

# Computational modelling of language acquisition: an introduction

One approach to studying how children acquire language is to simulate language acquisition through computational modelling. Computational models implement theories of language acquisition and simulation outcomes can then be tested against existing real-world data or in new empirical research. It is more than ten years ago that *Journal of Child Language* published a special issue on the topic, edited and introduced by Brian MacWhinney (MacWhinney, 2010). Now is thus a good time to take stock of recent developments by bringing together a collection of articles that explore recent research and insights from computational modelling of child language acquisition.

Contributions in the previous special issue tended to focus on a narrow range of the overall problem of language acquisition. A narrow focus has advantages in terms of detailing the exact learning mechanisms and pieces of input relevant for acquiring specific linguistic properties. However, the past decades of empirical research on language acquisition have demonstrated the relevance of taking a broader perspective and combining linguistic levels. Learning across different linguistic levels (i.e., phonetics, phonology, morphology, syntax, semantics, pragmatics) may help children to identify speech and language units, categorize and assign meaning to these units, and acquire linguistic structure. Children make, for example, use of prosodic information to unravel syntactic structure (see discussion in Morgan & Demuth, 1996), or use verb meaning to acquire syntax and vice versa (Pinker, 1984; Gleitman, 1990). Other research has suggested continuity between the lexical and (morpho-)syntactic levels, as children gradually combine lexical constructions to arrive at more abstract grammatical representations (Marchman & Bates, 1994; Theakston & Lieven, 2017). Such ideas and findings are incompatible with and not captured by a narrow focus on language acquisition. Broadening the computational approach by modelling acquisition across linguistic levels and perhaps also modalities (i.e., oral, written) would therefore be an important next step. In addition, taking into account *non-linguistic* resources and requirements is relevant, given (neuro-)physiological and cognitive factors that have been suggested and shown to play a significant role in shaping language development (e.g., working memory, Archibald, 2017; executive functioning, Shokrkon & Nicoladis, 2022).

The goal of this special issue is to provide an overview of the current state of computational modelling in child language acquisition, with a focus on broadening the perspective. The articles in the current issue represent a variety of perspectives on language acquisition, including different empirical phenomena and theoretical approaches. We invited authors who are actively publishing on computational models of child language and who integrate different linguistic levels and modalities, or consider non-linguistic aspects in their modelling work. Furthermore, there is some work involving more than one language and beyond a monolingual and typically developing populations. Below, we first provide a summary of each contribution, starting from phonetics-related studies, continuing onto syntax and the lexicon, and ending with the language-literacy interface. We then look more closely into what we learn from the six contributions combined, reflect on connections

between the modelling work and empirical literature, and outline joint future directions for computational and empirical approaches to child language acquisition.

### *Summary of each contribution*

De Seyssel et al. are included as the first contribution in this volume, as they reflect on how computational modelling (which they call learning simulations) compares and relates to other theoretical and statistical-modelling approaches. They specifically argue that learning simulations with realistic input and multiple linguistic levels can provide a proof of concept about the role of broad learning mechanisms in general language acquisition (i.e., not restricted to a specific linguistic level or phenomenon). This approach is illustrated with an AI-based learning simulation called STELA. This model shows that statistical learning from the raw, untranscribed audio signal, replicates infants' perceptual development as observed in phonetic and lexical experiments, suggesting that acquisition may take place simultaneously across levels and in the absence of explicit linguistic categories. In the context of this volume, this contribution also offers the complementary approach of taking a learning mechanism (statistical learning) rather than a linguistic phenomenon as the starting point of the investigation.

A different perspective on phonetic learning is provided in Meier and Guenther's overview of the neurocomputational modelling of speech motor control development in infancy with the DIVA model (speech sounds) and its GODIVA extension (speech sound combinations). DIVA implements the speech production system with multiple, connected, biophysically realistic artificial neural networks, each representing a cortical region or subcortical nucleus that is credited with a specific function in the speech production system. The GODIVA model is extended with a planning loop to model the sequential production of speech sounds. Meier and Guenther review how the DIVA model can account for the empirically observed gradual expansion of speech motor control in infancy, while the GODIVA extension suggests that the gradual automation of larger speech production programs can account for children's expanding production capacities. This line of work provides insight in the neural systems underlying speech production, in the emerging connections between auditory, somatosensory, and articulatory representations, and in the timing of their involvement in speech production development. As this volume's only contribution on neurocomputational modelling, it illustrates that computational modelling can elucidate how neural-level changes underlie and give rise to stages in language development.

The issue continues with two contributions on (morpho-)syntactic development. Pine, Freudenthal and Gobet provide a comprehensive review of their work on the modelling of children's verb marking errors with a learner called MOSAIC. Verb-marking errors have featured prominently in the literature on children's morphosyntactic development, especially in studies framed within generative theory. These studies typically explain children's omission of tense and agreement morphemes (e.g., *that go in there* instead of *that goes in there*) in terms of maturation of innate abstract grammatical structure, features or constraints. MOSAIC, an unsupervised learning algorithm that relies on co-occurrence statistics for representing syntactic rules, takes a different approach and learns to progressively produce longer utterances as a function of amount of input to which it is exposed. The empirical focus of research with MOSAIC has been on one specific phenomenon. The strength of this work, and its broad perspective, lies in the modelling of different dimensions of variation, such as the variation across sentence types within one language (declaratives versus wh-questions),

variation over time within children, cross-linguistic variation (Dutch, English, German, Spanish), and, more recently, atypical development. In this respect, MOSAIC has been quite successful, which, as pointed out by Pine and colleagues, raises important conceptual questions about the mechanisms that underlie children's learning of morphosyntactic properties.

Pearl delves deeper into syntax by describing three case studies that involve linking theories (i.e., theories that we, as adults, have to link thematic roles such as agent, theme, patient, goal or experiencer, which are specified by the verb's lexical semantics, to syntactic argument positions such as subject, direct object or indirect object, which are specified by that verb's syntactic frame), the passive, and pronoun interpretation. For each case, Pearl reviews the syntactic knowledge children need to acquire, the relevant aspects for acquisition theory that need to be implemented in the computational cognitive model, input to the model, the evaluation against behavioral data, and, importantly, what we learn from this. For example, the modelling results for linking indicate that learning syntax involves learning from syntactic contexts, as well as from non-syntactic sources such as animacy and thematic roles. Pearl dedicates a part of her contribution to outlining the relevant components of the acquisition process that a computational cognitive model should consider. In doing so, she identifies several directions for future modelling work. One angle follows from considering the implications of children's immature non-linguistic systems on generating output, extraction of information from input, and the use of information for learning. Simultaneous acquisition may be another fruitful angle: while De Seyssel et al. emphasize cross-level simultaneous acquisition, Pearl illustrates the potential of cross-structure simultaneous acquisition.

A model of vocabulary acquisition is presented by Alhama and colleagues, who set out to elucidate the helping or hindering role of distributional properties of the input, in particular the co-occurrences between words. Based on Vector Space Models to operationalize neighbourhood density (as a specific measure of word co-occurrences), the results show that words that share fewer contexts with other words are acquired earlier. This suggests that children may extract meaning from word co-occurrences and that co-occurrences may even be part of children's semantic representations. Interestingly, the results were substantially impacted by specific modelling choices, including the quantitative definition of context and the algorithm used to derive the vector representations. More generally, this work thus highlights how computational modelling can contribute to the specification and operationalization of concepts and processes in child language acquisition.

The last contribution in this special issue is from Monaghan, who focuses on the interface between language and literacy development. In the literature on literacy, a gap exists between theoretical models of word representations that connect written form, spoken form, and word meaning, on the one hand, and behavioural models that describe pathways among the different learning tasks such as decoding (sets of) letters, mapping them onto speech sounds, and comprehending oral language, on the other hand. Monaghan's contribution demonstrates how this gap is narrowed by incorporating oral language experience and reading training in computational modelling. Both quantity and quality of early oral language experiences turn out to impact on the model's reading performance, but in different ways: while quantity affects the fidelity of the representations and effective mapping of sound and meaning, quality provides opportunities for vocabulary expansion. Other modelling results demonstrate that timing of oral language experience matters, and that the impact of early

oral language skills depends on the exact literacy training system (e.g., sound-based vs meaning-based). As such, this article illustrates how learning across the oral and written modalities can be implemented in computational modelling and has the potential to increase our insight into the mechanisms that underlie the behavioural models.

### *What do we learn from modelling across linguistic levels and modalities?*

All contributions in this special issue extend beyond a single linguistic level, as MacWhinney called for in his reflection on the 2010 special issue on computational modelling. In the (GO)DIVA and STELA models, these extensions are achieved by changing the capabilities of the learner. The DIVA model of speech-sound production is turned into the GODIVA model of speech sound combinations, by providing the learner with a planning loop and ability to automate frequently repeated movement sequences. The STELA model is able to learn at both the phonetic and lexical levels, as it integrates an acoustic and a language learner, thereby bridging the gap between continuous and discrete representations. These contributions thus show that, in some cases, increasing complexity of the learner is sufficient to expand its learning scope.

Other contributions achieve their extension across linguistic levels by providing the learner with input from multiple linguistic levels, alongside suitable additional learning mechanisms. Alhama et al. model the speed of lexical acquisition on the basis of co-occurrences between lexical items (rather than from properties of individual items), which their Vector Space Models can capture. The MOSAIC model has no direct access to semantic, pragmatic or phonological information, but Pine et al. marked the input for semantic properties, which enabled modelling the semantic conditioning of morphosyntactic errors. Pearl focuses on models of syntactic acquisition that integrate information from multiple places including non-syntactic information, such as animacy of an event participant, semantic information about participant event roles, and components of lexical meaning. Finally, the computational models on reading described by Monaghan include a stepwise extension of exposure in input to spoken, meaning and written representations of words. These contributions show that language learners are more powerful, as well as possibly more correct models of language acquisition, if they are equipped to detect cross-level regularities and combine cross-modality information.

While these modelling approaches all illustrate that it is possible to model language acquisition across linguistic levels or modalities, they make connections between relatively closely-related linguistic phenomena. This keeps the progress tractable and informative, but does not address how to link linguistic levels that are less closely related. De Seyssel et al. specifically note that machine-learning or AI algorithms might be capable of learning across multiple levels, but that more detailed learning mechanisms will remain necessary for insight in the exact solution to the language learning problem. A next iteration of progress and expansion may be achieved by modelling acquisition at well-understood interfaces, such as the syntax-prosody interface (Morgan & Demuth, 1996), as another step forward to formally unravelling the complexities of the acquisition of language as a whole.

### *What do we learn from including non-linguistic aspects?*

Several contributions to this special issue also extend their modelling approach beyond speech and language by including aspects of the neural architecture, domain-general learning mechanisms, or domain-general cognition.

Meier and Guenther, De Seyssel et al., and Alhama et al., all use neural networks in their models, showing that this domain-general architecture can be used to account for several aspects of child language acquisition. The (GO)DIVA models presented in Meier and Guenther use neural networks to represent smaller and larger neural components, formally linking child language acquisition to biophysical aspects of the (developing) human body. This work also highlights how neurocomputational models of language acquisition are contingent on a solid understanding of the neural underpinnings of a given linguistic ability, such as is available for speech production. The neural-network implementations of De Seyssel et al. and Alhama et al. come from machine learning and related fields and share that they learn vector representations with prediction as a driving mechanism to achieve distributional or statistical learning. This work shows that domain-general learning mechanisms can provide the basis for the acquisition of (some aspects of) language.

General cognitive mechanisms are incorporated in the computational modelling work of Pine et al. and Pearl, on respectively morphosyntax and syntactic acquisition. However, the authors of both contributions are also cautious in making strong claims about exactly which aspects of cognition are modelled in their work, how these develop in children, and how they relate to language behavior. The MOSAIC model of Pine and colleagues includes input processing strategies and limitations through primacy and recency effects. These effects, which are underpinned by insights from the field of psychology, suggest that language-learning children are sensitive to both the beginning and end of an unfamiliar utterance and may reflect, respectively, rehearsal and processing limitations. Pearl discusses modelling of two potential consequences of cognitive immaturity, namely inaccurate representation of information and ignoring of accurate information. Results showed that modelled children matched empirical data on children's interpretation preferences best when either one of these two cognitive limitations was taken into account.

Language does not develop in isolation and these contributions illustrate the various ways in which non-linguistic aspects can be incorporated in the computational modelling of language acquisition to achieve better models, that is, models that better capture human behavior. The results presented in the articles in this special issue are promising and pave the way for future research. To yield interpretable and informative models and results, such research will require cross-disciplinary collaborations involving neurology, psychology, and linguistics.

#### *Connection between the contributions and the empirical literature*

All contributions implement hypothesized language-acquisition mechanisms and provide a proof of concept by building on a well-established empirical base of child-language acquisition data and child-directed speech. This has the advantage that models can be trained on sufficient data and evaluated against well-established human behaviour. This allows for the focused investigation of language-acquisition mechanisms that computational models uniquely afford.

A natural consequence of this strong empirical base is that several biases from the empirical literature are propagated into the computational work. Firstly, each paper models one or more well-studied phenomena, for example perceptual attunement (De Seyssel et al.) and the optional-infinite stage (Pine et al.). Secondly, all six contributions focus on monolingual language development, and four only consider an ‘average’ typically developing child (De Seyssel et al.; Pearl; Alahama et al.; Monaghan). The third bias evident in this special issue is the English preponderance in the language-development literature (Kidd & Garcia, 2022; Cristia et al., 2023), with three contributions modelling only the acquisition of English, two including other (Indo-European) languages in addition to English (DeSeyssel et al.; Pine et al.), and one contribution discussing general learning mechanisms that would presumably apply to all languages (Meier & Guenther, although examples are in English).

Future computational studies could, thus, substantially expand this current scope, with several articles in this volume already providing illustrations, directions, and potential challenges. Regarding the English bias, De Seyssel et al., Alhama, and Monaghan all mention cross-linguistic (and cross-orthography) research as an avenue to testing the universality of learning mechanisms. Monaghan furthermore suggests that such cross-linguistic and cross-orthography research could provide a fruitful basis for modelling transfer effects in cases of (sequential) bi- or multilingualism. However, Alahama and De Seyssel et al. both point out that the large amounts of data needed for (their) computational models are currently not available for most languages of the world.

As for the bias towards ‘average, typical’ language development, three contributions already illustrate what can be learned from modelling beyond this bias. Meier and Guenther’s brief description of the DIVA accounts of developmental speech disorders illustrates how an implemented model of an ‘average, typical’ language user can be used to hypothesize causes of disorders. Supplementing the MOSAIC learner with a frequency-based defaulting mechanism (MOSAIC+), Pine and colleagues were able to simulate the error profiles of children with Developmental Language Disorder and their frequent use of default forms, both within and across languages. They directly manipulated the default threshold and, therefore, no further insight is obtained in the mechanisms that underlie the use of default forms by language-impaired children. While Pearl does not explicitly discuss atypical development, the empirical observations she highlights could potentially serve as an inspiration for modelling the syntactic development of children with working memory difficulties or impaired cognitive inhibition.

A final way for computational modelling to move beyond the current biases in the empirical literature, is by offering novel predictions. While none of the contributions in this special issue do this, these existing models could (theoretically) be trained on input data from a not-yet-studied language or inspected for unpredicted or emergent behaviours. Such results could then provide the starting hypothesis of a new empirical cycle, further integrating computational models in the study of child language acquisition.

## *Conclusion*

Overall, the present volume of 6 papers illustrates that computational modelling of child language acquisition has moved forward substantially since the previous special issue on computational modelling, as recent advances provide a formal understanding of acquisition across linguistic levels and in connection with non-linguistic aspects of cognition.

Future work could aim to model more phenomena, integrate across linguistic levels that are further removed, clarify the neuro-physiological and domain-general underpinnings of more aspects of language acquisition, or clarify the impact of cognitive processing. Other advances can be found in modelling of children who acquire more than one language or whose language acquisition appears disordered, and in expanding current models beyond English and the Indo-European languages.

Practically speaking, we hope that this volume inspires empirical researchers to seek more synergy with computational researchers, for example by creating datasets in a manner that could be useful for computational work or using predictions from computational models as a starting point of rigorous empirical tests. Conversely, recent empirical advances will hopefully continue to contribute to computational studies. Such distributed but ultimately joint efforts will ultimately lead to a better understanding of the mechanisms underlying child language acquisition.

### References

- Archibald, L. M. (2017). Working memory and language learning: A review. *Child Language Teaching and Therapy*, 33(1), 5–17. <https://doi.org/10.1177/0265659016654206>
- Cristia, A., Foushee, R., Aravena-Bravo, P., Cychosz, M., Scaff, C., & Casillas, M. (2022). Combining observational and experimental approaches to the development of language and communication in rural samples: Opportunities and challenges. *Journal of Child Language*, 50(3), 495-517. <http://doi.org/10.1017/S0305000922000617>
- Gleitman, L. (1990). The structural sources of verb meaning. *Language Acquisition*, 1, 3–55. [https://doi.org/10.1207/s15327817la0101\\_2](https://doi.org/10.1207/s15327817la0101_2)
- Kidd, E., & Garcia, R. (2022). How diverse is child language acquisition research? *First Language*, 42(6), 703-735. <https://doi.org/10.1177/01427237211066405>
- MacWhinney, B. (2010). Computational models of child language learning: An introduction. *Journal of Child Language*, 37(3), 477-485. <https://doi.org/10.1017/S0305000910000139>
- Marchman, V. A., & Bates, E. (1994). Continuity in lexical and morphological development: A test of the critical mass hypothesis. *Journal of child language*, 21(2), 339-366. <https://doi.org/10.1017/S0305000900009302>
- Morgan, J. L., & Demuth, K. (1996). Signal to syntax: An overview. In J. L. Morgan & K. Demuth (Eds.), *Signal to syntax: Bootstrapping from speech to grammar in early acquisition* (pp.1 – 22). Mahwah, NJ: Lawrence Erlbaum Associates.
- Pinker, S. (1984). *Language Learnability & Language Development*. Harvard University Press.
- Theakston, A., & Lieven, E. (2017). Multiunit Sequences in First Language Acquisition. *Topics in cognitive science*, 9(3), 588–603. <https://doi.org/10.1111/tops.12268>

Shokrkon, A., & Nicoladis, E. (2022). The Directionality of the Relationship Between Executive Functions and Language Skills: A Literature Review. *Frontiers in psychology, 13*, 848696. <https://doi.org/10.3389/fpsyg.2022.848696>