context and generalizations based on the broader community and the larger pool of potential participants. Children in our sample grow up in an industrialized, urban Central-European context in a city with approximately 600,000 inhabitants. They often live in nuclear two-generational families with few household members. Information on socioeconomic status was not formally recorded, although the majority of families come from mixed, mainly mid to high socioeconomic backgrounds with high levels of parental education. The median individual monthly net income in the year 2021 was ~ 1,600€ for the city of Leipzig.

Adults were recruited via *Prolific* (Palan & Schitter, 2018). *Prolific* is an online participant recruitment service with a predominantly European and US-American subject pool. One hundred English speakers with an average age of 31.34 years (SD = 10.77, range = 18–63, 64 females) were included. Participants live in a variety of different countries: the UK, Italy, Spain, Poland, Netherlands, Canada, Australia, Ireland, South Africa, Norway, Portugal, France, Austria, Finland, Greece, Germany, the U.S., Mexico, Chile, Iceland, New Zealand, Czech Republic, Hungary, Latvia, and Switzerland. In this sample, most participants resided in the United Kingdom ( $n = 47$ ), South Africa ( $n = 8$ ), and Portugal ( $n = 6$ ). Additional detailed information can be found in the data set online. For completing the study, subjects were paid above the fixed minimum wage (on average £10.00 per hour; see Supplements for further detail).

## Procedure

Children in our in-person sample were tested on a tablet in a quiet room in their kindergarten. An experimenter guided the child through the study.

Children in the remote sample received a personalized link to the study website, and families could participate at any time or location. At the beginning of the online study, families were invited to enter our “virtual institute”. We welcomed them with a short introductory video of the study leader, describing the research background and

further procedure. Then, caregivers were informed about data security and were asked for their informed consent. They were asked to enable the sound and seat their child centrally in front of their device. Before the study started, families were instructed on how to set up their webcam and enable the recording permissions. We stressed that caregivers should not help their children. Study participation was video recorded whenever possible in order to ensure that the children themselves generated the answers. Depending on the participant’s device, the website automatically presented the hedge or box version of the study. For families that used a tablet with a touchscreen, the hedge version was shown. Here, children could directly click on the touchscreen to indicate where the target is. For families that used a computer without a touchscreen, the website presented the box version of the task. We assumed that younger children in our sample would not be acquainted with using a computer mouse. Therefore, we asked children to point to the screen, while caregivers were asked to act as the “digital finger” of their children and click on the indicated box.

All participants received 15 test trials. In the box version, we decided to adjust the task difficulty according to the sample: children were presented with five boxes, while adults were presented with eight boxes as possible target locations.

## Analysis

All test trials without voice-over descriptions were included in our analyses. We ran all analyses in R version 4.3.0 (2023-04-21) (R Core Team, 2022). Regression models were fitted as Bayesian generalized linear mixed models (GLMMs) with default priors for all analyses, using the function `brm` from the package `brms` (Bürkner, 2017, 2018).

To estimate the developmental trajectory of gaze understanding and the effect of data collection mode, we fit a GLMM predicting the task performance in each trial by age (in months, z-transformed) and data collection mode (reference category: in-person supervised). The model included random intercepts for each participant and symmetric target position, and a random slope for symmetric target position within participants (model notation in R: `performance ~ age + datacollection + symmetricPosition + trialNr + (1 + symmetricPosition + trialNr | subjID)`).<sup>1</sup>

<sup>1</sup> In the pre-registration (<https://osf.io/snju6>), we specified the following model structure: “All models will include a fixed effect of target centrality and age, a random intercept for ID and a random slope for trial number by ID. For both study versions, we will compare the above specified null models with a model including data source (live vs. online) as a fixed effect.” Or, in R model formula: `R: performance ~ target_centrality + age + (1 | ID) + (1 + trial | ID)`. In this paper, we added `symmetricPosition` (synonymous to `target_centrality`) as a random slope because we expected that this item effect could vary between participants. To be better able to interpret trial

Here, symmetricPosition refers to the absolute distance from the stimulus center (i.e., smaller value meaning more central target position). We expected that trials could differ in their difficulty depending on the target centrality and that these item effects could vary between participants.

For the hedge version, performance was defined as the absolute click distance between the target center and the click x coordinate, scaled according to target widths, and modeled by a lognormal distribution. For the box version, the model predicted correct responses (0/1) using a Bernoulli distribution with a logit link function. We inspected the posterior distribution (mean and 95% credible interval (CrI)) for the age and data collection estimates.

## Results

Children showed nearly perfect precision in the first training trial. As visual access to the target location decreased in the subsequent training trials, imprecision levels increased (see Supplements). Within test trials, children's imprecision levels did not vary as a function of trial number. We take this as evidence that (A) children were comfortable touching the screen, (B) children understood the task instructions insofar as they aimed at locating the target, and (C) our experimental design successfully manipulated task difficulty.

We found a strong developmental effect: with increasing age, participants got more accurate in locating the target. In the hedge version, children's click imprecision decreased with age, while in the box version, the proportion of correct responses increased (see Fig. 2A and F). Most participants in the box version performed above chance level. By the end of their sixth year of life, children came close to the adult's proficiency level. Most importantly, however, we found substantial inter-individual variation across study versions and age groups. For example, some 3-year-olds were more precise in their responses than some 5-year-olds. Even though variation is smaller, we could even find inter-individual differences in the adult sample.

As Fig. 2A and F show, our remotely collected child data resembled the data from the kindergarten sample. We found evidence that responses of children participating remotely were slightly more precise. This difference was mainly driven by the younger participants and was especially prominent in the box version of the task. It is conceivable that caregivers were especially prone to influence the behavior of younger children. In the box version, caregivers might

Footnote 1 (continued)

number effects, we decided to include it as a fixed effect. Data collection mode (formerly named “data source”) proved as a meaningful predictor and was accordingly added to the model.

have had more opportunities to interfere since they carried out the clicking for their children.<sup>2</sup>

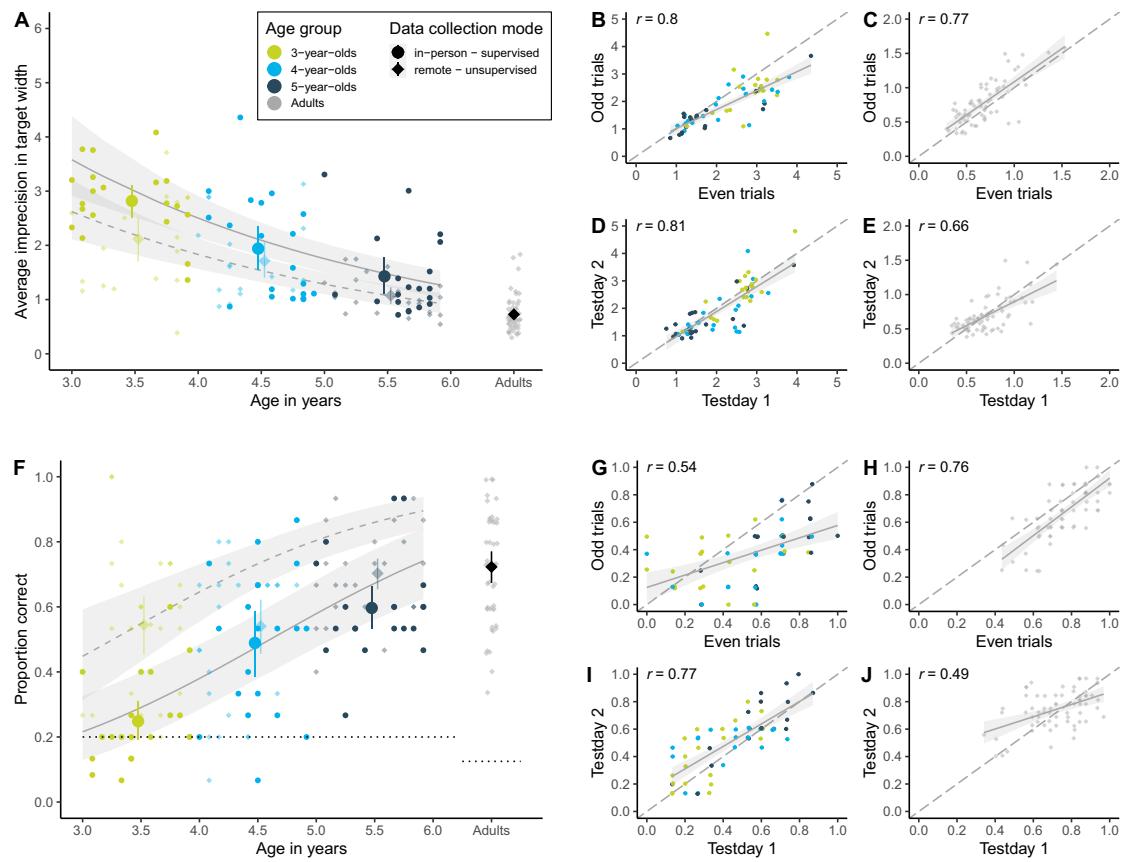
Our GLMM analysis corroborated the visual inspection of the data: in the hedge version, the estimates for age ( $\beta = -0.32$ ; 95% CrI [-0.41; -0.23]) and data collection mode  $-0.31$  (95% CrI [-0.48; -0.14]) were negative and reliably different from zero. In the box version, the estimate of age ( $\beta = 0.68$  (95% CrI [0.44; 0.93])) and the estimate of data collection mode ( $\beta = 1.10$  (95% CrI [0.66; 1.56])) were positive and reliably different from zero. Note that even though confidence intervals from the data collection estimates were wide, the effect was positive and reliably different from zero in that our remote sample performed more accurately than our in-person sample.

There was no effect of trial number (hedge version:  $\beta = 0.00$ ; 95% CrI [-0.02; 0.01]; box version:  $\beta = -0.02$ ; 95% CrI [-0.05; 0.01]). However, trials differed in difficulty depending on where the target landed (hedge version:  $\beta = 0.47$ ; 95% CrI [0.40; 0.54]; box version:  $\beta = -1.59$ ; 95% CrI [-1.88; -1.31]). When the target landed closer to the center of the screen, participants were more accurate in locating it.

## Discussion

Our task measured inter-individual differences in both children and adults; that is, we found substantial variation in individuals across age groups. For example, some 3-year-olds showed greater precision levels than some 5-year-olds. This holds across both study versions. However, due to the continuous study design, the hedge version was able to capture more fine-grained differences in individual performance. We see substantial developmental gains: with increasing age, participants became on average more and more precise in locating the target. The 5-year-olds reached a proficiency level close to the adults' level. For neither study version nor age group did we find any floor or ceiling effects. The presentation as a web app with cartoon-like features kept children interested and motivated throughout the 15 test trials. Furthermore, we found a comparable developmental trajectory for an unsupervised remote child sample. This illustrates the flexibility of the task design.

<sup>2</sup> In an exploratory analysis, we coded parental behavior and environmental factors during remote unsupervised testing. We focused on the subsample with the greatest performance difference between data collection modes: the 3-year-olds in the box version of the task ( $n = 16$ ). We reasoned that if parental interference cannot explain the greatest performance difference in our sample, the effects would be negligible in the remaining sample. Based on our model comparison, we conclude that there is no clear evidence of a stable effect of parental interference. See Supplements for further detail.



**Fig. 2** Measuring inter-individual variation. **A** Developmental trajectory in the continuous hedge version. Performance is measured as imprecision, i.e., the absolute distance between the target's center and the participant's click (averaged across trials). The unit of imprecision is counted in the width of the target, i.e., a participant with imprecision of 1 clicked on average one target width to the left or right of the true target center. **B** Internal consistency (odd-even split) in hedge child sample. **C** Internal consistency in hedge adult sample. **D** Test-retest reliability in hedge child sample. **E** Test-retest reliability in hedge adult sample. **F** Developmental trajectory in the discrete box version. Performance is measured as the proportion of correct responses, i.e., how many times the participant clicked on the box that contained the target. The dotted black line shows the level of performance expected by chance (for child sample 20%, i.e., one out of five boxes; for adult sample 12.5%, i.e., one out of eight boxes). **G** Internal consistency (odd-even split) in box child sample. **H** Internal

consistency in box adult sample. **I** Test-retest reliability in box child sample. **J** Test-retest reliability in box adult sample. For (A) and (F), regression lines show the predicted developmental trajectories (with 95% Crl) based on GLMMs, with the *line type* indicating the data collection mode. Large points with 95% CI (based on non-parametric bootstrap) represent performance means by age group (binned by year). Small points show the mean performance for each subject averaged across trials. For adult data in (A) and (F), we added minimal horizontal and vertical noise to avoid overplotting. The *shape of data points* represents data collection mode: opaque circles for in-person supervised data collection and translucent diamonds for remote unsupervised data collection. The *color of data points* denotes age group. For (B–E) and (G–J), regression lines with 95% CI show smooth conditional mean based on a linear model (generalized linear model for box version), with Pearson's correlation coefficient  $r$

## Internal consistency and test-retest reliability

As a next step, we aimed to investigate whether the variation that we captured with the TANGO is reliable. We assessed internal consistency (as split-half reliability) and test-retest reliability. Task procedure, data collection, and sample sizes

were pre-registered (<https://osf.io/xqm73> for the child sample and <https://osf.io/nu62m> for the adult sample). Participants were equally distributed across the two study versions. Data was collected between July 2021 and June 2022.

The study design and procedure obtained ethical clearance by the MPG Ethics commission Munich, Germany, falling under a packaged ethics application (Appl.

No. 2021\_45), and was approved by an internal ethics committee at the Max Planck Institute for Evolutionary Anthropology. The research adheres to the legal requirements of psychological research with children in Germany.

## Participants

Participants were recruited in the same way as in the previous study. The child sample consisted of 120 children, including 41 3-year-olds (mean = 42.34 months, SD = 3.10, range = 37–47, 20 girls), 41 4-year-olds (mean = 53.76 months, SD = 3.15, range = 48–59, 21 girls), and 38 5-year-olds (mean = 66.05 months, SD = 3.40, range = 60–71, 19 girls).

Additional 65 children were recruited but not included in the analysis due to absence on the second test day ( $n = 49$ ), canceled testing because of COVID-19 cases in the kindergarten ( $n = 7$ ), children did not want to participate a second time ( $n = 5$ ), children already participated in the first data collection round and were included in the above-mentioned *Individual Differences* sample ( $n = 3$ ), or children did not understand the task instructions ( $n = 1$ ; manifested in too early clicking in the training trials while the instructions were still playing, and no clicking by themselves in the test trials). Two additional children were recruited for the first day (as backup) in case another child would be absent on the second test day. Similar to our first study, we did not collect participant-specific demographics. For a community-based description of our participant pool, see Participant section of the first study.

As in our first study, adults were recruited via *Prolific* (Palan & Schitter, 2018). The adult sample included 136 English speakers with an average age of 25.73 years (SD = 8.09, range = 18–71, 87 females; see Supplements for further details). Most participants resided in South Africa ( $n = 48$ ), the United Kingdom ( $n = 19$ ), and the United States ( $n = 14$ ). See Supplements and the available online data set for more detailed information.

## Procedure

We applied the same procedure as in the first study, with the following differences. Participants completed the study twice, with a delay of  $14 \pm 3$  days. The target locations, as well as the succession of agents and target colors, were randomized once and then held constant across participants. The child sample received 15 test trials. In the hedge version, each bin occurred once, making up ten of the test trials. For the remaining five test trials, we repeated one out of two adjacent bins (i.e., randomly chose between bins 1 & 2, bins 3 & 4, etc.). In the box version, we ensured that each of the five boxes occurred exactly three times during

test trials. Adults in the hedge version received 30 test trials, each of the ten bins occurring exactly three times. Adults in the box version received 32 test trials, with each of the eight boxes occurring exactly four times. For the four training trials, we repeated a fixed order of random bins/boxes. For the adult sample, we decided to increase the number of trials in order to get more accurate reliability estimates. Trial numbers were multipliers of the possible target locations and therefore differed between hedge and box versions. For the child sample, we stuck to the same number of trials to not risk higher attrition rates.

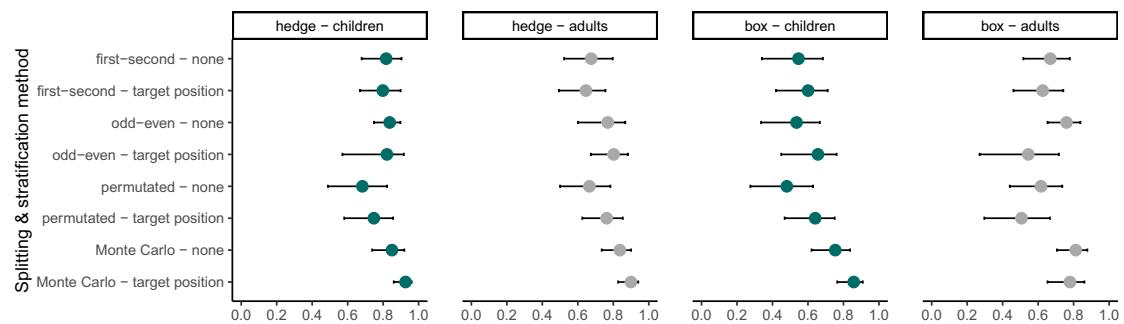
## Analysis

We assessed reliability in two ways. First, we focused on internal consistency by calculating split-half reliability coefficients.<sup>3</sup> For each subject, trials were split into odd and even trials. Performance was aggregated and then correlated using Pearson correlation coefficients. For this, we used the data of the first test day. Performance was defined according to each study version: in the hedge version, performance referred to the mean absolute difference between the target center and the click coordinate, scaled according to target widths; in the box version, we computed the mean proportion of correct choices.

Pronk et al., (2022) recently compared various methods for computing split-half reliability that differ in how the trials are split into parts and whether they are combined with stratification by task design. To compare our traditional approach of a simple odd-even split, we additionally calculated split-half reliability estimates using first-second, odd-even, permuted, and Monte Carlo splits without and with stratification by target position. First-second and odd-even splits belong to single sample methods since each participant has a single pair of performance scores, while permuted (without replacement) and Monte Carlo (with replacement) splits make use of resampling. Analyses were run using the function `by_split` from the `splithfr` package (Pronk et al., 2021).

Second, we assessed test-retest reliability. We calculated performance scores (depending on the study version as described above) for each participant in each test session and correlated them using Pearson correlation coefficients. Furthermore, for our child sample, we report an age-corrected correlation between the two test days using a GLMM-based approach (Rouder & Haaf, 2019). We fit trial-by-trial data with a fixed effect of age, a random intercept for each subject, and a random slope for test day (model notation

<sup>3</sup> The assessment of internal consistency was not pre-registered and was included as an additional measure of reliability.



**Fig. 3** Internal consistency. Reliability coefficients per splitting method, stratification level, study version, and age group. Error bars show the 95% confidence intervals of the coefficient estimates, calculated with the function by\_split from the splithalfr package (Pronk et al., 2021)

in R: performance ~ age + (0 + reiday | subID)). For the hedge version, performance was modeled by a lognormal distribution, while the model for the box version used a Bernoulli distribution with a logit link function. The model computes a correlation between the participant-specific estimates for each test day. This can be interpreted as the test-retest reliability. By using this approach, we do not need to compromise on data aggregation and, therefore, loss of information. Since the model uses hierarchical shrinkage, we obtain regularized, more accurate person-specific estimates. Most importantly, the model includes age as a fixed effect. The correlation between the two person-specific estimates is consequently the age-independent estimate for test-retest reliability. This rules out the possibility that a high correlation between test days arises from domain-general cognitive development instead of study-specific inter-individual differences. A high correlation between our participant-specific model estimates would indicate a high association between test days.

## Results

We found that the TANGO measured systematic variation: split-half and test-retest reliability was medium to high. For internal consistency, we show traditional odd-even splits on our data and the corresponding *Pearson* correlation coefficients in Fig. 2B, C, G, and H.

Figure 3 compares split-half reliability coefficients by splitting and stratification method (Pronk et al., 2021). In the hedge version, the split-half reliability coefficients ranged from 0.65 to 0.93. In the box version, split-half reliability coefficients ranged from 0.48 to 0.86. Similar to the results of Pronk et al. (2021), we found that more robust splitting methods that are less prone to task design or time confounds yielded higher reliability coefficients. In most cases, stratifying by target position led to similar or even higher estimates compared to no stratification. As expected, we found higher

coefficients for the samples with higher variation, i.e., for our continuous hedge version of the task.

For test-retest reliability, we show the association between raw performance scores of the two test days and corresponding *Pearson* correlation coefficients in Fig. 2D, E, I and J.<sup>4</sup> See Supplements for reliability estimates by age group.

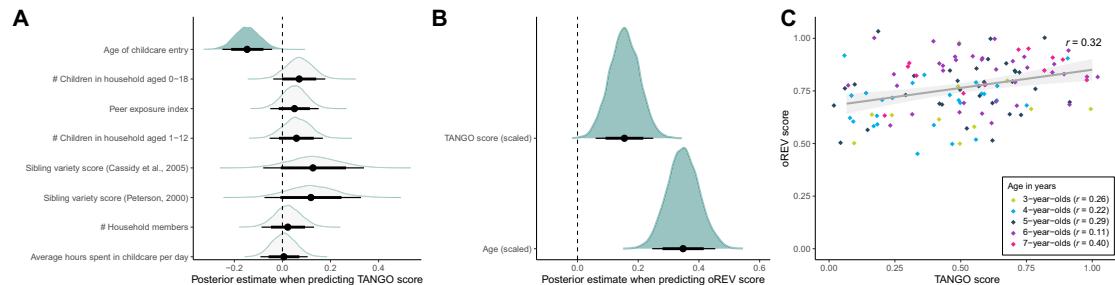
The age-corrected, GLMM-based retest reliabilities for children yielded similar results. In the hedge version, the correlation between test days was 0.89 (95% CrI [0.64;1.00]). In the box version, the correlation between test days was 0.91 (95% CrI [0.70;1.00]).

For both study versions, reliability estimates based on the GLMM approach were higher than the *Pearson* correlations. The GLMM-based estimates are less noisy due to the fact that the model uses all available information (e.g., participant age) and does not rely on data aggregation across trials.

## Discussion

Our results indicated that the measured variation was systematic. As expected, the continuous measure of the hedge version yielded higher reliability estimates than the discrete box version. For children, the model-based reliability estimates showed that the task did capture individual differences even when correcting for age. This corroborates what we already saw in Fig. 2: there was a clear overlap between age groups, indicating that age is predictive of performance for the mean but is not the main source of individual differences.

<sup>4</sup> In the hedge version, we excluded one 3-year-old, one 5-year-old, and two adults from the test-retest analysis. The performance of the mentioned participants was 3 standard deviations above/below the mean of each sample. Including the two children yielded a *Pearson* correlation coefficient of  $r = 0.88$ . Including the two adults yielded a *Pearson* correlation coefficient of  $r = 0.73$ .



**Fig. 4** Validity of the TANGO. **A** Influence of social-environmental factors on gaze understanding. **B** Influence of gaze understanding on receptive vocabulary. For (A) and (B), the graphs show the posterior distributions for the respective predictor of each model. *Filled green density curves* show that adding the respective predictor improved the model fit compared to the null model. *Black dots* represent means, *thicker black lines* 80% CrI and *thinner black lines* 95% CrI. The oREV score is the proportion of correctly selected pictures in

the receptive vocabulary task. Similarly, the TANGO score refers to proportion of correctly located targets (see Supplements for further detail). **C** Influence of gaze understanding on receptive vocabulary by age. The *regression line* with 95% CI shows a smooth conditional mean based on a generalized linear model, with Pearson's correlation coefficient  $r$ . *Dots* show the mean performance for each subject averaged across trials with minimal horizontal and vertical noise added to avoid overplotting. The *color of dots* denotes age group

## Validity

After having probed our new testing infrastructure and the psychometric properties of the TANGO, we aimed at establishing its validity. One way to assess validity is to correlate the social-cognitive ability in question to concepts that are thought to be theoretically related. Social cognition is often described as developing in response to social interaction (Devine & Hughes, 2018; Hughes & Leekam, 2004). It is assumed that opportunities to play, communicate and argue with peers help children to understand the human mind. Therefore, many studies link social cognition to opportunities for social interaction captured in demographic variables such as parent-child interaction quality and quantity, mental state talk, and center-based childcare (Bulgarelli & Molina, 2016; Dunn et al., 1991; Pavarini et al., 2013). In particular, family constellation, the number and age of siblings, and their interaction have been linked to social cognition (Cassidy et al., 2005; Dunn et al., 1991; Perner et al., 1994; Peterson, 2000; Zhang et al., 2021).

To assess such external validity for the TANGO, we handed out a brief demographic questionnaire to families of our kindergarten and online child sample and asked for (1) the total number of household members, (2) the number of children, (3) age of the other children, (4) whether the child was in daycare, and if yes, (5) since when and (6) for how long on an average day. 109 families filled out the questionnaire and were included in the analysis. We used parents' responses to construct different scores suggested in the literature (Cassidy et al., 2005; Peterson, 2000), capturing aspects of children's opportunities for social interaction with adults and peers. Only the predictor "age of childcare entry" improved the model fit compared to the null model (see Fig. 4A; for model comparisons, see Supplements): the older the children were when entering

childcare, the less likely they were to correctly use the available gaze cue. Figure 4A shows that all other predictor scores were positively linked to gaze understanding. Effect sizes were probably influenced by the lack of variance in the predictors: variables like household size and number of siblings typically vary very little among German households (see Supplements for distribution characteristics of the predictors). Albeit the effects were weak, they are consistent with the literature.

In addition, children's sensitivity to gaze has been linked to language acquisition (Brooks & Meltzoff, 2005; Del Bianco et al., 2019; Okumura et al., 2017). Discovering the attentional focus of your counterpart is thought to facilitate word learning, for example by identifying the referent of a new word (Tomasello, 2003). For 117 children, we also collected data with a receptive vocabulary test (oREV; Bohn et al., 2023) approximately 6 months (mean = 0.52 years, SD = 0.08, range = 0.06–0.80) after their participation in the TANGO. In the oREV task, children are shown four pictures (see Supplements for further detail) and hear a verbal prompt asking them to select one of the pictures. The oREV score is the proportion of correctly selected pictures. We found a substantial relationship between gaze understanding 6 months prior and receptive vocabulary, even when correcting for age (see Fig. 4B and C). Taken together, our newly developed task shows connections to external variables and psychological constructs that are characteristic of measures of social cognition.

## General discussion

We have presented a new experimental paradigm to study gaze understanding across the lifespan. This paper contributes to methodological advances in developmental

psychology in the following ways: first, we captured fine-grained individual differences in gaze understanding at different ages – from early childhood until adulthood. Individuals behaved consistently differently from one another (i.e., we found substantial variation between individuals across age groups). Second, our task showed satisfactory psychometric properties with respect to internal consistency and test-retest reliability estimates. Third, our new browser-based testing infrastructure ensures standardized, portable data collection at scale, both remotely as well as in person. In sum, the TANGO provides a step toward more robust and reliable research methods, especially with regard to measuring developmental change in a fundamental social-cognitive ability. The web app (<https://ccp-odc.eva.mpg.de/tango-demo/>) and its source code (<https://github.com/ccp-eva/tango-demo>) are freely accessible for use and modification.

Our continuous measure of children's gaze understanding moves away from treating a social-cognitive ability as an all-or-nothing matter (e.g., dichotomous measures in pass/fail situations) toward an ability on a continuum (Beaudoin et al., 2020; Hughes & Devine, 2015). Identifying variability in social-cognitive abilities is vital for accurately quantifying developmental change, revealing relations between different aspects of cognition and children's real-life social surroundings, and for meaningful comparisons across human cultures and across animal species. Dedicated measures of individual differences will help us to design meaningful interventions and progress in psychological theory building (Hedge et al., 2018).

Our continuous hedge version yields higher internal consistency estimates than the categorical box version. Both study versions exhibit high test-retest reliability, also when controlling for age. Therefore, when a sufficient number of trials is presented, the box version of the task can also yield reliable individual estimates (cf. Hughes et al. (2000); improved reliability through aggregation). When testing time is limited (and the number of trials might be low), we recommend using the continuous study version for higher internal consistency. However, the categorical box version demonstrates design features that might be preferable in some research contexts: for example, researchers could induce different levels of salience for each box. Our task could consequently be used to study bias, preferences, and diverse desires (e.g., matching the box appearance to some feature/behavioral characteristic of the agent).

In the split-half reliability calculations, the more accurately the statistical method represents the task structure, the higher the reliability estimates are. Therefore, we argue that future research should aim at implementing statistical analyses that mirror the complexity of the experimental design. Theoretically informed, computational cognitive models are a promising approach forward (Haines et al.,

2020). Computational models take advantage of all available information and model variation between and within individuals in an even more fine-grained and psychologically interpretable manner. Computational frameworks could also be used to model performance and their underlying cognitive processes across tasks. With nested hierarchical models, we could assess the systematic relation between various social-cognitive abilities and recover potentially shared structures between cognitive processes (Bohn et al., 2023).

The TANGO fulfills several demands that were proposed by Schaafsma et al. (2015)'s wish list: it measures proficiency on a continuum, avoids floor and ceiling effects, measures variation across age ranges, shows satisfactory reliability estimates, and has a high face value.

In addition to the new task design itself, we designed a new testing infrastructure. The TANGO is presented as an interactive web app. This enables presentation across devices without any prior installation. Stimuli presentation is achieved through the use of SVGs. This has several advantages: the aspect ratio and stimulus quality are kept constant no matter which size the web browser displays. The cartoon-like presentation makes the task engaging for children and adults alike. Most importantly, we can dynamically modify the stimulus details (e.g., target positions) on a trial-by-trial basis. Presented agents, voice-over instructions, and objects can be easily adapted for future task modifications or specific linguistic and cultural settings.

The browser-based implementation allows for different data collection modes: participants can be tested in person with supervision or remotely at home. Test instructions are standardized, and with prior informed consent, the webcam records study participation. This allows us to scale up data collection: testing is flexible, fast, and requires no further experimenter training. We compared children participating in-person and supervised in kindergartens with children who participated remotely at home. Our results suggest a comparable developmental trajectory of gaze understanding in both samples. Children in the remote sample were slightly more precise. This effect was most pronounced in the 3-year-olds in the box version (for an analysis of the webcam recordings, see Supplements). Therefore, we recommend using a tablet for remote data collection. Children can click for themselves, and caregivers have less chance to interfere. The design choices of the infrastructure underline how our study design can act as a versatile framework for addressing further research questions on social-cognitive development.

With respect to validity, we found that performance in the TANGO was related to relevant external variables and cognitive measures. Family-level variables, capturing a child's opportunity for social interaction, systematically influenced gaze understanding. Even though the effects

were small and confidence intervals were wide, it is remarkable that we were able to detect relationships between this fundamental social-cognitive ability and very distant, real-life variables. In addition, we assessed the influence of gaze understanding on receptive vocabulary. We found a substantial relationship between the two variables, even when correcting for age. Taken together, this speaks to the validity of the TANGO.

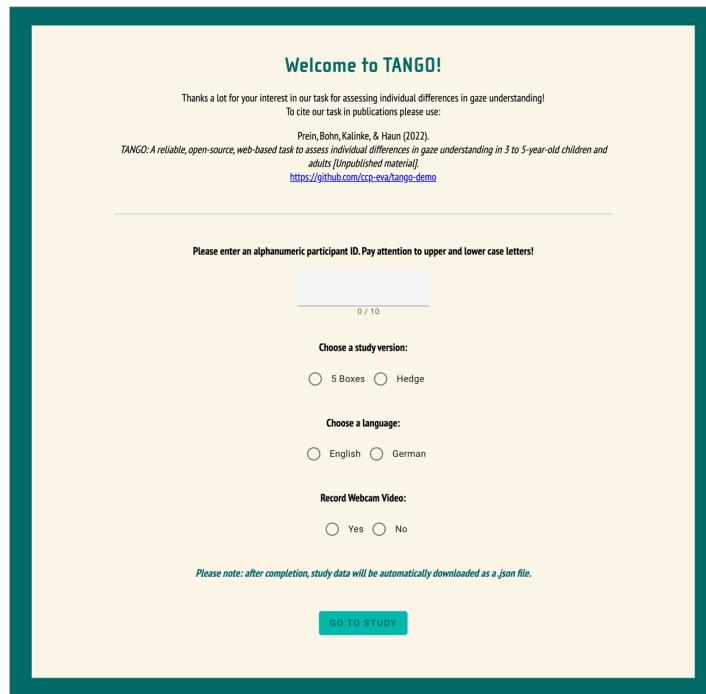
### Limitations

First, we want to address the scope and interpretation of the TANGO. We believe that solving the task requires locating the attentional focus of an agent as the gaze cues the target location. This speaks to the face validity of the TANGO and its focus on an inherently social stimulus. However, we do not want to claim that the TANGO does not also recruit other, domain-general processes. For example, we believe that a considerable part of gaze understanding relies on vector-following: not just in our task but also in real life. From that perspective, gaze understanding could be seen as a particular case of vector-following that is learned and used in social interactions. Future research could assess how

much variation of the gaze understanding task is shared with a physical vector-following task. In addition, computational cognitive models might prove helpful in defining children's behavior on a process-level and disentangling parameters that influence task performance (e.g., spatial acuity).

Second, the influence of testing modality requires further attention. Remote data collection loosens the standardization of the experimental procedure, as we cannot prevent caregivers from interfering. Steering the child's behavior becomes less possible when touchscreens are used, and the child can click on the screen directly. This is why we recommend using tablets for remote data collection. However, it should be noted that families' access to technological devices varies, both across socio-environmental as well as cultural settings.

Third, the children in our sample live in an industrialized, urban Central-European context. It is unclear how our results would generalize to different socio-cultural contexts. A related limitation is that we did not collect demographic information on a participant-level and, instead, had to rely on a community-level description of the sample. This is important to keep in mind when gauging the generalizability of our new measure.



**Fig. 5** TANGO demo website. We want to highlight that researchers are welcome to use and modify our task according to their needs. The number of training and test trials and the number of boxes can be

adjusted within the JavaScript code, while agents and targets can be exchanged within the HTML code

Finally, we utilized subtle gaze cues in order to increase difficulty and capture individual differences. However, in real-life settings, children could be more accustomed to a combination of head and eye orientation changes, and subtle gaze differences might be less common.

## Conclusions

We have presented a new experimental paradigm to study gaze understanding across the lifespan. The TANGO captures individual differences and shows highly satisfactory psychometric properties with respect to internal consistency and test-retest reliability. The browser-based testing infrastructure allows for standardized, portable data collection at scale, both remotely as well as in person. Associations with social-environmental factors and language skills illustrate the validity of the task. Ultimately, this work shows a promising way forward toward more precise measures of cognitive development. The data sets and the analysis code are freely available in the associated online repository (<https://github.com/ccp-eva/gazecues-methods>). A demo version of the task is available at the following website (see Fig. 5): <https://ccp-odc.eva.mpg.de/tango-demo/>. The code base and respective assets can be accessed in the following repository: <https://github.com/ccp-eva/tango-demo>. These resources allow interested researchers to use, extend and adapt the task.

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.3758/s13428-023-02159-5>.

**Author note** The authors made the following contributions. Julia Christin Prein: Conceptualization, Software, Formal Analysis, Writing—Original Draft Preparation, Writing - Review & Editing; Manuel Bohn: Conceptualization, Writing - Original Draft Preparation, Writing - Review & Editing; Steven Kalinke: Software, Writing - Review & Editing; Daniel B. M. Haun: Conceptualization, Writing - Review & Editing.

**Funding** Open Access funding enabled and organized by Projekt DEAL. This study was funded by the Max Planck Society for the Advancement of Science, a noncommercial, publicly financed scientific organization (no grant number). We thank all the children, caregivers, and adults who participated in the study. We thank Jana Jurkat for her help with data collection.

**Open practices statement** The web application (<https://ccp-odc.eva.mpg.de/tango-demo/>) described here is open-source (<https://github.com/ccp-eva/tango-demo>). The data sets generated during and/or analyzed during the current study are available in the [gazecues-methods] repository (<https://github.com/ccp-eva/gazecues-methods>). All experiments were pre-registered (<https://osf.io/zjhsc/>).

## Declarations

**Conflicts of interest** The authors declare that they have no conflict of interest.

**Consent to participate** Informed consent was obtained from all individual participants included in the study or their legal guardians.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Astor, K., Thiele, M., & Gredebäck, G. (2021). Gaze following emergence relies on both perceptual cues and social awareness. *Cognitive Development*, 60, 101121. <https://doi.org/10.1016/j.cogdev.2021.101121>
- Astor, K., Lindskog, M., Forssman, L., Kenward, B., Fransson, M., Skalkidou, A., Gredebäck, G. (2020). Social and emotional contexts predict the development of gaze following in early infancy. *Royal Society Open Science*, 7(9), 201178. <https://doi.org/10.1098/rsos.201178>
- Baron-Cohen, S., Leslie, A. M., & Frith, U. (1985). Does the autistic child have a “theory of mind”? *Cognition*, 21, 37–46. [https://doi.org/10.1016/0010-0277\(85\)90022-8](https://doi.org/10.1016/0010-0277(85)90022-8)
- Barrett, H. C., Broesch, T., Scott, R. M., He, Z., Baillargeon, R., Wu, D., Bolz, M., Henrich, J., Setoh, P., Wang, J., & Laurence, S. (2013). Early false-belief understanding in traditional non-Western societies. *Proceedings of the Royal Society B: Biological Sciences*, 280(1755), 20122654. <https://doi.org/10.1098/rspb.2012.2654>
- Beaudoin, C., Leblanc, É., Gagné, C., & Beauchamp, M. H. (2020). Systematic review and inventory of theory of mind measures for young children. *Frontiers in Psychology*, 10, 2905. <https://doi.org/10.3389/fpsyg.2019.02905>
- Begeer, S., Bernstein, D. M., van Wijhe, J., Scheeren, A. M., & Koot, H. M. (2012). A continuous false belief task reveals egocentric biases in children and adolescents with autism spectrum disorders. *Autism*, 16(4), 357–366. <https://doi.org/10.1177/1362361311434545>
- Benson, J. E., Sabbagh, M. A., Carlson, S. M., & Zelazo, P. D. (2013). Individual differences in executive functioning predict preschoolers’ improvement from theory-of-mind training. *Developmental Psychology*, 49(9), 1615–1627. <https://doi.org/10.1037/a0031056>
- Bernstein, D. M., Thornton, W. L., & Sommerville, J. A. (2011). Theory of mind through the ages: Older and middle-aged adults exhibit more errors than do younger adults on a continuous false belief task. *Experimental Aging Research*, 37(5), 481–502. <https://doi.org/10.1080/0361073X.2011.619466>
- Bohn, M., & Köyken, B. (2018). Common ground and development. *Child Development Perspectives*, 12(2), 104–108. <https://doi.org/10.1111/cdep.12269>
- Bohn, M., Le, K. N., Peloquin, B., Köyken, B., & Frank, M. C. (2021a). Children’s interpretation of ambiguous pronouns based on prior discourse. *Developmental Science*, 24(3), e13049. <https://doi.org/10.1111/desc.13049>

## Appendix A — Main Publications

### Behavior Research Methods

- Bohn, M., Tessler, M. H., Merrick, M., & Frank, M. C. (2021b). How young children integrate information sources to infer the meaning of words. *Nature Human Behaviour*, 5(8), 1046–1054. <https://doi.org/10.1038/s41562-021-01145-1>
- Bohn, M., Prein, J., Koch, T., Bee, R. M., Delikaya, B., Haun, D., & Gagarina, N. (2023). oREV: An item response theory-based open receptive vocabulary task for 3- to 8-year-old children. *Behavior Research Methods*. <https://doi.org/10.3758/s13428-023-02169-3>
- Bohn, M., Tessler, M. H., Kordt, C., Hausmann, T., & Frank, M. C. (2023). An individual differences perspective on pragmatic abilities in the preschool years. *Developmental Science*, e13401. Advance online publication. <https://doi.org/10.1111/desc.13401>
- Brooks, R., & Meltzoff, A. N. (2005). The development of gaze following and its relation to language. *Developmental Science*, 8(6), 535–543. <https://doi.org/10.1111/j.1467-7687.2005.00445.x>
- Bulgarelli, D., & Molina, P. (2016). Social cognition in preschoolers: effects of early experience and individual differences. *Frontiers in Psychology*, 7(1762). <https://doi.org/10.3389/fpsyg.2016.01762>
- Bürkner, P.-C. (2018). Advanced Bayesian Multilevel Modeling with the R Package brms. *The Royal Journal*, 10(1), 395. <https://doi.org/10.32614/RJ-2018-017>
- Bürkner, P.-C. (2017). brms: an R package for bayesian multilevel models using stan. *Journal of Statistical Software*, 80(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- Buttelmann, D., Kühn, K., & Zmyj, N. (2022). The relations among theory of mind, inhibitory control, and aggressive behavior in 4-year-old children – A multi-measure multi-informant approach. *Journal of Cognition and Development*, 23(1), 111–134. <https://doi.org/10.1080/15248372.2021.1987240>
- Byers-Heinlein, K., Tsui, R. K.-Y., van Renswoude, D., Black, A. K., Barr, R., Brown, A., ... Singh, L. (2021). The development of gaze following in monolingual and bilingual infants: A multi-laboratory study. *Infancy*, 26(1), 4–38. <https://doi.org/10.1111/infa.12360>
- Callaghan, T., Rochat, P., Lillard, A., Claux, M. L., Odden, H., Itakura, S., Tapanya, S., & Singh, S. (2005). Synchrony in the onset of mental-state reasoning: evidence from five cultures. *Psychological Science*, 16(5), 378–384. <https://doi.org/10.1111/j.0956-7976.2005.01544.x>
- Carlson, S. M., & Moses, L. J. (2001). Individual differences in inhibitory control and children's theory of mind. *Child Development*, 72(4), 1032–1053. <https://doi.org/10.1111/1467-8624.00333>
- Carlson, S., Moses, L. J., & Claxton, L. J. (2004). Individual differences in executive functioning and theory of mind: An investigation of inhibitory control and planning ability. *Journal of Experimental Child Psychology*, 87(4), 299–319. <https://doi.org/10.1016/j.jecp.2004.01.002>
- Cassidy, K. W., Fineberg, D. S., Brown, K., & Perkins, A. (2005). Theory of mind may be contagious, but you don't catch it from your twin. *Child Development*, 76(1), 97–106. <https://doi.org/10.1111/j.1467-8624.2005.00832.x>
- Coburn, P. I., Bernstein, D. M., & Begeer, S. (2015). A new paper and pencil task reveals adult false belief reasoning bias. *Psychological Research*, 79(5), 739–749. <https://doi.org/10.1007/s00426-014-0606-0>
- Coelho, E., George, N., Conty, L., Hugueville, L., & Tijus, C. (2006). Searching for asymmetries in the detection of gaze contact versus averted gaze under different head views: A behavioural study. *Spatial Vision*, 19(6), 529–545. <https://doi.org/10.1163/1568-6806779194026>
- Cutting, A. L., & Dunn, J. (1999). Theory of Mind, Emotion Understanding, Language, and Family Background: Individual Differences and Interrelations. *Child Development*, 70(4), 853–865. <https://doi.org/10.1111/1467-8624.00061>
- Del Bianco, T., Falck-Ytter, T., Thorup, E., & Gredebäck, G. (2019). The Developmental Origins of Gaze-Following in Human Infants. *Infancy*, 24(3), 433–454. <https://doi.org/10.1111/infa.12276>
- Devine, R. T., & Hughes, C. (2018). Family Correlates of False Belief Understanding in Early Childhood: A Meta-Analysis. *Child Development*, 89(3), 971–987. <https://doi.org/10.1111/cdev.12682>
- Dunn, J., Brown, J., Slomkowski, C., Tesla, C., & Youngblade, L. (1991). Young children's understanding of other people's feelings and beliefs: Individual differences and their antecedents. *Child Development*, 62(6), 1352–1366.
- Frank, M. C., Sugarman, E., Horowitz, A. C., Lewis, M. L., & Yurovsky, D. (2016). Using Tablets to Collect Data From Young Children. *Journal of Cognition and Development*, 17(1), 1–17. <https://doi.org/10.1080/15248372.2015.1061528>
- Frischen, A., Bayliss, A. P., & Tipper, S. P. (2007). Gaze cueing of attention: Visual attention, social cognition, and individual differences. *Psychological Bulletin*, 133(4), 694–724. <https://doi.org/10.1037/0033-2909.133.4.694>
- Gathercole, S. E., Pickering, S. J., Ambridge, B., & Wearing, H. (2004). The Structure of Working Memory From 4 to 15 Years of Age. *Developmental Psychology*, 40, 177–190. <https://doi.org/10.1037/0012-1649.40.2.177>
- Gola, A. A. H. (2012). Mental verb input for promoting children's theory of mind: A training study. *Cognitive Development*, 27(1), 64–76. <https://doi.org/10.1016/j.cogdev.2011.10.003>
- Gopnik, A., & Slaughter, V. (1991). Young Children's Understanding of Changes in Their Mental States. *Child Development*, 62(1), 98. <https://doi.org/10.2307/1130707>
- Haines, N., Kvam, P. D., Irving, L. H., Smith, C., Beauchaine, T. P., Pitt, M. A., Ahn, W.-Y., & Turner, B. (2020). *Theoretically Informed Generative Models Can Advance the Psychological and Brain Sciences: Lessons from the Reliability Paradox* [Preprint]. PsyArXiv. <https://doi.org/10.31234/osf.io/xr7y3>
- Happé, F., Cook, J. L., & Bird, G. (2017). The Structure of Social Cognition: In(ter)dependence of Sociocognitive Processes. *Annual Review of Psychology*, 68(1), 243–267. <https://doi.org/10.1146/annurev-psych-010416-044046>
- Hedge, C., Powell, G., & Sumner, P. (2018). The reliability paradox: Why robust cognitive tasks do not produce reliable individual differences. *Behavior Research Methods*, 50(3), 1166–1186. <https://doi.org/10.3758/s13428-017-0935-1>
- Hernik, M., & Broesch, T. (2019). Infant gaze following depends on communicative signals: An eye-tracking study of 5- to 7-month-olds in Vanuatu. *Developmental Science*, 22(4), e12779. <https://doi.org/10.1111/desc.12779>
- Herrmann, E., Hernández-Lloreda, M. V., Call, J., Hare, B., & Tomasello, M. (2010). The Structure of Individual Differences in the Cognitive Abilities of Children and Chimpanzees. *Psychological Science*, 21(1), 102–110. <https://doi.org/10.1177/0956797609356511>
- Hughes, C., & Devine, R. T. (2015). Individual Differences in Theory of Mind From Preschool to Adolescence: Achievements and Directions. *Child Development Perspectives*, 9(3), 149–153. <https://doi.org/10.1111/cdep.12124>
- Hughes, C., & Ensor, R. (2007). Executive function and theory of mind: Predictive relations from ages 2 to 4. *Developmental Psychology*, 43(6), 1447–1459. <https://doi.org/10.1037/0012-1649.43.6.1447>
- Hughes, C., & Leekam, S. (2004). What are the Links Between Theory of Mind and Social Relations? Review, Reflections and New Directions for Studies of Typical and Atypical Development. *Social Development*, 13(4), 590–619. <https://doi.org/10.1111/j.1467-9507.2004.00285.x>

- Hughes, C., Adlam, A., Happé, F., Jackson, J., Taylor, A., & Caspi, A. (2000). Good Test-Retest Reliability for Standard and Advanced False-Belief Tasks across a Wide Range of Abilities. *Journal of Child Psychology and Psychiatry*, 41(4), 483–490. <https://doi.org/10.1111/j.1469-7610.00633>
- Hughes, C., Ensor, R., & Marks, A. (2011). Individual differences in false belief understanding are stable from 3 to 6 years of age and predict children's mental state talk with school friends. *Journal of Experimental Child Psychology*, 108(1), 96–112. <https://doi.org/10.1016/j.jecp.2010.07.012>
- Imuta, K., Henry, J. D., Slaughter, V., Selcuk, B., & Ruffman, T. (2016). Theory of mind and prosocial behavior in childhood: A meta-analytic review. *Developmental Psychology*, 52(8), 1192–1205. <https://doi.org/10.1037/dev0000140>
- Itakura, S., & Tanaka, M. (1998). Use of experimenter-given cues during object-choice tasks by chimpanzees (*Pan troglodytes*), an orangutan (*Pongo pygmaeus*), and human infants (*Homo sapiens*). *Journal of Comparative Psychology*, 112(2), 119–126. <https://doi.org/10.1037/0735-7036.112.2.119>
- Kidd, E., Donnelly, S., & Christiansen, M. H. (2018). Individual Differences in Language Acquisition and Processing. *Trends in Cognitive Sciences*, 22(2), 154–169. <https://doi.org/10.1016/j.tics.2017.11.006>
- Lecce, S., Bianco, F., Devine, R. T., Hughes, C., & Banerjee, R. (2014). Promoting theory of mind during middle childhood: A training program. *Journal of Experimental Child Psychology*, 126, 52–67. <https://doi.org/10.1016/j.jecp.2014.03.002>
- Lee, K., Eskritt, M., Symons, L. A., & Muir, D. (1998). Children's use of triadic eye gaze information for "mind reading." *Developmental Psychology*, 34(3), 525–539. <https://doi.org/10.1037/0012-1649.34.3.525>
- Macdonald, R. G., & Tatler, B. W. (2013). Do as eye say: Gaze cueing and language in a real-world social interaction. *Journal of Vision*, 13(4), 6–6. <https://doi.org/10.1167/13.4.6>
- Mahy, C. E. V., Bernstein, D. M., Gerrard, L. D., & Atance, C. M. (2017). Testing the validity of a continuous false belief task in 3- to 7-year-old children. *Journal of Experimental Child Psychology*, 160, 50–66. <https://doi.org/10.1016/j.jecp.2017.03.010>
- Mayer, A., & Träuble, B. (2015). The weird world of cross-cultural false-belief research: A true- and false-belief study among samoan children based on commands. *Journal of Cognition and Development*, 16(4), 650–665. <https://doi.org/10.1080/15248372.2014.926273>
- Mayes, L. C., Klin, A., Tercyak, K. P., Cicchetti, D. V., & Cohen, D. J. (1996). Test-Retest Reliability for False-Belief Tasks. *Journal of Child Psychology and Psychiatry*, 37(3), 313–319. <https://doi.org/10.1111/j.1469-7610.1996.tb01408.x>
- McEwen, F., Happé, F., Bolton, P., Rijssdijk, F., Ronald, A., Dworzynski, K., & Plomin, R. (2007). Origins of individual differences in imitation: Links with language, pretend play, and socially insightful behavior in two-year-old twins. *Child Development*, 78(2), 474–492. <https://doi.org/10.1111/j.1467-8624.2007.01010.x>
- Milligan, K., Astington, J. W., & Dack, L. A. (2007). Language and Theory of Mind: Meta-Analysis of the Relation Between Language Ability and False-belief Understanding. *Child Development*, 78(2), 622–646. <https://doi.org/10.1111/j.1467-8624.2007.01018.x>
- Molleman, L., Kurvers, R. H. J. M., & van den Bos, W. (2019). Unleashing the BEAST: A brief measure of human social information use. *Evolution and Human Behavior*, 40(5), 492–499. <https://doi.org/10.1016/j.evolhumbehav.2019.06.005>
- Moore, C. (2008). The Development of Gaze Following. *Child Development Perspectives*, 2(2), 66–70. <https://doi.org/10.1111/j.1750-8606.2008.00052.x>
- Mundy, P., Block, J., Delgado, C., Pomares, Y., Van Hecke, A. V., & Parlade, M. V. (2007). Individual differences and the development of joint attention in infancy. *Child Development*, 78(3), 938–954. <https://doi.org/10.1111/j.1467-8624.2007.01042.x>
- Okumura, Y., Kanakogi, Y., Kobayashi, T., & Itakura, S. (2017). Individual differences in object-processing explain the relationship between early gaze-following and later language development. *Cognition*, 166, 418–424. <https://doi.org/10.1016/j.cognition.2017.06.005>
- Palan, S., & Schitter, C. (2018). Prolific.ac subject pool for online experiments. *Journal of Behavioral and Experimental Finance*, 17, 22–27. <https://doi.org/10.1016/j.jbef.2017.12.004>
- Pavarini, G., de Hollanda Souza, D., & Hawk, C. K. (2013). Parental Practices and Theory of Mind Development. *Journal of Child and Family Studies*, 22(6), 844–853. <https://doi.org/10.1007/s10826-012-9643-8>
- Perner, J., Ruffman, T., & Leekam, S. R. (1994). Theory of Mind Is Contagious: You Catch It from Your Sibs. *Child Development*, 65(4), 1228–1238. <https://doi.org/10.2307/1131316>
- Peterson. (2000). Kindred spirits: Influences of siblings' perspectives on theory of mind. *Cognitive Development*, 15(4), 435–455. [https://doi.org/10.1016/S0885-2014\(01\)00040-5](https://doi.org/10.1016/S0885-2014(01)00040-5)
- Peterson, & WellmanSlaughter, H. M. V. (2012). The mind behind the message: advancing theory-of-mind scales for typically developing children, and those with deafness, autism, or asperger syndrome: The mind behind the message. *Child Development*, 83(2), 469–485. <https://doi.org/10.1111/j.1467-8624.2011.01728.x>
- Poulin-Dubois, D., Goldman, E. J., Meltzer, A., & Psaradellis, E. (2023). Discontinuity from implicit to explicit theory of mind from infancy to preschool age. *Cognitive Development*, 65, 101273. <https://doi.org/10.1016/j.cogdev.2022.101273>
- Pronk, T., Molenaar, D., Wiers, R. W., & Murre, J. (2022). Methods to split cognitive task data for estimating split-half reliability: a comprehensive review and systematic assessment. *Psychonomic Bulletin & Review*, 29(1), 44–54. <https://doi.org/10.3758/s13423-021-01948-3>
- R Core Team. (2022). *R: A language and environment for statistical computing* [Manual]. R Foundation for Statistical Computing.
- Rakoczy, H. (2022). Foundations of theory of mind and its development in early childhood. *Nature Reviews Psychology*, 1(4), 223–235. <https://doi.org/10.1038/s44159-022-00037-z>
- Raviv, L., & Arnon, I. (2018). The developmental trajectory of children's auditory and visual statistical learning abilities: Modality-based differences in the effect of age. *Developmental Science*, 21(4), e12593. <https://doi.org/10.1111/desc.12593>
- Slaughter, V., & Repacholi, B. (2003). Introduction: individual differences in theory of mind. What are we investigating? *Individual differences in theory of mind: implications for typical and atypical development*. Psychology Press.
- Rizzo, M. T., & Killen, M. (2018). Theory of mind is related to children's resource allocations in gender stereotypic contexts. *Developmental Psychology*, 54(3), 510. <https://doi.org/10.1037/dev0000439>
- Rouder, J. N., & Haaf, J. M. (2019). A psychometrics of individual differences in experimental tasks. *Psychonomic Bulletin & Review*, 26(2), 452–467. <https://doi.org/10.3758/s13423-018-1558-y>
- Samuel, S., Legg, E. W., Lurz, R., & Clayton, N. S. (2018). Egocentric bias across mental and non-mental representations in the Sandbox Task. *Quarterly Journal of Experimental Psychology*, 71(11), 2395–2410. <https://doi.org/10.1177/1747021817742367>
- Samuel, S., Legg, E. W., Lurz, R., & Clayton, N. S. (2018). The unreliability of egocentric bias across self-other and memory-belief distinctions in the Sandbox Task. *Royal Society Open Science*, 5(11), 181355. <https://doi.org/10.1098/rsos.181355>
- Schaafsma, S. M., Pfaff, D. W., Spunt, R. P., & Adolphs, R. (2015). Deconstructing and reconstructing theory of mind. *Trends in Cognitive Sciences*, 19(2), 65–72. <https://doi.org/10.1016/j.tics.2014.11.007>

Behavior Research Methods

- Shepherd, S. (2010). Following gaze: gaze-following behavior as a window into social cognition. *Frontiers in Integrative Neuroscience*, 4(5). <https://doi.org/10.3389/fnint.2010.00005>
- Slaughter, V. (2015). Theory of Mind in Infants and Young Children: A Review. *Australian Psychologist*, 50(3), 169–172. <https://doi.org/10.1111/ap.12080>
- Sodian, B. (2023). Reply to Poulin-Dubois et al. (2023): Replication problems concerning both implicit and explicit false belief reasoning greatly reduced the chance of finding longitudinal correlations. *Cognitive Development*, 65, 101294. <https://doi.org/10.1016/j.cogdev.2022.101294>
- Sodian, B., Licata, M., Kristen-Antonow, S., Paulus, M., Killen, M., & Woodward, A. (2016). Understanding of Goals, Beliefs, and Desires Predicts Morally Relevant Theory of Mind: A Longitudinal Investigation. *Child Development*, 87(4), 1221–1232. <https://doi.org/10.1111/cdev.12533>
- Sommerville, J. A., Bernstein, D. M., & Meltzoff, A. N. (2013). Measuring Beliefs in Centimeters: Private Knowledge Biases Preschoolers' and Adults' Representation of Others' Beliefs. *Child Development*, 84(6), 1846–1854. <https://doi.org/10.1111/cdev.12110>
- Tomasello, M. (2003). *Constructing a language a usage based theory of language acquisition* (8th ed., p. 388). Harvard University Press.
- Tomasello, M., Hare, B., Lehmann, H., & Call, J. (2007). Reliance on head versus eyes in the gaze following of great apes and human infants: The cooperative eye hypothesis. *Journal of Human Evolution*, 52(3), 314–320. <https://doi.org/10.1016/j.jhevol.2006.10.001>
- Underwood, B. J. (1975). Individual differences as a crucible in theory construction. *American Psychologist*, 30(2), 128–134. <https://doi.org/10.1037/h0076759>
- Walker, S. (2005). Gender Differences in the Relationship Between Young Children's Peer-Related Social Competence and Individual Differences in Theory of Mind. *The Journal of Genetic Psychology*, 166(3), 297–312. <https://doi.org/10.3200/GNTP.166.3.297-312>
- Wellman, H. M. (2012). Theory of mind: Better methods, clearer findings, more development. *European Journal of Developmental Psychology*, 9(3), 313–330. <https://doi.org/10.1080/17405629.2012.680297>
- Wellman, H. M., & Liu, D. (2004). Scaling of Theory-of-Mind Tasks. *Child Development*, 75(2), 523–541. <https://doi.org/10.1111/j.1467-8624.2004.00691.x>
- Wellman, H. M., Cross, D., & Watson, J. (2001). Meta-analysis of theory-of-mind development: The truth about false belief. *Child Development*, 72, 655–684. <https://doi.org/10.1111/1467-8624.00304>
- Wimmer, H., & Perner, J. (1983). Beliefs about beliefs: representation and constraining function of wrong beliefs in young children's understanding of deception. *Cognition*, 13(1), 103–128. [https://doi.org/10.1016/0010-0277\(83\)90004-5](https://doi.org/10.1016/0010-0277(83)90004-5)
- Zhang, Z., Yu, H., Long, M., & Li, H. (2021). Worse Theory of Mind in Only-Children Compared to Children With Siblings and Its Intervention. *Frontiers in Psychology*, 12, 5073. <https://doi.org/10.3389/fpsyg.2021.754168>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## **Study II**

Running head: MODELING VARIATION IN GAZE FOLLOWING

1

Variation in gaze following across the life span: A process-level perspective

MODELING VARIATION IN GAZE FOLLOWING

2

**Variation in gaze following across the life span: A process-level perspective**

**Authors:** Julia Christin Prein<sup>1</sup>, Luke Maurits<sup>1</sup>, Annika Werwach<sup>1,3,4</sup>, Daniel B. M. Haun<sup>1,\*</sup>, Manuel Bohn<sup>1,2,\*</sup>

**Affiliations:** <sup>1</sup> Department of Comparative Cultural Psychology, Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany. <sup>2</sup> Institute of Psychology, Leuphana University Lüneburg, Germany. <sup>3</sup> Center for Lifespan Psychology, Max Planck Institute for Human Development, Berlin, Germany. <sup>4</sup> Max Planck School of Cognition, Leipzig, Germany. \* shared senior authorship

**ORCiD:** *Julia Christin Prein* <https://orcid.org/0000-0002-3154-6167>

**Conflicts of interest:** The authors declare that they have no conflict of interest.

**Data availability statement:** The gaze following task (<https://ccp-odc.eva.mpg.de/tango-demo/>) is open source (<https://github.com/ccp-eva/tango-demo>). The data sets generated during and/or analyzed during the current study are available in the following repository (<https://github.com/ccp-eva/gazecues-modeling>). All experiments and analyses were pre-registered prior to data collection (<https://osf.io/zjhsc/>).

**Acknowledgements:** We thank Jana Jurkat for her help with data collection and participant recruitment. We would also like to thank Steven Kalinke for his technical programming support. We thank all the children, caregivers, and adults who participated in the study.

**Funding:** This study was funded by the Max Planck Society for the Advancement of Science, a noncommercial, publicly financed scientific organization (no grant number). Manuel Bohn was supported by a Jacobs Foundation Research Fellowship (grant no. 2022-1484-00).

**Ethical statement:** The study obtained ethical clearance by the MPG Ethics commission Munich, Germany, falling under an umbrella ethics application (Appl. No. 2021\_45). Informed consent was obtained from all individual participants or their legal guardians. The research adhered to the legal requirements of psychological research with children in Germany.

MODELING VARIATION IN GAZE FOLLOWING

4

**Research highlights**

- Gaze following develops beyond infancy. Highest precision levels in localizing attentional foci are reached in young adulthood with a slight decrease towards old age.
- We present a computational model that describes gaze following as a process of estimating pupil angles and the corresponding gaze vectors.
- The model explains individual differences and recovers signature patterns in the data. To estimate the relation between gaze- and vector following, we designed a non-social vector following task.
- We found substantial correlations between gaze following and vector following, as well as Level 2 perspective-taking. Other Theory of Mind tasks did not correlate.

**Abstract**

Following eye gaze is fundamental for many social-cognitive abilities, for example, when judging what another agent can or cannot know. While the emergence of gaze following has been thoroughly studied on a group level, we know little about (a) the developmental trajectory beyond infancy and (b) the sources of individual differences. In Study 1, we examined gaze following across the lifespan ( $N = 471$  3- to 80-year-olds; children from a mid-sized German city; international remotely tested adults). We found a steep performance improvement during preschool years, in which children became more precise in locating the attentional focus of an agent. Precision levels then stayed comparably stable throughout adulthood with a minor decline toward old age. In Study 2, we formalized the process of gaze following in a computational cognitive model that allowed us to conceptualize individual differences in a psychologically meaningful way ( $N = 60$  3- to 5-year-olds, 50 adults). According to our model, participants estimate pupil angles with varying levels of precision based on observing the pupil location within the agent's eyes. In Study 3, we empirically tested how gaze following relates to vector following in non-social settings and perspective-taking abilities ( $N = 102$  4- to 5-year-olds). We found that gaze following is associated with both of these abilities but less so with other Theory of Mind tasks. This work illustrates how the combination of reliable measurement instruments and formal theoretical models allows us to explore the in(ter)dependence of core social-cognitive processes in greater detail.

*Keywords:* social-cognitive development, theory of mind, gaze following, individual differences, cognitive modeling, lifespan

*Word count:* 10469

## MODELING VARIATION IN GAZE FOLLOWING

6

### Introduction

Following the gaze of others is valuable for extracting information from the environment. It guides us to “informational hotspots” (Meltzoff et al., 2010, p. 1) and can be used to identify internal states such as intentions or emotions (Corkum & Moore, 1998; Pfeiffer et al., 2013). As one of the most fundamental social-cognitive abilities, gaze following is an integral part of almost every form of social interaction, including communication, collaboration, cultural and social learning (Bohn & Frank, 2019; Emery, 2000; Frith & Frith, 2012; Hessels, 2020; Moore, 2008; Rakoczy, 2022; Tomasello & Rakoczy, 2003), and has been extensively studied in infancy (for review, see Del Bianco et al., 2019). In the most commonly used paradigm (e.g., Astor et al., 2020; Byers-Heinlein et al., 2021; Gredebäck et al., 2010; Ishikawa et al., 2022), the experimenter looks directly at the infant before shifting their head and eyes to one of two objects. Infants’ looking times to the target or the proportion of choosing the target over the distractor are measured. Research in this tradition finds that infants as young as three to four months can follow the gaze of another agent (Astor et al., 2021; D’Entremont et al., 1997; D’Entremont, 2000; Del Bianco et al., 2019).

Previous research has shown a refinement of gaze following abilities in a child’s first and second year of life (e.g., Astor et al., 2021; Brooks & Meltzoff, 2002; Butterworth & Jarrett, 1991). At the end of their first year of life, infants can follow gaze to locations outside their current visual field and move themselves to gain proper perceptual access (Butterworth & Jarrett, 1991; Corkum & Moore, 1995; Deák et al., 2000; Moll & Tomasello, 2004). However, we do not know much about the developmental progression beyond these qualitative milestones. One possibility is that the ability to follow gaze does not improve beyond infancy. Yet, most, if not all, cognitive abilities continue to develop throughout childhood (e.g., Gathercole et al., 2004;

Gredebäck et al., 2010). It seems likely that children fine-tune their gaze following as they get older - presumably while using them in social interactions. To capture the development in gaze following beyond infancy, Study 1 included participants from preschool to old age.

Up until today, only a handful of gaze following studies differentiate between manipulating head and eye movement. Michel et al. (2021) found that gaze following in four-month-olds was likely driven by other's head- instead of eye movements. Corkum and Moore (1995), Lempers (1979), and Lempers et al. (1977) suggest that infants, at least until 19 months, struggle when eye and head direction diverge. From farther distance, body or face orientation can act as more salient cues to determine another's area of attention. However, eye direction indicates a more precise location of focus (Emery, 2000; Stiefelhagen & Zhu, 2002; however, see Loomis et al., 2008 for peripheral vision), and allows to anticipate likely future actions (Friesen & Rao, 2011; Zohary et al., 2022). The three studies included in this work, therefore, focused on subtle gaze cues and isolated eye movement alone.

Group-level analyses of gaze following abilities (e.g., average age at which children as a group reach an above-chance performance) may mask individual differences between children. Measuring individual differences in basic aspects of social cognition is important to understand the underlying processes and to quantify the impact of environmental influences and other cognitive abilities (Birch et al., 2017; Del Bianco et al., 2019). Across the three studies, we measured gaze following continuously by using a task that is designed to capture individual-level variation (Prein et al., 2023): the TANGO (Task for Assessing iNdividual differences in Gaze understanding - Open) avoids floor and ceiling effects in children and adults and is thus particularly suited to examine how gaze following changes with age.

## MODELING VARIATION IN GAZE FOLLOWING

8

A promising approach to interpreting individual differences in cognitive abilities is computational cognitive modeling. Existing computational models have described gaze following via reinforcement learning (Ishikawa et al., 2020), as a consequence of goal interference and mapping self-experience onto other agents (Friesen & Rao, 2011) or an interplay of object distances/saliencies and head poses (Jasso & Triesch, 2006; Lau & Triesch, 2004; Recasens et al., 2015). As such, these models focus on the motivation behind gaze following or on situations in which the context (e.g., head orientation, surrounding objects) offers cues to the gaze direction. Our goal, however, was to formulate a theory that models how people estimate gaze direction based on the eyes alone (i.e., “literal” gaze following) when no target objects are present<sup>1</sup> and the other’s head is frontally oriented.

To our knowledge, there are three views that focus on the eyes alone that conceptualize gaze as (1) a beam, (2) a cone, or (3) a line (alternatively called a vector, ray, or line-of-sight). Guterstam et al. (2019) proposed the idea of gaze as a force-carrying beam. This theory, however, focuses on how people’s implicit assumptions about the physical properties of an object changes when someone looks at it. The idea of gaze as a cone is mostly concerned with the question of how people determine whether someone else looks at them (Gamer & Hecht, 2007; Horstmann & Linke, 2021). The conception most relevant to the present study sees eye gaze as a line: Yaniv and Shatz (1990) proposed that children extrapolate an imaginary trajectory between the agent

---

<sup>1</sup> As Jasso and Triesch (2006) put it: “There is a general agreement that tests of gaze following in the absence of targets are more stringent than with targets, because that eliminates the possibility that infants are simply following the target’s saliency” (p. 2).

and the object to identify the focus of attention (similar to “geometrical” gaze following; (Butterworth & Jarrett, 1991)). Anecdotal evidence was already reported by Walker and Gollin (1977), who observed two children pointing their fingers into the air, drawing a line between an agent and an object, and saying, “He sees that” (p. 354). Michelon and Zacks (2006) found that response time in Level 1 visual perspective-taking tasks depends on the distance between the agent and the object (i.e., the length of the line-of-sight). Some researchers additionally highlight iris eccentricity, that is, the ratio of the visible sclera on each side of the pupil (S. M. Anstis et al., 1969; Symons et al., 2004; Todorović, 2006). Todorović (2006) defined gaze direction as “the vector positioned along the visual axis, pointing from the fovea of the looker through the center of the pupil to the gazed-at spot” (p. 3550). Symons et al. (2004) furthermore reasoned that “the perceiver must use the asymmetrical configuration of the dark-white contrast of another individual’s eyes, and trace along two invisible sight-lines to their convergent point, that is, the third part of the triad (e.g., an object or a person)” (p. 452). Based on their finding that adults’ gaze direction sensitivity decreases when only one eye is shown, Symons et al. (2004) conclude that information from both eyes must be integrated.

While the previously mentioned conceptualizations focus on the direction of eye gaze, they (A) cannot explain how people differ in their abilities to precisely estimate gaze direction, and (B) are not clearly expressed as formal, mathematical models with explicit assumptions and testable predictions. In the words of (Gamer & Hecht, 2007): “Given the social relevance of determining gaze direction, the psychophysics of gaze is underdeveloped” (p. 705).

Here, we propose a cognitive model of gaze following, which builds upon the notion of eye gaze as line-of-sight tracing and extends this by explicitly modeling individual differences. Our gaze model assumes participants infer the locus of someone’s attention to be where two

MODELING VARIATION IN GAZE FOLLOWING

10

estimated gaze vectors meet. Each of these gaze vectors results from connecting the center of the agent's eyeball and the center of the pupil. Because the center of the eyeball is not directly observable, vector estimation happens with a degree of uncertainty. Development of gaze following corresponds to a decrease in uncertainty. Individual differences correspond to systematic differences in uncertainty.

By focusing on individual differences, we can further address the relationship between gaze following and other cognitive abilities. A longstanding question has been whether gaze following is related to Theory of Mind (ToM) (e.g., Brooks & Meltzoff, 2015). Moll and Meltzoff (2011) have suggested that joint attention (including gaze following) might be seen as "Level 0 perspective-taking", which provides the foundation for later-emerging, more complex perspective-taking abilities. On the other hand, an alignment of infants' visual attention to another's gaze does not necessarily indicate understanding the intentions of the agent (Aslin, 2007). Infants could simply align their orientation without processing what exactly the other is seeing (Butterworth & Jarrett, 1991). In fact, one might question if such an alignment reflects an understanding of visual perspectives at all because the "target" or "object of representation" is not necessarily specified (Perner et al., 2003, p. 358). Consequently, Astor and Gredebäck (2022) have listed the relationship between gaze following and perspective-taking as one of their five big open questions in gaze following research. Therefore, in Study 3, we assessed how gaze following relates to ToM abilities, especially visual perspective-taking.

Taken together, the present study had three main goals: first, we studied the development of gaze following beyond infancy (Study 1). Instead of capturing the youngest age at which children follow gaze, we examined how this ability changes with age. Our second goal was to provide a process-level theory of gaze following – and, most importantly, individual differences

therein. We proposed a computational cognitive model, which formalized gaze following as a form of vector following, and tested whether our model explained empirical data (Study 2). Third, we examined which (social-)cognitive components comprise gaze following (Study 3). Based on our model, we predicted that gaze following should be related to non-social vector following. Additionally, we assessed the link between gaze following and ToM measures, with a particular focus on visual perspective-taking.

### **Study 1: Gaze following across the lifespan**

The study was pre-registered prior to data collection: <https://osf.io/snju6> (child sample) and <https://osf.io/6yjz3> (adult sample). The study obtained ethical clearance by the MPG Ethics commission in Munich, Germany, falling under an umbrella ethics application (Appl. No. 2021\_45). Data was collected between May 2021 and April 2023.

### **Participants**

We collected data online from 3- to 80- year-olds (see Supplements for further details). The child sample consisted of 471 participants and was recruited via an internal database of families in Leipzig, Germany, who volunteered to participate in child development studies. Participants came from ethnically homogeneous, mixed socioeconomic backgrounds with mid to high parental education levels. They lived in an industrialized, urban Central-European context in a mid-size German city (approx. 600,000 inhabitants; median individual monthly net income approx. 1,600€ as of 2021). Most were raised monolingually in a nuclear two-generational family setting. Information on demographics and socioeconomic status was not formally recorded on a participant level.

## MODELING VARIATION IN GAZE FOLLOWING

12

Adults were recruited via *Prolific* (Palan & Schitter, 2018). *Prolific* is an online participant recruitment tool from the University of Oxford with predominantly European and US-American subjects. Participants consisted of 240 English-speaking adults who reported to have normal or corrected-to-normal vision. For completing the study, subjects were paid above the fixed minimum wage (~£10.00/hour).

### Materials

We used the continuous version of the TANGO (Prein et al., 2023). The task was presented as a web application (demo <https://ccp-odc.eva.mpg.de/tango-demo/>; source code <https://github.com/ccp-eva/tango-demo>). The TANGO showed satisfactory internal consistency and retest reliability (*Pearson's r* from .7 to .8; Prein et al. (2023)) and no floor or ceiling effects for children and adults.

### Procedure

Children and teenagers received a personalized link to the study website. Caregivers were asked to provide technical support, while explicitly being reminded not to help in responding. Webcam videos were recorded whenever consented and technically feasible in order to monitor whether children and teenagers responded on their own. Adults completed the online study unsupervised.

Each trial presented an agent standing in a window, watching a balloon (*i.e.*, target) falling to the ground (see Figure 4A; however, Study 1 presented animal agents). The target fell behind a hedge while the agent's gaze followed the target's trajectory. In test trials, a hedge covered the target's position. Participants were asked to touch or click where they estimated the target to be based on the agent's gaze. Four familiarization trials ensured participants understood the task and

felt comfortable with the response format. Then, 15 test trials followed. Completing 19 trials took 5-10 minutes.

We measured imprecision, defined as the absolute difference between the target center and the x coordinate of the participant's click. The screen width was divided into ten bins. Within each bin, exact target coordinates were randomly generated. Each target bin, agent, and target color occurred equally often and did not appear in more than two consecutive trials.

### **Analysis**

We ran all analyses in R version 4.3.3 (2024-02-29) (R Core Team, 2022). Regression models were fit as Bayesian generalized linear mixed models (GLMMs) with default priors using the function `brm` from the package `brms` (Bürkner, 2017, 2018).

We fit GLMMs that make different assumptions about the developmental trajectory, modeling the relationship between age and performance as linear, quadratic, or cubic. In addition, we fit a Gaussian Process model (Bürkner, 2017), which assumes a smooth relationship but avoids enforcing a particular shape. Per individual, imprecision was aggregated across trials and modeled as a lognormal distribution.<sup>2</sup> The unit of imprecision was counted in target widths, i.e.,

---

<sup>2</sup> Originally, we fit our models on a trial-by-trial basis with the following structure: `performance ~ age + symmetricPosition + trialNr + (1 + symmetricPosition + trialNr | subjID)`. However, the Gaussian Process model was computationally heavy. Thus, we simplified the model structure, aggregated data on a subject level, and included only age as an effect. See

## MODELING VARIATION IN GAZE FOLLOWING

14

an imprecision of 1 meant clicking one balloon width to the left or right of the true target center.

We inspected the posterior distributions (mean and 95% Credible Interval (CrI)) for the age estimates and compared models via model weights and the difference in expected log pointwise predictive density (ELPD) estimated using leave-one-out cross-validation (LOO) (Vehtari et al., 2017).

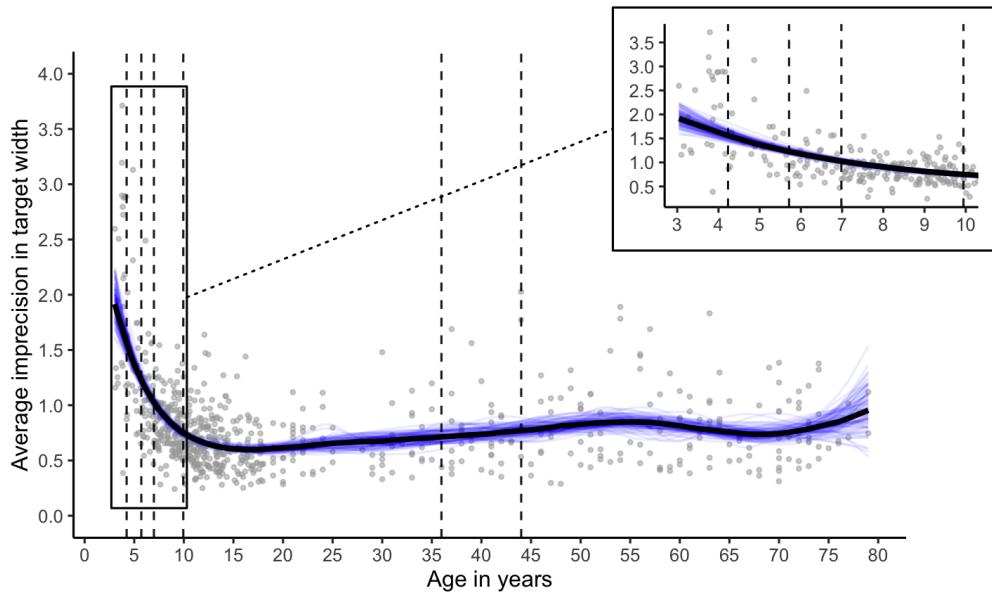
To obtain a concise but principled characterization of the developmental trajectory, we additionally performed a Bayesian change point analysis, using the package **RBeast** (Zhao et al., 2019). We sought the most likely change points in our data, assuming a constant mean (i.e., a flat line, zero-degree polynomial) within each segment. To avoid “overreactions” to outlying data points, we constrained the model to have minimally 10 data points between consecutive change points (= half of the data points collected per adult decade). We inspected the posterior probability of different numbers of change points and the locations of these change points (mean and 95% CrI).<sup>3</sup>

---

Supplements for a comparison between the original and the here-reported model structures. The model predictions did not differ notably.

<sup>3</sup> In a supplementary analysis, we varied the parameters of our change point analysis and modified the number of allowed change points, the minimum number of data points between change points, and the polynomial order. When we allowed more explorative room, the models became more sensitive and added more fine-grained change points. The exact location of the change points varied slightly. However, the overall interpretation stayed the same, fitting our initial visual inspection: while early childhood was characterized by much change, adults showed

## Result



**Figure 1: Developmental trajectory of gaze following across the lifespan.** Grey dots show the mean level of imprecision (ie., absolute distance between the target’s center and the participant’s click) for each subject averaged across trials. The unit of imprecision is counted in target width, i.e., a participant with imprecision of 1 clicked on average one target width to the left or right of the true target center. Blue lines show 100 draws from the expectation of the posterior predictive distribution of the Gaussian Process model, with its mean predicted developmental trajectory as a

---

a relatively stable level of imprecision, with a minor decay toward elderly adulthood. See Supplements for further detail.

## MODELING VARIATION IN GAZE FOLLOWING

16

solid black line. Vertical, black, dashed lines show the locations of the changes with the highest posterior probability according to our Bayesian change point analysis.

High levels of variation pointed to substantial individual differences in all age groups (overall imprecision mean = 0.81, sd = 0.82, range = [0 - 10.73]). We found substantial evidence for a non-linear development in gaze following across the lifespan. The Gaussian process model was clearly preferred over the polynomial models due to the highest predictive accuracy according to the LOO ELPD estimates: elpd\_diff between Gaussian process and cubic model = -33.75 (SE = 8.83); elpd\_diff between Gaussian process and quadratic model = -95.07 (SE = 15.19); elpd\_diff between Gaussian process and linear model = -127.17 (SE = 18.18); all in favor of the Gaussian process model. Moreover, the Gaussian Process model showed the greatest model weight (approximating 1). For the imprecision in gaze following, the standard deviation of the Gaussian process (SD = 1.52, 95% CrI [0.25; 5.19]) indicated nonlinearity.

The Bayesian change point analysis revealed 6 major shifts in gaze following during the lifespan (MAP estimate with 23.49% probability). The change points occurred at 4.23 years (95% CrI [4.13; 4.33], mean imprecision until change point = 2.15, SD = 0.85); 5.71 years (95% CrI [5.38; 5.90], mean = 1.36, SD = 0.53); 6.98 years (95% CrI [6.64; 8.04], mean = 1.06, SD = 0.40); 9.94 years (95% CrI [8.66; 12.55], mean = 0.82, SD = 0.23); 35.97 years (95% CrI [27.30; 39.69], mean = 0.65, SD = 0.23); and 44.01 years (95% CrI [40.37; 44.43], mean = 0.79, SD = 0.40). In short: we found a rapid initial improvement in gaze following in early childhood, followed by a long period of minor, very slow change with slightly increasing levels of imprecision toward old age.

**Discussion**

We investigated the shape of change in gaze following across the lifespan and found a non-linear developmental trajectory, in which young children quickly enhanced their level of proficiency. Performance peaked around early adulthood, while there was a minor decay in later adulthood. These results support the idea that humans fine-tune their existing gaze following ability after the first emergence in infancy. Furthermore, we observed substantial individual differences in all age groups. While variation was highest in the three- and four-year-olds, it remained relatively stable across the lifespan.

Previous studies found that already four-month-olds demonstrate basic gaze following abilities (Del Bianco et al., 2019). Since we measured an active location choice on a touchscreen, we could not collect data from infants. In our sample, three-year-olds were still rather imprecise in their gaze following ability (average imprecision was approx. two target widths). How can we explain this divergence? First, we used subtle eye movements as cues. Many existing studies let the agents move eye and head in parallel (Behne et al., 2005; Povinelli et al., 1997), establishing a confound with the more salient head movement. Relying exclusively on eye movements might be more difficult for children than presenting them with a combined eye and head orientation (Carpenter et al., 1998). Silverstein et al. (2021) used a similar manipulation of gaze cues without head rotation and found that 6- to 18-month-olds were around or just above chance for gaze following. The authors argue that infants might fixate on another's face most of the time, while eye movement alone might not be strong enough to guide their attention. Furthermore, our study required participants to (1) precisely follow an agent's gaze, (2) interpret this as a cue, and (3) use the cue to guide their own behavior. It is conceivable that three-year-olds followed the agent's gaze but did not translate this into precise, active behavior. Moll and Kadipasaoglu (2013) argue

## MODELING VARIATION IN GAZE FOLLOWING

18

that social forms of perspective-taking evolve prior to visual perspective-taking, which only emerges within the third year of life. Young children might simply not be interested in a differential, spatial representation of the surrounding objects. Taken together, this might explain why our younger participants located the agent's gaze rather imprecisely.

Regarding our sample of elderly adults, we expect a sampling bias (Bethlehem, 2010; Gosling et al., 2004; Remillard et al., 2014). First, certainly not all older people have a high-speed internet connection or are knowledgeable in its use. Second, the elderly adults participating in *Prolific* studies might show greater cognitive flexibility compared to their offline counterparts. Therefore, a representative sample may show a greater age decline in gaze following compared to our reported sample. In addition, older people might be more likely to suffer from visual impairments. Even though we filtered participants to only include normal- to correct-to-normal vision, we cannot guarantee that our participants showed no symptoms of reduced vision.

### **Study 2: Computational cognitive model**

Our lifespan study showed that gaze following develops throughout childhood, and variation between individuals appears in all age groups. The TANGO has previously been shown to reliably capture inter-individual differences in gaze following (Prein et al., 2023). The variation between participants was thus likely genuine and not due to random noise. In Study 2, we aimed to understand the developmental change and individual differences on a process level. We present a theory of gaze following that explains how participants process the available gaze information and trace a line-of-sight to identify the agent's focus. We formalized this inference process in a computational cognitive model that replicates a schematic representation of how participants make inferences in the task's context (i.e., a model of the task and not the data).

The study design and procedure obtained ethical clearance in the same way as Study 1.

The study and the model were pre-registered prior to data collection: <https://osf.io/r3bhn>. Data were collected between May and August 2021.

### **Computational model**

Our model quantifies a participant's cognitive ability to follow gaze by inverting a probabilistic process that generates the participant's clicks from observing the eyes of the agent.

It is formally defined as:

$$P(\theta|x_c, \alpha_l, \alpha_r) \propto P(x_c|\alpha_l, \alpha_r, \theta)P(\theta)$$

where  $\theta$  is an individual's cognitive ability to locate the focus of the agent's attention,  $x_c$  is the coordinate the participant clicked, and  $\alpha_l$  and  $\alpha_r$  are the pupil angles for the left and right eye, respectively. The pupil angle  $\alpha$  is defined as the angle between a line connecting the center of the eye to the pupil and a line extended vertically downward from the center of the eye (see Figure 2A)<sup>4</sup>. Please note that in our case, the center of the pupil is simultaneously the center of the iris (see S. Anstis (2018) for the influence of moving irises, pupils, and corneal reflexes).

---

<sup>4</sup> This model mirrors the logic of the TANGO programming code. In the online experiment, we read out the center point coordinates of the target and the agent's eyeball (i.e., the SVG coordinates), and then calculated a line between these two points: this was our gaze vector (acting in the functionally same way as a pupil angle). Knowing the eyeball radius, we calculated the point of intersection at which the gaze vector met the eyeball boundary. Finally, the agent's pupil moved from the center of the eyeball along the gaze vector to the intersection point. This way, the

## MODELING VARIATION IN GAZE FOLLOWING

20

Based on our verbal task instructions, we assumed that participants (1) expected the agent’s looks to be directed at the target, and (2) to click on the coordinate they estimated the agent to look at. Consequently, we did not assume that participants’ clicks were noisy in any way but that they clicked on the screen location where they genuinely thought the target was (and that the agent was looking at).

The true eye angles ( $\alpha_l$  and  $\alpha_r$ ) cannot be directly observed and have to be estimated based on the position of the pupils within the eyes, resulting in approximate values ( $\widehat{\alpha}_l$  and  $\widehat{\alpha}_r$ ). We presumed this estimation to be a noisy process. Thus, we conceptualized the development of the cognitive ability to follow gaze as a reduction of noise in the estimates (i.e., an increased certainty about the pupil angles).

Any clicked value of  $x_c$  implied a “matched pair” of the estimated pupil angles  $\widehat{\alpha}_l$  and  $\widehat{\alpha}_r$ , with the property that lines extended along those two angles met at the precise location of the target. As a consequence, we can rewrite the likelihood function of the model above:

$$P(x_c | \alpha_l, \alpha_r, \theta) \propto P(\widehat{\alpha}_l, \widehat{\alpha}_r | \alpha_l, \alpha_r, \theta) P(x_c)$$

$P(x_c)$  is a prior over potential target locations, which we assumed to be skewed towards the screen center: We anticipated that participants have an a priori expectation that the target will land close to the middle, because the target was last visible in the screen center before disappearing behind the hedge and because the agent was located centrally on the screen. We

agent was animated to “look at” the target. In the gaze model, we assumed participants go through these steps in reverse order.

estimated the strength of this center bias (i.e., the standard deviation of a Normal distribution around the screen center) based on the data:  $P(x_c) \sim \mathcal{N}(960, \sigma^p)$ .

The width of this distribution is defined by  $\sigma^p$ . For children, we assumed that the center bias changed with age and estimated  $\sigma^p$  via a linear regression as a function of the child's age ( $age_i$ ):  $\sigma^p = \beta_0^{\sigma^p} + age_i \cdot \beta_1^{\sigma^p}$ . Therefore, the participant-specific distribution for  $P(x_c)$  was constrained by the performance in the TANGO and the child's age. For the adults,  $\sigma^p$  was not age-specific.

The main inferential task for the participant lay in estimating the pupil angles, i.e., sampling from the first term of the right-hand side equation above,  $P(\hat{\alpha}_l, \hat{\alpha}_r | \alpha_l, \alpha_r, \theta)$ . For this, we assumed that the pair of estimated pupil angles were sampled from a probability distribution which is the product of two Normal distributions of equal variance,  $\sigma_v$ , centered on the true pupil angles:

$$P(\hat{\alpha}_l, \hat{\alpha}_r | \alpha_l, \alpha_r, \theta) \propto \phi(\hat{\alpha}_l; \alpha_l, \sigma_v) \phi(\hat{\alpha}_r; \alpha_r, \sigma_v),$$

As  $\sigma_v$  determined the level of accuracy with which participants estimated the pupil angles, it is the component of the model that defines  $\theta$ . When  $\sigma_v$  is very small (i.e., the distribution around the pupil angle is narrow), clicks far away from the target are unlikely, as these would require estimated pupil angles very different from the true pupil angles. When  $\sigma_v$  is very large (i.e., the distribution around the pupil angle is wide), almost any pupil angles may be sampled, corresponding to a roughly uniform distribution over click coordinates. We expected  $\sigma_v$  to vary between individuals. Consequently, individuals differed in the level of precision with which they can locate the target based on observing the agent's eyes.

## MODELING VARIATION IN GAZE FOLLOWING

22

The shape of the  $P(x_c | \alpha_l, \alpha_r, \theta)$  distribution leads to a testable group-level prediction. As the pupil location varies, a fixed amount of uncertainty around the pupil angle corresponds to a varying degree of uncertainty in the estimated target location (see Figure 2B & C). When the agent directs their gaze toward the very left or right side, the distribution around the target location from which participants sample is comparatively wider than when the agent gazes centrally to the ground. For illustrative purposes, imagine a similar phenomenon: pointing a torch light to a flat surface on the ground. When one points the light cone directly at the surface, the light beam is concentrated in a clearly defined, small, symmetric area. When one points the light cone further away from oneself (shining at an angle), the light from one half of the cone must travel further to reach the surface than the light from the other half, resulting in an asymmetric light pattern. As the angle increases, the light is spread over a wider area, and the surface is illuminated less evenly. Consequently, for the same  $\sigma_v$ , the further out a target coordinate lies, the wider and less symmetric the distribution. This increases both the variance and the bias in a participant's estimate of the agent's attentional focus, resulting in a decreased performance in the task. As  $\sigma_v$  decreases and the cone narrows, the extent to which performance varies at different angles decreases.

Our gaze model consequently predicted that TANGO trials vary in difficulty (see Figure 2B and C): participants should be more imprecise in locating the target the further out it lands, resulting in a U-shaped pattern<sup>5</sup>. If our data matched the pattern of this model prediction, this

---

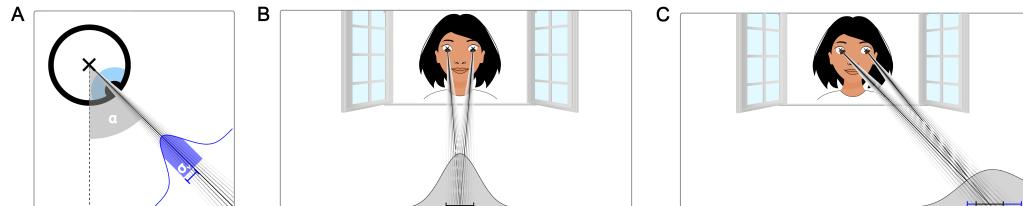
<sup>5</sup> In our screen-based study, this effect should decrease again towards the most outward sides.

Since the computer screen has a natural border, trials in which the target lands furthest out to the

could act as evidence for the gaze model. Therefore, our gaze model provides a quantitative theory of gaze following. In the following, we tested these predictions in children and adults.

---

left/right become slightly easier again. In these cases, the uncertainty about the pupil angle faces practically only the inner side (facing the center) of the screen, since the natural border of the screen limits where participants can click. In another adult sample with more trials, we could recover this pattern. For further elaboration, see Supplements.



**Figure 2: Gaze model** (A) Visualization of the gaze model (simplified, for one eye). Participants are assumed to observe the pupil location and estimate the center of the agent's eye. Connecting these two point estimates as a line yields the unique vector that extends from the center of the agent's eyeball through the center of the pupil to the attentional focal point (line-of-sight). The angle between this vector and a line pointing vertically to the ground (black dashed line) is the pupil angle ( $\alpha$ ). Participants are assumed to sample (grey lines) from Normal distributions (blue line) centered around the true pupil angle. The variance of the Normal distribution ( $\sigma_v$ ) is expected to vary between participants. (B & C) Geometrical features of the gaze model. As the pupil angle varies, a fixed amount of uncertainty in the angle corresponds to a varying degree of uncertainty in the estimated target location. The distribution around the target location from which participants sample is wider when the agent gazes toward the side than when she gazes centrally. The blue line on the ground shows the added level of uncertainty in the estimated target position for the target location further outward (C).

## Participants

The sample consisted of 60 children, including 20 three-year-olds (mean age = 3.47 years, SD = 0.34, range = 3.07 - 3.97, 11 girls), 20 four-year-olds (mean age = 4.61 years, SD = 0.26, range = 4.09 - 4.98, 10 girls), 20 five-year-olds (mean age = 5.66 years, SD = 0.24, range = 5.01 - 5.96, 12 girls). Children were recruited via an internal database, where each parent previously

consented to child development studies, and data was collected in kindergartens in Leipzig, Germany.

In addition, we included 50 adults from Study 1 (mean age = 31.92 years, SD = 12.15, range = 18 - 63, 36 female). Since developmental change was minimal in our adult sample (see Study 1) and the cognitive models were computationally heavy, we decided to only include the first 50 adults who had completed the study.

### Procedure

We applied the same procedure as in Study 1. Children were tested in a quiet room in their kindergarten, while an experimenter guided the child through the study on a tablet. Adults participated online.

### Analysis

We quantified how well our gaze model explained the gaze following process in two ways. First, we aggregated the model predictions and data for each target bin and age group (3-, 4-, 5-years-olds, adults), and computed correlation to quantify how well the model was able to recover the data. Second, we compared the predictions of our gaze model to two simple alternative models that assume participants do not rely on the agent's gaze at all: a random guessing model and a center bias model. The random guessing model assumed participants randomly clicked on the screen and was implemented as sampling from a uniform distribution over all possible coordinates,  $\mathcal{U}(0,1920)$ . The center bias model assumed participants always clicked near the screen center and was implemented as sampling from a Normal distribution with the screen center as the mean, and one balloon width as the variance,  $\mathcal{N}(960,160)$ . Note that the center bias model also predicted imprecision should be higher for targets further out on the

## MODELING VARIATION IN GAZE FOLLOWING

26

screen. However, compared to the gaze model, it predicted a steep effect towards the sides, resulting in a V-shaped pattern (imprecision as the distance between the target location and the screen center). All cognitive models were implemented in WebPPL (Goodman & Stuhlmüller, 2014).<sup>6</sup>

We compared models via the marginal likelihood of the data under each model. The pairwise ratio of marginal likelihoods for two models is also known as the Bayes Factor, which quantifies the quality of a model’s predictions by averaging over the possible values of the model’s parameters weighted by the prior probabilities of those parameter values. It can be used to estimate how much more likely the data under one model are compared to the other. Bayes Factors implicitly consider model complexity: models with more parameters often have broader prior distributions over parameters, which might weaken potential gains in predictive accuracy.

## Results

We found very clear support for our gaze model, both in children as well as adults. A strong correlation between the data mean and the gaze model estimate ( $\sigma_v$ ) showed that the data mean is suitable to quantify individual differences: child sample  $r = 0.95$ , 95%CI [0.92, 0.97];

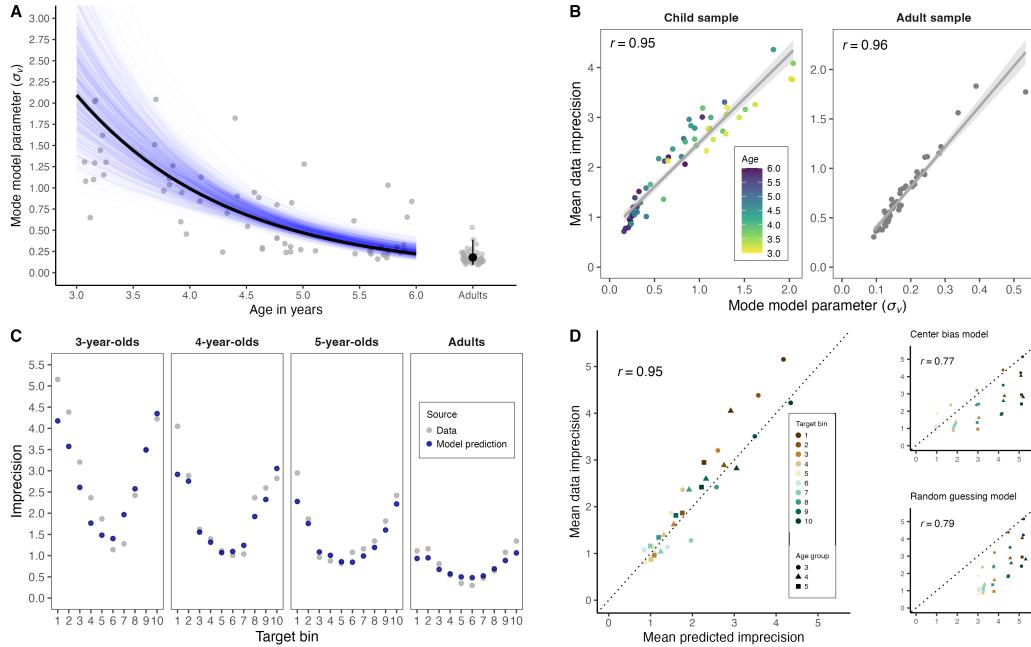
---

<sup>6</sup> In an exploratory analysis, we simulated data for two more alternatives: a line-of-sight tracing model that assumed no inferential noise in the participants’ gaze following ability, and a model building up on the line-of-sight tracing with added motor noise. Please note that these did not predict a U-shape pattern, since all target locations would be influenced by the motor noise equally. For further details, see Supplements.

adult sample  $r = 0.96$ , 95%CI [0.93, 0.98]). The gaze model predicted a U-shaped pattern which we also observed in our data (see Figure 3C). The model comparison strongly favored our gaze model over the center bias model (child sample  $\log BF_{10} = 1,015.33$ ; adult sample  $\log BF_{10} = 2,575.75$ ) and the random guessing model (child sample  $\log BF_{10} = 388.98$ ; adult sample  $\log BF_{10} = 919.03$ ). For the child sample, we went one step further into the analysis. When correlating the observed data across all target positions with the predictions of the three models, we found a high similarity for the gaze model:  $r = 0.95$ , 95%CI [0.90, 0.98], while the correlations with the alternative models were smaller (center bias model:  $r = 0.77$ , 95%CI [0.57, 0.89]; random guessing model:  $r = 0.78$ , 95%CI [0.58, 0.89]). The gaze model's prior over potential target locations assumed that participants' clicks would be skewed towards the screen center. For children, we estimated the width of the prior's distribution based on the participants' age. Interestingly, we found that the differential age effects on the center bias prior were minor (intercept = 300.14, 95%HDI [277.19; 321.28]; slope = 1.18, 95%HDI [0.00; 18.78]). The slope of the prior indicated that older children were only marginally less drawn towards the screen center compared to the younger children.

## MODELING VARIATION IN GAZE FOLLOWING

28



**Figure 3: Gaze model** (A) Developmental trajectory of the estimated model parameter. Grey dots show individual level parameter values. The black line shows the maximum a posteriori (MAP) estimate; blue lines show 1000 draws from it. The large black point with 95% HDI shows the mean of the model parameter for the adult sample. We added minimal horizontal and vertical noise to the adult individual level parameter values to avoid overplotting. (B) Correlation between estimated mode of the model parameter and data mean per individual, faceted by age group. In the child sample, color denotes age in years. Grey regression lines with 95% CI show smooth conditional means based on linear models, with Pearson's correlation coefficient  $r$ . (C) Pattern recovery. Imprecision in target width for each target bin by age group. Model predictions in blue; data in grey. (D) Correlation between the observed data and the predictions of the three models by target position and age group (across individuals; only for the child sample).

## Discussion

We presented a formal cognitive model of gaze following to describe how gaze following develops with age and varies between individuals. We modeled gaze following as a process in which participants estimate pupil angles based on the pupil location within the eye. By following the resulting gaze vector, they consequently arrive at the attentional focus of the agent. Individual differences can be explained as varying levels of imprecision in the pupil angle estimation. We assume the basic process of gaze following to be the same across the lifespan, though individuals become increasingly precise with age. While we did not include an infancy sample, we believe the process of gaze following should operate similarly — if not other cues such as the object saliences or head directions are followed instead. By conducting model comparisons, we ruled out simpler explanations of the data.

In addition, we observed differences in performance depending on gaze location: participant data showed that precision levels dropped as the agent's gaze moved further away from the center. Our gaze model predictions recovered this “signature pattern” in the data. Future research could use this signature in the data as evidence of whether diverse communities employ the same inferential mechanism to solve the task, speaking for a shared cognitive architecture.

Interestingly, the U-shaped pattern in the TANGO task can be conceptually compared to the result patterns of Michelon and Zacks (2006): in their Level 1 perspective-taking task, an increased distance between the agent and the target decreased performance in adults (i.e., reaction times). Targets closer to the midline were more easily traced than ones further away from the agent. The authors concluded that visual acuity is generally higher for locations on the vertical than the diagonal axis (“oblique effect”; (Appelle, 1972; Heeley et al., 1997; Mikellidou et al., 2015)). Similarly, Symons et al. (2004) have found that adults’ acuity for gaze direction depends

MODELING VARIATION IN GAZE FOLLOWING

30

on the target location. Not only is the increased imprecision for TANGO trials in which the target lands further out consistent with this finding, but our gaze model poses a viable explanation for this effect.

A limitation of our model is that we cannot disentangle how much of the participants' uncertainty comes from a noisy estimate of the agent's attentional focus and how much is due to imprecise clicking (e.g., experiencing motor issues, adding random noise to the click). However, we believe imprecise clicking to be of minor concern, since the children in our sample seemed determined in where they clicked and to not have issues with aiming or motor control (see precision in training trials in Supplements).

A critical feature of our model is that it assumes gaze following to rely on vector following: subjects are modeled to calculate pupil angles which serve as gaze vectors to point to the attentional focus point of an agent. Even though this vector following component is a geometrical calculation, one must first interpret the agent's eyes as a relevant social stimulus. Therefore, our model describes gaze following as a particular form of vector following in a social context.

**Study 3: Components of gaze following**

Study 3 examined the components of gaze following and whether it can be fully reduced to physical vector following. The positive link between TANGO and social-environmental factors like age of childcare entry (Prein et al., 2023) underline how social interaction is integral to gaze following. Nevertheless, it is unclear how gaze following relates to other forms of perspective-taking (Astor & Gredebäck, 2022). Therefore, we investigated associations between gaze following and other measures of social-cognitive abilities.

First, we experimentally isolated the vector following component of the TANGO. We designed a new non-social vector following task that shared all crucial design features of the TANGO. Second, we assessed children's social-cognitive abilities by administering a ToM task battery, comprising four tasks from the ToM scale by Wellman and Liu (Wellman & Liu, 2004) and two additional perspective-taking tasks (Flavell, Flavell, et al., 1981; Flavell, Everett, et al., 1981). We reasoned that the TANGO shares task demands with the non-social vector following task while it shares its social context with the ToM tasks. We could, therefore, imagine both an absence or presence of relationship between gaze following and ToM. As stated in our pre-registration, we further assessed whether the two perspective-taking tasks related to gaze following. Our reasoning was that similar underlying mechanisms might be needed to solve these tasks since they require participants to take into account another person's point of view.

The study design and procedure obtained ethical clearance in the same way as Study 1. The study was pre-registered prior to data collection: (<https://osf.io/xsqkt>). Data collection took place in Leipzig, Germany, between February and March 2023.

### **Participants**

The sample consisted of 102 children (mean age = 4.54 years, SD = 0.31, range = 3.99 - 5.03, 54 girls). Information on individual socio-economic status was not formally recorded.

### **Procedure**

Children were tested in a quiet room in their kindergarten. An experimenter guided the child through the study. For maximum control of extraneous participant variables, we employed a within-subjects study design. Participants performed the tasks in this order: (1) non-social vector following task, (2) ToM task battery, (3) TANGO. We decided on a fixed order to compare

MODELING VARIATION IN GAZE FOLLOWING

32

participants' performance straight-forwardly with each other. To increase engagement and decrease fatigue or fuzziness, we switched between tablet tasks and tasks with personal interaction. We presented the non-social vector following task before the TANGO so that participants would not be biased to interpret the stimuli as "agent-like".

**Non-social vector following.** Modeling the structure of the TANGO, we designed a non-social vector following task. This task was presented on a tablet and used the concept of magnetism. On the upper part of the screen, there was a tube with a circular window, containing a gearwheel. On the floor, there was a magnet. The magnet got switched on (with a cartoon-like sound), whereupon the gearwheel moved towards the magnet. The gearwheel moved so that its center aligned with the magnet center while staying inside the circular window. Participants were then asked to locate the magnet. Access to the magnet's true location was manipulated by a wooden wall: participants either had full, partial, or no visual access to the true magnet location. Compared to the TANGO, the circular window acted functionally similar to the agent's eyeball, while the gearwheel acted similar to the pupil. Participants were expected to estimate a vector from the center of the circular window to the gearwheel and extend this as a line toward the ground to locate the magnet. We deliberately decided against displaying an arrow: we aimed to keep the mechanistic functions of the TANGO and magnet stimuli as similar as possible. In both cases, the starting point of the vector needs to be estimated by the participant. With an arrow, we would have drastically reduced the level of uncertainty, since the arrow already displays all information (arrow tip as the "gaze" direction). Furthermore, we wanted to avoid referential or iconic stimuli.

Children received 19 trials with one full visual access trial, two partial visual access trials, and 16 test trials. The first trial of each type comprised a voice-over description of the presented

events. We conducted our analysis with 15 test trials (excluding the voice-over trial). The outcome variable was imprecision, defined as the absolute distance between the magnet's x coordinate and the x coordinate of the participant's click. Magnet coordinates were randomized: The full width of the screen was divided into ten bins; each bin occurred equally often, while the same bin could occur in two consecutive trials; and exact coordinates within each bin were randomly generated.

**ToM task battery.** We administered four tasks from the Wellman and Liu (2004) ToM scale (see Supplements for further detail). We excluded three tasks: the Diverse Desires task to avoid ceiling effects; and both tasks involving emotions (Belief Emotion and Real-Apparent Emotion), as we aimed at assessing the “cold, cognitive” (vs. “emotional”) aspects of social cognition. We added two perspective-taking level-2 tasks (Flavell, Flavell, et al. (1981); Flavell, Everett, et al. (1981); where children were asked whether a turtle appeared to be on its back or feet / a worm lay on a red or blue blanket from the experimenter’s point of view) with the aim of increasing the variability we can capture between individuals, and since we hypothesized that perspective-taking would rely on similar mechanisms than gaze following, both relying on another’s person egocentric frame of reference.

**Gaze following.** As in Study 1 and 2, we presented children with the TANGO (Prein et al., 2023). To accentuate the social aspect of the TANGO, we exchanged the animal agents (used in the previous two studies) with human faces, which were modeled after the local population in appearance (already created for another project on cross-cultural similarities in gaze

## MODELING VARIATION IN GAZE FOLLOWING

following (<https://osf.io/tdsvc>). This further highlighted the contrast (i.e., social vs. non-social context) to the non-social vector following task.<sup>7</sup>

**Analysis**

By design, the TANGO and the non-social vector following task involved vector following. Based on our computational model, we expected children's performance in both tasks to correlate with each other. For each task, we calculated the mean level of imprecision for each subject and correlated them using *Pearson's* correlation coefficient.

For the ToM battery, we aggregated the score of all solved tasks. Please note that the ToM score acted as an umbrella term and included the two perspective-taking tasks. Regarding the relationship between the two vector following tasks and the ToM measures, we could imagine two possible scenarios: (A) If gaze following recruited a social-cognitive ability beyond geometric vector following, we expected that ToM measures would correlate more strongly with the gaze following task than with the non-social vector following task. (B) If gaze following relied purely on task-specific geometric processes, then the correlation between gaze following and ToM measures would be comparable to the correlation between non-social vector following and the ToM measures. For the association between the aggregate ToM scores and the gaze following / non-social vector following tasks, we used *Spearman's* rank correlation coefficients.

---

<sup>7</sup> In an exploratory analysis, we compared children's imprecision levels in the TANGO task with animal vs. human agents. Based on a GLMM analysis, we conclude that there was no evidence of a stable effect of stimulus choice. See Supplements for further detail.

We compared the correlation between gaze following and ToM measures and the correlation between non-social vector following and ToM measures by using the Williams' test from the function `cocor.dep.groups.overlap` (designed for two dependent overlapping correlations) from the package `cocor` (Diedenhofen & Musch, 2015).

To estimate which components best explain the gaze following score, we conducted a model comparison with GLMMs predicting the mean imprecision in gaze following by age, imprecision in non-social vector following, the ToM aggregate score, or the aggregate of the two perspective-taking tasks (subset of ToM battery; example of model notation in R: `tango_mean ~ age_centered + magnet_scaled + perspective_scaled`). The outcome variable was modeled by a lognormal distribution. We wanted to assess whether the ToM aggregate score or the singled-out perspective-taking score added additional explanatory value when predicting the gaze following score. We hypothesized that perspective-taking seemed most closely theoretically related to gaze following as in both cases the participant was asked to judge another person's point of view.

## Results

Children performed similarly well in the social and non-social vector following tasks (gaze following mean = 1.92,  $sd = 0.77$ , range = [0.57 — 4.49]; vector following mean = 1.90,  $sd = 1.01$ , range = [0.57 — 5.83]). The mean aggregate ToM score was 3.46 ( $sd = 1.53$ , range = [0 — 6]), while the mean aggregate perspective-taking score was 1.38 ( $sd = 0.81$ , range = [0 — 2]).

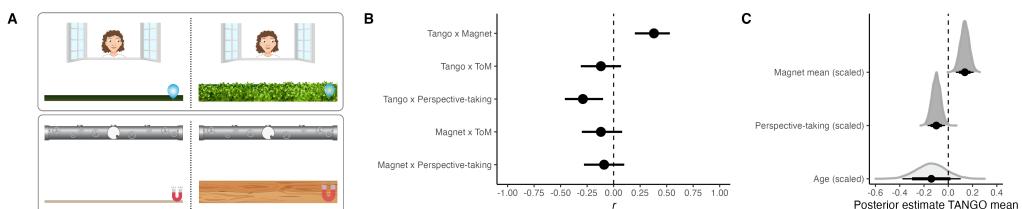
Gaze following substantially correlated with the non-social vector following task,  $r = 0.38$ , 95%CI [0.20, 0.53] (see Figure 4B). While the tasks highly overlap in task demands, their measures shared some, but not all of their variance.

MODELING VARIATION IN GAZE FOLLOWING

ToM abilities did not correlate with gaze following ( $\rho = -0.12$ , 95%CI [-0.31, 0.07]) or non-social vector following ( $\rho = -0.12$ , 95%CI [-0.30, 0.08]), and the correlations did not differ from each other, Williams' test  $t(99) = 0$ ,  $p = 1$ .

Interestingly, gaze following and perspective-taking correlated with each other,  $\rho = -0.29$ , 95%CI [-0.46, -0.10]. Please note that the TANGO quantifies imprecision in gaze following. Therefore, a negative correlation suggests that imprecision in gaze following corresponds to less perspective-taking. Non-social vector following and perspective-taking did not correlate,  $\rho = -0.09$ , 95%CI [-0.28, 0.10]. However, according to the Williams' test, the two correlations did not differ significantly from each other,  $t(99) = -1.86$ ,  $p = 0.07$ .

Our model comparison revealed that gaze following was best predicted by a model including non-social vector following ( $\beta = 0.14$ , 95% CrI [0.06; 0.21]) and perspective-taking ( $\beta = -0.10$ ; 95% CrI [-0.17, -0.03]), even when controlling for age ( $\beta = -0.14$ , 95% CrI [-0.38, 0.10]) (see Figure 4C and Supplements for the model comparison).



**Figure 4: Components of gaze following.** (A) Study procedures. Top: TANGO (i.e., gaze following; social vector following). Bottom: Magnet (i.e., non-social vector following). The left-hand side shows screenshots of the familiarization phase; right-hand side shows screenshots of the test phase. For illustrative purposes, translucent targets are shown in the test phase. Participants cannot see them. (B) Correlations between gaze following, non-social vector

following, ToM, and perspective-taking. Dots show the correlation coefficients, error bars represent 95% CIs. Please note that the ToM score acts as an umbrella term and also includes two perspective-taking tasks. (C) Influence of age, perspective-taking and non-social vector following on gaze following. The graph shows the posterior distributions for the respective predictor. Black dots represent means, thicker black lines 80% CrI, and thinner black lines 95% CrI.

## Discussion

Our gaze model assumed vector calculations on a process level. By isolating non-social vector following experimentally, we could show that gaze following does indeed, to a certain degree, rely on this component. However, non-social vector following alone did not suffice to explain variation in gaze following. Additionally, perspective-taking proved to be a relevant social-cognitive ability when predicting the performance in the TANGO. The overall ToM aggregate score did not add explanatory power.

The TANGO and the two perspective-taking tasks can be seen as instances of visuospatial perspective-taking. Often, researchers distinguish between Level 1 and Level 2 perspective-taking tasks. Level 1 perspective-taking is concerned with the visibility of objects from a particular viewpoint (Surtees et al., 2013), and can usually be solved independently from another's agent frame of reference and without mental rotation [not matching the “mentalizing criterion”; Quesque and Rossetti (2020)]. In contrast, Level 2 perspective-taking tasks ask participants to judge how (visually or conceptually) the world looks for another person. Solving these tasks presumably relies on a simulation of mentally transforming one's own body schema into the space of another agent (Erle & Topolinski, 2017). Participants must understand that different viewpoints lead to different perspectives (Flavell, Everett, et al., 1981): two agents can see the same object “in different, incompatible, fine-grained ways” (Rakoczy, 2022). In short,

## MODELING VARIATION IN GAZE FOLLOWING

38

Level 1 perspective-taking addresses *what* an agent can see, while Level 2 perspective-taking addresses *how or where* they see it (Moll & Tomasello, 2006). The TANGO falls into the first category - however, in contrast to existing research, on a continuous scale - while the perspective-taking tasks presented within the ToM scale fall into the second category.

In the TANGO, participants locate the target by assuming that another agent perceives and looks at the it. The here-applied Level 2 perspective-taking tasks add another layer by asking how one's perspective differs from another's and how exactly the other person sees an object.

Michelon and Zacks (2006) assessed the different processes that participants used to master Level 1 and Level 2 perspective-taking tasks and concluded that participants rely on line-of-sight tracing in the former and perspective transformation in the latter case. Line-of-sight tracing consists of locating the avatar and the target, and drawing a line between them. Note how our gaze model in Study 2 shares the underlying idea of connecting points in space via a line. Level 2 perspective-taking requires an additional step in which participants “perform a perspective transformation so that one’s imagined position matches the position of the avatar” (Michelon & Zacks, 2006). Therefore, the processes to solve Level 1 perspective-taking tasks might be computationally lighter since there is no need to adapt a new reference frame and there is no potential conflict between one’s own and the other person’s reference frame. Still, the assumed processes overlap, which could explain the correlation between the TANGO and the administered perspective-taking tasks.

Surtees et al. (2013) further differentiated between visual and spatial perspective-taking. While the former helps to judge if and how an agent sees an object, the latter involves judging the relative spatial locations of an agent and the object. Spatial perspective-taking does not necessitate mental states since computing a line of sight does not demand the presence of another

agent (Michelon & Zacks, 2006) and can, therefore, be applied to non-agentive objects with a front (Surtees et al., 2013). This could explain how our participants solved the non-social vector following task.

Interestingly, we found weaker correlations between gaze following and the other ToM tasks. Similarly, in a longitudinal study, Brooks and Meltzoff (2015) found no direct association between gaze following at 10.5 months and explicit ToM at 4.5 years. While the ToM tasks and the TANGO shared the social context, the cognitive processes needed to solve each task might vary. As Rakoczy (2022) reflected, perception-goal psychology (which includes gaze following) comprises understanding that others see different objects or pursue different goals. However, this ability does not necessarily entail understanding more complex meta-representational aspects; for example, understanding that mental states can be false or involve aspectual information.

In previous work, we established that the TANGO is suited to capture meaningful variability across individuals (Prein et al., 2023), which is a crucial task feature when we are interested in revealing the relationship between different cognitive abilities. Importantly, the tasks we used to measure ToM abilities were not designed to capture individual differences: the aggregate score of few dichotomous items is of limited use when it comes to quantifying genuine differences between individuals. However, since these tasks are the gold standard in the social-cognitive literature (Bialecka-Pikul et al., 2021; Byom & Mutlu, 2013; Poulin-Dubois et al., 2023; Rakoczy, 2022; Wellman, 2018), and measures with satisfying psychometric properties are, to the best of our knowledge, still scarce (e.g., Beaudoin et al., 2020; Mayes et al., 1996), we nonetheless relied on them in this study. Thus, lower correlations between ToM abilities and gaze following may reflect poor measurement characteristics on the side of ToM tasks rather than a

MODELING VARIATION IN GAZE FOLLOWING

40

genuine absence of association. We would like to point out that we already stated this concern in our pre-registration (<https://osf.io/xsqkt>).

If the reliability of the tasks at hand is known, one can estimate the “true” correlation between the latent constructs by applying an attenuation formula or structural equation models (Metsämuuronen, 2022; Trafimow, 2016). Adjusting for the measurement error would increase the so-called true correlation. While we can estimate the split-half reliability for the TANGO (Prein et al., 2023) and the non-social vector following task, we do not have reliability estimates of the ToM measures and cannot apply said approaches. This, in turn, underlines the importance of reporting the psychometric properties of a task. The development of new measures to capture individual differences in social-cognitive abilities seems essential to move this research further.

**General discussion**

In three studies, we shed light on the cognitive process underlying gaze following and its developmental trajectory across the lifespan. Study 1 focused on how gaze following changes with age. We found a steep performance improvement in the preschool years in which children became more precise in locating the attentional focus of an agent. During teenage years and early adulthood, participants reached their peak performance. Precision levels then stayed comparably stable, with a minor decay toward older adulthood. Beyond these aggregated developmental patterns, we found that individual differences exist throughout the lifespan. In Study 2, we proposed a computational cognitive model that described gaze following at a process level. We modeled gaze following as a process in which participants use the pupil location within the eyes to estimate pupil angles. To locate the attentional focus of an agent (and find a target), they extend the resulting gaze vector towards the ground. Individuals vary in their levels of

uncertainty around these pupil angles. Our gaze model outperformed two alternative models, which assumed participants solved the task via a center bias or random guessing. Knowing the TANGO to be a reliable individual differences measure, we investigated potential components of gaze following in Study 3. A fundamental assumption of our computational model was that gaze following relies on a vector following process. We experimentally isolated this component by designing a non-social vector following task. Furthermore, we assessed the relationship between gaze following and traditional ToM tasks. We found that gaze following does, indeed, share a substantial part of its variance with the non-social counterpart of vector following. Additionally, perspective-taking correlated with gaze following, whereas the other ToM measures (focusing on diverse desires, knowledge access, and false beliefs) did not.

The developmental trajectory seen in Study 1 shows how abilities that emerge in infancy can continue to develop throughout childhood. While previous research established that four-month-olds can follow gaze toward one out of two objects, this is not the end point of development. By studying gaze following on a continuum, we assessed not just *whether*, but *how precisely* children locate others' attentional focus. In previous work, we have shown that these individual differences are meaningful (e.g., connected to theoretically related constructs, and showing high split half and retest reliability; (Prein et al., 2023)). Capturing individual variation is crucial when we study development and the improvement in social-cognitive abilities.

Preschool children increased their precision level to locate an agent's attentional focus, which then stayed comparably stable across adulthood. Older adults decreased slightly in their precision levels. This developmental trajectory of a first emergence with a rapid improvement, followed by a plateau and slight decline toward older age, might be representative of many cognitive processes.

MODELING VARIATION IN GAZE FOLLOWING

42

In Study 2, we proposed a theoretical framework to interpret the development and individual differences in gaze following. Our computational cognitive model assumes that participants estimate a pupil angle (i.e., the angle between a line extending vertically downwards from the pupil center and a line connecting the pupil and eye center; line-of-sight; see Figure 2A). The model parameter estimates a participant’s latent ability to follow gaze and can explain why individuals differ in their precision to locate an agent’s attentional focus. The model proposes that development in gaze following equals a reduction in noise when estimating the agent’s pupil angles. We found strong evidence for the proposed gaze model when comparing it against two alternative models and correlating its predictions with the observed data (see Figure 3D), both for children and adults. Notably, the model recovers signature patterns in the data (see Figure 3C).

In Study 3, we tested the relation between gaze following and non-social vector following. As implied by our gaze model, we found that gaze following relates to the ability to estimate vectors in space. Already Butterworth and Jarrett (1991) have argued that, as the scene of actions, space is the commonality between different minds. The spatial vector following ability might be helpful in several social-cognitive tasks, for example, action prediction (Friesen & Rao, 2011), and intention understanding. Predicting which object another agent likely wants to grasp or calculating their movement pathway could rely on similar vector following abilities.

Gaze as displayed in the TANGO versus as in real-life social interactions differs with regard to which information is available. In the TANGO, the agent’s eyes are big and round, with uniform colors, white sclera, fully visible iris and pupils, and continuous eye movement. Eye gaze in natural social interactions often provides less information and is more ambiguous. Even though the TANGO and our gaze model are designed within a 2D world, we believe the mechanisms can be extrapolated into the 3D world. The processes of understanding gaze in daily

life likely rely on the same principles as proposed in this paper. In Study 3, we presented the first evidence that this might be the case. We administered Level 2 perspective-taking tasks in which participants needed to adapt another person's frame of reference in a real-world social interaction. The correlation between this task and the TANGO speaks toward a unified mechanism behind these two visual perspective-taking tasks, regardless of the testing setup or stimulus features. Clearly, our real-life environment is visually more cluttered and diverse than the one presented in our tablet task. Here, however, additional informational sources are available to infer where others are looking; for example, body or head orientation or common ground (Bohn & Köyメン, 2018; Moll & Kadipasaoglu, 2013; Osborne-Crowley, 2020). From a modeling perspective, a shared interaction history or diverse desires might be represented as non-uniform prior distributions over locations in the visual scenery. This way, our gaze model could be expanded to include more complex processes like mental state reasoning.

Theories of gaze following differ in whether they illustrate why children pay attention to gaze in the first place versus how they identify the exact location of gaze. While we described the process behind children's increasing precision in gaze following, we still need to further explore the driving forces behind this development. Existing theories broadly vary in how much importance they place on (A) experience and social environment, and (B) social awareness vs. domain-general learning mechanisms (see categorization by Astor and Gredebäck (2022)). For example, it has been hypothesized that infants might be equipped with an innate gaze module and special neural mechanisms to detect eyes (Baron-Cohen, 1995; Batki et al., 2000); that children identify contingencies in social interactions and, in these, get reinforced to follow gaze (Corkum & Moore, 1998; Silverstein et al., 2021; Triesch et al., 2006); or that children are intrinsically socially motivated and simulate their own experiences to understand others (Astor et al., 2020; Friesen & Rao, 2011; Ishikawa et al., 2020; Meltzoff, 2007; Tomasello, 1999). As seen

## MODELING VARIATION IN GAZE FOLLOWING

44

in Prein et al. (2023), precision in gaze following is linked to receptive vocabulary and opportunities for social interaction (e.g., number of siblings and age when entering childcare). Which exact kind of interactions are most helpful to improve precision in gaze following remains unknown.

### Limitations

In this paper, we have focused on studying variation across ages and individuals. However, our findings rely on participants from Western, Educated, Industrialized, Rich, and Democratic (WEIRD) backgrounds (Henrich et al., 2010). Cultural variation has been found in many foundational aspects of cognition and socialization: for example, in parent-child interaction and communication (Nielsen et al., 2017). First evidence suggests that cultural variation in face-to-face interactions does not influence infants' gaze shifts (Hernik & Broesch, 2019). While we cannot generalize our here-reported developmental trajectory to different socio-environmental settings, we predict that our presented process-model of gaze following holds true across communities. Analyzing cross-cultural variation and checking for predicted "signature patterns" in the data will inform our modeling work and theory building in further detail (Amir & McAuliffe, 2020).

In Study 1, we recruited older participants online which might have selected a particular subgroup of this age range. Seventy-year-olds who have working Wi-Fi connections, know how to use a computer, and are registered on *Prolific* might not be representative of their age group. We can imagine that results from a more diverse, in-person data collection might show different developmental trajectories toward old age.

Our computational model of gaze following estimates one person-specific parameter for how accurately participants locate another person's attentional focus. The model assumes no motor imprecision in this estimation. However, younger children could have located the agent's focus at one particular point but clicked somewhere slightly off for motor control reasons. This would blur the model's estimation of the inferential component. However, in the first training trial, in which children were simply asked to touch the balloon, we found nearly perfect precision levels (cf. Prein et al., 2023). Motor issues and inaccurate aiming, resulting in falsely wide estimations in the model's inferential component, seem unlikely.

In Study 3, we matched the non-social vector following task as closely as possible to the TANGO. However, the starting positions differ: the magnet never appeared in the center of the screen. The starting point of the balloon might be especially important when interpreting the U-shaped pattern in Study 2. Furthermore, in the TANGO, two eyes are presented and information of the two (matching) cues needs to be integrated to infer the target's location. In the non-social vector following task, only one circular window with a gearwheel inside is presented as a directional cue, and there is no need to integrate two different information sources. In addition, we want to mention that the gaze presented in the TANGO might be more prominent compared to real-life social interactions (e.g., perfectly round and visible sclera; see discussion above). Future research should investigate how factors like self-propelled movement, spatial layout, and number of information sources influence the mechanisms of gaze following.

### **Conclusion**

In three studies, we have illuminated the lifelong development of precisely estimating another's gaze direction, and the psycho-physical process behind this. We have shown that gaze

MODELING VARIATION IN GAZE FOLLOWING

46

following continues to develop beyond infancy, and that individuals differ in their precision levels to localize the gaze direction of an agent. Our proposed process-level theory of gaze following modeled individual differences in precision as varying levels of uncertainty in the estimated gaze vectors. Consequently, we found that imprecision in gaze following relates to non-social vector following, as proposed by the model. Additionally, gaze following was linked to visual perspective-taking but no other aspects of ToM. The present research shows how precise and reliable measures and process models jointly inform each other and lead to a more comprehensive understanding of the psychological phenomenon in question.

**References**

- Amir, D., & McAuliffe, K. (2020). Cross-cultural, developmental psychology: Integrating approaches and key insights. *Evolution and Human Behavior*, 41(5), 430–444.  
<https://doi.org/10.1016/j.evolhumbehav.2020.06.006>
- Anstis, S. (2018). The Role of the Pupil, Corneal Reflex, and Iris in Determining the Perceived Direction of Gaze. *I-Perception*, 9(4), 2041669518765852.  
<https://doi.org/10.1177/2041669518765852>
- Anstis, S. M., Mayhew, J. W., & Morley, T. (1969). The Perception of Where a Face or Television 'Portrait' Is Looking. *The American Journal of Psychology*, 82(4), 474–489.  
<https://doi.org/10.2307/1420441>
- Appelle, S. (1972). Perception and discrimination as a function of stimulus orientation: The "oblique effect" in man and animals. *Psychological Bulletin*, 78(4), 266–278.  
<https://doi.org/10.1037/h0033117>
- Aslin, R. N. (2007). What's in a look? *Developmental Science*, 10(1), 48–53.  
<https://doi.org/10.1111/j.1467-7687.2007.00563.x>
- Astor, K., & Gredebäck, G. (2022). Gaze following in infancy: Five big questions that the field should answer. In *Advances in Child Development and Behavior* (p. S0065240722000192). Elsevier. <https://doi.org/10.1016/bs.acdb.2022.04.003>
- Astor, K., Lindskog, M., Forssman, L., Kenward, B., Fransson, M., Skalkidou, A., Tharner, A., Cassé, J., & Gredebäck, G. (2020). Social and emotional contexts predict the development

MODELING VARIATION IN GAZE FOLLOWING

48

of gaze following in early infancy. *Royal Society Open Science*, 7(9), 201178.

<https://doi.org/10.1098/rsos.201178>

Astor, K., Thiele, M., & Gredebäck, G. (2021). Gaze following emergence relies on both perceptual cues and social awareness. *Cognitive Development*, 60, 101121.

<https://doi.org/10.1016/j.cogdev.2021.101121>

Baron-Cohen, S. (1995). *Mindblindness: An essay on autism and theory of mind*. The MIT Press.

Batki, A., Baron-Cohen, S., Wheelwright, S., Connellan, J., & Ahluwalia, J. (2000). Is there an innate gaze module? Evidence from human neonates. *Infant Behavior and Development*, 23(2), 223–229. [https://doi.org/10.1016/S0163-6383\(01\)00037-6](https://doi.org/10.1016/S0163-6383(01)00037-6)

Beaudoin, C., Leblanc, É., Gagner, C., & Beauchamp, M. H. (2020). Systematic Review and Inventory of Theory of Mind Measures for Young Children. *Frontiers in Psychology*, 10, 2905. <https://doi.org/10.3389/fpsyg.2019.02905>

Behne, T., Carpenter, M., & Tomasello, M. (2005). One-year-olds comprehend the communicative intentions behind gestures in a hiding game. *Developmental Science*, 8(6), 492–499. <https://doi.org/10.1111/j.1467-7687.2005.00440.x>

Bethlehem, J. (2010). Selection Bias in Web Surveys. *International Statistical Review*, 78(2), 161–188. <https://doi.org/10.1111/j.1751-5823.2010.00112.x>

Białecka-Pikul, M., Białek, A., Kosno, M., Stępień-Nycz, M., Blukacz, M., & Zubek, J. (2021). Early mindreading scale: From joint attention to false-belief understanding. *European Journal of Developmental Psychology*, 0(0), 1–18. <https://doi.org/10.1080/17405629.2021.1911799>

- Birch, S. A. J., Li, V., Haddock, T., Ghrear, S. E., Brosseau-Liard, P., Baimel, A., & Whyte, M. (2017). Perspectives on Perspective Taking. In *Advances in Child Development and Behavior* (Vol. 52, pp. 185–226). Elsevier. <https://doi.org/10.1016/bs.acdb.2016.10.005>
- Bohn, M., & Frank, M. C. (2019). The Pervasive Role of Pragmatics in Early Language. *Annual Review of Developmental Psychology*, 1(1), 223–249. <https://doi.org/10.1146/annurev-devpsych-121318-085037>
- Bohn, M., & Köymen, B. (2018). Common Ground and Development. *Child Development Perspectives*, 12(2), 104–108. <https://doi.org/10.1111/cdep.12269>
- Brooks, R., & Meltzoff, A. N. (2002). The importance of eyes: How infants interpret adult looking behavior. *Developmental Psychology*, 38(6), 958–966. <https://doi.org/10.1037/0012-1649.38.6.958>
- Brooks, R., & Meltzoff, A. N. (2015). Connecting the dots from infancy to childhood: A longitudinal study connecting gaze following, language, and explicit theory of mind. *Journal of Experimental Child Psychology*, 130, 67–78. <https://doi.org/10.1016/j.jecp.2014.09.010>
- Bürkner, P.-C. (2017). Brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal of Statistical Software*, 80(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- Bürkner, P.-C. (2018). Advanced Bayesian Multilevel Modeling with the R Package brms. *The R Journal*, 10(1), 395. <https://doi.org/10.32614/RJ-2018-017>

MODELING VARIATION IN GAZE FOLLOWING

50

Butterworth, G., & Jarrett, N. (1991). What minds have in common is space: Spatial mechanisms serving joint visual attention in infancy. *British Journal of Developmental Psychology*, 9(1), 55–72. <https://doi.org/10.1111/j.2044-835X.1991.tb00862.x>

Byers-Heinlein, K., Tsui, R. K.-Y., van Renswoude, D., Black, A. K., Barr, R., Brown, A., Colomer, M., Durrant, S., Gampe, A., Gonzalez-Gomez, N., Hay, J. F., Hernik, M., Jartó, M., Kovács, Á. M., Laoun-Rubenstein, A., Lew-Williams, C., Liszkowski, U., Liu, L., Noble, C., ... Singh, L. (2021). The development of gaze following in monolingual and bilingual infants: A multi-laboratory study. *Infancy*, 26(1), 4–38. <https://doi.org/10.1111/infa.12360>

Byom, L., & Mutlu, B. (2013). Theory of mind: Mechanisms, methods, and new directions. *Frontiers in Human Neuroscience*, 7.

Carpenter, M., Nagell, K., & Tomasello, M. (1998). Social cognition, joint attention, and communicative competence from 9 to 15 months of age. *Monographs of the Society for Research in Child Development*, 63(4), i–vi, 1–143.

Corkum, V., & Moore, C. (1995). Development of joint visual attention in infants. In *Joint attention: Its origins and role in development* (pp. 61–83). Lawrence Erlbaum Associates, Inc.

Corkum, V., & Moore, C. (1998). The origins of joint visual attention in infants. *Developmental Psychology*, 34(1), 28–38. <https://doi.org/10.1037/0012-1649.34.1.28>

D'Entremont, B. (2000). A perceptual–attentional explanation of gaze following in 3- and 6-month-olds. *Developmental Science*, 3(3), 302–311. <https://doi.org/10.1111/1467-7687.00124>

D'Entremont, B., Hains, S. M. J., & Muir, D. W. (1997). A demonstration of gaze following in 3- to 6-month-olds. *Infant Behavior and Development*, 20(4), 569–572. [https://doi.org/10.1016/S0163-6383\(97\)90048-5](https://doi.org/10.1016/S0163-6383(97)90048-5)

Deák, G. O., Flom, R. A., & Pick, A. D. (2000). Effects of gesture and target on 12- and 18-month-olds' joint visual attention to objects in front of or behind them. *Developmental Psychology*, 36(4), 511–523.

Del Bianco, T., Falck-Ytter, T., Thorup, E., & Gredebäck, G. (2019). The Developmental Origins of Gaze-Following in Human Infants. *Infancy*, 24(3), 433–454. <https://doi.org/10.1111/infa.12276>

Diedenhofen, B., & Musch, J. (2015). Cocor: A Comprehensive Solution for the Statistical Comparison of Correlations. *PLoS ONE*, 10(4), e0121945. <https://doi.org/10.1371/journal.pone.0121945>

Emery, N. J. (2000). The eyes have it: The neuroethology, function and evolution of social gaze. *Neuroscience & Biobehavioral Reviews*, 24(6), 581–604. [https://doi.org/10.1016/S0149-7634\(00\)00025-7](https://doi.org/10.1016/S0149-7634(00)00025-7)

Erle, T. M., & Topolinski, S. (2017). The grounded nature of psychological perspective-taking. *Journal of Personality and Social Psychology*, 112(5), 683–695. <https://doi.org/10.1037/pspa0000081>

MODELING VARIATION IN GAZE FOLLOWING 52

Flavell, J. H., Everett, B. A., Croft, K., & Flavell, E. R. (1981). Young children's knowledge about visual perception: Further evidence for the Level 1–Level 2 distinction.

*Developmental Psychology, 17*, 99–103. <https://doi.org/10.1037/0012-1649.17.1.99>

Flavell, J. H., Flavell, E. R., Green, F. L., & Wilcox, S. A. (1981). The Development of Three Spatial Perspective-Taking Rules. *Child Development, 52*(1), 356–358.

<https://doi.org/10.2307/1129250>

Friesen, A. L., & Rao, R. P. N. (2011). Gaze Following as Goal Inference: A Bayesian Model.

*Proceedings of the Annual Meeting of the Cognitive Science Society, 33*(33).

Frith, C. D., & Frith, U. (2012). Mechanisms of Social Cognition. *Annual Review of Psychology, 63*(1), 287–313. <https://doi.org/10.1146/annurev-psych-120710-100449>

Gamer, M., & Hecht, H. (2007). Are you looking at me? Measuring the cone of gaze. *Journal of Experimental Psychology: Human Perception and Performance, 33*(3), 705–715.

<https://doi.org/10.1037/0096-1523.33.3.705>

Gathercole, S. E., Pickering, S. J., Ambridge, B., & Wearing, H. (2004). The Structure of Working Memory From 4 to 15 Years of Age. *Developmental Psychology, 40*, 177–190.

<https://doi.org/10.1037/0012-1649.40.2.177>

Goodman, N. D., & Stuhlmüller, A. (2014). *The Design and Implementation of Probabilistic Programming Languages*.

Gosling, S. D., Vazire, S., Srivastava, S., & John, O. P. (2004). Should We Trust Web-Based Studies? A Comparative Analysis of Six Preconceptions About Internet Questionnaires.

*American Psychologist, 59*(2), 93–104. <https://doi.org/10.1037/0003-066X.59.2.93>

Gredebäck, G., Fikke, L., & Melinder, A. (2010). The development of joint visual attention: A longitudinal study of gaze following during interactions with mothers and strangers.

*Developmental Science*, 13(6), 839–848. <https://doi.org/10.1111/j.1467-7687.2009.00945.x>

Guterstam, A., Kean, H. H., Webb, T. W., Kean, F. S., & Graziano, M. S. A. (2019). Implicit model of other people's visual attention as an invisible, force-carrying beam projecting from the eyes. *Proceedings of the National Academy of Sciences of the United States of America*, 116(1), 328–333. <https://doi.org/10.1073/pnas.1816581115>

Heeley, D. W., Buchanan-Smith, H. M., Cromwell, J. A., & Wright, J. S. (1997). The oblique effect in orientation acuity. *Vision Research*, 37(2), 235–242.

[https://doi.org/10.1016/S0042-6989\(96\)00097-1](https://doi.org/10.1016/S0042-6989(96)00097-1)

Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world? *The Behavioral and Brain Sciences*, 33(2-3), 61-83; discussion 83-135.

<https://doi.org/10.1017/S0140525X0999152X>

Hernik, M., & Broesch, T. (2019). Infant gaze following depends on communicative signals: An eye-tracking study of 5- to 7-month-olds in Vanuatu. *Developmental Science*, 22(4), e12779. <https://doi.org/10.1111/desc.12779>

Hessels, R. S. (2020). How does gaze to faces support face-to-face interaction? A review and perspective. *Psychonomic Bulletin & Review*, 27(5), 856–881.

<https://doi.org/10.3758/s13423-020-01715-w>

MODELING VARIATION IN GAZE FOLLOWING

54

Horstmann, G., & Linke, L. (2021). Examining Gaze Cone Shape and Size. *Perception*, 50(12), 1056–1065. <https://doi.org/10.1177/03010066211059930>

Ishikawa, M., Senju, A., & Itakura, S. (2020). Learning Process of Gaze Following: Computational Modeling Based on Reinforcement Learning. *Frontiers in Psychology*, 11.

Ishikawa, M., Senju, A., Kato, M., & Itakura, S. (2022). Physiological arousal explains infant gaze following in various social contexts. *Royal Society Open Science*, 9(8), 220592. <https://doi.org/10.1098/rsos.220592>

Jasso, H., & Triesch, J. (2006). Using eye direction cues for gaze following – A developmental model. In.

Lau, B., & Triesch, J. (2004). Learning gaze following in space: A computational model. In *Proceedings of the Third International Conference on Development and Learning*.

Lempers, J. D. (1979). Young Children's Production and Comprehension of Nonverbal Deictic Behaviors. *The Journal of Genetic Psychology*, 135(1), 93–102. <https://doi.org/10.1080/00221325.1979.10533420>

Lempers, J. D., Flavell, E. R., & Flavell, J. H. (1977). The development in very young children of tacit knowledge concerning visual perception. *Genetic Psychology Monographs*, 95(1), 3–53.

Loomis, J. M., Kelly, J. W., Pusch, M., Bailenson, J. N., & Beall, A. C. (2008). Psychophysics of perceiving eye-gaze and head direction with peripheral vision: Implications for the dynamics of eye-gaze behavior. *Perception*, 37(9), 1443–1457. <https://doi.org/10.1080/p5896>

- Mayes, L. C., Klin, A., Tercyak, K. P., Cicchetti, D. V., & Cohen, D. J. (1996). Test-Retest Reliability for False-Belief Tasks. *Journal of Child Psychology and Psychiatry*, 37(3), 313–319. <https://doi.org/10.1111/j.1469-7610.1996.tb01408.x>
- Meltzoff, A. N. (2007). “Like me”: A foundation for social cognition. *Developmental Science*, 10(1), 126–134. <https://doi.org/10.1111/j.1467-7687.2007.00574.x>
- Meltzoff, A. N., Brooks, R., Shon, A. P., & Rao, R. P. N. (2010). “Social” robots are psychological agents for infants: A test of gaze following. *Neural Networks*, 23(8), 966–972. <https://doi.org/10.1016/j.neunet.2010.09.005>
- Metsämuuronen, J. (2022). Attenuation-Corrected Estimators of Reliability. *Applied Psychological Measurement*, 46(8), 720–737.  
<https://doi.org/10.1177/01466216221108131>
- Michel, C., Kayhan, E., Pauen, S., & Hoehl, S. (2021). Effects of Reinforcement Learning on Gaze Following of Gaze and Head Direction in Early Infancy: An Interactive Eye-Tracking Study. *Child Development*, 92(4), e364–e382.  
<https://doi.org/10.1111/cdev.13497>
- Michelon, P., & Zacks, J. M. (2006). Two kinds of visual perspective taking. *Perception & Psychophysics*, 68(2), 327–337. <https://doi.org/10.3758/BF03193680>
- Mikellidou, K., Cicchini, G. M., Thompson, P. G., & Burr, D. C. (2015). The oblique effect is both allocentric and egocentric. *Journal of Vision*, 15(8), 24.  
<https://doi.org/10.1167/15.8.24>

MODELING VARIATION IN GAZE FOLLOWING

56

Moll, H., & Kadipasaoglu, D. (2013). The primacy of social over visual perspective-taking.

*Frontiers in Human Neuroscience*, 7.

Moll, H., & Meltzoff, A. (2011). Perspective-Taking and its Foundation in Joint Attention. In J.

Roessler, H. Lerman, & N. Eilan (Eds.), *Perception, Causation, and Objectivity* (p. 0).

Oxford University Press. <https://doi.org/10.1093/acprof:oso/9780199692040.003.0016>

Moll, H., & Tomasello, M. (2004). 12- and 18-month-old infants follow gaze to spaces behind

barriers. *Developmental Science*, 7(1), F1–F9. <https://doi.org/10.1111/j.1467-7687.2004.00315.x>

Moll, H., & Tomasello, M. (2006). Level 1 perspective-taking at 24 months of age. *British*

*Journal of Developmental Psychology*, 24(3), 603–613.

<https://doi.org/10.1348/026151005X55370>

Moore, C. (2008). The Development of Gaze Following. *Child Development Perspectives*, 2(2),

66–70. <https://doi.org/10.1111/j.1750-8606.2008.00052.x>

Nielsen, M., Haun, D., Kärtner, J., & Legare, C. H. (2017). The persistent sampling bias in

developmental psychology: A call to action. *Journal of Experimental Child Psychology*,

162, 31–38. <https://doi.org/10.1016/j.jecp.2017.04.017>

Osborne-Crowley, K. (2020). Social Cognition in the Real World: Reconnecting the Study of

Social Cognition With Social Reality. *Review of General Psychology*, 24(2), 144–158.

<https://doi.org/10.1177/1089268020906483>

Palan, S., & Schitter, C. (2018). Prolific.ac—A subject pool for online experiments. *Journal of Behavioral and Experimental Finance*, 17, 22–27.

<https://doi.org/10.1016/j.jbef.2017.12.004>

Perner, J., Brandl, J. L., Garnham, A., & Peter Lang. (2003). What is a Perspective Problem? Developmental Issues in Belief Ascription and Dual Identity. *Facta Philosophica*, 5(2), 355–378. <https://doi.org/10.5840/factaphil20035220>

Pfeiffer, U. J., Vogeley, K., & Schilbach, L. (2013). From gaze cueing to dual eye-tracking: Novel approaches to investigate the neural correlates of gaze in social interaction. *Neuroscience & Biobehavioral Reviews*, 37(10, Part 2), 2516–2528.

<https://doi.org/10.1016/j.neubiorev.2013.07.017>

Poulin-Dubois, D., Goldman, E. J., Meltzer, A., & Psaradellis, E. (2023). Discontinuity from implicit to explicit theory of mind from infancy to preschool age. *Cognitive Development*, 65, 101273. <https://doi.org/10.1016/j.cogdev.2022.101273>

Povinelli, D. J., Reaux, J. E., Bierschwale, D. T., Allain, A. D., & Simon, B. B. (1997). Exploitation of pointing as a referential gesture in young children, but not adolescent chimpanzees. *Cognitive Development*, 12(4), 423–461. [https://doi.org/10.1016/S0885-2014\(97\)90017-4](https://doi.org/10.1016/S0885-2014(97)90017-4)

Prein, J. C., Kalinke, S., Haun, D. B. M., & Bohn, M. (2023). TANGO: A reliable, open-source, browser-based task to assess individual differences in gaze understanding in 3 to 5-year-old children and adults. *Behavior Research Methods*. <https://doi.org/10.3758/s13428-023-02159-5>

MODELING VARIATION IN GAZE FOLLOWING

58

Quesque, F., & Rossetti, Y. (2020). What Do Theory-of-Mind Tasks Actually Measure? Theory and Practice. *Perspectives on Psychological Science*, 15(2), 384–396.

<https://doi.org/10.1177/1745691619896607>

R Core Team. (2022). *R: A language and environment for statistical computing* [Manual]. R Foundation for Statistical Computing.

Rakoczy, H. (2022). Foundations of theory of mind and its development in early childhood.

*Nature Reviews Psychology*, 1(4), 223–235. <https://doi.org/10.1038/s44159-022-00037-z>

Recasens, A., Khosla, A., Vondrick, C., & Torralba, A. (2015). Where are they looking? *Advances in Neural Information Processing Systems*, 28.

Remillard, M. L., Mazor, K. M., Cutrona, S. L., Gurwitz, J. H., & Tjia, J. (2014). Systematic Review of the Use of Online Questionnaires of Older Adults. *Journal of the American Geriatrics Society*, 62(4), 696–705. <https://doi.org/10.1111/jgs.12747>

Silverstein, P., Feng, J., Westermann, G., Parise, E., & Twomey, K. E. (2021). Infants Learn to Follow Gaze in Stages: Evidence Confirming a Robotic Prediction. *Open Mind*, 5, 174–188. [https://doi.org/10.1162/opmi\\_a\\_00049](https://doi.org/10.1162/opmi_a_00049)

Stiefelhagen, R., & Zhu, J. (2002). Head orientation and gaze direction in meetings. *CHI '02 Extended Abstracts on Human Factors in Computing Systems*, 858–859. <https://doi.org/10.1145/506443.506634>

Surtees, A., Apperly, I., & Samson, D. (2013). Similarities and differences in visual and spatial perspective-taking processes. *Cognition*, 129(2), 426–438. <https://doi.org/10.1016/j.cognition.2013.06.008>

- Symons, L. A., Lee, K., Cedrone, C. C., & Nishimura, M. (2004). What are you looking at? Acuity for triadic eye gaze. *The Journal of General Psychology*, 131(4), 451–469.
- Todorović, D. (2006). Geometrical basis of perception of gaze direction. *Vision Research*, 46(21), 3549–3562. <https://doi.org/10.1016/j.visres.2006.04.011>
- Tomasello, M. (1999). Social Cognition Before the Revolution. In *Early Social Cognition*. Psychology Press.
- Tomasello, M., & Rakoczy, H. (2003). What Makes Human Cognition Unique? From Individual to Shared to Collective Intentionality. *Mind and Language*, 18(2), 121–147. <https://doi.org/10.1111/1468-0017.00217>
- Trafimow, D. (2016). The attenuation of correlation coefficients: A statistical literacy issue. *Teaching Statistics*, 38(1), 25–28. <https://doi.org/10.1111/test.12087>
- Triesch, J., Teuscher, C., Deák, G. O., & Carlson, E. (2006). Gaze following: Why (not) learn it? *Developmental Science*, 9(2), 125–147. <https://doi.org/10.1111/j.1467-7687.2006.00470.x>
- Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing*, 27(5), 1413–1432. <https://doi.org/10.1007/s11222-016-9696-4>
- Walker, L. D., & Gollin, E. S. (1977). Perspective role-taking in young children. *Journal of Experimental Child Psychology*, 24(2), 343–357. [https://doi.org/10.1016/0022-0965\(77\)90012-1](https://doi.org/10.1016/0022-0965(77)90012-1)

MODELING VARIATION IN GAZE FOLLOWING

60

Wellman, H. M. (2018). Theory of mind: The state of the art. *European Journal of Developmental Psychology*, 15(6), 728–755.

<https://doi.org/10.1080/17405629.2018.1435413>

Wellman, H. M., & Liu, D. (2004). Scaling of Theory-of-Mind Tasks. *Child Development*, 75(2), 523–541. <https://doi.org/10.1111/j.1467-8624.2004.00691.x>

Yaniv, I., & Shatz, M. (1990). Heuristics of reasoning and analogy in children's visual perspective taking. *Child Development*, 61(5), 1491–1501.

<https://doi.org/10.2307/1130758>

Zhao, K., Wulder, M. A., Hu, T., Bright, R., Wu, Q., Qin, H., Li, Y., Toman, E., Mallick, B., Zhang, X., & Brown, M. (2019). Detecting change-point, trend, and seasonality in satellite time series data to track abrupt changes and nonlinear dynamics: A Bayesian ensemble algorithm. *Remote Sensing of Environment*, 232, 111181.

<https://doi.org/10.1016/j.rse.2019.04.034>

Zohary, E., Harari, D., Ullman, S., Ben-Zion, I., Doron, R., Attias, S., Porat, Y., Sklar, A. Y., & Mckyton, A. (2022). Gaze following requires early visual experience. *Proceedings of the National Academy of Sciences*, 119(20), e2117184119.

<https://doi.org/10.1073/pnas.2117184119>

## Study III

Running head: GAZE-FOLLOWING ACROSS 17 COMMUNITIES

1

1 A universal of human social cognition: Children from 17 communities process gaze in  
2 similar ways

3 Manuel Bohn<sup>1,2,\*</sup>, Julia Prein<sup>1,2,\*</sup>, Agnes Ayikoru<sup>3</sup>, Florian M. Bednarski<sup>4</sup>, Ardain  
4 Dzabatou<sup>5</sup>, Michael C. Frank<sup>6</sup>, Annette M. E. Henderson<sup>4</sup>, Joan Isabella<sup>3</sup>, Josefine  
5 Kalbitz<sup>2</sup>, Patricia Kanngiesser<sup>7</sup>, Dilara Keşşafoglu<sup>8</sup>, Bahar Köyメン<sup>9</sup>, Maira V.  
6 Manrique-Hernandez<sup>2</sup>, Shirley Magazi<sup>10</sup>, Lizbeth Mújica-Manrique<sup>2</sup>, Julia Ohlendorf<sup>2</sup>,  
7 Damilola Olaoba<sup>2</sup>, Wesley R. Pieters<sup>10</sup>, Sarah Pope-Caldwell<sup>2</sup>, Katie Slocombe<sup>11</sup>, Robert Z.  
8 Sparks<sup>6</sup>, Jahnavi Sunderarajan<sup>2</sup>, Wilson Vieira<sup>2</sup>, Zhen Zhang<sup>12</sup>, Yufei Zong<sup>12</sup>, Roman  
9 Stengelin<sup>2,10,+</sup>, & Daniel B. M. Haun<sup>2,+</sup>

10 <sup>1</sup> Institute of Psychology in Education, Leuphana University Lüneburg

11 <sup>2</sup> Department of Comparative Cultural Psychology, Max Planck Institute for Evolutionary  
12 Anthropology

13 <sup>3</sup> Budongo Conservation Field Station

14 <sup>4</sup> School of Psychology, University of Auckland

15 <sup>5</sup> Université Marien Ngouabi

16 <sup>6</sup> Department of Psychology, Stanford University

17 <sup>7</sup> School of Psychology, University of Plymouth

18 <sup>8</sup> Department of Psychology, Koç University

19 <sup>9</sup> Division of Psychology, Communication, and Human Neuroscience, University of  
20 Manchester

*Appendix A — Main Publications*

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

2

<sup>21</sup> <sup>10</sup> Department of Psychology and Social Work, University of Namibia

<sup>22</sup> <sup>11</sup> Department of Psychology, University of York

<sup>23</sup> <sup>12</sup> CAS Key Laboratory of Behavioral Science, Institute of Psychology, Chinese Academy of  
<sup>24</sup> Sciences

<sup>25</sup> \* joint first author

<sup>26</sup> + joint last author

<sup>27</sup> Author Note

28        The authors would like to thank Luke Maurits for statistical advice. Manuel Bohn  
29    was supported by a Jacobs Foundation Research Fellowship (2022-1484-00). We are  
30    grateful to thank all children and caregivers for participating in the study. We thank the  
31    Max Planck Society for the Advancement of Science.

32        The authors made the following contributions. Manuel Bohn: Conceptualization,  
33    Methodology, Formal Analysis, Writing - Original Draft Preparation, Writing - Review &  
34    Editing; Julia Prein: Conceptualization, Methodology, Software, Investigation, Writing -  
35    Review & Editing; Agnes Ayikoru: Investigation; Florian M. Bednarski: Investigation,  
36    Writing - Review & Editing; Arda Dzabatou: Investigation; Michael C. Frank:  
37    Investigation, Writing - Review & Editing; Annette M. E. Henderson: Investigation,  
38    Writing - Review & Editing; Joan Isabella: Investigation; Josefina Kalbitz: Investigation,  
39    Writing - Review & Editing; Patricia Kanngiesser: Investigation, Writing - Review &  
40    Editing; Dilara Keşşahoğlu: Investigation, Writing - Review & Editing; Bahar Köyメン:  
41    Investigation, Writing - Review & Editing; Maira V. Manrique-Hernandez: Investigation;  
42    Shirley Magazi: Investigation; Lizbeth Mújica-Manrique: Investigation, Writing - Review  
43    & Editing; Julia Ohlendorf: Investigation; Damilola Olaoba: Investigation; Wesley R.  
44    Pieters: Investigation, Writing - Review & Editing; Sarah Pope-Caldwell: Investigation;  
45    Katie Slocombe: Investigation, Writing - Review & Editing; Robert Z. Sparks:  
46    Investigation; Jahnavi Sunderajan: Investigation; Wilson Vieira: Investigation; Zhen  
47    Zhang: Investigation, Writing - Review & Editing; Yufei Zong: Investigation; Roman  
48    Stengelin: Conceptualization, Methodology, Investigation, Writing - Review & Editing;  
49    Daniel B. M. Haun: Conceptualization, Funding acquisition, Writing - Review & Editing.

50        Correspondence concerning this article should be addressed to Manuel Bohn,  
51    Universitätsallee 1, 21335 Lüneburg, Germany. E-mail: manuel.bohn@leuphana.de

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

4

52

Abstract

53 Theoretical accounts assume that key features of human social cognition are universal.  
54 Here we focus on gaze-following, the bedrock of social interactions and coordinated  
55 activities, to test this claim. In a comprehensive cross-cultural study spanning five  
56 continents and 17 distinct cultural communities, we examined the development of  
57 gaze-following in early childhood. We identified key processing signatures through a  
58 computational model that assumes that participants follow an individual's gaze by  
59 estimating a vector emanating from the eye-center through the pupil. We found these  
60 signatures in all communities, suggesting that children worldwide processed gaze in highly  
61 similar ways. Absolute differences between groups were accounted for by a cross-culturally  
62 consistent relationship between children's exposure to touchscreens and their performance  
63 in the task. These results provide strong evidence for a universal process underlying a  
64 foundational socio-cognitive ability in humans that can be reliably inferred even in the  
65 presence of cultural variation in overt behavior.

66 A universal of human social cognition: Children from 17 communities process gaze in  
67 similar ways

68 **Research Transparency Statement**

69 All authors declare no conflicts of interest. Preregistration: The hypotheses, methods  
70 and parts of the analysis plan were preregistered ( [Whttps://osf.io/tdsvc](https://osf.io/tdsvc)) on March 12th,  
71 2022, prior to data collection which began on on March 18th, 2022. Additional analysis  
72 and deviations from the preregistration are reported in the Supplementary Material.

73 Materials: All study materials are publicly available  
74 (<https://ccp-odc.eva.mpg.de/tango-cc/>). Data: All primary data are publicly available  
75 (<https://github.com/ccp-eva/gafo-cc-analysis/>). Analysis scripts: All analysis scripts are  
76 publicly available (<https://github.com/ccp-eva/gafo-cc-analysis/>).

77 **Introduction**

78 Human socio-cognitive skills enable unique forms of communication and cooperation  
79 that provide a bedrock for cumulative culture and the formation of complex societies  
80 (Henrich, 2016; Heyes, 2018; Laland & Seed, 2021; Legare, 2019; Tomasello, 2020;  
81 Tomasello & Rakoczy, 2003; Wellman, 2014). The eyes are the proverbial “window to the  
82 mind” and eye gaze is essential for many social reasoning processes (Doherty, 2006; Emery,  
83 2000; Shepherd, 2010). Others’ eye gaze is used to infer their focus of visual attention,  
84 which is a critical aspect of coordinated activities, including communication and  
85 cooperation (Langton, Watt, & Bruce, 2000; Richardson & Dale, 2005; Rossano, 2012;  
86 Scaife & Bruner, 1975; Sebanz, Bekkering, & Knoblich, 2006; Tomasello, Hare, Lehmann,  
87 & Call, 2007).

88 The ability to follow gaze emerges early in development (Byers-Heinlein et al., 2021;  
89 Del Bianco, Falck-Ytter, Thorup, & Gredebäck, 2019; Gredebäck, Fikke, & Melinder, 2010;  
90 Tang, Gonzalez, & Deák, 2023). The earliest signs of gaze-following have been found in

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

6

91 infants as young as four months (Astor, Thiele, & Gredebäck, 2021; D'Entremont, Hains,  
92 & Muir, 1997). Initially, infants rely more on head direction than actual gaze direction  
93 (Lempers, Flavell, & Flavell, 1977; Michel, Kayhan, Pauen, & Hoehl, 2021). Throughout  
94 the first two years of life, children refine their abilities: they interpret gaze in mentalistic  
95 terms, for example, they follow gaze to locations outside their own visual field by moving  
96 around barriers (Moll & Tomasello, 2004). Importantly, individual differences in children's  
97 gaze-following abilities predict later life outcomes, most notably communicative abilities  
98 (Carpenter, Nagell, Tomasello, Butterworth, & Moore, 1998). For example, gaze-following  
99 at 10 months predicts language scores at 18 months of age (Brooks & Meltzoff, 2005).  
100 Difficulties with gaze-following have been linked to developmental disorders, including  
101 Autism (Itier & Batty, 2009; Thorup et al., 2016, 2018). This work highlights the  
102 importance of gaze-following as a foundational building block of human social interaction  
103 and its central place in theorizing.

104 A central assumption in the theoretical and empirical work discussed above is that,  
105 despite substantial variation in developmental contexts, gaze-following works and develops  
106 in the same way across human societies (Tomasello, 2019). This assumption – despite  
107 being central to many developmental theories – is currently not supported by evidence. On  
108 the contrary, cross-cultural studies have revealed substantial diversity in socio-cognitive  
109 development (Dixson, Komugabe-Dixson, Dixson, & Low, 2018; Mayer & Träuble, 2013;  
110 Miller, Wice, & Goyal, 2018; Taumoepeau, Sadeghi, & Nobilo, 2019; Wellman, 2014). One  
111 of the very few cross-cultural studies also found differences in the likelihood to follow gaze  
112 between communities (Callaghan et al., 2011).

113 One potential source for this paradox lies in the reliance on aggregated measures in  
114 cross-cultural studies. Absolute differences in mean performance across communities are  
115 interpreted as a signal of different underlying cognitive processes. In the present study, we  
116 resolve this paradox by instead focusing on processing signatures that can be investigated  
117 independently of absolute community-level differences. This allows us to directly evaluate

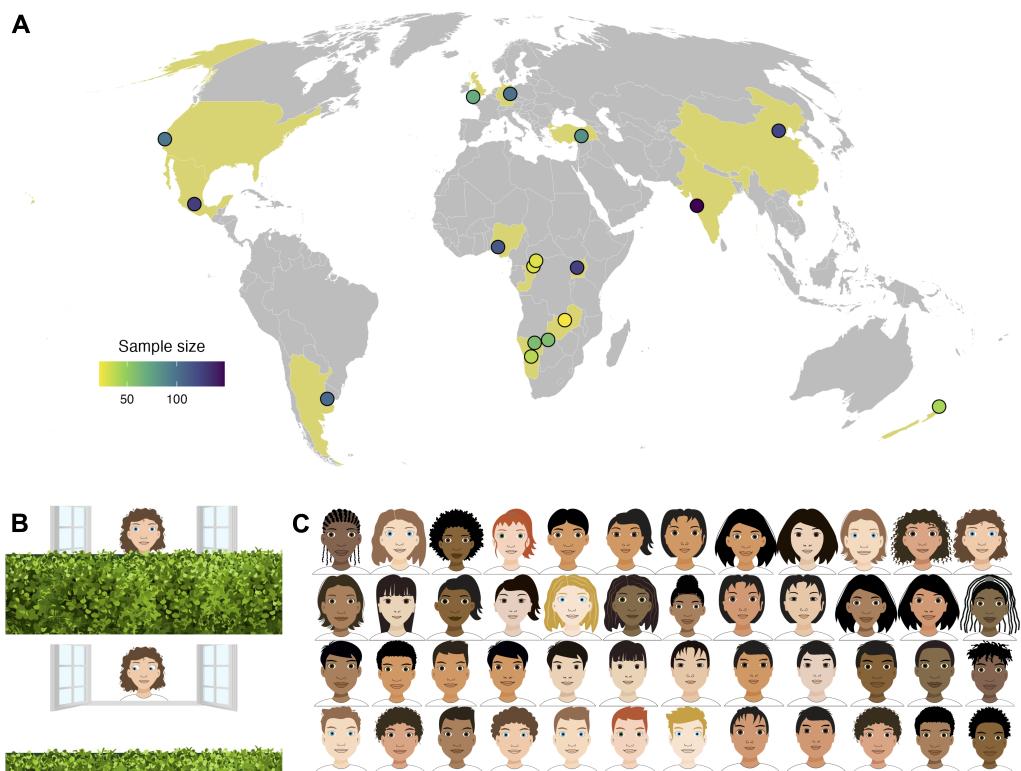
118 the empirical foundation of claims about universal features of human social cognition. To  
119 this end, we conducted a pre-registered, large-scale, cross-cultural study on the  
120 development of gaze-following abilities to study potentially universal processing signatures.

121 The processing signatures were derived from a computational model that assumes  
122 that participants follow gaze by estimating a vector emanating from the eye center through  
123 the pupil (Prein, Maurits, Werwach, Haun, & Bohn, 2023). The key innovation of the  
124 model is that it explains how individuals may use the same cognitive process but still differ  
125 in their measured abilities. The process always involves estimating a vector but also  
126 involves a degree of uncertainty because the eye center is not directly observable.  
127 Individuals are assumed to differ in their level of uncertainty with which they estimate the  
128 vector which causes differences in their observable behavior. Importantly, the assumed  
129 process leaves a key signature in the data that is observable independent of the absolute  
130 level of performance. In the present study, we therefor focus on this signature instead of  
131 absolute levels of performance when evaluating the claim whether there is evidence for a  
132 universal cognitive mechanism underlying gaze-following.

133 The 1377 participants who took part in the study lived in 17 different communities  
134 across 14 countries and five continents (Fig. 1A, Tab. 1). These countries represent ~46%  
135 of the world's population. Communities covered a broad spectrum of geographical  
136 locations, social and political systems, languages, and subsistence styles (see  
137 Supplementary Material). This diversity allowed us to overcome the common pitfall of  
138 cross-cultural studies that compare urban communities from the global north to rural  
139 communities from the global south (Barrett, 2020).

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

8



*Figure 1.* (A) Data collection sites. Points show the approximate geographical location of the data collection sites, coloring shows the sample sizes. (B) Screenshots from the task. Screenshots from the task. The upper scene depicts the start and the lower the choice phase in a test trial. Participants had to use the gaze of the agent to locate the balloon and touch the location on the hedge where they thought the balloon was. Agents, audio recordings and backgrounds were adapted to each cultural setting. (C) Drawings used as agents across cultural settings.

Table 1

*Participant demographics.*

Continent	Country	Community	N(male)	Age (range)	Language	Touchscreen exposure1
Americas	Argentina	Buenos Aires	105 (53)	4.72 (3.00 - 6.96)	Spanish (Rio-platense)	0.90
		Ocuilan	127 (63)	4.96 (2.57 - 6.95)	Spanish (Mexican)	0.77
	USA	Stanford	98 (54)	4.99 (2.52 - 7.90)	English (American)	0.98
Africa	Namibia	Hai  om	60 (38)	5.85 (2.74 - 8.34)	Hai  om	0.05
		Khwe	59 (24)	5.84 (3.38 - 8.63)	Khwedam	0.19
		Windhoek	39 (17)	5.69 (2.66 - 8.66)	English (Nigerian)2	0.95
Nigeria	Rep. Congo	Akure	114 (54)	5.07 (2.57 - 7.33)	English (Nigerian)	0.91
		BaYaka	29 (13)	7.80 (3.94 - 10.56)	BaYaka	0.00
		Bandongo	30 (11)	7.45 (3.50 - 10.95)	Lingala	0.00
Uganda	Nyabyeya			5.94 (2.67 - 8.92)	Kiswahili	0.34

Table 1 continued

Continent	Country	Community	N(male)	Age (range)	Language	Touchscreen exposure <sup>1</sup>
	Zambia	Chimfunshi	22 (5)	5.98 (2.88 - 8.00)	Bemba	0.14
Europe	Germany	Leipzig	100 (48)	4.88 (2.53 - 6.95)	German	0.89
	UK	Plymouth	70 (30)	6.02 (2.38 - 8.94)	English (British)	0.99
Asia	China	Beijing	123 (62)	5.47 (2.69 - 8.48)	Mandarin	0.95
	India	Pune	148 (73)	6.14 (3.06 - 8.83)	English (Indian) / Marathi	0.93
	Türkiye	Malatya	85 (40)	5.02 (2.75 - 7.12)	Turkish	1.00
Oceania	New Zealand	Auckland	43 (19)	5.14 (2.81 - 8.75)	English (New Zealand)	0.95

*Note.* 1 Proportion of participants who have access to touchscreens according to parental questionnaire. 2 Local collaborators and piloting suggested that Nigerian English is suitable for Windhoek as well.

<sup>140</sup>

<sup>141</sup> We used an animated picture book tablet task in which participants had to locate a hidden object based on observing an agent's gaze. Children watched a balloon disappear

<sup>142</sup>

<sup>143</sup> behind a hedge. An agent followed the trajectory of the balloon with their eyes (Fig. 1B).  
<sup>144</sup> The key dependent variable was the (im)precision with which children located the agent's  
<sup>145</sup> focus of attention, that is, the deviation between where the agent looked (where the  
<sup>146</sup> balloon was) and the child's response. We adapted visuals and audio instructions  
<sup>147</sup> specifically for each of the 17 communities. Previous work demonstrated excellent  
<sup>148</sup> individual-level measurement properties for this task in a German sample (Prein, Kalinke,  
<sup>149</sup> Haun, & Bohn, 2023).

<sup>150</sup> **Methods**

<sup>151</sup> **Participants**

<sup>152</sup> A total of 1377 children between 2.38 and 10.95 provided data for the study. Children  
<sup>153</sup> lived in 17 different communities, located in 14 different countries. Table 1 gives the  
<sup>154</sup> sample size per community together with some basic demographic information. The  
<sup>155</sup> recruitment strategy for each community is reported in the respective site descriptions. For  
<sup>156</sup> some children, the exact birthday was unknown. In such cases, we set the birthday to the  
<sup>157</sup> 30th of June of the year that would make them fall into the reported age category.

<sup>158</sup> Data from children was only included in the study when they contributed at least  
<sup>159</sup> four valid test trials. We also excluded the data from children with a diagnosed  
<sup>160</sup> developmental disorder. In sum, in addition to the sample size reported above, 74  
<sup>161</sup> additional children participated in the study but did not contribute data. The main  
<sup>162</sup> reasons for exclusion were: contribution of less than four valid test trials, technical failures,  
<sup>163</sup> and missing or implausible demographic information (e.g., when the number of children  
<sup>164</sup> living in the household was reported to be larger than the household itself or when the  
<sup>165</sup> number of children reported to live in the household equaled the number of children  
<sup>166</sup> younger than the child being tested). We did not exclude any participants for performance  
<sup>167</sup> reasons. A detailed description of each data collection site and the way children were

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

12

<sup>168</sup> recruited can be found in the Supplementary Material.

<sup>169</sup> **Setup and Procedure**

<sup>170</sup> The task was implemented as a browser-based interactive picture book using HTML  
<sup>171</sup> and JavaScript. Participants saw animated agents on a touch screen device, listened to  
<sup>172</sup> pre-recorded audio instructions and responded by touching the screen. In all communities,  
<sup>173</sup> a research assistant, fluent in the local language(s), guided the child through the task.

<sup>174</sup> Figure 1B shows a screenshot from the task. The task was introduced verbally by the  
<sup>175</sup> assistant as the balloon game in which the participant would play with other children to  
<sup>176</sup> find a balloon. On each trial, participants saw an agent located in a window in the center  
<sup>177</sup> of the screen. A balloon fell down from its starting position just below the agent. The  
<sup>178</sup> agent's gaze followed the trajectory of the balloon. That is, the pupils and the iris were  
<sup>179</sup> programmed to align with the center of the balloon. Once the balloon had landed on the  
<sup>180</sup> ground, the agent was instructed to locate it, that is, to touch the location on the screen  
<sup>181</sup> where they thought the balloon was. On each trial, we recorded the exact x-coordinate of  
<sup>182</sup> the participant's touch.

<sup>183</sup> There were two types of familiarization trials. In fam1 trials, the balloon fell down  
<sup>184</sup> and landed in plain sight. Participants simply had to touch the visible balloon. In fam2  
<sup>185</sup> trials, the trajectory of the balloon was visible but it landed behind a small barrier (a  
<sup>186</sup> hedge - see Figure 1B). Thus, participants needed to touch the hedge where they saw the  
<sup>187</sup> balloon land. Next came test trials. Here, the barrier moved up and covered the balloon's  
<sup>188</sup> trajectory. That is, participants only saw the agent's eyes move, but not the balloon. They  
<sup>189</sup> had to infer the location of the balloon based on the agent's gaze direction. During fam1,  
<sup>190</sup> fam2 and the first test trial, children heard voice overs commenting what happened on the  
<sup>191</sup> screen. Critically, the agent was described as wanting to help the child and always looking  
<sup>192</sup> at the balloon.

193 Children completed one fam1 trial, two fam2 trials and 16 test trials. We excluded  
194 the first test trial from the analysis because of the voice-over. Thus, 15 test trials were used  
195 in the analysis below. Each child saw eight different agents, four male, four female. The  
196 agent changed from trial to trial, with alternating genders. A coin toss before the first trial  
197 decided whether the first agent was male or female. The order in which agents were shown  
198 was randomized with the constraint that all agents had to be shown once until an agent  
199 was shown again. The color of the balloon also changed from trial to trial in a random  
200 order, also with the constraint that all colors appeared once before any one was repeated.

201 The location (x-coordinate) where the balloon landed was determined in the following  
202 way: The screen was divided in ten equally sized bins. On each trial, one of the bins was  
203 randomly selected and the exact x-coordinate was randomly chosen within that bin.  
204 Constraints were that the balloon landed in each bin equally often and the same bin  
205 appeared no more than twice in a row.

206 All children were tested with a touchscreen device with a size between 11 and 13 inch  
207 equipped with a webcam. The data was either stored locally or sent to a server. In  
208 addition to the behavioral data, we stored the webcam recording of the session for  
209 verification purposes. Culture-specific adaptations were made by changing the visuals and  
210 the audio instructions (see Supplementary Material for details).

211 In addition to the gaze-following task, caregivers responded to a short questionnaire  
212 about children's access to screens and touchscreens (binary answer) as well as the number  
213 of people, children and children younger than the focal child living in the household  
214 (numeric; see Supplementary Material for details). The numeric variables were scaled  
215 within cultural setting prior to inclusion into the regression models.

216

## Analysis

### 217 Regression models

218 We used Bayesian Regression models fit in R (R Core Team, 2023) using the package  
219 `brms` (Bürkner, 2017) for all analyses except the cognitive models (see below). We used  
220 default priors built into `brms`. The dependent variable in all regression models was  
221 imprecision, that is, the absolute distance between the true location of the balloon  
222 (x-coordinate of its center) and the location where the participant touched the screen. We  
223 used a Log-normal distribution to model the data because the natural lower bound for  
224 imprecision is zero and the data was right skewed with a long tail. Numeric predictors that  
225 entered the models were scaled to have a mean of zero and a standard deviation of 1.

226 To analyse cross-cultural variation in performance, we used a cross-validation  
227 procedure (see e.g., Stengelin, Ball, Maurits, Kanngiesser, & Haun, 2023). In the  
228 Supplementary Material we give a detailed justification of this approach. For each cultural  
229 setting, we randomly sampled a data set that was 5/6 the size of the full data set (training  
230 data). Then, we fit the model to this training data and used the estimated model  
231 parameters to predict the remaining 1/6 of the data (testing data). We then compared the  
232 model predictions from the different models by computing the mean difference between the  
233 true and predicted imprecision, over all trials in the testing data set. We repeated the  
234 cross-validation procedure 100 times and computed the percentage of cases in which one  
235 model outperformed the other. We compared three models: a null model assuming no  
236 systematic community-level variation, a model assuming variation between communities  
237 and a model assuming variation between communities and in developmental trajectories  
238 (see Supplementary Material for model equations).

239 To evaluate the processing signatures predicted by the cognitive model that trials in  
240 which the balloon lands further away from the center lead to larger imprecision (see next  
241 section for details), we fit a model predicting imprecision by age and target centrality

<sup>242</sup> (distance of the landing position form the center in pixel) with random intercepts for  
<sup>243</sup> participant and cultural setting and random slopes for target centrality within participant  
<sup>244</sup> and cultural setting (`brms` notation: `age + target_centrality + (target_centrality`  
<sup>245</sup> `| participant) + (age + target_centrality | culture)`).

<sup>246</sup> **Cognitive model**

<sup>247</sup> Recent computational work modeled gaze-following as social vector estimation (Prein,  
<sup>248</sup> Maurits, et al., 2023). When following gaze, onlookers observe the location of the pupil  
<sup>249</sup> within the eye and estimate a vector emanating from the center of the eye through the  
<sup>250</sup> pupil. The focus of attention is the location where the estimated vectors from both eyes hit  
<sup>251</sup> a surface (Fig. 3). It is assumed that this estimation process has some uncertainty because  
<sup>252</sup> the center of the eye is not directly observable and that individuals vary in their level of  
<sup>253</sup> uncertainty. As a consequence, even though individuals use the same general process, they  
<sup>254</sup> might differ in their absolute levels of precision. Crucially, this process model predicts a  
<sup>255</sup> clear performance signature in our gaze-following task: Trials in which the agent looks  
<sup>256</sup> further away from the center should result in lower levels of precision compared to trials in  
<sup>257</sup> which the agent looks closer to the center. This prediction is best understood by  
<sup>258</sup> considering a similar phenomenon: pointing a torch light to a flat surface. The width of the  
<sup>259</sup> light beam represents each individual's level of uncertainty in vector estimation. When the  
<sup>260</sup> torch is directed straight down, the light beam is concentrated in a relatively small area.  
<sup>261</sup> When the torch is rotated to the side, the light from one half of the cone must travel  
<sup>262</sup> further than the light from the other half to reach the surface. As a consequence, the light  
<sup>263</sup> is spread over a wider area (see Fig. 3).

<sup>264</sup> The model inversely models the process generating touches on the screen based on  
<sup>265</sup> observed eye movements and is defined as:

$$P(\theta|x_c, \alpha_l, \alpha_r) \propto P(x_c|\alpha_l, \alpha_r, \theta)P(\theta) \quad (1)$$

266 Here,  $\theta$  represents an individual's cognitive ability to locate the focus of the agent's  
 267 attention,  $x_c$  represents the touched coordinate, and  $\alpha_l$  and  $\alpha_r$  correspond to the left and  
 268 right pupil angles (each defined as the angle between a line connecting the center of the eye  
 269 to the pupil and a line extended vertically downward from the center of the eye).

270 The basic assumption in this model is that participants touch on the screen location  
 271 where they think the agent is looking. The true eye angles ( $\alpha_l$  and  $\alpha_r$ ) are not directly  
 272 observable and are estimated with noise, yielding  $\hat{\alpha}_l$  and  $\hat{\alpha}_r$ .

273 Each touch  $x_c$  implies a “matched pair” of estimated pupil angles  $\hat{\alpha}_l$  and  $\hat{\alpha}_r$ , with the  
 274 constraint that the lines extended along those two angles meet at the precise location of  
 275 where the target is believed to be. As a consequence, we can rewrite the likelihood function  
 276 of the model as:

$$P(x_c|\alpha_l, \alpha_r, \theta) \propto P(\hat{\alpha}_l, \hat{\alpha}_r|\alpha_l, \alpha_r, \theta)P(x_c) \quad (2)$$

277  $P(x_c)$  is a prior over potential target locations. Because the target was last visible in  
 278 the screen and because the agent was located in the center, we assumed that participants  
 279 have an a priori expectation that the target will land close to the middle. We estimated the  
 280 strength of this center bias (i.e., the standard deviation of a Normal distribution around  
 281 the screen center) based on the data:  $P(x_c) \sim \mathcal{N}(960, \sigma^p)$ .

282 The primary inferential task for participants is therefore to estimate the pupil angles  
 283 ( $\hat{\alpha}_l$  and  $\hat{\alpha}_r$ ), i.e., to sample from the term  $P(\hat{\alpha}_l, \hat{\alpha}_r|\alpha_l, \alpha_r, \theta)$ . Here, we assumed that the  
 284 pair of estimated pupil angles were sampled from a probability distribution which is the  
 285 product of two Normal distributions of equal variance,  $\sigma_v$ , centered on the true pupil angles:

$$P(\hat{\alpha}_l, \hat{\alpha}_r | \alpha_l, \alpha_r, \theta) \propto \phi(\hat{\alpha}_l; \alpha_l, \sigma_v) \phi(\hat{\alpha}_r; \alpha_r, \sigma_v), \quad (3)$$

286 Here,  $\sigma_v$  determines the level of accuracy with which participants estimated the pupil  
 287 angles, and it is thus the component of the model that defines  $\theta$ . Smaller values of  $\sigma_v$  result  
 288 in a narrow distribution around the pupil angle, making touches far away from the target  
 289 less likely. Conversely, larger values for  $\sigma_v$  lead to a wider distribution, making touches far  
 290 away from the target more likely. To circle back to the analogy introduced above,  $\sigma_v$   
 291 corresponds to the width of the light beam. Thus, the goal of the model was to estimate  
 292 participant-specific values for  $\sigma_v$ :  $\sigma_{v_i}$ . For more details on how  $\sigma_{v_i}$  was estimated, see the  
 293 Supplementary Material.

294 To summarize, the model assumes that participant's touches are generated by a  
 295 process that relies on noisy estimates of the agent's gaze direction. The precision, with  
 296 which the gaze direction is estimated, varies between participants and increases with  
 297 development.

298 As stated above, the key signature prediction of the model is that precision decreases  
 299 when the balloon lands further away from the center. This pattern, however, also arises  
 300 when participants ignore the agent's gaze completely and instead follow simple heuristics.  
 301 We implemented these heuristics as alternative models and directly compared them to the  
 302 focal model. According to the center bias model, they always try to touch in the center of  
 303 the screen:  $P(x_c) \sim \mathcal{N}(960, 160)$ . (960 is the x-coordinate of the center and 160 is the  
 304 width of the balloon). According to the random guessing model, they randomly touch  
 305 coordinates on the screen:  $P(x_c) \sim \mathcal{U}(0, 1920)$ .

306 The cognitive models were implemented in the probabilistic programming language  
 307 **webpp1** (Goodman & Stuhlmüller, 2014). All models were run separately for each cultural  
 308 setting. Information on the prior distributions for all model parameters can be found in the  
 309 associated online repository. We compared models based on the marginal likelihood of the

310 data for each model, which represents the likelihood of the data while averaging over the  
311 prior distribution on parameters. The pair-wise ratio of marginal likelihoods for two models  
312 is known as the Bayes Factor. Bayes Factors are a quantitative measure of the predictive  
313 quality of a model, taking into account the possible values of the model parameters  
314 weighted by their prior probabilities. The incorporation of the prior distribution over  
315 parameters in the averaging process implicitly considers model complexity: models with  
316 more parameters typically exhibit broader prior distributions over parameter values and  
317 broader prior distribution can attenuate the potential gains in predictive accuracy that a  
318 model with more parameters might otherwise achieve (Lee & Wagenmakers, 2014).

319 **Results**

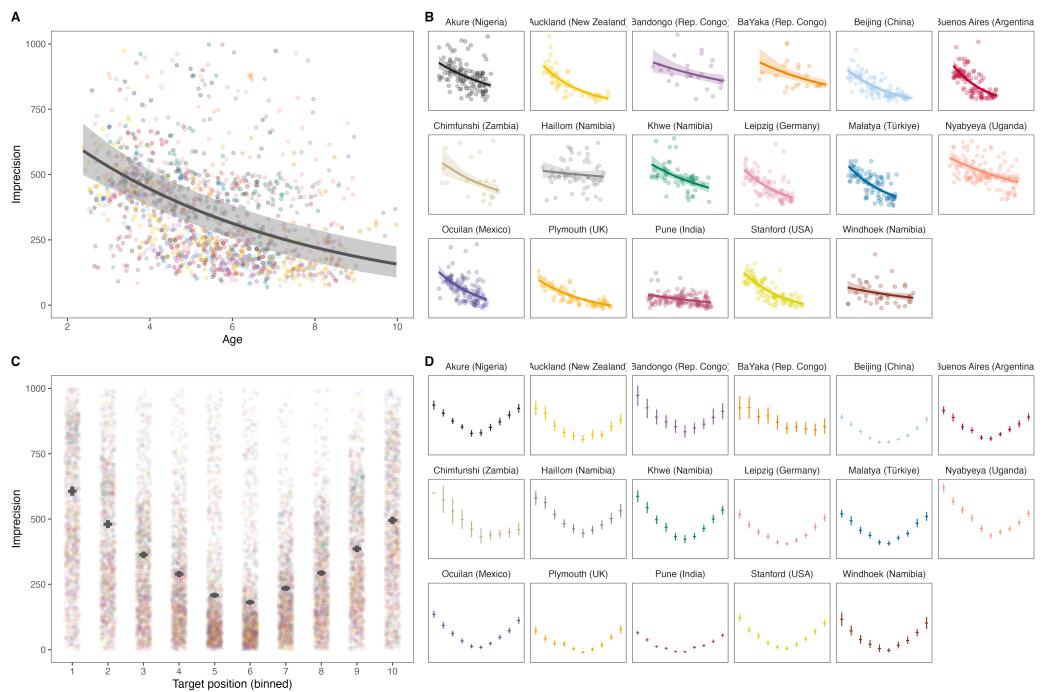
320 **Cross-cultural variation in development**

321 As the first step, we investigated developmental improvements, that is, how children  
322 become more precise at estimating the target location with age. Across all 17 communities,  
323 we found a substantial increase in average levels of precision with age (fixed effect of age:  $\beta$   
324 = -0.30, 95% Credible Interval (CrI) (-0.40 - -0.21); range of community-level (random)  
325 effects:  $\beta_{min} = -0.06$ , 95% CrI (-0.18 - 0.05) to  $\beta_{max} = -0.59$ , 95% CrI (-0.71 - -0.48)).

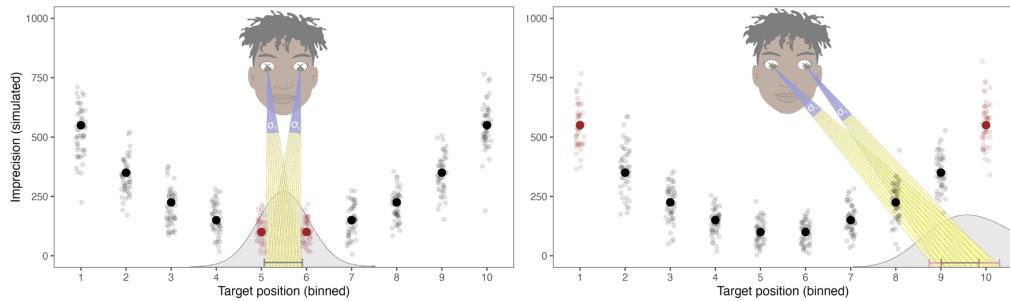
326 Nevertheless, there were also marked differences between communities (see Fig. 2A).  
327 The cross-validation procedure found that a model assuming cross-cultural variation in  
328 average performance as well as cross-cultural variation in developmental trajectories  
329 outperformed simpler models – assuming no variation in the shape of developmental  
330 trajectories or no variation between settings at all – in 98% of cases.

331 Average differences in precision between communities were small compared to  
332 differences between individuals: communities did not form homogeneous clusters but  
333 largely overlapping distributions in that some individuals from communities with a lower  
334 average level of precision performed better compared to some individuals from a setting

with a very high average level of precision. Similarly, in all communities, some 4-year-olds outperformed children two years older than them (see Fig. 2A). The lack of adequate individual-level measurement instruments in previous large-scale developmental cross-cultural studies made it impossible to contrast these perspectives.



*Figure 2.* A) Developmental trajectory across and B) by community. The developmental trajectories are predicted based on a model of the data aggregated for each participant. C) Performance by target location on the screen across, and D) by community. Each bin covers 1/10th of the screen. Points show means, and error bars 95% confidence intervals for the data within that bin aggregated across participants. Transparent dots in A) and C) show aggregated data for each individual.



*Figure 3.* Graphical illustration of the cognitive model. Individuals infer the target of an agent’s attention by estimating a vector based on the position of the pupils within the eyes. This process is noisy, illustrated by the different vectors (transparent lines). Individuals differ in their level of precision (indicated by sigma). For a given level of precision, the further the target lands from the centre of the screen, the less precise the model predicts individuals to be. Solid and transparent dots show simulated means and individual data points to illustrate the predicted effect of target position.

<sup>339</sup> **Universal processing signatures**

<sup>340</sup> The key processing signature predicted by the cognitive model was that precision  
<sup>341</sup> should decrease when the balloon landed further away from the center. This signature was  
<sup>342</sup> clearly visible across all 17 communities (fixed effect for target centrality:  $\beta = 0.47$ , 95%  
<sup>343</sup> CrI (0.40 - 0.54); range of community-level (random) effects:  $\beta_{min} = 0.58$ , 95% CrI (0.51 -  
<sup>344</sup> 0.66) to  $\beta_{max} = 0.16$ , 95% CrI (-0.01 - 0.33)). Visualization of the data showed the  
<sup>345</sup> predicted u-shaped pattern in all communities (see Fig. 2B). When we compared the focal  
<sup>346</sup> vector-based gaze estimation model described above to the alternative center-bias and  
<sup>347</sup> random guessing models, we found overwhelming support for the gaze estimation model  
<sup>348</sup> ( $\min BF_{10} > 100\,000$  for comparisons with both alternative models, see Supplementary  
<sup>349</sup> Materials) in every community.

350 **Predictors of variation**

351 We used the caregiver questionnaire to explain community- and individual-level  
352 variation. On an individual level, we found that children with access to touchscreen devices  
353 had higher levels of precision ( $\beta = -0.14$ , SE = 0.04, 95% CrI = -0.21 - -0.07). This effect  
354 was consistent across communities in that allowing the effect of access to touchscreens to  
355 vary across communities did not improve model fit (see Supplementary Material). On a  
356 community level, we also saw that average performance was lowest in communities in which  
357 touchscreen devices were the least frequent (community-level correlation between  
358 age-corrected imprecision and proportion of children with access to touchscreens:  $r =$   
359 -0.90, 95% CI = -0.96 - -0.74). Thus, familiarity with the device used for data collection  
360 likely explains variation between communities. Children with more touchscreen experience  
361 were probably better at task handling and thus more likely to precisely touch the location  
362 they inferred the agent to look at.

363 However, there was substantial variation between individuals that could not be  
364 explained by differential exposures to touchscreens alone. For example, in Malatya  
365 (Turkeyie) where 100% of children had access to touchscreens there was still substantial  
366 variation between individuals (see Fig.1B). This strongly indicates that other factors likely  
367 contributed to individual differences. Social interaction has been highlighted as an  
368 important driver of social-cognitive development (Barresi & Moore, 1996; Carpendale &  
369 Lewis, 2020; Perner, Ruffman, & Leekam, 1994; Rakoczy, 2022; e.g., Tomasello, 2019) and  
370 thus we hypothesized (and pre-registered) that more opportunities for social interaction –  
371 approximated by living in larger households with more children – would be associated with  
372 higher levels of precision. When predicting performance by relative opportunities for social  
373 interactions within a community – while accounting for absolute differences and the  
374 prevalence of touchscreens – we found no strong associations between any of the  
375 demographic indicators and performance (see Supplementary Material).

376

## Discussion

377 Following and understanding gaze is a foundational building block of human social  
378 cognition (Langton et al., 2000; Richardson & Dale, 2005; Rossano, 2012; Scaife & Bruner,  
379 1975; Sebanz et al., 2006; Tomasello et al., 2007). A substantial body of work has explored  
380 the developmental onset of gaze-following in a few selected cultural communities  
381 (Byers-Heinlein et al., 2021; Gredebäck et al., 2010; Moore, 2008; Tang et al., 2023). The  
382 data reported here provides strong evidence that children from a large and diverse set of  
383 communities process others' gaze in similar ways. We found key performance signatures  
384 predicted by a model treating gaze-following as a form of social vector estimation across all  
385 17 communities. With the focus on individual-level processing signatures, the study goes  
386 beyond previous studies on gaze-following – focused on the onset of gaze-following in  
387 infancy (Callaghan et al., 2011; Hernik & Broesch, 2019) – as well as comprehensive  
388 cross-cultural studies that compared average developmental trajectories (Blake et al., 2015;  
389 House et al., 2020; Kanngiesser et al., 2022; Van Leeuwen et al., 2018).

390

The cognitive processes underlying gaze-following might be rooted in humans'  
391 evolved cognitive architecture, which is – presumably – later refined during social  
392 interaction (Astor et al., 2020; Movellan & Watson, 2002; Senju et al., 2015). The  
393 phylogenetic roots of these processes might possibly lie much deeper as primates from a  
394 wide range of species follow gaze (Itakura, 2004; Kano & Call, 2014; Rosati & Hare, 2009;  
395 Tomasello, Call, & Hare, 1998). Yet, similarities in overt behavior do not imply the same  
396 underlying cognitive processes. The present study defines clear performance signatures that  
397 can be explored in other species to test such evolutionary hypotheses.

398

Our study combined precise individual-level cognitive measurement and  
399 individual-level assessment of experience (here: touchscreen exposure) in a large and  
400 diverse sample to directly investigate the impact of specific cultural experiences on  
401 developmental outcomes. Instead of establishing universality by maximizing the cultural

402 distance between two or three tested communities (Norenzayan & Heine, 2005), this  
403 large-scale cross-cultural approach treats children's cultural experience at scale, shedding  
404 light on the big "middle ground" of children's cultural experience (Barrett, 2020).

405 The study has important limitations. The fact that performance in the task was  
406 correlated with exposure to touchscreens might have overshadowed other sources of  
407 variation. However, we think it is an important innovation that we were able to account  
408 for this effect. Most developmental cross-cultural studies do not even question the  
409 portability of their measurement instruments. Importantly, the key result that the  
410 processing signatures were seen in all cultural settings, is immune to this finding. The  
411 potential that lies in the otherwise precise individual-level measurement that our task  
412 achieves is largely unexploited. The questionnaire items only offer a very coarse picture  
413 into children's actual lived experiences. Whilst household size was a useful proxy for  
414 regular social interaction opportunities, the measure does not directly measure the factors  
415 that previous work has suggested to be related to the development of gaze-following in  
416 younger children, such as attachment quality or the use of gaze in early communicative  
417 interactions (Astor et al., 2020; Movellan & Watson, 2002; Senju et al., 2015). Future work  
418 could increase the resolution with which everyday experiences in children from diverse  
419 communities are recorded to compare the drivers behind social-cognitive development as  
420 we observe it. Recent work in the field of language acquisition has shown how technological  
421 innovations allowed for direct recording of social interactions across communities which can  
422 be used to close this explanatory gap (Bergelson et al., 2023; Donnelly & Kidd, 2021).

423 In sum, our work pioneers an approach that introduces computational modeling and  
424 precise individual-level measurement to the cross-cultural study of cognitive development.  
425 This approach allowed us to test for universals in the human cognitive architecture rather  
426 than just overt behavior. As such, it can serve as a blueprint for future research on a broad  
427 spectrum of cognitive abilities and offers a much-needed empirical foundation for theories  
428 on the nature of the human mind. Children from diverse cultures deploy similar cognitive

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

24

<sup>429</sup> processes in interpreting gaze, pointing to a universal foundation of basic social cognition,

<sup>430</sup> which is refined during development.

431

## References

- 432 Astor, K., Lindskog, M., Forssman, L., Kenward, B., Fransson, M., Skalkidou, A., ...
- 433 Gredebäck, G. (2020). Social and emotional contexts predict the development of gaze
- 434 following in early infancy. *Royal Society Open Science*, 7(9), 201178.
- 435 Astor, K., Thiele, M., & Gredebäck, G. (2021). Gaze following emergence relies on both
- 436 perceptual cues and social awareness. *Cognitive Development*, 60, 101121.
- 437 Barresi, J., & Moore, C. (1996). Intentional relations and social understanding. *Behavioral*
- 438 *and Brain Sciences*, 19(1), 107–122.
- 439 Barrett, H. C. (2020). Towards a cognitive science of the human: Cross-cultural
- 440 approaches and their urgency. *Trends in Cognitive Sciences*, 24(8), 620–638.
- 441 Bergelson, E., Soderstrom, M., Schwarz, I.-C., Rowland, C. F., Ramirez-Esparza, N., R.
- 442 Hamrick, L., et al.others. (2023). Everyday language input and production in 1,001
- 443 children from six continents. *Proceedings of the National Academy of Sciences*, 120(52),
- 444 e2300671120.
- 445 Blake, P. R., McAuliffe, K., Corbit, J., Callaghan, T. C., Barry, O., Bowie, A., et al.others.
- 446 (2015). The ontogeny of fairness in seven societies. *Nature*, 528(7581), 258–261.
- 447 Brooks, R., & Meltzoff, A. N. (2005). The development of gaze following and its relation to
- 448 language. *Developmental Science*, 8(6), 535–543.
- 449 Bürkner, P.-C. (2017). Brms: An r package for bayesian multilevel models using stan.
- 450 *Journal of Statistical Software*, 80(1), 1–28.
- 451 Byers-Heinlein, K., Tsui, R. K.-Y., Van Renswoude, D., Black, A. K., Barr, R., Brown, A.,
- 452 et al.others. (2021). The development of gaze following in monolingual and bilingual
- 453 infants: A multi-laboratory study. *Infancy*, 26(1), 4–38.
- 454 Callaghan, T., Moll, H., Rakoczy, H., Warneken, F., Liszkowski, U., Behne, T., ... Collins,
- 455 W. A. (2011). Early social cognition in three cultural contexts. *Monographs of the*
- 456 *Society for Research in Child Development*, i–142.
- 457 Carpendale, J., & Lewis, C. (2020). *What makes us human: How minds develop through*

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

26

- 458        *social interactions*. Routledge.
- 459    Carpenter, M., Nagell, K., Tomasello, M., Butterworth, G., & Moore, C. (1998). Social  
460        cognition, joint attention, and communicative competence from 9 to 15 months of age.  
461        *Monographs of the Society for Research in Child Development*, i–174.
- 462    D'Entremont, B., Hains, S. M., & Muir, D. W. (1997). A demonstration of gaze following  
463        in 3-to 6-month-olds. *Infant Behavior and Development*, 20(4), 569–572.
- 464    Del Bianco, T., Falck-Ytter, T., Thorup, E., & Gredebäck, G. (2019). The developmental  
465        origins of gaze-following in human infants. *Infancy*, 24(3), 433–454.
- 466    Dixson, H. G., Komugabe-Dixson, A. F., Dixson, B. J., & Low, J. (2018). Scaling theory of  
467        mind in a small-scale society: A case study from vanuatu. *Child Development*, 89(6),  
468        2157–2175.
- 469    Doherty, M. J. (2006). The development of mentalistic gaze understanding. *Infant and  
470        Child Development*, 15(2), 179–186.
- 471    Donnelly, S., & Kidd, E. (2021). The longitudinal relationship between conversational  
472        turn-taking and vocabulary growth in early language development. *Child Development*,  
473        92(2), 609–625.
- 474    Emery, N. J. (2000). The eyes have it: The neuroethology, function and evolution of social  
475        gaze. *Neuroscience & Biobehavioral Reviews*, 24(6), 581–604.
- 476    Goodman, N. D., & Stuhlmüller, A. (2014). *The design and implementation of probabilistic  
477        programming languages*. <http://dippl.org>.
- 478    Gredebäck, G., Fikke, L., & Melinder, A. (2010). The development of joint visual  
479        attention: A longitudinal study of gaze following during interactions with mothers and  
480        strangers. *Developmental Science*, 13(6), 839–848.
- 481    Henrich, J. (2016). *The secret of our success: How culture is driving human evolution,  
482        domesticating our species, and making us smarter*. princeton University press.
- 483    Hernik, M., & Broesch, T. (2019). Infant gaze following depends on communicative signals:  
484        An eye-tracking study of 5-to 7-month-olds in vanuatu. *Developmental Science*, 22(4),

- 485 e12779.
- 486 Heyes, C. (2018). *Cognitive gadgets*. Harvard University Press.
- 487 House, B. R., Kanngiesser, P., Barrett, H. C., Broesch, T., Cebioglu, S., Crittenden, A. N.,  
488 et al.others. (2020). Universal norm psychology leads to societal diversity in prosocial  
489 behaviour and development. *Nature Human Behaviour*, 4(1), 36–44.
- 490 Itakura, S. (2004). Gaze-following and joint visual attention in nonhuman animals. Wiley  
491 Online Library.
- 492 Itier, R. J., & Batty, M. (2009). Neural bases of eye and gaze processing: The core of social  
493 cognition. *Neuroscience & Biobehavioral Reviews*, 33(6), 843–863.
- 494 Kanngiesser, P., Schäfer, M., Herrmann, E., Zeidler, H., Haun, D., & Tomasello, M. (2022).  
495 Children across societies enforce conventional norms but in culturally variable ways.  
496 *Proceedings of the National Academy of Sciences*, 119(1), e2112521118.
- 497 Kano, F., & Call, J. (2014). Cross-species variation in gaze following and conspecific  
498 preference among great apes, human infants and adults. *Animal Behaviour*, 91,  
499 137–150.
- 500 Laland, K., & Seed, A. (2021). Understanding human cognitive uniqueness. *Annual Review  
501 of Psychology*, 72, 689–716.
- 502 Langton, S. R., Watt, R. J., & Bruce, V. (2000). Do the eyes have it? Cues to the  
503 direction of social attention. *Trends in Cognitive Sciences*, 4(2), 50–59.
- 504 Lee, M. D., & Wagenmakers, E.-J. (2014). *Bayesian cognitive modeling: A practical course*.  
505 Cambridge University Press.
- 506 Legare, C. H. (2019). The development of cumulative cultural learning. *Annual Review of  
507 Developmental Psychology*, 1, 119–147.
- 508 Lempers, J. D., Flavell, E. R., & Flavell, J. H. (1977). The development in very young  
509 children of tacit knowledge concerning visual perception. *Genetic Psychology  
510 Monographs*, 95(1), 3–53.
- 511 Mayer, A., & Träuble, B. E. (2013). Synchrony in the onset of mental state understanding

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

28

- 512 across cultures? A study among children in samoa. *International Journal of Behavioral*  
513 *Development*, 37(1), 21–28.
- 514 Michel, C., Kayhan, E., Pauen, S., & Hoehl, S. (2021). Effects of reinforcement learning on  
515 gaze following of gaze and head direction in early infancy: An interactive eye-tracking  
516 study. *Child Development*, 92(4), e364–e382.
- 517 Miller, J. G., Wice, M., & Goyal, N. (2018). Contributions and challenges of cultural  
518 research on the development of social cognition. *Developmental Review*, 50, 65–76.
- 519 Moll, H., & Tomasello, M. (2004). 12-and 18-month-old infants follow gaze to spaces  
520 behind barriers. *Developmental Science*, 7(1), F1–F9.
- 521 Moore, C. (2008). The development of gaze following. *Child Development Perspectives*,  
522 2(2), 66–70.
- 523 Movellan, J. R., & Watson, J. S. (2002). The development of gaze following as a bayesian  
524 systems identification problem. *Proceedings 2nd International Conference on*  
525 *Development and Learning. ICDL 2002*, 34–40. IEEE.
- 526 Norenzayan, A., & Heine, S. J. (2005). Psychological universals: What are they and how  
527 can we know? *Psychological Bulletin*, 131(5), 763.
- 528 Perner, J., Ruffman, T., & Leekam, S. R. (1994). Theory of mind is contagious: You catch  
529 it from your sibs. *Child Development*, 65(4), 1228–1238.
- 530 Prein, J. C., Kalinke, S., Haun, D. B., & Bohn, M. (2023). TANGO: A reliable,  
531 open-source, browser-based task to assess individual differences in gaze understanding  
532 in 3 to 5-year-old children and adults. *Behavior Research Methods*, 1–17.
- 533 Prein, J. C., Maurits, L., Werwach, A., Haun, D. B. M., & Bohn, M. (2023). Variation in  
534 gaze understanding across the life span: A process-level perspective. *PsyArXiv*.  
535 <https://doi.org/10.31234/osf.io/dy73a>
- 536 R Core Team. (2023). *R: A language and environment for statistical computing*. Vienna,  
537 Austria: R Foundation for Statistical Computing. Retrieved from  
538 <https://www.R-project.org/>

- 539 Rakoczy, H. (2022). Foundations of theory of mind and its development in early childhood.
- 540 *Nature Reviews Psychology*, 1(4), 223–235.
- 541 Richardson, D. C., & Dale, R. (2005). Looking to understand: The coupling between
- 542 speakers' and listeners' eye movements and its relationship to discourse comprehension.
- 543 *Cognitive Science*, 29(6), 1045–1060.
- 544 Rosati, A. G., & Hare, B. (2009). Looking past the model species: Diversity in
- 545 gaze-following skills across primates. *Current Opinion in Neurobiology*, 19(1), 45–51.
- 546 Rossano, F. (2012). Gaze in conversation. *The Handbook of Conversation Analysis*,
- 547 308–329.
- 548 Scaife, M., & Bruner, J. S. (1975). The capacity for joint visual attention in the infant.
- 549 *Nature*, 253(5489), 265–266.
- 550 Sebanz, N., Bekkering, H., & Knoblich, G. (2006). Joint action: Bodies and minds moving
- 551 together. *Trends in Cognitive Sciences*, 10(2), 70–76.
- 552 Senju, A., Vernetti, A., Ganea, N., Hudry, K., Tucker, L., Charman, T., & Johnson, M. H.
- 553 (2015). Early social experience affects the development of eye gaze processing. *Current*
- 554 *Biology*, 25(23), 3086–3091.
- 555 Shepherd, S. V. (2010). Following gaze: Gaze-following behavior as a window into social
- 556 cognition. *Frontiers in Integrative Neuroscience*, 4, 5.
- 557 Stengelin, R., Ball, R., Maurits, L., Kanngiesser, P., & Haun, D. B. (2023). Children
- 558 over-imitate adults and peers more than puppets. *Developmental Science*, 26(2),
- 559 e13303.
- 560 Tang, Y., Gonzalez, M. R., & Deák, G. O. (2023). The slow emergence of gaze-and
- 561 point-following: A longitudinal study of infants from 4 to 12 months. *Developmental*
- 562 *Science*, e13457.
- 563 Taumoepeau, M., Sadeghi, S., & Nobilo, A. (2019). Cross-cultural differences in children's
- 564 theory of mind in iran and new zealand: The role of caregiver mental state talk.
- 565 *Cognitive Development*, 51, 32–45.

GAZE-FOLLOWING ACROSS 17 COMMUNITIES

30

- 566 Thorup, E., Nyström, P., Gredebäck, G., Bölte, S., Falck-Ytter, T., & Team, E. (2016).  
567 Altered gaze following during live interaction in infants at risk for autism: An eye  
568 tracking study. *Molecular Autism*, 7, 1–10.
- 569 Thorup, E., Nyström, P., Gredebäck, G., Bölte, S., Falck-Ytter, T., & Team, E. (2018).  
570 Reduced alternating gaze during social interaction in infancy is associated with elevated  
571 symptoms of autism in toddlerhood. *Journal of Abnormal Child Psychology*, 46,  
572 1547–1561.
- 573 Tomasello, M. (2019). *Becoming human: A theory of ontogeny*. Harvard University Press.
- 574 Tomasello, M. (2020). The adaptive origins of uniquely human sociality. *Philosophical  
Transactions of the Royal Society B*, 375(1803), 20190493.
- 575 Tomasello, M., Call, J., & Hare, B. (1998). Five primate species follow the visual gaze of  
577 conspecifics. *Animal Behaviour*, 55(4), 1063–1069.
- 578 Tomasello, M., Hare, B., Lehmann, H., & Call, J. (2007). Reliance on head versus eyes in  
579 the gaze following of great apes and human infants: The cooperative eye hypothesis.  
580 *Journal of Human Evolution*, 52(3), 314–320.
- 581 Tomasello, M., & Rakoczy, H. (2003). What makes human cognition unique? From  
582 individual to shared to collective intentionality. *Mind & Language*, 18(2), 121–147.
- 583 Van Leeuwen, E. J., Cohen, E., Collier-Baker, E., Rapold, C. J., Schäfer, M., Schütte, S., &  
584 Haun, D. B. (2018). The development of human social learning across seven societies.  
585 *Nature Communications*, 9(1), 2076.
- 586 Wellman, H. M. (2014). *Making minds: How theory of mind develops*. Oxford University  
587 Press.

## **Study IV**

# MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

1

# 1 Measuring variation in gaze following across communities, ages, and individuals – a 2 showcase of the TANGO-CC

Julia Christin Prein (ORCID: 0000-0002-3154-6167)<sup>1,2</sup>, Florian M. Bednarski (ORCID: 0000-0003-4384-4791)<sup>3</sup>, Ardain Dzabatou<sup>4</sup>, Michael C. Frank (ORCID: 0000-0002-7551-4378)<sup>5</sup>, Annette M. E. Henderson (ORCID: 0000-0003-4384-4791)<sup>3</sup>, Josefine Kalbitz<sup>2</sup>, Patricia Kanngiesser (ORCID: 0000-0003-1068-3725)<sup>6</sup>, Dilara Keşifoğlu (ORCID: 0000-0002-7356-0733)<sup>7</sup>, Bahar Köyメン (ORCID: 0000-0001-5126-8240)<sup>8</sup>, Maira V. Manrique-Hernandez<sup>2</sup>, Shirley Magazi (ORCID: 0009-0006-0479-9800)<sup>9</sup>, Lizbeth Mújica-Manrique<sup>2</sup>, Julia Ohlendorf<sup>2</sup>, Damilola Olaoba<sup>2</sup>, Wesley R. Pieters (ORCID: 0000-0002-6152-249X)<sup>9</sup>, Sarah Pope-Caldwell<sup>2</sup>, Umay Sen (ORCID: 0000-0001-9488-0851)<sup>10</sup>, Katie Slocombe (ORCID: 0000-0002-7310-1887)<sup>11</sup>, Robert Z. Sparks (ORCID: 0000-0001-7545-0522)<sup>5</sup>, Roman Stengelin (ORCID: 0000-0003-2212-4613)<sup>2,9</sup>, Jahnavi Sunderarajan<sup>2</sup>, Kirsten Sutherland<sup>2</sup>, Florence Tusiime<sup>12</sup>, Wilson Vieira (ORCID: 0009-0001-9400-6328)<sup>2</sup>, Zhen Zhang (ORCID: 0000-0001-9300-0920)<sup>13</sup>, Yufei Zong (ORCID: 0009-0000-5012-0244)<sup>13</sup>, Daniel B. M. Haun (ORCID: 0000-0002-3262-645X)<sup>2,+</sup>, and & Manuel Bohn (ORCID: 0000-0001-6006-1348)<sup>1,2,+</sup>

<sup>16</sup> Institute of Psychology in Education  
<sup>17</sup> Leuphana University Lüneburg

## MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

2



MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

4

58

**Author Note**

59

60       The authors made the following contributions. Julia Christin Prein (ORCID:  
61       0000-0002-3154-6167): Conceptualization, Formal Analysis, Methodology, Project  
62       administration, Resources, Software, Visualization, Writing - Original Draft Preparation,  
63       Writing - Review & Editing; Florian M. Bednarski (ORCID: 0000-0003-4384-4791): Methodology,  
64       Resources, Writing - Review & Editing; Ardain Dzabatou: Methodology, Resources, Writing -  
65       Review & Editing; Michael C. Frank (ORCID: 0000-0002-7551-4378): Methodology, Resources,  
66       Writing - Review & Editing; Annette M. E. Henderson (ORCID: 0000-0003-4384-4791):  
67       Methodology, Resources, Writing - Review & Editing; Josefine Kalbitz: Methodology, Resources,  
68       Writing - Review & Editing; Patricia Kanngiesser (ORCID:0000-0003-1068-3725): Methodology,  
69       Resources, Writing - Review & Editing; Dilara Keşsafoğlu (ORCID: 0000-0002-7356-0733):  
70       Methodology, Resources, Writing - Review & Editing; Bahar Köyメン (ORCID:  
71       0000-0001-5126-8240): Methodology, Resources, Writing - Review & Editing; Maira V.  
72       Manrique-Hernandez: Methodology, Resources, Writing - Review & Editing; Shirley Magazi  
73       (ORCID: 0009-0006-0479-9800): Methodology, Resources, Writing - Review & Editing; Lizbeth  
74       Mújica-Manrique: Methodology, Resources, Writing - Review & Editing; Julia Ohlendorf:  
75       Methodology, Resources, Writing - Review & Editing; Damilola Olaoba: Methodology,  
76       Resources, Writing - Review & Editing; Wesley R. Pieters (ORCID:0000-0002-6152-249X):  
77       Methodology, Resources, Writing - Review & Editing; Sarah Pope-Caldwell: Methodology,  
78       Resources, Writing - Review & Editing; Umay Sen (ORCID: 0000-0001-9488-0851): Methodology,  
79       Resources, Writing - Review & Editing; Katie Slocombe (ORCID: 0000-0002-7310-1887):  
80       Methodology, Resources, Writing - Review & Editing; Robert Z. Sparks (ORCID:  
81       0000-0001-7545-0522): Methodology, Resources, Writing - Review & Editing; Roman Stengelin  
82       (ORCID: 0000-0003-2212-4613): Conceptualization, Methodology, Resources, Writing - Review  
83       & Editing; Jahnavi Sunderarajan: Methodology, Resources, Writing - Review & Editing; Kirsten

<sup>84</sup> Sutherland: Methodology, Resources, Writing - Review & Editing; Florence Tusiime:  
<sup>85</sup> Methodology, Resources, Writing - Review & Editing; Wilson Vieira (ORCID:  
<sup>86</sup> 0009-0001-9400-6328): Methodology, Resources, Writing - Review & Editing; Zhen Zhang  
<sup>87</sup> (ORCID: 0000-0001-9300-0920): Methodology, Resources, Writing - Review & Editing; Yufei  
<sup>88</sup> Zong (ORCID: 0009-0000-5012-0244): Methodology, Resources, Writing - Review & Editing;  
<sup>89</sup> Daniel B. M. Haun (ORCID: 0000-0002-3262-645X): Conceptualization, Funding acquisition,  
<sup>90</sup> Methodology, Writing - Review & Editing; Manuel Bohn (ORCID: 0000-0001-6006-1348):  
<sup>91</sup> Conceptualization, Funding acquisition, Methodology, Project administration, Supervision,  
<sup>92</sup> Writing - Review & Editing.

<sup>93</sup> Correspondence concerning this article should be addressed to Julia Christin Prein  
<sup>94</sup> (ORCID: 0000-0002-3154-6167), Universitätsallee 1, 21335 Lüneburg, Germany. E-mail:  
<sup>95</sup> julia.prein@leuphana.de

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

6

96

**Abstract**

97 Cross-cultural studies are crucial for investigating the cultural variability and universality of  
98 cognitive developmental processes. However, cross-cultural assessment tools in cognition  
99 across languages and communities are limited. This paper describes a gaze following task  
100 designed to measure basic social cognition across individuals, ages, and communities  
101 (TANGO-CC). The task was developed and psychometrically assessed in one cultural setting  
102 and, with input of local collaborators, adapted for cross-cultural data collection. Minimal  
103 language demands and the web-app implementation allow fast and easy contextual adaptations  
104 to each community. The TANGO-CC captures individual- and community-level variation and  
105 shows good internal consistency in a data set from 2.5- to 11-year-old children from 17 diverse  
106 communities. Within-community variation outweighed between-community variation. We  
107 provide an open-source website for researchers to customize and use the task  
108 (<https://ccp-odc.eva.mpg.de/tango-cc>). The TANGO-CC can be used to assess basic social  
109 cognition in diverse communities and provides a roadmap for researching community-level and  
110 individual-level differences across cultures.

111       *Keywords:* cross-cultural psychology, social cognition, gaze following, individual  
112 differences, reliability

113 Word count: XXX

114 Measuring variation in gaze following across communities, ages, and individuals – a  
115 showcase of the TANGO-CC

## Introduction

For decades, researchers have advocated for more diverse samples in psychological research and cautioned against relying solely on convenience samples from the Global North (Arnett, 2008; Henrich et al., 2010; Lillard, 1998). Despite numerous calls for change, the samples reported in major psychology journals still lack diversity (Apicella et al., 2020; Gutchess & Rajaram, 2023; Thalmayer et al., 2021). This hinders progress in theory building and testing: we cannot draw inferences about universal and variable aspects of the human cognitive system from data collected exclusively in one single community (Krys et al., 2024). While this sampling bias is often discussed within adult psychology, it is equally relevant to developmental psychology (Nielsen et al., 2017). Early experiences shape the way children think about and interact with the world and an ontogenetic perspective is needed to explore the foundational aspects of human behavioral diversity (Amir & McAuliffe, 2020; Broesch et al., 2023; Liebal & Haun, 2018; Torréns et al., 2023).

129 There are numerous challenges with collecting cross-cultural, developmental data (Amir  
130 & McAuliffe, 2020; Broesch et al., 2023). Cross-cultural studies need reliable and valid measures  
131 to capture variation between communities and/or individuals systematically. Even though this  
132 applies to all areas of cognitive development, we focus on social cognition in this paper.

Social cognition refers to how an individual processes information in social situations which allows them to understand and predict others' behavior (Adolphs, 1999; Decety, 2020; Frith & Frith, 2007; Zeigler-Hill et al., 2015). If, in theory, stimuli used in social cognition tasks should relate to people's everyday experiences, then tasks themselves should be tuned to the features of specific communities. Indeed, task performance can be diminished when stimuli do not reflect the characteristics of the participants' communities (Peña, 2007). For example, Elfenbein and Ambady (2002) found better emotion recognition for members of the same

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

8

<sup>140</sup> national, ethnic, or regional group. Selcuk et al. (2023) concluded that children often attribute  
<sup>141</sup> mental states more accurately and more frequently to individuals from the same community.

<sup>142</sup> From a psychometric perspective, the situation looks dire. Studies on social cognition  
<sup>143</sup> with US-American and European samples rarely report psychometric information (for a review,  
<sup>144</sup> see Beaudoin et al., 2020), and the picture further deteriorates when we look at cross-cultural  
<sup>145</sup> social cognition tasks (Bourdage et al., 2023; Hajdúk et al., 2020; Waschl & Chen, 2022). Thus, it  
<sup>146</sup> is already challenging to find reliable and valid tasks that have measurement sensitivity to  
<sup>147</sup> detect individual differences within one community, let alone tasks that do so across different  
<sup>148</sup> communities. In this paper, we describe the construction and psychometric evaluation of a  
<sup>149</sup> cross-cultural measure of basic social cognition (gaze following) in children as a concrete  
<sup>150</sup> example of how to address this problem.

<sup>151</sup> The approaches that researchers can take to collect cross-cultural data lie on a  
<sup>152</sup> continuum: the decision for a specific method partly depends on whether researchers aim to  
<sup>153</sup> increase the *depth* (culture specificity) or *breadth* (standardization across multiple communities)  
<sup>154</sup> of their work (Amir & McAuliffe, 2020). At one extreme, researchers translate the psychological  
<sup>155</sup> construct into a separate design or task for each community (termed “assembly”; He and Vijver  
<sup>156</sup> (2012), Waschl and Chen (2022)). While this approach allows greater flexibility and sensitivity to  
<sup>157</sup> cultural differences, it might not be feasible to study a multitude of communities as it becomes  
<sup>158</sup> too demanding and time-consuming. Furthermore, the results are limited to each community  
<sup>159</sup> and absolute task scores might not be comparable across communities. A study following this  
<sup>160</sup> approach is Wefers et al. (2023) who investigated how cultural variations in parenting styles  
<sup>161</sup> modulated infants’ responses to disruptions in social interactions. While studies in the Global  
<sup>162</sup> North often apply the Still-Face Paradigm to assess infants’ reactions to unresponsive partners,  
<sup>163</sup> Wefers et al. (2023) reasoned that this paradigm might not capture infants’ everyday interaction  
<sup>164</sup> routines in communities with proximal (*i.e.*, emphasis on body stimulation) parenting styles. By  
<sup>165</sup> developing the novel No-Touch Paradigm, they found that indeed infants’ responses to

<sup>166</sup> unresponsive partners were modulated by the cultural context in which they grew up: Kichwa  
<sup>167</sup> infants from rural Ecuador showed stronger reactions to unresponsive partners in the No-Touch  
<sup>168</sup> Paradigm compared to the Still-Face Paradigm, while reactions of urban German infants  
<sup>169</sup> differed less in both paradigms.

<sup>170</sup> At the other extreme, researchers use the same standardized procedure across diverse  
<sup>171</sup> communities, potentially providing a simple translation or modification of stimuli to ensure  
<sup>172</sup> they are culturally appropriate (termed “adoption” and “adaptation”, respectively; Waschl &  
<sup>173</sup> Chen (2022)). This approach is less sensitive to each community’s unique characteristics but  
<sup>174</sup> renders quantitative comparisons of data more feasible. An example following this approach is  
<sup>175</sup> the Multilingual Assessment Instrument for Narratives (MAIN; (Gagarina et al., 2012)), which  
<sup>176</sup> assesses narrative abilities in mono- and multilingual children. Extensive piloting and  
<sup>177</sup> adaptation of MAIN materials ensured that the instrument is culturally appropriate, robust, and  
<sup>178</sup> suitable for cross-linguistic comparisons (Gagarina et al., 2012), and new and revised language  
<sup>179</sup> versions are continuously added to the MAIN database (Gagarina & Lindgren, 2020).

<sup>180</sup> The present paper aims to describe the development and psychometric properties of a  
<sup>181</sup> social cognition task that can be adapted to diverse communities. On the continuum described  
<sup>182</sup> above, our task lies more toward standardized approaches but allows for some customization of  
<sup>183</sup> the stimuli to each local community. The task focuses on one of the most fundamental  
<sup>184</sup> social-cognitive abilities: gaze following, that is, the ability to identify the attentional focus of  
<sup>185</sup> another agent. Gaze following develops early in infancy (Del Bianco et al., 2019; Tang et al.,  
<sup>186</sup> 2024) and contributes to social learning, communication, and collaboration (Bohn & Köymen,  
<sup>187</sup> 2018; Hernik & Broesch, 2019; Shepherd, 2010; Tomasello et al., 2007). While the question of  
<sup>188</sup> how social-environmental and cultural factors impact gaze following was recently posed as one  
<sup>189</sup> of the big open questions in gaze following research (Astor & Gredebäck, 2022), studies focusing  
<sup>190</sup> on cultural variations of gaze following are rare. Callaghan et al. (2011) investigated gaze  
<sup>191</sup> following behind barriers in 12- and 17-month-olds from rural Canada (n = 35), Peru (n = 38),

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

10

192 and India ( $n = 65$ ). In their setup, an agent looked at a toy behind (experimental condition) or a  
193 sticker in front of (control condition) a barrier, and children's crawling toward the barrier to  
194 follow the agent's gaze was assessed. While the absolute crawling rates differed across  
195 communities, children in all three communities crawled more often to gain visual access when  
196 the agent looked at an object behind the barrier than in front of it. Hernik and Broesch (2019)  
197 studied 22 infants between 5 to 7 months of age from Vanuatu and used an eye-tracking  
198 procedure displaying a local actor looking at one of two objects. Even though face-to-face  
199 interactions are less common in parent-child interactions in Vanuatu than in Western  
200 communities, the result patterns of gaze following in Ni-Vanuatu infants resembled those of  
201 their Western counterparts: infants from Vanuatu followed the agent's gaze to the target object  
202 only when preceded by infant-directed, and not adult-directed, speech. While these two studies  
203 point towards potential cross-cultural stability of gaze following, the lack of psychometrically  
204 evaluated tasks and the small number of communities studied limit the generalizability of these  
205 findings and – more importantly – does not allow for studying individual-differences. A more  
206 comprehensive cross-cultural study on gaze following was recently conducted by Bohn et al.  
207 (2024). The researchers tested the universality of gaze following by studying 17 different  
208 communities on five continents and found evidence for a similar processing mechanism across  
209 communities.

210 The task presented in this manuscript was developed for the study by Bohn et al. (2024)  
211 and is based on a previously established gaze-following task called 'TANGO' (Task for  
212 Assessing iNdividual differences in Gaze understanding - Open) by Prein et al. (2023). The  
213 TANGO measures how precisely participants locate an agent's attentional focus. It reliably  
214 measured individual differences in a German child sample and an English-speaking remote  
215 adult sample (Prein et al., 2023). However, we cannot claim the task's generalizability and  
216 reliability based on a mono-cultural sample. This paper showcases the TANGO-CC (TANGO –  
217 Cross-Cultural), a standardized gaze following task that has been adapted to 13 languages and  
218 even more communities, and evaluates its psychometric properties by leveraging a large and

<sup>219</sup> diverse data set of 2.5- to 11-year-olds from 17 diverse communities (Bohn et al., 2024). We  
<sup>220</sup> describe the task's development and provide a tutorial for the open-source website  
<sup>221</sup> (<https://ccp-odc.eva.mpg.de/tango-cc/>).

<sup>222</sup> **Task development**

<sup>223</sup> **Approach**

<sup>224</sup> The TANGO-CC was implemented in Leipzig, Germany, and thoroughly assessed in  
<sup>225</sup> terms of reliability and validity (Prein et al., 2023). During this process, the cross-cultural  
<sup>226</sup> adaptation of the task was prepared by a team of cross-cultural psychologists and cognitive  
<sup>227</sup> scientists. In this paper, we assess the TANGO-CC's measurement quality (*i.e.*, variability and  
<sup>228</sup> reliability) across 17 diverse communities by analyzing the data set from Bohn et al. (2024). In  
<sup>229</sup> the following, we describe the different steps in detail.

<sup>230</sup> The TANGO-CC is a screen-based task that measures the imprecision with which  
<sup>231</sup> participants locate a balloon by following an agent's gaze (see Figure 1). Participants click or  
<sup>232</sup> touch the location on the screen where they believe the balloon to be. Precision is measured as  
<sup>233</sup> the distance between the participant's click on the screen and the balloon's real position.

<sup>234</sup> During the task development, we decided to implement the task's main functionality  
<sup>235</sup> independently of the task's appearance. We programmed a function that calculates the x and y  
<sup>236</sup> coordinates of where the agent's pupil and iris should move to follow the balloon, given the  
<sup>237</sup> eyes' and balloon's original positions and measures. As the measures of the eyes and balloon  
<sup>238</sup> are read out dynamically from the image on screen, stimuli can be easily adapted and  
<sup>239</sup> exchanged (*i.e.*, no coordination values for animation are hard-coded into the task's source  
<sup>240</sup> code). After having programmed this "backbone" functionality of the task (*i.e.*, animate the eyes  
<sup>241</sup> so that they follow the balloon), we added the task's audio instructions and superficial  
<sup>242</sup> appearance (e.g., background scene, hedge, agent faces).

<sup>243</sup> This basic version of the TANGO was psychometrically evaluated in a German child

## MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

12

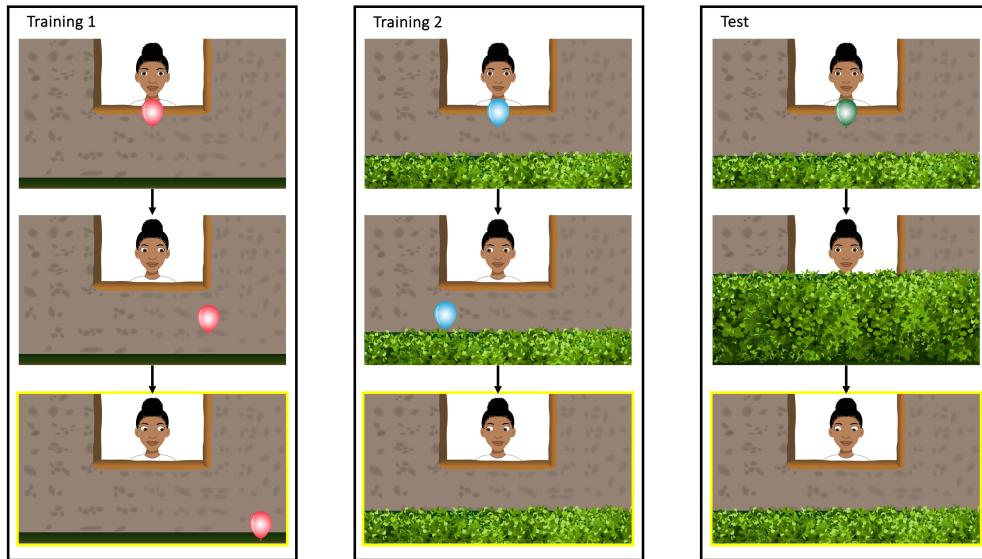
244 sample and an English-speaking remote adult sample and was found to be highly reliable and  
245 valid (Prein et al., 2023). While children became more precise in locating the agent's attentional  
246 focus with age, individuals differed across all age groups and showed no floor- or ceiling effects.  
247 In the Leipzig sample, performance in the TANGO was weakly related to factors of children's  
248 daily social environment and could predict children's receptive vocabulary 6 months later. In a  
249 computational cognitive model, Prein et al. (2024) described gaze following as a form of social  
250 vector following and empirically found that performance in the TANGO was related to  
251 children's non-social vector following and visual perspective-taking abilities. These connections  
252 to related constructs indicate the task's convergent validity in the German child sample.

253 To adapt the task for cross-cultural data collection, we generated a set of human cartoon  
254 faces and background scenes with input from local researchers and research assistants. The  
255 stimulus pool was adjusted and expanded until the researchers and research assistants from  
256 each target community judged the selected stimuli to be representative of the local population  
257 and typical accommodation (see Figure 2). Audio instructions were translated from English or  
258 German into the corresponding local language(s). By back-translating these instructions, we  
259 ensured the original meaning did not change. Sometimes, specific words were slightly modified  
260 in the target language (e.g., "bush" instead of "hedge") to ensure that all participants understood  
261 the instructions well. Based on these adaptations, the TANGO-CC could be applied in 17  
262 communities and 13 different languages (Bohn et al., 2024). In the following, we describe how  
263 researchers can use and customize the TANGO-CC in more detail.

### 264 **Features of the TANGO-CC**

#### 265 **Trials**

266 The task consists of three different trial types: training 1, training 2, and test trials (see  
267 Figure 1). In every trial, participants see an agent (boy or girl) looking out of a house with a  
268 balloon (red, blue, green, or yellow) in front of them. The balloon falls down to the ground. The  
269 eyes of the agent follow the movement of the balloon in a way that the balloon center and the



**Figure 1**

**Screenshots of the trials.** In training 1, an agent looks at a balloon that falls to the ground, and participants have to respond by clicking/touching the balloon. In training 2, the balloon falls behind the hedge while its flight is still visible. Participants respond by clicking the hedge where they think the balloon is. In test trials, the balloon's movement and final position are covered by a hedge, and participants respond by clicking the hedge. In the task, all movements are smoothly animated (no still pictures). Yellow frames indicate the time point when participants respond (only illustrative, not shown during the task).

270 pupil center always align. Depending on the trial type, participants have different visual access  
 271 to the balloon's position. In training 1, participants see the full trajectory of the balloon and  
 272 directly have to click the balloon itself. In training 2, participants see most of the balloon's  
 273 movement, but a hedge covers the final location. In test trials, a hedge grows at the beginning  
 274 of the trial and participants see neither the movement nor the final position of the balloon. The  
 275 first trial of each type contains an audio description of the presented events (see Supplements).  
 276 Notably, the instructions explicitly state that the agent is looking at the balloon.

277 The outcome variable is the distance between the participant's click and the balloon's  
 278 center. Trials can be completed quickly and efficiently so that children can complete 15 trials  
 279 within 10 minutes, and few children fail to complete the task. By using self-explanatory

<sup>280</sup> animations, language demands are kept to a minimum. The task uses simple audio instructions,  
<sup>281</sup> which makes the task accessible to children from different age groups, and no reading skills are  
<sup>282</sup> required. There is no feedback during the task to prevent learning effects across trials.

<sup>283</sup> ***Randomization***

<sup>284</sup> The order of the agents, balloon colors (red, yellow, green, blue), and balloon positions  
<sup>285</sup> are each randomized independently. For the balloon positions, the entire width of the screen  
<sup>286</sup> (1920 in “SVG units”) is divided into ten bins. Exact coordinates (value between 0 for the far left  
<sup>287</sup> and 1920 for the far right) within each bin are then randomly generated. The number of  
<sup>288</sup> repetitions for each agent, balloon color, and balloon bin is calculated based on the total number  
<sup>289</sup> of trials and the number of unique agents, balloon colors, and bins, respectively. All agents,  
<sup>290</sup> balloon colors, and bins appear equally often and are not repeated in more than two consecutive  
<sup>291</sup> trials. If the total number of trials is not divisible by the number of unique elements, some  
<sup>292</sup> elements (*i.e.*, some agents, balloon colors, bins) are randomly repeated to make up for the  
<sup>293</sup> remainder.

<sup>294</sup> ***Cross-cultural customization***

<sup>295</sup> When researchers visit the TANGO–CC’s website  
<sup>296</sup> (<https://ccp-odc.eva.mpg.de/tango-cc/>), they can select the language for audio instructions  
<sup>297</sup> which are currently available for 13 different languages and five more dialects (see Table 1). To  
<sup>298</sup> add a new language, researchers have two options: (1) for using their own audio instructions in  
<sup>299</sup> the offline version of the task, researchers can download the task, exchange the audio  
<sup>300</sup> instructions in the dist folder (in the folder sounds > custom) and select “Custom” in the  
<sup>301</sup> language drop-down menu. For detailed instructions, see the TANGO–CC’s manual  
<sup>302</sup> (<https://ccp-odc.eva.mpg.de/tango-cc/manual.html>). (2) For adding a new language in the online  
<sup>303</sup> version of the task, researchers can contact the first author of this paper. Please note that this  
<sup>304</sup> option requires new audio recordings by the interested researchers, which will then be openly  
<sup>305</sup> available for all users of the task. All written instructions in the task are solely for the research  
<sup>306</sup> assistant to help them guide participants through the task; these instructions are solely available

307 in English. The task can either be started with the default settings or further customized by  
308 adapting the number of trials, agents, and background scenes. The default settings use the  
309 version applied in Bohn et al. (2024) based on the selected language (see Supplements).

310 If researchers choose to customize the task (see Figure 2), they can adjust the number of  
311 trials for each trial type, but not their sequence. Specifically, trial types build on each other and  
312 participants need to complete each trial type (without skipping any) to understand the structure  
313 of the task. The minimum number of trials per type is 1; the maximum is 100. Furthermore,  
314 researchers can customize backgrounds by selecting one of four different backgrounds. Finally,  
315 researchers can choose from 50 diverse cartoon-like human faces (50% female, 50% male) and  
316 freely select how many different faces to include (min 1, max 50). Once all the settings are  
317 selected, the customized task is compiled. To save the selected settings, researchers can  
318 bookmark the URL to easily access the customized task.

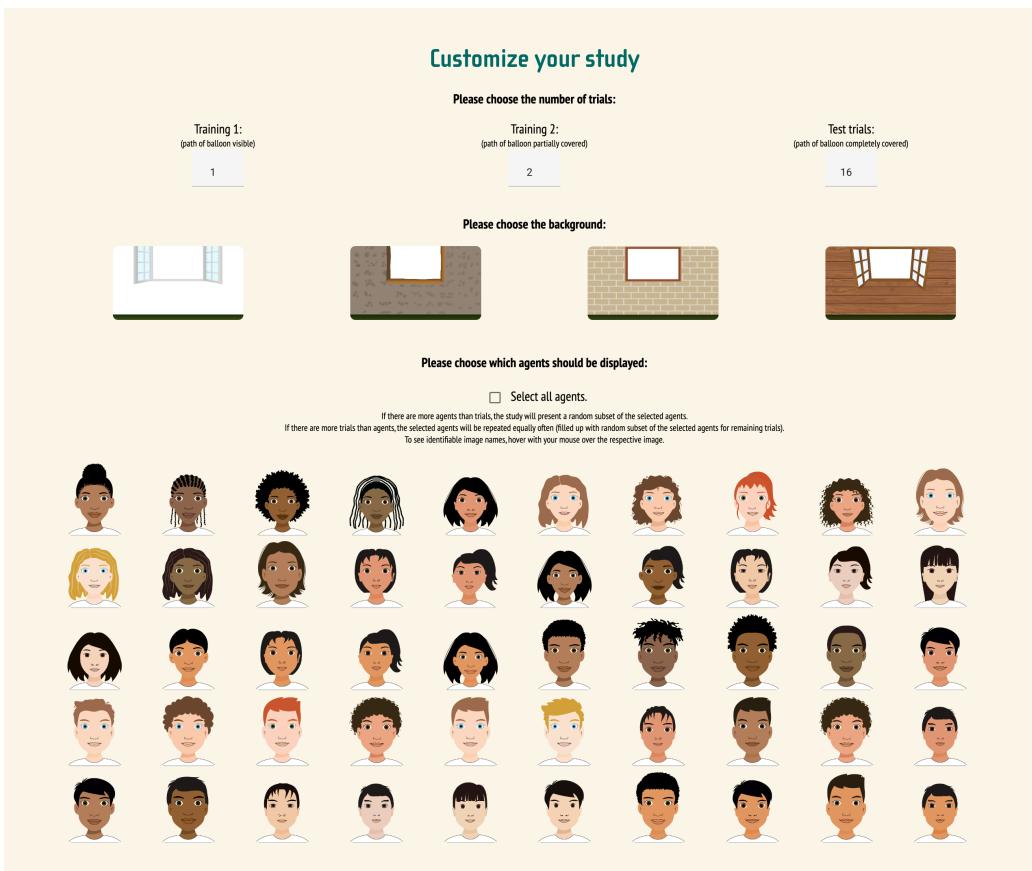
319 In the last step, researchers can enter an alphanumeric participant identifier (1 - 8  
320 characters) and enable a webcam recording of the participant, if needed. A webcam recording  
321 might prove especially helpful for unsupervised online data collection to ensure that the  
322 participant is alone during the task and no help is provided. The participant identifier and  
323 webcam choice have to be provided every time the task is run.

324 The source code of the task is openly available on GitHub  
325 (<https://github.com/ccp-eva/tango-cc>). By directly editing the HTML and JavaScript code,  
326 researchers can further modify the task as needed.

327 We created a public OSF page (<https://doi.org/10.17605/OSF.IO/P2EGU>) on which we  
328 plan to collect data sets that used the TANGO-CC. Researchers who have collected data using  
329 the TANGO-CC can share their data with the community by contacting the first author of this  
330 paper or visiting the OSF repository.

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

16



**Figure 2**

**Screenshot of the customizable components of the TANGO-CC.** Researchers can select the language of the audio instructions (see Table 1), the number of trials per trial type, the background, and the agent's face.

**Table 1***Current language options available for the audio instructions in the TANGO-CC*

Languages	Language family	Speaker's country of origin
Bemba	Bantu	Zambia
Chinese	Sino-Tibetan	China
English	Indo-European	USA / UK / India / Nigeria / New Zealand
German	Indo-European	Germany
(#Akhoe) Hai  om	Khoe-Kwadi	Namibia
Khwedam	Khoe-Kwadi	Namibia
Lingala	Bantu	Rep. Congo
Marathi	Indo-European	India
Shona *	Bantu	Zimbabwe
Spanish	Indo-European	Argentina (Rioplatense Spanish) / Mexico (Mexiquense Spanish)
Kiswahili	Bantu	Uganda
Turkish	Turkic	Türkiye
Yaka	Bantu	Rep. Congo

*Note.* In cases where more than one speaker's country of origin is listed, the audio instructions were recorded multiple times by different speakers. For example, the English instructions are available in five different versions. \* Please note that audio instructions are available in Shona but no data of this version is included in the present data set.

<sup>331</sup> **Task implementation**

<sup>332</sup> The task was implemented in JavaScript, HTML, and CSS and is presented as a  
<sup>333</sup> web app. It can be accessed on any modern web browser on any device (e.g., computer or tablet)  
<sup>334</sup> and does not require prior installation (though please note that configurations of browsers and  
<sup>335</sup> JavaScript may change in the future). Participant's responses can be recorded on a touchscreen  
<sup>336</sup> or with a mouse or trackpad. The online version of the task can be used for unsupervised data  
<sup>337</sup> collection (for example, using online platforms like *Prolific*; see Prein et al. (2023)). The task can  
<sup>338</sup> be shared easily internationally by providing the URL. Importantly, the web app  
<sup>339</sup> implementation does not require a working WIFI connection: An offline version of the task can  
<sup>340</sup> be downloaded and quickly set up for devices that support Node.js (<https://nodejs.org/en>). This  
<sup>341</sup> is an especially useful feature for researchers working in locations with limited internet access.

<sup>342</sup> The stimuli are embedded as Scalable Vector Graphics (SVG; an image format that stores

## MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

18

343 image elements via mathematical formulas based on points and lines on a grid). SVGs ensure  
 344 that the picture quality, aspect ratio, and relative object positioning are constant. Furthermore,  
 345 stimuli are added as individual components to the image scene which allows for an easy  
 346 adaptation of the task's elements (in contrast to other image formats that consist of only one  
 347 combined layer that would need entire replacement). The task is programmed so that responses  
 348 are only registered when the participant clicks the relevant part of the screen (*i.e.*, in test trials,  
 349 when they click the hedge). Furthermore, clicks are only registered after the voice recordings  
 350 stop playing. An audio reminder is played again if no click is registered within 5 seconds.

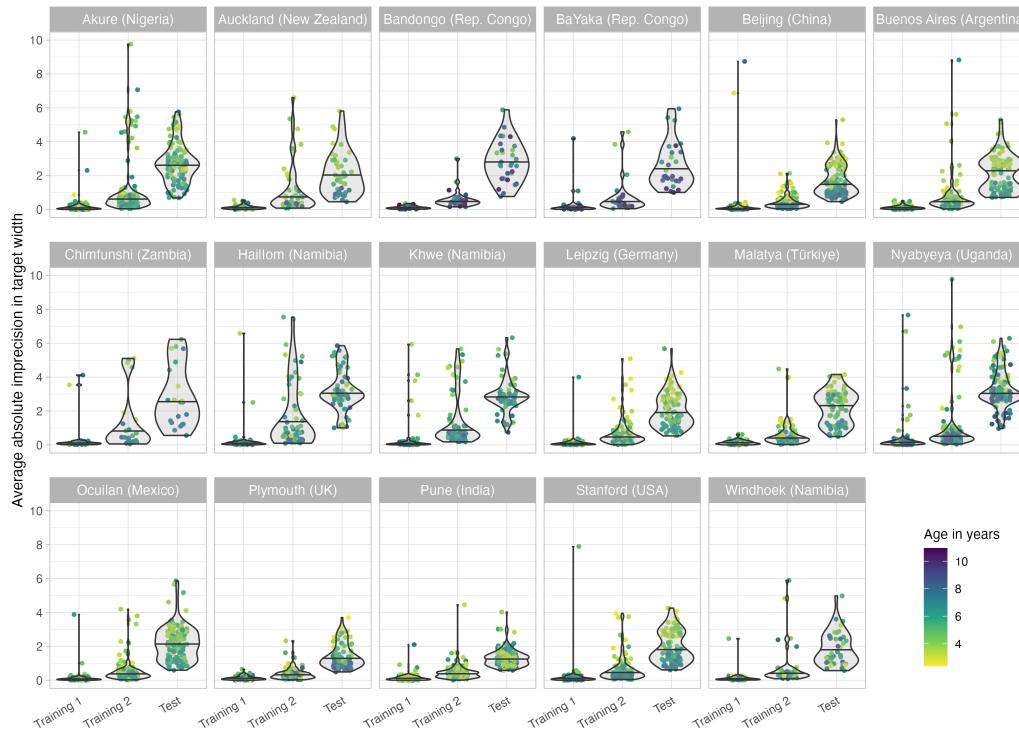
351 The website does not use cookies, nor does it upload any data to servers; that is, the data  
 352 is only stored locally on the device. The output of the task is a CSV file (and WEBM file if a  
 353 webcam recording was selected) that contains the participants' responses and can be easily  
 354 imported into statistical software for further analysis. The file will be stored in the device's  
 355 downloads folder and is named after the following pattern:  
 356 “tangoCC-participantID-YYYY-MM-DD\_hh\_mm\_ss”. To modify the storage location on the  
 357 device, researchers can change the designated downloads folder in their browser settings.

358 **Psychometric evaluation**

359 **Data set**

360 We used the data set from Bohn et al. (2024) for the psychometric evaluation of the  
 361 TANGO-CC. The data set contains a sample of  $N = 1377$  children, aged 2.5 to 11 years.  
 362 Participants came from 17 communities on five continents, in rural and urban settings, with  
 363 varying degrees of market integration and technology exposure. Bohn et al. (2024) carried out  
 364 19 trials (1 training 1, 2 training 2, and 16 test trials, of which the first of each type had audio  
 365 instructions) on a touchscreen device. Faces, backgrounds, and languages were chosen by  
 366 researchers and assistants with experience in the specific community. For further details on the  
 367 communities, participant information, and data collection procedures, see the supplements of  
 368 Bohn et al. (2024).

<sup>369</sup> **Individual differences**



**Figure 3**

**Variability measured by the TANGO-CC.** Mean imprecision in locating the agent's attentional focus by community (alphabetically) and trial type. Imprecision is defined as the distance between the participant's click and the balloon's center in units of balloon width. For a depiction of each trial's procedure, see Figure 1.

<sup>370</sup> All analyses and the data set can be accessed on GitHub

<sup>371</sup> (<https://github.com/ccp-eva/tango-cc-methods/>). First, we inspected the mean and standard <sup>372</sup> deviations by community and compared performance in each trial type (training 1, training 2, <sup>373</sup> test trials). Performance was defined as the absolute distance between the target center and the <sup>374</sup> x coordinate of the participant's click (measured in balloon widths). Across communities, <sup>375</sup> children performed best in training 1 (mean = 0.19, sd = 0.63), followed by training 2 (mean = <sup>376</sup> 0.79, sd = 1.44) and test trials (mean = 2.21, sd = 2.03; see Figure 3).

377 To formally estimate the effect of trial type on performance in the TANGO-CC, we  
 378 fitted a generalized linear mixed model (GLMM) predicting the task performance by trial type  
 379 (reference category: test trials). All analyses were run in R version 4.4.0 (2024-04-24) (R Core  
 380 Team, 2024). GLMMs were fitted with default priors using the function `brm` from the package  
 381 `brms` (Bürkner, 2017, 2018). The model included random effects for trial type by community  
 382 (model notation in R: `imprecision ~ trialtypes + (triaitype |`  
 383 `community)`), and imprecision was modeled by a lognormal distribution. We inspected  
 384 the posterior distribution (mean and 95% Credible Interval (CrI)) for the trial type estimates.

385 Our GLMM analysis supported the visual inspection of the data: the fixed-effect  
 386 estimates for training 1 ( $\beta = -3.26$ ; 95% CrI [-3.41; -3.10]) and training 2 ( $\beta = -1.47$ ; 95% CrI [-1.58;  
 387 -1.35]) were negative and reliably different from zero.<sup>1</sup> This effect was found across all  
 388 communities (random effects of trial type within community: minimum estimate for training 1  
 389 = -2.87; 95% CrI [-3.11; -2.60]; minimum estimate for training 2 = -1.27; 95% CrI [-1.51; -0.98]). The  
 390 almost perfect performance in training trials indicated that children understood the task and  
 391 were able to correctly indicate the location of the balloon when its path was (mostly) visible. In  
 392 test trials, children's imprecision was higher, indicating that the task was more challenging. All  
 393 communities showed substantial individual variation and overlapped in their imprecision levels  
 394 (see Figure 3).

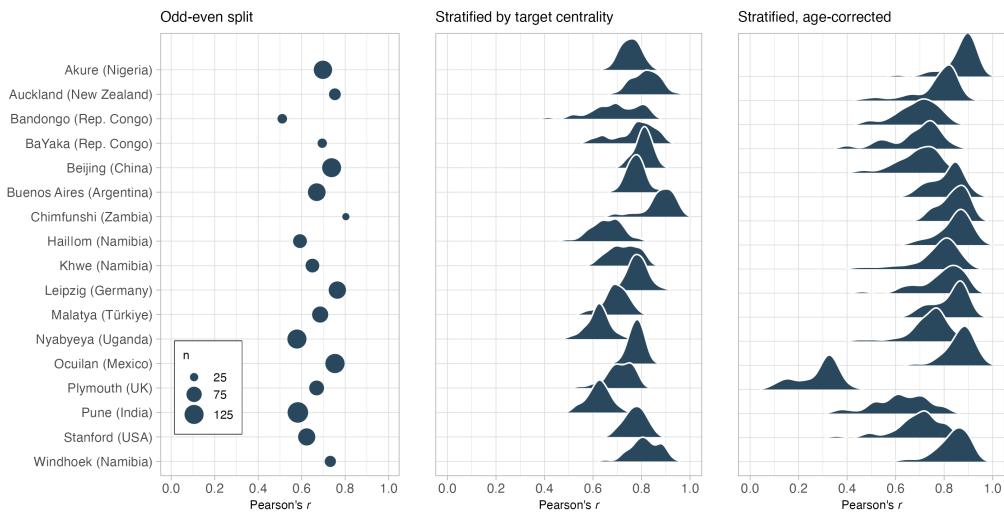
395 To identify the sources of variation, we computed intraclass correlations (ICC). The  
 396 variation in children's imprecision within communities was substantially larger than the  
 397 variation between the communities. The mean within-community variance was 1.28, ranging  
 398 from 0.24 (in Pune, India) to 3.46 (in Chimfunshi, Zambia). Between-community variance was  
 399 0.34. The ICC, representing the proportion of between-community variance relative to the total  
 400 variance (sum of within- and between-community variance), was 0.02. This indicated that only

---

<sup>1</sup> Please note that the TANGO-CC measures imprecision in gaze following. Therefore, a negative sign indicates that children showed less imprecision (*i.e.*, were more precise) in the training trials than in the test trials.

401 2% of the total variability in the data could be attributed to differences between communities,  
 402 while the remaining 98% were attributed to differences within communities (Kusano et al., 2024).

403 **Reliability**



**Figure 4**

**Reliability of the TANGO-CC by community.** Internal consistency estimates by community, following three different approaches. In the odd-even split, the size of points reflects the sample size in each community. In the stratified approach with and without age correction, density curves show the posterior distributions of the GLMM.

404 To assess reliability, we estimated internal consistency in each community in three  
 405 different ways. First, data of each participant was split into odd and even trials and a Pearson  
 406 correlation was calculated between the aggregated scores of the two halves. Second, using the  
 407 function `by_split` from the `splithalfr` package (Pronk et al., 2022), data was stratified  
 408 by target centrality (capturing trial difficulty), and a Pearson correlation was calculated between  
 409 the matched halves. Third, a data set was generated with stratified test halves by target  
 410 centrality and we applied the GLMM approach introduced by Rouder and Haaf (2019). A GLMM  
 411 was fitted with the mean imprecision as the outcome, age as the predictor, and test half and  
 412 participant id as random effects (model notation: `imprecision ~ age + (0 + half`  
 413 `| subj_id)`). The model estimates correlations between participant-specific estimates for

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

22

414 each test half. The hierarchical shrinkage of the model enables accurate person-specific  
415 estimates. By incorporating age as a fixed effect, the correlation between the two  
416 person-specific estimates represents the age-independent estimate for internal consistency. This  
417 removes the possibility that a good internal consistency estimate results from general cognitive  
418 development rather than task-specific inter-individual differences. Because the process of  
419 generating stratified data sets is partly random, the model was fitted 50 times for each  
420 community. The posterior estimate of the correlation between the two person-specific estimates  
421 was taken as the age-independent estimate for internal consistency.

422 The results are shown in Figure 4. Across communities, internal consistency estimates  
423 ranged from 0.51 to 0.80 for the odd-even split, 0.62 to 0.89 for the stratified internal consistency,  
424 and 0.62 to 0.87 for the age-corrected approach (Plymouth, UK, being an outlier with 0.28).  
425 Following Cohen's suggestions (Cohen, 1988, 1992), these correlations constitute large effects ( $r$   
426  $> .50$ ), and indicate good internal consistency.<sup>2</sup> The results are comparable to the internal  
427 consistency estimates found in the original TANGO study (Prein et al., 2023), and also resemble  
428 reliability estimates of classical false belief tasks (Hughes et al., 2000).

429 In an exploratory analysis, we found that communities with larger individual variation  
430 showed higher internal consistency estimates (Pearson's  $r = 0.46$ , 95%CI [-0.03; 0.77]). This  
431 suggests that the less variation a task can capture within a community, the lower the reliability.  
432 However, please note that this correlation could be influenced by outliers and that the sample  
433 size here ( $N = 17$  communities) is too small to make substantial claims.

**434 Discussion**

435 The TANGO-CC measures imprecision in gaze following across individuals, ages, and  
436 communities. Children's imprecision in gaze following showed highly similar result patterns  
437 across communities: children performed better in the training than the test trials, and

---

<sup>2</sup> Note that for scale reliability and Cronbach's  $\alpha$ , values of .7 to .8 have been suggested to be acceptable (Field et al., 2012; Kline, 1999). However, Kline (1999) suggested that values below .7 could be realistic for psychological constructs due to their variable nature.

438 within-community variation greatly exceeded between-community variation. Furthermore, the  
439 task showed satisfactory to high reliability across all communities. Therefore, the TANGO-CC  
440 is a suitable task to capture individual differences in social-cognitive development in diverse  
441 communities.

442 The TANGO-CC's design process lays out a much-needed pragmatic approach to  
443 studying community-level and individual-level differences across cultures: While we performed  
444 a detailed psychometric evaluation of the task in a German setting, we collaborated with local  
445 researchers for the cross-cultural stimulus development and selection. Importantly, we  
446 re-assessed the TANGO-CC's psychometric properties in a large and diverse data set. While we  
447 cannot generalize our findings to all communities worldwide, we found that the TANGO-CC  
448 captured reliable individual variation in all 17 communities studied by Bohn et al. (2024). We  
449 hope that not just the TANGO-CC but also our pragmatic approach to constructing it will be  
450 helpful to other researchers. We recommend that researchers consider generalizability concerns  
451 and cross-cultural applications of their tasks and collaborate with local researchers at the early  
452 stages of task development (Torréns et al., 2023). Using the TANGO-CC (or any other task) in a  
453 new community requires sensitivity to the specific context, piloting, and, most importantly, the  
454 involvement of researchers or research assistants from the specific community.

455 Bourdage et al. (2023) pointed out a major challenge with adapting social cognition tasks  
456 to diverse communities: the number of world cultures is vast, and communities are constantly  
457 changing. Therefore, a promising approach might be to provide tasks with a modular system  
458 where components can be modified (*i.e.*, building block structure). In the case of the  
459 TANGO-CC, the task cannot only be adapted to different languages, cartoon faces, and  
460 backgrounds (see Figure 2) but also updated with new stimuli. Unlike studies that present  
461 sequential, hand-painted pictures that are difficult to adapt (Mehta et al., 2011), the TANGO-CC  
462 uses SVGs (Scalable Vector Graphics) that can be easily exchanged.

463 Compared to one of the most commonly used social cognition measures – the

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

24

464 change-of-location false belief task (Baron-Cohen et al., 1985; Wimmer & Perner, 1983) – the  
465 TANGO–CC has several advantages: a continuous outcome measure that can capture individual  
466 differences from three years to adulthood, a short task duration that allows for more trials per  
467 child, stimuli that can be easily adapted, and known psychometric properties across 17  
468 communities. The task is presented as a web app that enables efficient data collection with large  
469 sample sizes, while it can also be used to collect data offline in locations without a reliable  
470 internet connection. The TANGO–CC follows a standardized procedure, which leaves no room  
471 for rater errors and greatly simplifies online training of research assistants. Furthermore,  
472 minimal language demands and an engaging, playful design increase the task’s usability and  
473 reduce non-completion rates.

474 The TANGO–CC is a screen-based task. Bohn et al. (2024) have shown that children  
475 with no prior touchscreen exposure were less precise in the TANGO–CC than children with  
476 prior experience. However, individual differences were also found in communities with 100%  
477 touchscreen exposure, showing that this factor alone could not explain children’s performance  
478 in the task (Bohn et al., 2024). Notably, even though the touchscreen experience caused absolute  
479 differences in task performance, all communities showed the same processing signature. A  
480 recent computational cognitive model described gaze following as a process of estimating pupil  
481 angles and the corresponding gaze vectors (Prein et al., 2024). Bohn et al. (2024) found clear  
482 support for this model in every community studied, suggesting that children all over the world  
483 process gaze in a similar way. Nevertheless, the mode of stimulus presentation needs to be kept  
484 in mind when administering the TANGO–CC, especially in communities with little technology  
485 exposure. Additional touchscreen training (e.g., more trials of training 1) might prove helpful in  
486 these cases.

487 Schilbach et al. (2013) pointed out that witnessing social interactions as an observer  
488 undoubtedly differs from actively participating in social interactions. First evidence suggests  
489 that the TANGO–CC indeed taps into social cognition as utilized in real life: Prein et al. (2024)

490 found that children's perspective-taking abilities in a personal social interaction were linked to  
491 performance in the TANGO, but less so to a matched, non-social vector following task.  
492 However, this study exclusively relied on a German sample and future research should  
493 investigate whether the relationship between the TANGO-CC and perspective-taking abilities  
494 holds across communities.

495 We have reported reliability estimates for each community by calculating internal  
496 consistency. Ideally, we would have additionally evaluated the task's test-retest reliability in  
497 each community and checked for relationships with theoretically related constructs to assess  
498 validity. Unfortunately, this might not always be feasible in large-scale cross-cultural studies  
499 due to organizational and financial constraints. An example of assessing the TANGO's  
500 predictive validity is a study conducted in Leipzig, Germany, which used the TANGO to predict  
501 children's receptive vocabulary 6 months later (Prein et al., 2023). Future cross-cultural studies  
502 could investigate the TANGO-CC's predictive validity and its relationship to other  
503 social-cognitive abilities (e.g., Theory of Mind, language development) in diverse communities.

504 Measurement invariance (*i.e.*, measuring the same construct across different  
505 communities) is often seen as a requirement for a "fair" cross-cultural comparison: it is  
506 important that any group differences are not the result of the task unintentionally tapping into  
507 different underlying constructs. As Kusano et al. (2024) put it: "The research challenge is to  
508 achieve a balance between ensuring methodological "fairness" at the individual level while also  
509 recognizing and capturing genuine sociocultural variability" (p. 34). We argue that the  
510 TANGO-CC measures a fundamental social-cognitive ability that is likely similar across  
511 communities. Selcuk et al. (2023) pointed out that researchers should study both within- and  
512 between-culture variability in the development of social cognition since sometimes  
513 within-culture differences exceed between-culture differences. Indeed, we found that  
514 within-group variability was greater than between-group variability. While we believe that the  
515 TANGO-CC can be used to compare mean differences across communities, we would

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

26

516 recommend using it to study individual differences within communities.

517 For years, researchers have called for more diverse sampling and culturally valid  
518 measures of cognitive development (Matsumoto & Yoo, 2006; e.g., Mehta et al., 2011; Nielsen et  
519 al., 2017). As Hajdúk et al. (2020) said, “using large samples and multisite approaches will align  
520 with efforts to improve reproducibility [replicability] and will clarify both the type and extent of  
521 cultural influences on social cognition” (p. 463). Similarly, Elson et al. (2023) have called for  
522 standardized, psychometrically evaluated measures that can be re-used by other researchers in  
523 order to “build a cumulative evidence base in psychology” (p. 2). This underlines how efforts to  
524 improve replicability can be combined with the goal of increasing the generalisability of  
525 psychological research findings (Li et al., 2024; Syed, 2021). Li et al. (2024) have argued that  
526 replicable and generalizable results rely on stimulus sets with slight variations, more diverse  
527 samples, and data collection at a greater scale, which are indeed all steps the TANGO-CC has  
528 taken. Openly sharing the TANGO-CC’s materials will allow other researchers to (hopefully)  
529 replicate the results and deepen our cumulative understanding of social-cognitive development  
530 across diverse communities.

531 **Conclusion**

532 The TANGO-CC captures individual differences in social-cognitive development across  
533 diverse communities. The task’s customizability, minimal language demands, and its efficient  
534 data collection method make it a valuable tool for cross-cultural research. The task showed  
535 satisfactory to high reliability (internal consistency) in a large data set including 17 diverse  
536 communities on five continents. We hope that the TANGO-CC – and its pragmatic construction  
537 process – will provide a roadmap for future cross-cultural studies on cognitive development.

538

## References

- 539 Adolphs, R. (1999). Social cognition and the human brain. *Trends in Cognitive Sciences*, 3(12),  
 540 469–479. [https://doi.org/10.1016/S1364-6613\(99\)01399-6](https://doi.org/10.1016/S1364-6613(99)01399-6)
- 541 Amir, D., & McAuliffe, K. (2020). Cross-cultural, developmental psychology: Integrating  
 542 approaches and key insights. *Evolution and Human Behavior*, 41(5), 430–444.  
 543 <https://doi.org/10.1016/j.evolhumbehav.2020.06.006>
- 544 Apicella, C., Norenzayan, A., & Henrich, J. (2020). Beyond WEIRD: A review of the last decade  
 545 and a look ahead to the global laboratory of the future. *Evolution and Human Behavior*, 41(5),  
 546 319–329. <https://doi.org/10.1016/j.evolhumbehav.2020.07.015>
- 547 Arnett, J. J. (2008). The neglected 95%: Why American psychology needs to become less  
 548 American. *American Psychologist*, 63(7), 602–614. <https://doi.org/10.1037/0003-066X.63.7.602>
- 549 Astor, K., & Gredebäck, G. (2022). Gaze following in infancy: Five big questions that the field  
 550 should answer. In *Advances in Child Development and Behavior* (p. S0065240722000192).  
 551 Elsevier. <https://doi.org/10.1016/bs.acdb.2022.04.003>
- 552 Baron-Cohen, S., Leslie, A. M., & Frith, U. (1985). Does the autistic child have a "theory of mind"  
 553 ? *Cognition*, 21, 37–46. [https://doi.org/10.1016/0010-0277\(85\)90022-8](https://doi.org/10.1016/0010-0277(85)90022-8)
- 554 Beaudoin, C., Leblanc, É., Gagner, C., & Beauchamp, M. H. (2020). Systematic Review and  
 555 Inventory of Theory of Mind Measures for Young Children. *Frontiers in Psychology*, 10, 2905.  
 556 <https://doi.org/10.3389/fpsyg.2019.02905>
- 557 Bohn, M., & Köymen, B. (2018). Common Ground and Development. *Child Development  
 558 Perspectives*, 12(2), 104–108. <https://doi.org/10.1111/cdep.12269>
- 559 Bohn, M., Prein, J. C., Ayikoru, A., Bednarski, F. M., Dzabatou, A., Frank, M. C., Henderson, A. M.  
 560 E., Isabella, J., Kalbitz, J., Kanngiesser, P., Keşsafoğlu, D., Koymen, B., Manrique-Hernandez,  
 561 M., Magazi, S., Mújica-Manrique, L., Ohlendorf, J., Olaoba, D., Pieters, W., Pope-Caldwell, S.,  
 562 ... Haun, D. (2024). *A universal of human social cognition: Children from 17 communities  
 563 process gaze in similar ways*. PsyArXiv. <https://doi.org/10.31234/osf.io/z3ahv>
- 564 Bourdage, R., Narme, P., Neeskens, R., Papma, J., & Franzen, S. (2023). An Evaluation of

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

28

- 565      Cross-Cultural Adaptations of Social Cognition Testing: A Systematic Review.
- 566      *Neuropsychology Review*. <https://doi.org/10.1007/s11065-023-09616-0>
- 567      Broesch, T., Lew-Levy, S., Kärtner, J., Kanngiesser, P., & Kline, M. (2023). A roadmap to doing  
568      culturally grounded developmental science. *Review of Philosophy and Psychology*, 14(2),  
569      587–609. <https://doi.org/10.1007/s13164-022-00636-y>
- 570      Bürkner, P.-C. (2017). Brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal*  
571      *of Statistical Software*, 80(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- 572      Bürkner, P.-C. (2018). Advanced Bayesian Multilevel Modeling with the R Package brms. *The R*  
573      *Journal*, 10(1), 395. <https://doi.org/10.32614/RJ-2018-017>
- 574      Callaghan, T., Moll, H., Rakoczy, H., Warneken, F., Liszkowski, U., Behne, T., & Tomasello, M.  
575      (2011). Early social cognition in three cultural contexts. *Monographs of the Society for*  
576      *Research in Child Development*, 76(2), vii–viii, 1–142.  
577      <https://doi.org/10.1111/j.1540-5834.2011.00603.x>
- 578      Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed). L. Erlbaum  
579      Associates.
- 580      Cohen, J. (1992). A power primer. *Psychological Bulletin*, 112(1), 155–159.  
581      <https://doi.org/10.1037/0033-2909.112.1.155>
- 582      Decety, J. (Ed.). (2020). *The social brain: A developmental perspective* (pp. xii, 426). The MIT  
583      Press. <https://doi.org/10.7551/mitpress/11970.001.0001>
- 584      Del Bianco, T., Falck-Ytter, T., Thorup, E., & Gredebäck, G. (2019). The Developmental Origins of  
585      Gaze-Following in Human Infants. *Infancy*, 24(3), 433–454.  
586      <https://doi.org/10.1111/infa.12276>
- 587      Elfenbein, H. A., & Ambady, N. (2002). On the universality and cultural specificity of emotion  
588      recognition: A meta-analysis. *Psychological Bulletin*, 128(2), 203–235.  
589      <https://doi.org/10.1037/0033-2909.128.2.203>
- 590      Elson, M., Hussey, I., Alsalti, T., & Arslan, R. C. (2023). Psychological measures aren't  
591      toothbrushes. *Communications Psychology*, 1(1), 1–4.

- 592        <https://doi.org/10.1038/s44271-023-00026-9>
- 593    Field, A., Miles, J., & Field, Z. (2012). *Discovering statistics using R*. Sage.
- 594    Frith, C. D., & Frith, U. (2007). Social Cognition in Humans. *Current Biology*, 17(16), R724–R732.
- 595        <https://doi.org/10.1016/j.cub.2007.05.068>
- 596    Gagarina, N., Klop, D., Kunnari, S., Tantale, K., Välimaa, T., Balčiūnienė, I., Bohnacker, U., &
- 597        Walters, J. (2012). *MAIN: Multilingual Assessment Instrument for Narratives*. Zentrum für
- 598        Allgemeine Sprachwissenschaft.
- 599    Gagarina, N., & Lindgren, J. (2020). New language versions of MAIN: Multilingual Assessment
- 600        Instrument for Narratives – Revised. *ZAS Papers in Linguistics*, 64, 274.
- 601        <https://doi.org/10.21248/zaspil.64.2020.543>
- 602    Gutchess, A., & Rajaram, S. (2023). Consideration of culture in cognition: How we can enrich
- 603        methodology and theory. *Psychonomic Bulletin & Review*, 30(3), 914–931.
- 604        <https://doi.org/10.3758/s13423-022-02227-5>
- 605    Hajdúk, M., Achim, A. M., Brunet – Gouet, E., Mehta, U. M., & Pinkham, A. E. (2020). How to
- 606        move forward in social cognition research? Put it into an international perspective.
- 607        *Schizophrenia Research*, 215, 463–464. <https://doi.org/10.1016/j.schres.2019.10.001>
- 608    He, J., & Vijver, F. van de. (2012). Bias and Equivalence in Cross-Cultural Research. *Online*
- 609        *Readings in Psychology and Culture*, 2(2). <https://doi.org/10.9707/2307-0919.1111>
- 610    Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world? *The*
- 611        *Behavioral and Brain Sciences*, 33(2-3), 61-83; discussion 83-135.
- 612        <https://doi.org/10.1017/S0140525X0999152X>
- 613    Hernik, M., & Broesch, T. (2019). Infant gaze following depends on communicative signals: An
- 614        eye-tracking study of 5- to 7-month-olds in Vanuatu. *Developmental Science*, 22(4), e12779.
- 615        <https://doi.org/10.1111/desc.12779>
- 616    Hughes, C., Adlam, A., Happé, F., Jackson, J., Taylor, A., & Caspi, A. (2000). Good Test-Retest
- 617        Reliability for Standard and Advanced False-Belief Tasks across a Wide Range of Abilities.
- 618        *Journal of Child Psychology and Psychiatry*, 41(4), 483–490.

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

30

- 619        <https://doi.org/10.1111/1469-7610.00633>
- 620    Kline, P. (1999). *The Handbook of Psychological Testing* (2nd ed.). Routledge.
- 621    Krys, K., De Almeida, I., Wasiel, A., & Vignoles, V. L. (2024). WEIRD–Confucian comparisons:  
622        Ongoing cultural biases in psychology’s evidence base and some recommendations for  
623        improving global representation. *American Psychologist*. <https://doi.org/10.1037/amp0001298>
- 624    Kusano, K., Napier, J., & Jost, J. (2024). *The Mismeasure of Culture: When Measurement Invariance  
625        Requirements Hinder Cross-Cultural Research in Psychology*. OSF.  
<https://doi.org/10.31234/osf.io/9qe2k>
- 626    Li, W., Germine, L. T., Mehr, S. A., Srinivasan, M., & Hartshorne, J. (2024). Developmental  
627        psychologists should adopt citizen science to improve generalization and reproducibility.  
628        *Infant and Child Development*, 33(1), e2348. <https://doi.org/10.1002/icd.2348>
- 629    Liebal, K., & Haun, D. B. M. (2018). Why Cross-Cultural Psychology Is Incomplete Without  
630        Comparative and Developmental Perspectives. *Journal of Cross-Cultural Psychology*, 49(5),  
631        751–763. <https://doi.org/10.1177/0022022117738085>
- 632    Lillard, A. (1998). Ethnopsychologies: Cultural variations in theories of mind. *Psychological  
633        Bulletin*, 123(1), 3–32. <https://doi.org/10.1037/0033-2909.123.1.3>
- 634    Matsumoto, D., & Yoo, S. H. (2006). Toward a New Generation of Cross-Cultural Research.  
635        *Perspectives on Psychological Science*, 1(3), 234–250.  
<https://doi.org/10.1111/j.1745-6916.2006.00014.x>
- 636    Mehta, U. M., Thirthalli, J., Naveen Kumar, C., Mahadevaiah, M., Rao, K., Subbakrishna, D. K.,  
637        Gangadhar, B. N., & Keshavan, M. S. (2011). Validation of Social Cognition Rating Tools in  
638        Indian Setting (SOCRATIS): A new test-battery to assess social cognition. *Asian Journal of  
639        Psychiatry*, 4(3), 203–209. <https://doi.org/10.1016/j.ajp.2011.05.014>
- 640    Nielsen, M., Haun, D., Kärtner, J., & Legare, C. H. (2017). The persistent sampling bias in  
641        developmental psychology: A call to action. *Journal of Experimental Child Psychology*, 162,  
642        31–38. <https://doi.org/10.1016/j.jecp.2017.04.017>
- 643    Peña, E. D. (2007). Lost in Translation: Methodological Considerations in Cross-Cultural  
644        *Psychology*.

- 646 Research. *Child Development*, 78(4), 1255–1264. <https://www.jstor.org/stable/4620701>
- 647 Prein, J. C., Kalinke, S., Haun, D. B. M., & Bohn, M. (2023). TANGO: A reliable, open-source,  
648 browser-based task to assess individual differences in gaze understanding in 3 to 5-year-old  
649 children and adults. *Behavior Research Methods*, 56(3), 2469–2485.  
650 <https://doi.org/10.3758/s13428-023-02159-5>
- 651 Prein, J. C., Maurits, L., Werwach, A., Haun, D. B. M., & Bohn, M. (2024). *Variation in gaze  
652 following across the life span: A process-level perspective*. PsyArXiv.  
653 <https://doi.org/10.31234/osf.io/dy73a>
- 654 Pronk, T., Molenaar, D., Wiers, R. W., & Murre, J. (2022). Methods to split cognitive task data for  
655 estimating split-half reliability: A comprehensive review and systematic assessment.  
656 *Psychonomic Bulletin & Review*, 29(1), 44–54. <https://doi.org/10.3758/s13423-021-01948-3>
- 657 R Core Team. (2024). *R: A language and environment for statistical computing* [Manual]. R  
658 Foundation for Statistical Computing.
- 659 Rouder, J. N., & Haaf, J. M. (2019). A psychometrics of individual differences in experimental  
660 tasks. *Psychonomic Bulletin & Review*, 26(2), 452–467.  
661 <https://doi.org/10.3758/s13423-018-1558-y>
- 662 Schilbach, L., Timmermans, B., Reddy, V., Costall, A., Bente, G., Schlicht, T., & Vogeley, K. (2013).  
663 Toward a second-person neuroscience. *The Behavioral and Brain Sciences*, 36(4), 393–414.  
664 <https://doi.org/10.1017/S0140525X12000660>
- 665 Selcuk, B., Gonultas, S., & Ekerim-Akulut, M. (2023). Development and use of theory of mind  
666 in social and cultural context. *Child Development Perspectives*, 17(1), 39–45.  
667 <https://doi.org/10.1111/cdep.12473>
- 668 Shepherd, S. (2010). Following Gaze: Gaze-Following Behavior as a Window into Social  
669 Cognition. *Frontiers in Integrative Neuroscience*, 4(5). <https://doi.org/10.3389/fnint.2010.00005>
- 670 Syed, M. (2021). *Reproducibility, Diversity, and the Crisis of Inference in Psychology*. OSF.  
671 <https://doi.org/10.31234/osf.io/89buj>
- 672 Tang, Y., Gonzalez, M. R., & Deák, G. O. (2024). The slow emergence of gaze- and

MEASURING GAZE FOLLOWING ACROSS COMMUNITIES

32

- 673 point-following: A longitudinal study of infants from 4 to 12 months. *Developmental Science*,  
674 27(3), e13457. <https://doi.org/10.1111/desc.13457>
- 675 Thalmayer, A. G., Toscanelli, C., & Arnett, J. J. (2021). The neglected 95% revisited: Is American  
676 psychology becoming less American? *American Psychologist*, 76(1), 116–129.  
677 <https://doi.org/10.1037/amp0000622>
- 678 Tomasello, M., Hare, B., Lehmann, H., & Call, J. (2007). Reliance on head versus eyes in the gaze  
679 following of great apes and human infants: The cooperative eye hypothesis. *Journal of*  
680 *Human Evolution*, 52(3), 314–320. <https://doi.org/10.1016/j.jhevol.2006.10.001>
- 681 Torréns, M. G., Wefers, H., & Kärtner, J. (2023). Der Einfluss von Kultur auf die Entwicklung in  
682 den ersten drei Lebensjahren. *Zeitschrift für Entwicklungspsychologie Und Pädagogische*  
683 *Psychologie*, 55(1), 14–18. <https://doi.org/10.1026/0049-8637/a000271>
- 684 Waschl, N., & Chen, M. (2022). Cross-Cultural Considerations for Adapting Valid  
685 Psychoeducational Assessments. In O. S. Tan, K. K. Poon, B. A. O'Brien, & A. Rifkin-Graboi  
686 (Eds.), *Early Childhood Development and Education in Singapore* (pp. 113–140). Springer.  
687 [https://doi.org/10.1007/978-981-16-7405-1\\_7](https://doi.org/10.1007/978-981-16-7405-1_7)
- 688 Wefers, H., Schuhmacher, N., Chacón, L. H., & Kärtner, J. (2023). Universality without  
689 uniformity – infants' reactions to unresponsive partners in urban Germany and rural  
690 Ecuador. *Memory & Cognition*, 51(3), 807–823. <https://doi.org/10.3758/s13421-022-01318-x>
- 691 Wimmer, H., & Perner, J. (1983). Beliefs about beliefs: Representation and constraining function  
692 of wrong beliefs in young children's understanding of deception. *Cognition*, 13(1), 103–128.  
693 [https://doi.org/10.1016/0010-0277\(83\)90004-5](https://doi.org/10.1016/0010-0277(83)90004-5)
- 694 Zeigler-Hill, V., Welling, L. L., & Shackelford, T. K. (2015). How can an understanding of  
695 evolutionary psychology contribute to social psychology? In *Evolutionary perspectives on*  
696 *social psychology* (pp. 3–12). Springer, Cham.

# Appendix B — Further Publications

The following Appendix B contains other publications that were written in the context of the dissertation but not included in the main text, with their respective abstracts.

## Action anticipation based on an agent's epistemic state in toddlers and adults

**Citation:** Schuwerk, T., Kampis, D., Baillargeon, R., Biro, S., Bohn, M., Byers-Heinlein, K., Dörrenberg, S., Fisher, C., Franchin, L., Fulcher, T., Garbisch, I., Geraci, A., Grosse Wiesmann, C., Hamlin, K., Haun, D. B. M., Hepach, R., Hunnius, S., Hyde, D. C., Karman, P., ..., Prein, J., ... Rakoczy, H. (2021). *Action anticipation based on an agent's epistemic state in toddlers and adults*. *Child Development* (In-Principle Acceptance of Registered Report Stage 1: Study Design). PsyArXiv. <https://doi.org/10.31234/osf.io/x4jbm>

**Abstract:** Do toddlers and adults engage in spontaneous Theory of Mind (ToM)? Evidence from anticipatory looking (AL) studies suggests that they do. But a growing body of failed replication studies raised questions about the paradigm's suitability. In this multi-lab collaboration, we test the robustness of spontaneous ToM measures. We examine whether 18- to 27-month-olds' and adults' anticipatory looks distinguish between two basic forms of an agent's epistemic states: knowledge and ignorance. In toddlers [ANTICIPATED n = 520 50% FEMALE] and adults [ANTICIPATED n = 408, 50% FEMALE] from diverse ethnic backgrounds, we found [SUPPORT/NO SUPPORT] for epistemic state-based action anticipation. Future research can probe whether this conclusion extends to more complex kinds of epistemic states, such as true and false beliefs.

*Please note that this abstract was written for the Registered Report and does not entail results yet. Text in square brackets indicates placeholder text to be filled in after data collection.*

## **PREVIC: An adaptive parent report measure of expressive vocabulary in children between 3 and 8 years of age**

**Citation:** Bohn, M., Prein, J. C., Engicht, J., Haun, D., Gagarina, N., & Koch, T. (2023). *PREVIC: An adaptive parent report measure of expressive vocabulary in children between 3 and 8 years of age*. PsyArXiv. <https://doi.org/10.31234/osf.io/hvncp>

**Abstract:** Parent report measures have proven to be a valuable research tool to study early language development. Caregivers are given a list of words and are asked which of them their child has already used. However, most available measures are not suited for children beyond infancy, come with substantial licensing costs or lack a clear psychometric foundation. Here we present the PREVIC (Parent Report of Expressive Vocabulary in Children), an open access, high quality vocabulary checklist for German-speaking children between three and eight years of age. The PREVIC was constructed leveraging the advantages of Item Response Theory: we designed a large initial item pool of 379 words and collected data from  $N = 1190$  caregivers of children between three and eight years of age. Based on this data, we computed a range of fit indices for each item (word) and used an automated item selection algorithm to compile a final pool that contains items that a) vary in difficulty and b) fit the Rasch (one-parameter logistic) model. The resulting task is highly reliable and shows convergent validity. The IRT-based construction allowed us to design an adaptive version of the task, which substantially reduces the duration of the task while retaining measurement precision. The task – including the adaptive version – was implemented as a website and is freely accessible online (<https://ccp-odc.eva.mpg.de/previc-demo/>). The PREVIC fills an important gap in the toolkit of researchers interested in language development and provides an ideal starting point for the development of converging measures in other languages.

## **oREV: An item response theory-based open receptive vocabulary task for 3- to 8-year-old children**

**Citation:** Bohn, M.\*, Prein, J.\*, Koch, T., Bee, R. M., Delikaya, B., Haun, D., & Gagarina, N. (2024). oREV: An item response theory-based open receptive vocabulary task for 3- to 8-year-old children. *Behavior Research Methods*, 56(3), 2595–2605. <https://doi.org/10.3758/s13428-023-02169-3>

**Abstract:** Individual differences in early language abilities are an important predictor of later life outcomes. High-quality, easy-access measures of language abilities are rare, especially in the preschool and primary school years. The present study describes the construction of a new receptive vocabulary task for children between 3 and 8 years of age. The task was implemented as a browser-based web application, allowing for both in-person and remote data collection via the internet. Based on data from  $N = 581$  German-speaking children, we estimated the psychometric properties of each item in a larger initial item pool via item response modeling. We then applied an automated item selection procedure to select an optimal subset of items based on item difficulty and discrimination. The so-constructed task has 22 items and shows excellent psychometric properties with respect to reliability, stability, and convergent and discriminant validity. The construction, implementation, and item selection process described here makes it easy to extend the task or adapt it to different languages. All materials and code are freely accessible to interested researchers. The task can be used via the following website: <https://ccp-odc.eva.mpg.de/orev-demo>.

## **Validation of an open source, remote web-based eye-tracking method (WebGazer) for research in early childhood**

**Citation:** Steffan, A., Zimmer, L., Arias-Trejo, N., Bohn, M., Dal Ben, R., Flores-Coronado, M. A., Franchin, L., Garbisch, I., Grosse Wiesmann, C., Hamlin, J. K., Havron, N., Hay, J. F., Hermansen, T. K., Jakobsen, K. V., Kalinke, S., Ko, E.-S., Kulke, L., Mayor, J., Meristo, M., ..., Prein, J., ..., Schuwerk, T. (2024). Validation of an open source, remote web-based eye-tracking method (WebGazer) for research in early childhood. *Infancy*, 29(1), 31–55. <https://doi.org/10.1111/infa.12564>

**Abstract:** Measuring eye movements remotely via the participant’s webcam promises to be an attractive methodological addition to in-person eye-tracking in the lab. However, there is a lack of systematic research comparing remote web-based eye-tracking with in-lab eye-tracking in young children. We report a multi-lab study that compared these two measures in an anticipatory looking task with toddlers using WebGazer.js and jsPsych. Results of our remotely tested sample of 18-27-month-old toddlers ( $N = 125$ ) revealed that web-based eye-tracking successfully captured goal-based action predictions, although the proportion of the goal-directed anticipatory looking was lower compared to the in-lab sample ( $N = 70$ ). As expected, attrition rate was substantially higher in the web-based (42%) than the in-lab sample (10%). Excluding trials based on visual inspection of the match of time-locked gaze coordinates and the participant’s webcam video overlayed on the stimuli was an important preprocessing step to reduce noise in the data. We discuss the use of this remote web-based method in comparison with other current methodological innovations. Our study demonstrates that remote web-based eye-tracking can be a useful tool for testing toddlers, facilitating recruitment of larger and more diverse samples; a caveat to consider is the larger drop-out rate.

# Appendix C — Social Cognition Survey

In Autumn/Winter 2020, we conducted a short online expert survey on defining social cognition, as reported in the Introduction. Below you find the full survey, including the questions and the answers of the experts.

## Documentation of the Expert Online Survey “Defining Social Cognition”

### Mail that we sent out:

Dear all,

At the Max Planck Institute for Evolutionary Anthropology, we are planning to design a test battery that reliably measures individual differences in social cognition in children between 2 and 5 years of age. As a first step in this process, we are reaching out to the community of developmental psychologists. Your input will help us to decide on which aspects of social cognition to focus on.

We designed a short (5-minutes) survey and we would very much appreciate your input. Here's the link: [https://www.soscisurvey.de/soc\\_cog/](https://www.soscisurvey.de/soc_cog/) [PLEASE NOTE: LINK IS NOT ACTIVE ANYMORE]

In addition to the survey, we are searching the literature for studies that looked at individual differences in children's social cognition. This will point us to the established tasks that "work" (i.e., produce meaningful variation).

Thank you very much for your time. If you have any feedback, we would be very happy to receive it.

Best,

Julia Prein, Manuel Bohn and Daniel Haun

**Data collection period:** November – December 2020

**Screenshots of the online survey:**

Page 1:

**MAX PLANCK INSTITUTE**  
FOR EVOLUTIONARY ANTHROPOLOGY

**Thank you for your interest in our survey on "Defining Social Cognition".**

The goal of this survey is to understand how scientists define "Social Cognition". The survey takes approximately 5 minutes.

There are no known risks associated with the content of this survey; however, as with any online related activity, the risk of a breach of confidentiality is always possible. Your participation and your answers to the questions in this survey are completely voluntary and you can withdraw at any time by closing the browser window or tab. Personal information will be kept confidential.

This survey is conducted by Julia Prein, Manuel Bohn and Daniel Haun from the Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany. If you have questions about this project, please contact Julia Prein ([julia\\_prein@eva.mpg.de](mailto:julia_prein@eva.mpg.de); see imprint).

By clicking "Yes, I agree" below, you are indicating that you are at least 18 years old, have read and understood this consent form, and agree to participate in this survey. Please print a copy of this page for your records.

No, I do not agree (do not participate in this survey).  
 Yes, I agree (participate in this survey).

**Next**

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020

0% completed

Page 2:

**MAX PLANCK INSTITUTE**  
FOR EVOLUTIONARY ANTHROPOLOGY

**Welcome to our online survey!**

Social Cognition is often used as an umbrella term to refer to a wide range of cognitive processes and abilities. In this survey, we are particularly interested in Social Cognition from a developmental and individual differences perspective. First, we will ask some questions about how you define Social Cognition as a psychological construct. Next, we will ask your opinion on which aspects of Social Cognition are likely to vary between children of the same age.

We want to use this survey to inform the design of a new set of tasks to measure individual differences in Social Cognition in children. Your expert opinion will help us to select the most important aspects of Social Cognition to focus on.

Thank you for taking the time to help us!

**Next**

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020

14% completed

Page 3:

**MAX PLANCK INSTITUTE  
FOR EVOLUTIONARY ANTHROPOLOGY**



We searched the literature for definitions of Social Cognition. Below you can find a (non-exhaustive) list.\* Which definition do you most agree with?

Social cognition...

is the process by which actors, at individual or collective levels, decode and encode their social world, using mental models, knowledge structures and cultural understandings to process information, extract meaning and determine appropriate action [1]

is the ability to attribute mental states to oneself and others [2]

is the sum of those processes that allow individuals of the same species (conspecifics) to interact with one another [3]

encompasses all the information-processing mechanisms that underlie how people capture, process, store, and apply information about others to navigate social situations [4]

concerns the various psychological processes that enable individuals to take advantage of being part of a social group [5]

is the ability to construct representations of the relations between oneself and others, and to use those representations flexibly to guide social behavior [6]

concerns learning about what matters in the social world [7]

is concerned with the study of the thought processes, both implicit and explicit, through which humans attain understanding of self, others, and their environment [8]

can be constructed as the process by which individuals develop the ability to monitor, control, and predict the behavior of others [9]

refers to the skills we use to think about others and ourselves in psychological terms [10]

other:

\* You can find the references for these definitions on the last page of the survey.

[Next](#)

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020      29% completed

Page 4:



**MAX PLANCK INSTITUTE**  
FOR EVOLUTIONARY ANTHROPOLOGY

Social Cognition is a multi-faceted construct with many dimensions. Here is a (non-exhaustive) list of the dimensions that reoccurred in our literature search. Please note that we do not want to focus on the areas in which you can use Social Cognition. Rather, we are interested in what Social Cognition is (i.e., which different psychological subprocesses make up Social Cognition).

Please think about situations in which you discuss or write about Social Cognition. Which dimensions do you mention most often (choose five)?

beliefs  
 knowledge  
 emotions  
 goals  
 desires  
 reasoning  
 perspective-taking  
 attention  
 pretense  
 intentions  
 Please enter what you think is missing  
 Please enter what you think is missing  
 Please enter what you think is missing

**Next**

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020      43% completed

**MAX PLANCK INSTITUTE  
FOR EVOLUTIONARY ANTHROPOLOGY**



We are primarily interested in individual differences in the development of Social Cognition. After an extensive literature review, we created the following list of aspects of Social Cognition.  
Which aspects would you expect children of the same age to vary?

The left side of the scale means that the presented term would not vary at all, while the right side of the scale means that the aspects would vary a lot between individuals. You are free to skip any item for which you don't have a strong opinion.

Children of the same age vary in how likely they are to...

	little variation	lots of variation
take another's perspective	○ ○ ○ ○ ○ ○	█ █ █ █ █ █
understand the subjectivity of knowledge states	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
learn from others	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
manipulate emotions	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
understand intentions	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
distinguish goal-directed from non-purposeful behavior	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
direct attention	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
recognise others as agents	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
play pretense	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
deceive others	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
simulate others' reasoning processes	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
attribute beliefs to others	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
remember one's own previous knowledge states	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
read facial expressions	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
recognise goals	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
follow gaze	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
understand diverse desires	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
share attention	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
recognise emotions	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
attribute knowledge and ignorance	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○
imitate actions	○ ○ ○ ○ ○ ○	○ ○ ○ ○ ○ ○

Did we miss something? Can you think of other aspects of Social Cognition that are likely to vary between children of the same age? Please enter them below.  
If you have a reference in mind for the aspect(s) that we missed, we would love to have it. Thank you!

[Next](#)

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020      57% completed

Page 6:

**MAX PLANCK INSTITUTE**  
FOR EVOLUTIONARY ANTHROPOLOGY



If you have any comments, thoughts or suggestions, we would love to hear them!

[Next](#)

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020      71% completed

Page 7:

**MAX PLANCK INSTITUTE**  
FOR EVOLUTIONARY ANTHROPOLOGY



Finally, we would like to ask for some personal information. Answering these questions is voluntary.

You are a ...	[Please choose]
You work in ...	[Please choose]
Do you conduct studies with children?	[Please choose]

You currently live in the following country ...

If you feel comfortable, you can tell us your name.\*

\* This will allow us to link your responses to your name. We would like to have your name to read up on your views. Your name will be kept strictly confidential and will not be published in any form.

[Next](#)

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020      86% completed

## MAX PLANCK INSTITUTE FOR EVOLUTIONARY ANTHROPOLOGY



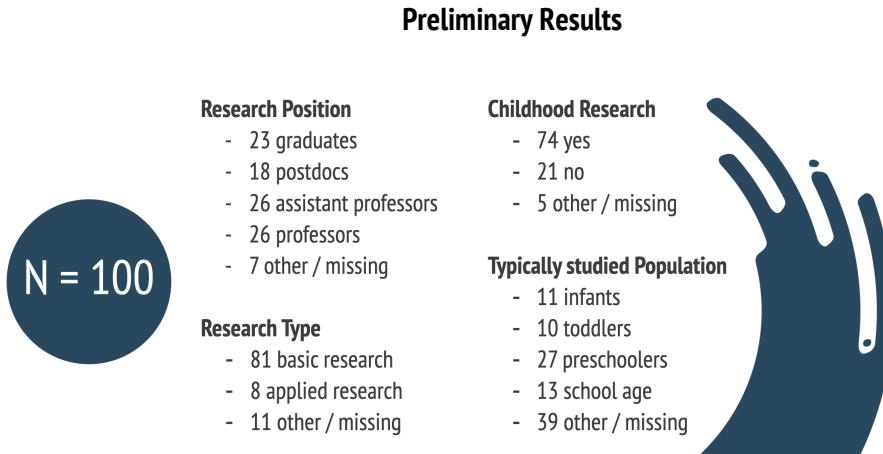
### Thank you for completing this questionnaire!

Your answers were transmitted, you may close the browser window or tab now.

In our first question of the survey, we listed some definitions of social cognition. Here is the list of references.

- [1] is the process by which actors, at individual or collective levels, decode and encode their social world, using mental models, knowledge structures and cultural understandings to process information, extract meaning and determine appropriate action. Glynn, M.A., & Watkiss, L. (2016). Social cognition. In M. Augier, D. Teece (Eds.), *The Palgrave Encyclopedia of Strategic Management*. Palgrave Macmillan, London. [https://doi.org/10.1057/978-1-349-94848-2\\_614-1](https://doi.org/10.1057/978-1-349-94848-2_614-1)
- [2] is the ability to attribute mental states to oneself and others. Dzibek, I., Fleck, S., Kalbe, E., Rogers, K., Hassenstab, J., Brand, M., . . . & Convit, A. (2006). Introducing MASC: A movie for the assessment of social cognition. *Journal of Autism and Developmental Disorders*, 36, 623–636. <https://doi.org/10.1007/s10803-006-0107-0>
- [3] is the sum of those processes that allow individuals of the same species (conspecifics) to interact with one another. Frith, C. D., & Frith, U. (2007). Social cognition in humans. *Current Biology*, 17 (16), R724–R732. <https://doi.org/10.1016/j.cub.2007.05.068>
- [4] encompasses all the information-processing mechanisms that underlie how people capture, process, store, and apply information about others to navigate social situations. Decety, J. (Ed.). (2020). *The Social Brain: A Developmental Perspective*. MIT Press.
- [5] concerns the various psychological processes that enable individuals to take advantage of being part of a social group. Frith, C. D. (2008). Social cognition. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences*, 363(1499), 2033–2039. <https://doi.org/10.1098/rstb.2008.0005>
- [6] is the ability to construct representations of the relations between oneself and others, and to use those representations flexibly to guide social behavior. Adolphs, R. (2001). The neurobiology of social cognition. *Current Opinion in Neurobiology*, (11), 231–239. [https://doi.org/10.1016/S0959-4388\(00\)00202-6](https://doi.org/10.1016/S0959-4388(00)00202-6)
- [7] concerns learning about what matters in the social world. Higgins, E. T. (2000). Social cognition: Learning about what matters in the social world. *European Journal of Social Psychology*, 30 (1), 3–39. [https://doi.org/10.1002/\(SICI\)1099-0992\(200001/02\)30:1::AID-EJSP987>3.0.CO;2-1](https://doi.org/10.1002/(SICI)1099-0992(200001/02)30:1::AID-EJSP987>3.0.CO;2-1)
- [8] is concerned with the study of the thought processes, both implicit and explicit, through which humans attain understanding of self, others, and their environment. Moskowitz, G.B. (2013). Social Cognition. *obo in Psychology*. <https://doi.org/10.1093/obo/9780199828340-0099>
- [9] can be constructed as the process by which individuals develop the ability to monitor, control, and predict the behavior of others. Rochat, P. and Striano, T. (1999). Social cognitive development in the first year. In Rochat, P. (Ed.), *Early Social Cognition: Understanding others in the first months of life*, pp. 3–34, Lawrence Erlbaum Associates Publishers.
- [10] refers to the skills we use to think about others and ourselves in psychological terms. Lewis, C., & Carpendale, J. (2014). Social cognition. In P. K. Smith & C. H. Hart (Eds.), *Wiley Blackwell handbooks of developmental psychology. The Wiley Blackwell handbook of childhood social development* (p. 531–548). Wiley Blackwell

Julia Prein, Max Planck Institute for Evolutionary Anthropology – 2020



Comparative  
Cultural Psychology

## Social Cognition...

33%

“...is the process by which actors, at individual or collective levels, decode and encode their social world, using mental models, knowledge structures and cultural understandings to process information, extract meaning and determine appropriate action.”

Glynn & Watkiss (2016)

21%

“...encompasses all the information-processing mechanisms that underlie how people capture, process, store, and apply information about others to navigate social situations.”

Decety (2020)

10%

“...is concerned with the study of the thought processes, both implicit and explicit, through which humans attain understanding of self, others, and their environment.”

Moskowitz (2013)

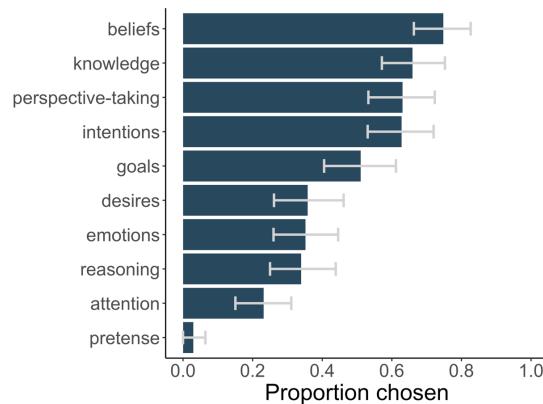
N = 100, 10 definitions (+ possibility for own definition)

## Participants' own definitions of Social Cognition

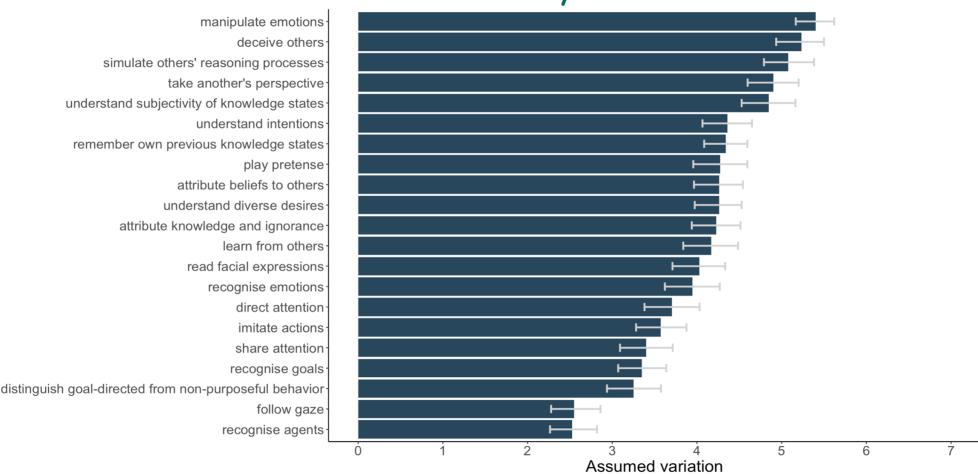
- “Cognition about social matters ranging from individuals (others and one self) and their subjective perspectives, feelings, etc. via groups to institutions”
- “Is the set of computations and representations specifically devoted to deal with conspecifics and their interactions”
- “Social cognition is the process by which we adopt or recognise a shared or non-shared perspective with other people”



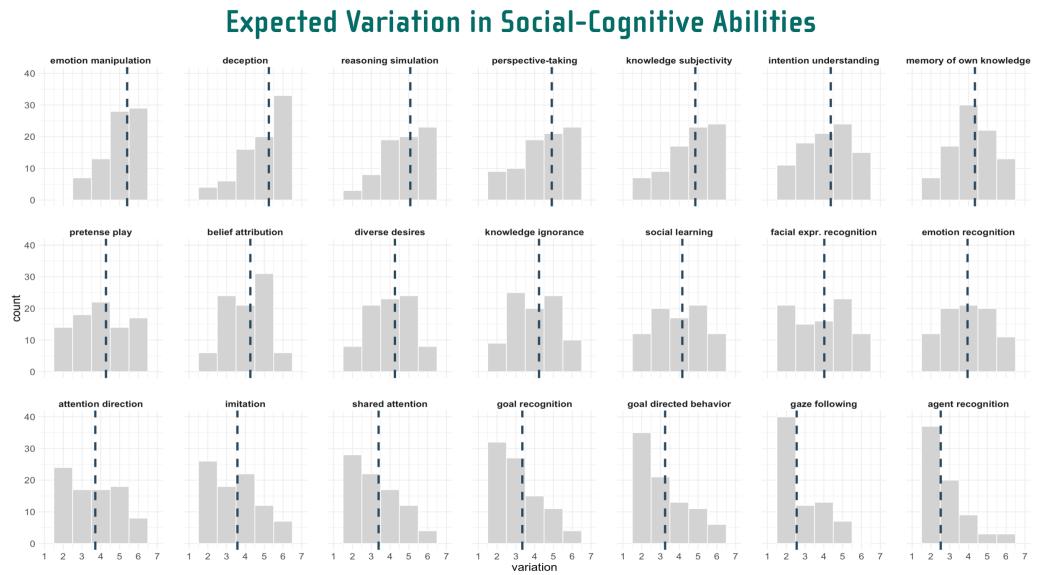
Please think about situations in which you discuss or write about Social Cognition. Which dimensions do you mention most often (choose five)?



Children of the same age vary in how likely they are to...



## Appendix A — Main Publications



# Selbstständigkeitserklärung

Julia Christin Prein  
[Straße Hausnummer]  
[PLZ Ort]  
[Telefon]  
[Email]

Hiermit erkläre ich, dass ich mich noch keiner Doktorprüfung unterzogen oder mich um Zulassung zu einer solchen beworben habe.

Ich versichere, dass die Dissertation *Rethinking Variation in Social Cognition: Gaze Following across Individuals, Ages, and Communities* in der gegenwärtigen oder einer anderen Fassung noch keiner anderen Hochschule zur Begutachtung vorgelegen hat.

Ich versichere an Eides statt, dass ich die eingereichte Dissertation *Rethinking Variation in Social Cognition: Gaze Following across Individuals, Ages, and Communities* selbstständig und ohne zulässige fremde Hilfe verfasst habe. Anderer als der von mir angegebenen Hilfsmittel und Schriften habe ich mich nicht bedient. Alle wörtlich oder sinngemäß anderen Schriften entnommenen Stellen habe ich kenntlich gemacht. Über die strafrechtlichen Folgen gemäß § 156 Strafgesetzbuch wurde ich in Kenntnis gesetzt.

Hamburg, [Datum]

[Unterschrift]