

Emergent syntax: the unremitting value of computational modeling for understanding the origins of complex language

Willem H. Zuidema
Artificial Intelligence Laboratory
Vrije Universiteit Brussel
Pleinlaan 2, 1050 Brussels, Belgium
jelle@arti.vub.ac.be
<http://arti.vub.ac.be/~jelle>

AI-MEMO 01-04, AI-Lab, Vrije Universiteit Brussel
Please, do not quote.

Abstract. In the debate on the origins of language “verbal” theories and mathematical models have drawn most attention. In this essay, we argue that A-life models add an important new perspective in this debate. We discuss some shortcomings of existing theories, in particular the lack of appreciation of the frequency dependency of language evolution and the role of selforganization. We show that A-life models can help to evaluate the validity of language evolution scenario’s, and help to adapt and extend them by showing both *restrictions* and *opportunities* that are deemed to be overlooked in both verbal and mathematical theorizing. Unlike many accounts that put evolution and selforganization in opposition, we adopt the view of selforganization as the *substrate* of evolution. We give an interpretation of an existing mathematical model that is in accordance with this view, and consider several important extensions. Because the extensions yield more complex behaviors, they are studied in computational (A-life) models. These models show good examples of the possible roles of selforganization in the origins of grammar.

1 Introduction

The debate on the origins of language has been dominated by “verbal” theories, both in scientific publications (see e.g. [9]) and in popular, best-selling books (e.g. [19, 7]). Recently also mathematical models of the evolution of language, especially those of Martin Nowak et al., have received much attention (e.g. [16, 17, 15]). These models are sometimes seen as a validation of the earlier verbal theories. Steven Pinker, e.g., writes in the accompanying newsstory of [17] that the paper shows “*the evolvability of [one of] the most striking features of language*”, i.e. its compositionality.

Although we appreciate the major contributions in these books and papers, we still observe many shortcomings in the proposed theories. Both the verbal

and the mathematical accounts tend to overlook many crucial details. Verbal theories often underestimate the intricacies of the evolutionary dynamics and take “evolution” too much as a general problem solver. The mathematical models often make crucial simplifications that are linguistically poorly motivated.

In particular, both types of theories have shown little appreciation for the importance of the frequency dependency of language evolution and the role of selforganization there-in. “Frequency dependent selection” is the type of evolutionary process where the fitness of a certain trait depends not only on its intrinsic quality, but also on its relative frequency in the population. Language evolution is necessarily frequency dependent because “it takes two to talk”. Selforganization is here taken to be the phenomenon that complex patterns can result from many simple interactions. We adopt a pragmatic approach and speak of selforganization if a complex pattern arises that would not be directly expected from the local interaction underlying it (i.e. there is no “blueprint” that specifies all the details of the pattern).

A-life models have shed light on both the intricacies of the dynamics of language evolution and the explanatory role of selforganization. In this paper we will first discuss some of the exposed shortcomings of verbal and mathematical theories. Next we will explore how the combination of a mathematical model (from [15]) and A-life simulations can help to adapt and extend the existing language evolution scenarios. We believe that such an approach can eventually both avoid the problematic simplifications of mathematical models, and the *ad hoc-ness* that often surrounds A-life models. In some sense, this paper thus aims to contribute to an emerging and needed “synthetic” methodology for understanding the origins of linguistic structure.

2 Linguistics & the origins of syntax

The origins of language are a heavily debated topic, with many rivaling theories. For sake of clarity, we will restrict ourselves here to what is probably the most well-known theory on the origins of human language: Steven Pinker’s book “The language instinct” [19], based on an earlier paper by Pinker and Paul Bloom [20].

2.1 The language instinct

Pinker & Bloom’s work is based on the basic assumptions of generative linguistics: (i) Syntax is independent of semantics, as one can observe in syntactic correct but meaningless sentences; (ii) Underlying this syntax must be a productive, formal system, to obtain the infiniteness of human language; (iii) This system is acquired by infants so fast and with so little data, and all human languages share so many characteristics, that there must be a shared, innate component: the universal grammar.

Pinker & Bloom original contribution, is the claim that such a system cannot arise spontaneously or as a side effect of other (cognitive) developments. For that, language is too intricate and complex. Moreover, language clearly shows

signs of adaptiveness: it allows humans to communicate infinitely many messages, including messages that refer to events that happened at other times and places than the present, messages that convey conditional and causal information or complete narratives that allow one to draw lessons from someone else’s experience.

Each of these features can be beneficial in itself, and thus be an intermediate step in the evolution. Inevitably, Pinker & Bloom argue, one should conclude that genetic evolution is the explanation for human’s remarkable linguistic talent. The innate blueprint for an individual’s “internal language”, the universal grammar, has been gradually updated to get from a non-grammatical protolanguage to the grammatical complexity of present human language. They summarize their argument as follows:

Our conclusion is based on two facts that we would think would be entirely uncontroversial: language shows signs of complex design for the communication of propositional structures, and the only explanation for the origin of organs with complex design is the process of natural selection. [20]

The authors thus argue that language is adaptive and take as a fact that evolution leads to such adaptive solutions. The work therefore fits in the tradition of “adaptationism”, although they do consider in some detail the way early language could have functioned and be advantageous, and although they discuss some possible intermediate steps from an extensive “animal” communication system, to the human linguistic abilities.

2.2 Why such accounts are unsatisfactory

There are many reasons why this type of explanations — extensive and well-documented as they may be — should not stop A-life researchers to continue modeling language origins. First of all, although these theories are based in part on formal models of language competence, the scenario for language evolution is purely verbal. A-life models can provide a more formal approach to the “evolution” part of the theory. Second, the theories of language evolution leave little room for *selforganization* as a component in the explanation. We believe that selforganization — the phenomenon that “many simple interactions can lead to complex patterns” (probably A-life’s best established wisdom) — is very likely to play a role, because many different dynamical processes must have interacted in the evolution of language.

A-life can thus help to evaluate the validity of language evolution scenario’s, and help to adapt and extend it, both by showing *restrictions* (e.g. the frequency dependency of language evolution) and by showing *opportunities* (e.g. selforganization of linguistic structure) that are deemed to be overlooked in verbal theorizing.

2.3 Novel restrictions

The main reason why “verbal” accounts of the origins of language tend to be unsatisfactory is that they seem hardly restricted by empirical or theoretical bounds. *Computational modeling* offers a novel approach to these issues, because such models are at least restricted by whether or not the *combination* of assumptions implemented in the model yield the hypothesized outcome: syntactic language. Scenario’s that seem perfectly plausible in words, often do not work in simulations without some crucial modifications or additional assumptions.

In particular, “adaptationist” explanations have turned out to be not very informative because evolution cannot be assumed to implement everything that can be useful. Even when a global fitness criterion is assumed, evolutionary dynamics can follow very different trajectories, which lead to very different, more or less adequate solutions. Good explanations should therefore describe the *evolutionary dynamics* in detail, and specify the genetic encoding, the selection pressures and the sequence of mutations that brings the system from one state to the other. Computational models — as anyone with some experience with evolutionary algorithms can confirm — immediately show that evolution is not a trivial route to the top; it is instead an intricate and open-ended process.

Evolutionary explanations of the origins of language, however, face some particular other difficulties as well. Language has two aspects that are particularly important: (i) language is transmitted, at least in part, culturally, and learned by one individual from the other; (ii) language is a group phenomenon, that occurs only between individuals and has no apparent value for an individual in isolation. These aspects make that the fitness of individual is not a function of its language acquisition system alone, but is dependent on the cultural dynamics and the composition of the group it is in as well. In other words, the “language instinct” did not evolve in isolation with respect to some *objective quality measure*, but instead co-evolved with changing language [7] and co-evolved with the changing language systems of other individuals (“frequency dependent selection”, as e.g. in [4]). The evolutionary dynamics in such a system might lead to very different results than one would expect in an “adaptationist” scenario. Figure 1 shows an example of the unexpected effects of frequency dependent selection in a model of the evolution of syntax in groups of agents [25].

In the mathematical models that we considered this fact is only partly acknowledged. In [16, 17, 15] the fitness of individuals are frequently dependent; however, for deriving the evolutionary dynamics only the average fitness in the population is considered. This averaging completely obscures the real problem of language evolution. The “difficulties in imagining how language could have arisen by darwinian evolution”, that the Nowak et al. claim to solve, have mainly to do with the problem of imagining how a syntactic individual can be successful in a non-syntactic population.

2.4 Novel opportunities

Pinker & Bloom’s work is symptomatic for the popular conviction that one can only choose between two types of explanations for the origins of human, gram-

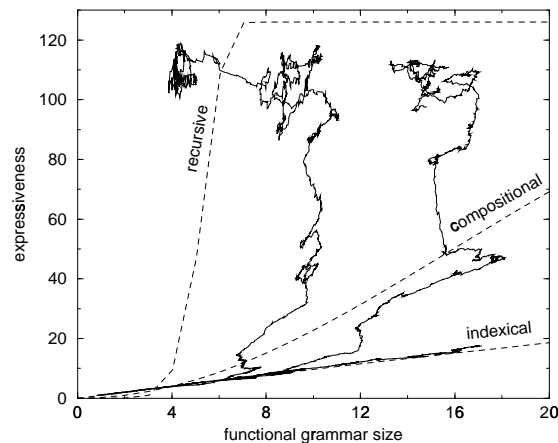


Fig. 1. Results from a model of the evolution of language in a population of agents, based on [8]. Agents’ language capabilities are represented with context-free grammars. Agents get scores for successfully communicating with each other. After some time the population is replaced with mutants of the present generation, where successful agents get more offspring than unsuccessful ones. The frequency dependent selection restricts the evolutionary dynamics in an unexpected manner. Frequently, agents with better than average language capabilities can nevertheless not persist in the population. We find that there exist distinct dynamical regimes, that lead to qualitatively different grammars (labeled “indexical”, “compositional” and “recursive”). The graph shows trajectories through a phase space, with horizontally the size of the grammar and vertically the number of distinct strings the grammars can parse. The three trajectories are obtained with the same parameters and just differ in the “random seed”; the dashed curves are theoretical predictions (see [23, 25]).

mathematical language. One refers to the powerful learning abilities of humans, and assumes that a human child, born as a *tabula rasa*, rapidly learns the regularities in its linguistic environment. Language specifics originate in the cultural transmission of language knowledge from one individual to the other.

The other type of explanation assumes a large innate component, e.g. the “universal grammar”, that specifies the principles of human languages. This innate component reduces the load on learning to merely identifying the “parameters” of a particular language’s grammar and building a lexicon of word–meaning associations. The essential, universal aspects of language originate in the process of natural selection: genetic evolution has updated our “language organ” to meet our adaptive demands.

Both types of explanation of course acknowledge that there is something innate and something learned about language: Chinese babies learn perfect English when brought up in England, and chimpanzees never learn to speak a human language fluently, no matter how carefully trained. However, they usually ignore the fact that every biological trait arises from the *interaction* between genetic coding and an environment. Discussions of linguistic “nurture & nature”, hinder

the understanding of such interactions by repeatedly suggesting that language knowledge is either learned, or *explicitly* coded for in our innate material. For instance:

Any aspect of language that the speaker knows must either be learnable from positive evidence, that is to say, through exposure to sentences of the language, or be part of the innate equipment of the human mind [5].

Underlying such reasoning is a strong intuition that the patterns observed in human language are too complicated to arise “spontaneously”. However, an impressive amount of examples — not in the last place from A-life research — show that intuitions about the complexity of underlying mechanisms are often flawed. These examples, ranging from Alan Turing work on pattern formation and morphogenesis, to Chris Langton’s “Ant”, are usually described as selforganization. Mechanisms of such spontaneous pattern formation in linguistics remain largely unexplored, although some recent studies suggest that its impact can be large and thus its misunderstanding unfortunate.

With “spontaneous” we do not mean “accidental”, though. Pinker & Bloom’s metaphor, where they compare the “appearance of design” of language to that of television sets, makes a valid point:

It would be vanishingly unlikely for something that was not designed as a television set to display television programs; the engineering demands are simply too complex. [20]

However, we believe that by putting selforganization and evolution in opposition as Pinker & Bloom do (but also proponents of selforganization [22]) they exclude the most important type of explanation. Evolution needs a substrate to operate on (i.e. the parameters that are under evolutionary control), and selforganization needs a mechanism to set the right parameters. Boerlijst & Hogeweg [3] have therefore proposed to view *selforganization* as the *substrate for evolution*. In such a view the “design” of language is neither accidental, nor hard-wired in an innate blue-print: innate are only the parameters of a selforganizing process. Note that such a view differs fundamentally from the simplistic “some parts are innate, some parts learned” explanations: cultural dynamics and evolutionary dynamics fundamentally interact.

Applying this view to language, offers a fresh perspective on some of the recurring problems in linguistics (see e.g. [13, 14]). It might very well be that children use grammatical rules in their speech without ever having encountered them. But such rules don’t need to be hard-wired in an infant’s genome, if they are a consequence of the interaction between the infant’s brain structures, its perceptual and motoric machinery, and its physical and cultural environment. Consequences of such interactions are often *trivial* under the simplifications of mathematical analysis, but *counterintuitive* in more complex computational models. A-life models have e.g. shown — contrary to conventional wisdom — that language universals concerning vowel systems or syntactic parameters emerge from articulatory and acoustic constraints [6] and processing constraints [13] respectively.

3 A-life & the origins of language

Recent work that studied the evolution of language in computational models has explored both the dynamics of language evolution and the selforganization of linguistic patterns, and has produced a wealth of new hypotheses and insights. Such models are *relatively precise* implementations of the underlying set of assumptions, and allow one to evaluate the internal coherence of such a set. Moreover, they are *productive*, in the sense that they often show unexpected behaviors that help to generate new hypotheses and concepts. And although they are necessarily simplified representations, the fact that their behavior can be experimentally evaluated makes it possible to study more complex phenomena than with analytical methods alone. Computational models therefore pre-eminently can make tractable systems with many variables and interactions.

Of course, we cannot expect these models to be informative about specifics of natural syntax, such as the position of auxiliary verbs in an English sentence. Rather, these models can give insights in the origins of some general but fundamental aspects of natural language. E.g., the facts that human language is (infinitely) expressive; for a large part specific [10, 21, 18] and distinctive [6]; compositional [1, 16, 2, 12]; recursive [8, 11, 25]; diverse on a global scale, but uniform on a local scale; dynamic, constantly subject to innovations. And that languages share universal tendencies [13, 6], and are used for very diverse purposes, including information exchange (communication), but also expression, manipulation, intimidation and social cohesion [25].

How can we best appreciate the contributions that each of these models bring to our understanding of the origins of complex language? Here we see an important role for formalization. In the following we present a mathematical model to describe language evolution, and shortly consider how some of the mentioned models fit in.

3.1 A mathematical framework

Nowak et al. use in [15] an elegant formalism¹ to describe both the cultural dynamics of language, and the evolutionary dynamics that operate on the parameters of the cultural process. If the cultural process would exhibit selforganization (which it can, with the extensions discussed below) this model thus implements the “selforganization as a substrate for evolution” approach that we outlined above. We will discuss here only the model for cultural dynamics; the analysis of the evolutionary dynamics in [15] is less convincing, as it fails — among other things — to make explicit which is the unit of selection.

Cultural dynamics Assume that there is a finite number of states (grammar types) that an individual can be in. Further, assume that newcomers (infants) learn their grammar from the population, where more successful grammars have

¹ Nowak et al. adapted this formalism in turn from Eigen & Schuster's quasi-species theory

a higher probability to be learned and mistakes are made in learning. The system can now be described in terms of the relative frequency of each grammar type in the population. The change of relative frequencies is a function of the frequencies, the success-measure of each grammar and the mistakes in learning:

$$\dot{x}_i = \sum_j^N x_j f_j Q_{ji} - \phi x_i \quad (1)$$

The components of this equation have the following interpretation:

- x_i is the fraction of individuals in the population that have a grammar of type i . i and j are indices that range from 1 to N , the number of different grammar types. \dot{x}_i describes the rate of change (the derivative) of x_i . The equation thus is an *ordinary differential equation*.
- f_i is the *relative fitness* (quality) of grammars of type i and equals $f_i = \sum_j x_j F_{ij}$, where F_{ij} is the expected communicative success from an interaction between an individual of type i and an individual of type j . The relative fitness f of a grammar thus depends on the frequencies of all grammar types, hence it is *frequency dependent*. The proper way to choose F depends on the characteristics of *language use* (production and interpretation).
- Q_{ij} is the probability that a child learning from a parent of type i , will end up with grammar of type j . The probability that the child ends up with the same grammar, Q_{ii} , is defined as q , the copying fidelity. The proper way to choose Q depends on the characteristics of *language acquisition* (learning and development).
- ϕ is the average fitness in the population and equals $\phi = \sum_i x_i f_i$. This term is needed to keep the sum of all fractions at 1.

The main result that Nowak et al. obtain is a “coherence threshold”: they show mathematically that there is a minimum value for q to keep coherence in the population. If q is lower than this value, all possible grammar types are equally frequent in the population and the communicative success is minimal. If q is higher than this value, one grammar type is dominant; the communicative success is much higher than before and reaches 100% if $q = 1$. This result is repeated in figure 2 (right panel).

3.2 Extensions

In order to make mathematical analysis possible, Nowak et al. make several crucial simplifications. Most importantly, they assume that all grammars have the same distance from each other. Consequently, the characteristics of *language use* are such that the communicative success between two agents is either maximal or minimal. The characteristics of *language acquisition* are such that an agent has either learned the right grammar or learned a random grammar. There is no information on whether an agent has moved closer to the target grammar; Nowak et al. have thus assumed the *worst case scenario* for language use and acquisition.

This assumption is clearly wrong: one only has to consider the fact that the grammatical similarity between English, German and Dutch is much stronger than between English and French or Japanese. Although a similarity metric remains difficult, the existence of some sort of a “grammar space” (and not *random* distances, such as they used later in the paper) is uncontroversial.

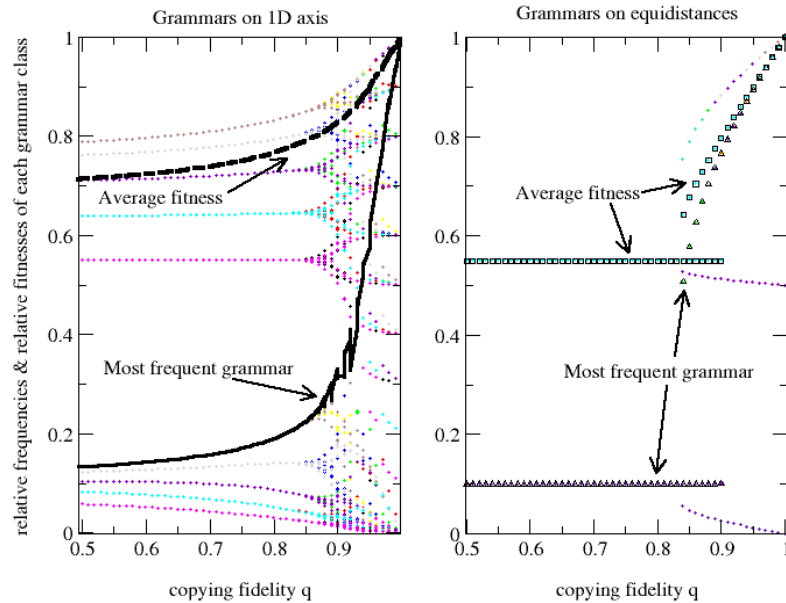


Fig. 2.

We studied and are studying several variants of the basic model, and many of the existing A-life models can be considered variants of it as well. Although these models do not always follow the same “silent” assumption of the mathematical model discussed here (e.g. that time is continuous and population size infinite), the main differences are in the characteristics of language use and language acquisition.

Language use (F) Nowak et al. assume that all grammars are equally expressive, and are all equally similar to each other. We considered several alternatives. First, we considered the other extreme case, where grammars vary only on one axis and where the similarity between grammars is determined by the distance on that axis. Under these circumstances, there is still a “coherence threshold”

but the dynamics show quite different characteristics. Some results are shown in the left panel of figure 2.

In models such as [12, 25] also the expressiveness of the grammar types differs. The language ability in these models is represented with context-free grammars, which can both implement non-syntactic idiosyncratic languages, and syntactic, compositional languages. These models shed light on the emergence of compositionality in the population.

Language acquisition (Q) Nowak et al. consider two extreme possibilities for the learning algorithm, and thus derive a lower and a upper bound on the number of training samples that a learning algorithm needs to reach the coherence threshold. However, for these results they have not taken into account that the choice of the grammar that a child has to learn is biased by how well previous generations have been able to learn and maintain it. We can now show that the lower bound that Nowak et al. derive is in fact not valid if one does take this into account [24].

The learning algorithms implemented in e.g. [2, 12] show interesting alternatives to the values for Q that Nowak et al. consider. Batali's model [2] shows a true example of selforganization: here both the linguistic abilities of the individual and the language of the population emerge from many interacting components and are not "blueprinted" in the linguistic representation or learning algorithm used.

4 Conclusions

We have reviewed critically some existing verbal and mathematical accounts of the origins of complex language. Using insights from A-life models, we have identified some novel restrictions for plausible scenario's that have to do with specifying the evolutionary dynamics, the units of selection and dealing with frequency dependent selection and the interaction with cultural dynamics. A-life models have also identified some opportunities: they have shown the selforganization of sharedness, specificity and distinctiveness in a population, and the compositionality and recursion of grammars. We interpret these models as support for the view of selforganization as a substrate for evolution. We conclude that A-life models help in evaluating and extending language evolution scenario's. Finally, we conclude that mathematical formalisms can help to put results from many A-life models into a coherent perspective.

Notes and Comments. This paper is based in part on [23], written under supervision of Paulien Hogeweg and strongly influenced by her work & ideas.

References

1. John Batali. Computational simulations of the emergence of grammar. In J. Hurford and M. Studdert-Kennedy, editors, *Approaches to the evolution of language: social and cognitive bases*. Cambridge University Press, 1998.

2. John Batali. The negotiation and acquisition of recursive grammars as a result of competition among exemplars. In T. Briscoe, editor, *Linguistic evolution through language acquisition: formal and computational models*. Cambridge University Press, 2000.
3. Maarten C. Boerlijst and Paulien Hogeweg. Self-structuring and selection. In C.G. Langton, C. Tayler, J.D. Farmer, and S. Rasmussen, editors, *Artificial Life II*, pages 255–276, 1991.
4. L.L. Cavalli-Sforza and M.W. Feldman. Paradox of the evolution of communication and of social interactivity. *Proc. Nat. Acad. Sci. USA*, 80:2017–2021, 1983.
5. Vivian Cook. *Linguistics and second language acquisition*. Macmillan, 1993.
6. Bart De Boer. *Self-Organisation in Vowel Systems*. PhD thesis, Vrije Universiteit Brussel AI-lab, 1999.
7. Terrence Deacon. *Symbolic species, the co-evolution of language and the human brain*. The Penguin Press, 1997.
8. T. Hashimoto and T. Ikegami. The emergence of a net-grammar in communicating agents. *BioSystems*, 38:1–14, 1996.
9. J. Hurford, M. Studdert-Kennedy, and C. Knight, editors. *Approaches to the evolution of language*. Cambridge University Press, 1998.
10. James Hurford. Biological evolution of the saussurean sign as a component of the language acquisition device. *Lingua*, 77(2):187–222, 1989.
11. S. Kirby. Learning, bottlenecks and the evolution of recursive syntax. In T. Briscoe, editor, *Linguistic evolution through language acquisition: formal and computational models*. Cambridge University Press, 2000.
12. S. Kirby. Syntax without natural selection: How compositionality emerges from vocabulary in a population of learners. In C. Knight, J. Hurford, and M. Studdert-Kennedy, editors, *The Evolutionary Emergence of Language: Social function and the origins of linguistic form*. Cambridge University Press, 2000.
13. Simon Kirby. *Function, selection and innateness: The emergence of language universals*. Oxford University Press, 1999.
14. Brian MacWhinney, editor. *The emergence of language*. Lawrence Erlbaum Associates, 1999.
15. Martin A. Nowak, Natalia Komarova, and Partha Niyogi. Evolution of universal grammar. *Science*, 291:114–118, 2001.
16. Martin A. Nowak and David C. Krakauer. The evolution of language. *Proc. Nat. Acad. Sci. USA*, 96:8028–8033, 1999.
17. Martin A. Nowak, Joshua B. Plotkin, and Vincent A.A. Jansen. The evolution of syntactic communication. *Nature*, 404:495–498, 2000.
18. M. Oliphant and J. Batali. Learning and the emergence of coordinated communication. *Center for research on language newsletter*, 11(1), 1996.
19. Steven Pinker. *The language instinct, how the mind creates language*. Harper Perennial, 1994.
20. Steven Pinker and Paul Bloom. Natural language and natural selection. *Behavioral and brain sciences*, 13:707–784, 1990.
21. Luc Steels. Self-organizing vocabularies. In Chris Langton, editor, *Proceedings of Alife V*, 1996.
22. Luc Steels. The synthetic modeling of language origins. *Evolution of Communication*, 1:1–35, 1997.
23. Willem H. Zuidema. Evolution of syntax in groups of agents. Master's thesis, Utrecht University, Theoretical Biology, january 2000.
24. Willem H. Zuidema. Iterated grammar acquisition. Technical Report AI-MEMO, AI-Lab Vrije Universiteit Brussel, 2001.

25. Willem H. Zuidema and Paulien Hogeweg. Selective advantages of syntactic language: a model study. In *Proceedings of the 22nd Annual Meeting of the Cognitive Science Society*, pages 577–582. Lawrence Erlbaum Associates, 2000.