



Universiteit
Leiden
The Netherlands

Theory of mind in language, minds, and machines: a multidisciplinary approach

Dijk, B.M.A. van

Citation

Dijk, B. M. A. van. (2025, January 17). *Theory of mind in language, minds, and machines: a multidisciplinary approach*. Retrieved from <https://hdl.handle.net/1887/4176419>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/4176419>

Note: To cite this publication please use the final published version (if applicable).

Chapter 9

Conclusions

This dissertation aimed to deepen our understanding of the relation between Theory of Mind (ToM) and language by combining computational, qualitative, and experimental methods. It proposed to study ToM and language through a new language resource consisting of children’s freely told narratives, and offered empirical studies on ToM and language in both humans and computational models of language and cognition. This dissertation’s empirical work was accompanied by broader reflection on how we can understand ToM and language in the context of modern Artificial Intelligence.

In this concluding chapter, we answer the Main Research Question (MRQ) in Section 9.1 by discussing answers to all Research Questions (RQs) as presented in Section 1.2.1 and placing them in a broader context. We end by providing a reflection on this dissertation in Section 9.2, which includes a discussion of its limitations (Section 9.2.1), and an outlook on future work (Section 9.2.2).

9.1 Answers to Research Questions

RQ1 (Chapter 2)

How can we predict the mental complexity of story characters with computational tools?

To answer this question, this chapter introduced and explored the use of narratives freely told by children (4-10y) of different ages as language data. A set of 51 stories was collected and annotated for Character Depth (CD) as proxy for ToM. The chapter

9.1. Answers to Research Questions

explored the extraction of linguistic features from narratives to predict CD. We found that higher lexical complexity of narratives predicted the occurrence of mentally more developed story characters, besides age as another story-external predictor, and the effect of lexical complexity was the same for both younger and older children. We argued that a more complex vocabulary allows a child to better organise and represent mental aspects of the story world. Syntax, however, was not associated with CD.

From a broader perspective, children’s vocabularies play a key role in development: studying children’s vocabularies is studying how they represent the world (Alexander Pan, 2011). Developmental differences are more pronounced in children’s vocabulary compared to syntax, as vocabulary is more sensitive to language exposure and experience, hence varies more (Hoff, 2006; Alexander Pan, 2011). Vocabulary has also been robustly linked to other variables such as academic achievement and reading comprehension (Griffin et al., 1998), but this is less obvious for syntax, except that mastery of specific syntactic structures has been shown to drive socio-cognitive development (Lohmann and Tomasello, 2003).

RQ2 (Chapter 3)

What is the contribution of narrative language data to research in (social) cognition and (computational) linguistics?

In this chapter we presented ChiSCor (Children’s Story Corpus) as a new natural language resource consisting of 619 fantasy stories told by 442 children aged 4-12y in social contexts. In addition, three case studies showcased ChiSCor’s potential for future work and they together answer RQ2.

The first case study analysed syntactic complexity in stories. We found that overall dependency distance in stories was similar to that reported for adult language use, and stable over the primary school age range (4-12y), which is surprising given that storytelling is cognitively demanding. Our finding aligns with work that shows that children master syntax rapidly (McNeill, 1966), and with the idea that storytelling solicits ‘maximal behaviour’ in challenging children’s ToM, linguistic and cognitive competences (Frizelle et al., 2018; Southwood and Russell, 2004).

The second case study benchmarked ChiSCor’s token frequencies as approximately Zipfian, and closer to Zipf’s law than BasiScript as reference corpus of Dutch children’s written essays (Tellings et al., 2018a). We argued that this finding testifies to the needs of speaker and receiver that are balanced in live, oral storytelling. This case study fuels further work on cognitive pressures in free speech, which hith-

erto received less attention than written language, but may also impact word length, syllable duration, and syntactic structures, which can all be studied with ChiSCor’s narratives (De Palma et al., 2021; Regier et al., 2015).

The third case study showed that with ChiSCor as relatively small dataset ($\pm 74k$ tokens), lemma vectors can be trained that are as informative as vectors trained on reference corpus BasiScript which is 46 times larger (Tellings et al., 2018a). Although Word2Vec is a stepping stone towards LLMs and not a language model itself, our finding aligns with work that shows that with smaller, curated datasets, LLMs that perform well can be trained (Samuel et al., 2023). There is also evidence that narratives form particularly effective training data for LLMs (Eldan and Li, 2023), and we argue that this is because narratives are often self-contained worlds, populated with (fictional) characters and their experiences, that may act as ‘surrogate groundings’ for learning to model different types of factual and social information.

RQ3 (Chapter 4)

How can a text classification task complement existing experimental work on the relation between Theory of Mind and language in children?

We found in 442 narratives from 442 children (4-12y) in ChiSCor overlap between ToM and linguistic complexity. By drawing on theory-motivated feature engineering, manual labelling of Character Depth as proxy for ToM, and a text classifier, our overall finding was that stories with mentally more sophisticated characters were also more linguistically complex. This finding is mostly in line with Chapter 2 regarding the importance of lexical features, although the effects of syntax are much more pronounced in this chapter compared to Chapter 2, possibly due to the larger sample and variation in syntactic features included. Our finding is relevant to work in NLP that classifies different levels of ToM or perspective representation in natural language, but does so in less explainable ways (Lee et al., 2021; Kovatchev et al., 2020; Sharma et al., 2020).

Our finding supports the idea that narratives challenge children to show ‘maximal behaviour’ regarding their language use in creating mentally complex characters. We are however aware that children’s narratives provide a window on, but not necessarily a full image of, ToM and language competence. Interestingly a similar argument runs for controlled tests of ToM (e.g. Beaudoin et al., 2020; Bloom and German, 2000) and language (e.g. Dockrell and Marshall, 2015). Hence, we see our approach in this chapter as *complementing* experimental work, as ToM and language are multi-faceted

9.1. Answers to Research Questions

phenomena that should be studied in both experimental and social contexts.

RQ4 (Chapter 5)

What different types of Character Perspective Representation occur in ChiSCor's narratives and what is their relation to children's age and language use?

We annotated different types of Character Perspective Representation (CPR) in a sample of 150 stories from 150 children (4-12y) in ChiSCor, and found that children employ almost the full CPR typology provided by Leech and Short (2007). We also found a relation with age in that CPR types that provide more direct access to character minds, are more often found in older children, although we also found that children of all ages draw about equally on basic types of CPR that provide less direct access to character minds.

Contrary to our expectations, we found no clear overlap between more advanced types of CPR and lexically and syntactically more complex contexts of CPR use, which would have been in line with findings of Chapter 2 and Chapter 4. Possibly, both the relatively sparse occurrences of some CPR types and the fact that contexts were single utterances (i.e. limited contexts), barred reliable estimation of linguistic properties.

RQ5 (Chapter 6)

In what way can we meaningfully employ Language Models in studying children's language development?

Meaningfully employing LLMs in the context of developmental research involves several steps. First, taking measures to prevent data contamination as much as possible. Second, for using a LLM as a representation of mature language use as typical use case in developmental contexts, verifying that the LLM has a decent amount of knowledge of the domain at issue. Third, choosing a LLM that is modest regarding volume of training data and size, which will improve generalisability of findings (Warstadt and Bowman, 2022). Fourth, using models zero-shot, so that the linguistic knowledge of the LLM can be directly employed, and using straightforward metrics like surprisal.

Following these steps, we employed a Dutch LM to predict masked instances of the perception verb *see* in a sample of 90 narratives from 68 children (4-12y) in ChiSCor and unravelled children's semantic development. We found that children's

use of *see* is close to mature use and manifests various complex attentive and cognitive meanings. This finding supports the idea that perception verbs can be linguistic devices for children to learn to represent information about characters' attentional and cognitive states (Johnson, 1999; Sweetser, 1990), bearing on the relation between ToM and language. Our finding contrasts with work arguing that in children's language use, denotational meanings of PVs are initially dominant (Adricula and Narasimhan, 2009; Davis, 2020; Davis and Landau, 2021; Elli et al., 2021; Landau and Gleitman, 2009, *inter alia*), but aligns with work arguing that complex meanings may be present early already because of the social situatedness of language (e.g. Enfield, 2023; San Roque and Schieffelin, 2019).

RQ6 (Chapter 7)

To what extent do Large Language Models show behaviour that is consistent with having Theory of Mind-like competence?

To answer this question, we employed ToM tests from developmental psychology that present narratives in which characters' beliefs, desires and intentions are crucial for understanding the story. Both children (n=73, 7-12y) and LLMs (n=11, base and instruction-tuned) answered narrative comprehension questions to test their ability to reason about character mind states, i.e. ToM. We found that overall, only the largest commercial LLMs (GPT-4, GPT-3.5, PaLM2-chat) performed at or above child level, but we also found generalisation issues, in that these LLMs excel in figurative language use but struggle with more advanced situation modelling, and are sometimes sensitive to reformulations of original ToM tests. LLMs fine-tuned to follow instructions outperformed their counterparts that were only pre-trained.

We explained our findings by linking fine-tuning LLMs to follow instructions to the reward for cooperative communication in human language evolution and development. Successful cooperation depends for a large part on coordinating with an interaction partner's perspective (Clark, 1996; Grice, 1975; Tomasello, 2008), which aligns with ToM and is what instruction-tuning rewards.

Overall, ToM-like ability differs in humans and LLMs in that in LLMs it seems less robust. Yet, we noted that ToM tests are in the child context foremost *behavioural* probes that correlate with various other linguistic, social and academic skills. For LLMs these further correlations are hitherto unexplored, which leaves LLM success or failure on ToM tests an open question.

The discussion regarding ToM in LLMs has parallels to animal cognition, where

9.1. Answers to Research Questions

the issue is whether from observational evidence of primate behaviour, we should make the inference that primates represent beliefs, or adopt simpler explanations that primates simply respond to behavioural regularities which is not ToM (Trott et al., 2023). This point also concerns the internal validity of ToM tests and the way they operationalise ToM through behaviour. However, validity as concept has multiple dimensions, including one that relates to practical value (Van Haastrecht et al., 2023), which recurs in our answer on the next research question.

RQ7 (Chapter 8)

What are the implications of Large Language Models' complex behaviour for studying human language understanding and cognition?

We identified and critically discussed three claims that often recur in debates about LLM behaviour in relation to language understanding and cognition.

Regarding claim 1) that LLMs are mere autocomplete devices, we argued that this overlooks much of the complexity of their internal representations learned during training. With respect to claim 2) that LLMs have formal but not functional language competence, we pointed out flaws in the current methodologies of assessing functional language competence in LLMs: among other things, the reliance on simplistic tasks, and double standards in evaluating reliance on heuristics in humans and LLMs. We concluded that at this point we cannot jump to conclusions about LLMs lacking or having functional competence. Concerning claim 3) that LLMs are not relevant in understanding human language acquisition, we argued that LLMs can be useful distributional models that show which linguistic and cognitive phenomena are *in principle* learnable from word co-occurrence statistics. This could be ToM (Kosinski, 2024; Trott et al., 2023), specific grammatical phenomena (Wilcox et al., 2023), or cognitive heuristics and biases (Suri et al., 2024).

Building forth on the position we developed in response to the claims mentioned above, we concluded by setting out a pragmatic philosophical framework regarding language understanding and intentionality in LLMs. We argued that the often recurring critique that humans tend to 'anthropomorphise' LLMs, overlooks the *pragmatic value* of ascribing mental states to LLMs. Pragmatic value means being able to efficiently predict and explain behaviour in a way that abstracts away from underlying complex mechanisms, which is essentially what we do for other humans all the time.

MRQ

How can we unravel the relation between Theory of Mind and language using computational methods and narrative?

Turning to our main research question, we have disclosed various *patterns* connecting ToM and language. Stories employing mentally more complex story characters have specific linguistic properties, some straightforward like the presence of mental state verbs, others more subtle such as the lexical complexity of a story or the presence of pragmatic markers. Statistical modelling, feature engineering and text classification as *computational* methods proved fruitful in extracting these patterns from children's narratives. We believe that the finding that narratives that manifest specific ToM levels have specific linguistic profiles, links children's ToM to their language/storytelling competences in a more direct and informative way than prediction through age or socioeconomic story-external features only, which was the focus of earlier work. It suggests that the way children use language in a context *that is relevant to them*, carries relevant information about their socio-cognitive ability. This opens up a richer set of methods and activities to study and engage children's ToM that likely aligns better with children's experiences and interests, and with ToM as a social tool to communicate and coordinate other humans' intentions and behaviour: from pretend play, to sharing experiences in group circle settings, to more focused storytelling in spoken or written form. We deem all such activities worthwhile for children's development, even though they do not yield immediate quantifiable results for a teacher. Though narratives have linguistic profiles that can be linked to ToM, this did not apply to the linguistic contexts of Character Perspective Representation (CPR), as related form of ToM.

We presented ChiSCor as new resource of children's *narratives* for research in (social) cognition and (computational) linguistics. We have shown that narrating in an interactive setting leaves social 'fingerprints' in how speaker and receiver needs are balanced in children's story language. In doing so, we attempted to lift the veil of other possible interesting properties of children's natural language use, which is underexplored in developmental research. It is true that experimental contexts allow focusing on specific (socio-)cognitive and linguistic phenomena and some degree of control. At the same time, the salience of a story that is narrated live by a peer may engage children's competences differently compared to the same story that is presented in a controlled setting. Early in language ontogeny already, children learn that language is for doing things in the world: directing attention, surprising, joking,

