# Computational Models of Language Evolution: Challenges and Future Perspectives

**Fernando Ferri**
(Istituto di Ricerca sulla Popolazione e le Politiche Sociali, Consiglio Nazionale delle Ricerche
Rome, Italy
fernando.ferri@irpps.cnr.it)

**Arianna D'Ulizia \***
(Istituto di Ricerca sulla Popolazione e le Politiche Sociali, Consiglio Nazionale delle Ricerche
Rome, Italy
**\*Corresponding author**
arianna.dulizia@irpps.cnr.it)

**Patrizia Grifoni**
(Istituto di Ricerca sulla Popolazione e le Politiche Sociali, Consiglio Nazionale delle Ricerche
Rome, Italy
patrizia.grifoni@irpps.cnr.it)

**Abstract:** This paper provides an analysis of the trends of the scientific production of language evolution models discussing the current developments and outlines the most promising future perspectives of this research field. A hybrid evaluation methodology has been applied in this study that integrates bibliometric and social research techniques to gain both quantitative and qualitative evidence of the research impact of language evolution models. Due to the ongoing interest in this research topic, the results of the analysis are valuable to many researchers to reveal the developments in the field and to plan future research directions.

**Keywords:** Language evolution, Grammatical evolution, Evolutionary computation, Agent-based models, Game theoretic models
**Categories:** I.6, F.4, F.4.2

## 1 Introduction

Human language continuously evolves as the study of human civilization reveals [Soules, 2002]. Analyzing, for instance, the differences between the seventeenth-century English and that of the twenty-first, it is evident how language can vary, ranging from phonological through orthographic, lexical, and syntactic change up to pragmatic changes [Juola, 2003].

The study of language evolution is a topic of increasing interest for the scientific community. It implies interdisciplinary competencies dealing with the investigation of "*how language evolves over time across multiple generations*" [Singh, 2005]. These competencies are necessary to cope with the complexity of the phenomenon of language evolution. Language, indeed, can be considered as a complex and non-linear dynamic system [Steels, 1997]. Providing a formal representation of the dynamics of processes occurring during language evolution is not a trivial task due to its

complexity and non-linearity. To this end, computational modeling has become fundamental for investigating and simulating the behavior and long-term dynamics of language [de Boer, 2006] [Christiansen and Kirby, 2003]. The main reason is that computational modeling allows theoretical models to be simulated and the results to be compared with empirical observations by allowing to validate the theory [Vogt, 2009]. We can use as a metaphor what happens in meteorology with computational climate change models. The earth's atmosphere and the oceanic masses are part of a complex dynamic system, which can be better understood with the help of computational modeling. Analogously, several linguistic phenomena, such as the lexicon emergence, the syntax acquisition, the symbol grounding (i.e. how words get their meanings), the emergence of compositionality (i.e. the property of systematically deriving the meaning of composite expressions from the meanings of their parts and the way in which these parts are combined), etc., are hard to be explained without the use of computational language evolution models.

This has led many researchers to apply computational modeling, giving rise to several relevant language evolution models. These models have been surveyed by several authors [Grifoni et al., 2016] [Jaeger et al., 2009] [Vogt, 2009]. In order to advance the field of language evolution modeling, it is useful to collect information and describe the developments in this field by carrying out a bibliometric analysis of the scientific publications from 2000 to 2015 and integrating it with a short interview with some authors of these papers. Due to the ongoing interest in this research topic, we think that such an analysis is valuable to many researchers to reveal the developments in the field and to plan future research directions.

Therefore, in this article, we provide a bibliometric study of language evolution models gathered from two relevant search engines (Web of Science and Scopus) and developed in the years 2000-2015. By using various measures, namely the temporal distributions and citation counts, this study evaluates the productivity and research impact of computational methods and linguistic representations used in the language evolution models. This evaluation is integrated with the analysis of the responses to a brief interview conducted with some authors of the models regarding the evolution of their research interests. Therefore, a hybrid methodology has been applied in this study that integrates bibliometric and social research techniques to gain both quantitative and qualitative evidence of the research impact of language evolution models in order to outline the most promising perspectives in this research field.

Compared to the previous surveys on language evolution, the research contribution of this article is shifted from an overview of language evolution models toward a bibliometric analysis of language evolution models and a discussion of possible research directions.

The remainder of the paper is organized as follows. Section 2 provides a discussion about current research challenges in language evolution and the motivation of our study. Section 3 gives some information about previous surveys on language evolution models. In Section 4, the analysis of scientific production based on bibliometric data is provided. In Section 5, a discussion about future challenges and perspectives of language evolution models is given. Section 6 concludes the paper.

## 2    Research challenges and motivation

Modeling has played a crucial role in language evolution research and a great deal of literature on models of language evolution has been provided over the years. However, despite the variety of modeling solutions proposed, still some computational and representational issues remain to be solved. Specifically, the main research challenges (RC) faced by language evolution modelers in recent years can be summarized as follows:

- RC1: the need for empirical evidence. Empirical evidence is information acquired by observation or experimentation. Defining language evolution models based on empirical evidence comes from the necessity to go beyond idealizations and approximations of the language evolution phenomenon. Without empirical evidence, indeed, the language evolution process is only numerically determined and, consequently, can lead to an unrealistic representation of the phenomenon. First language evolution models developed in the early 90's had this problem since they were based on theories stated vaguely and impossible to test empirically. As argued by Dediu and de Boer [Dediu and de Boer, 2016], "*hard evidence is scarce as it* (language evolution) *deals with events from the remote past*,, and spreads over many disciplines. To tackle this lack, several authors encouraged the development of new methods, tools and paradigms [Dediu and de Boer, 2016; Hauser et al., 2014; Vogt and de Boer, 2010] that investigate language evolution empirically and base theories on actual data. In light of that, most recent language evolution models have started to rely on various computational methods (e.g. agent-based, game theoretical, evolutionary computing) as a means to verify and test theoretical hypotheses on the language evolution process. In this regard, the article intends to investigate which computational methods are more/less used by recent language evolution models and to what extent these methods fulfill the need of empirical evidence;

- RC2: semantically-enriched linguistic representations. Embedding semantics is a challenge that arises from the need of modeling increasingly complex communication that better simulates what happens in reality. Several authors [Jackendoff, 1999; Luuk and Luuk, 2014] have investigated the evolutionary stages that give rise to the emergence of semantic structures. Jackendoff [Jackendoff, 1999] envisaged four steps from an unstructured use of symbols without grammar, to the introduction of a phonological structure, the concatenation of symbols, and the emergence of syntax with complex semantic relations. Similarly, [Luuk and Luuk, 2014] argued that language evolved following the sequence (elements, concatenation, embedding), meaning that starting from a limited set of symbols, it expands first by concatenating and then by embedding semantics. Following this steps, first language evolution models used linguistic representations based on a lexicon of few words and simple syntactic rules by focusing more on the first steps of the evolution process and overlooking the last step of semantic embedding. To overcome this oversimplification, most recent language evolution models have scaled up to semantically-enriched linguistic representations and a lexicon of thousands of words. In this regard, the article intends to investigate which linguistic representations are more/less used by recent language evolution models and to

what extent these formalisms fulfill the need of representing semantic aspects of the language;

- RC3: adaptable grammatical formalisms. Another main challenge concerns the ability to extend the grammar in order to cope with new, incoming changes [Steels, 2010]. Adaptable grammatical formalisms [Christiansen, 1990] are equipped with some formal means for modifying its own grammar rules while they are being used. A grammar, indeed, is said to be adaptable if it can be modified while it is being used [Ortega et al., 2007]. For example, if new concepts (words, signs, etc.) emerge, new grammatical rules have to be formulated and added to the grammar on the fly. The process of grammar adaptation is fundamental to efficiently model the dynamically changing behavior of language. Using static grammatical representation, indeed, means to erroneously consider language as an immutable artifact. Thereby, evolvable and extensible grammar formalisms have to be addressed. In this regard, the article intends to investigate whether the grammatical formalisms used in language evolution models are adaptable.

Therefore, this study aims to understand how language evolution modelers have faced these research challenges and with which results; on this basis, we aim to extract some conclusions about current developments and open challenges that still remain to be solved. To achieve that, a bibliometric analysis of several publications proposing language evolution models has been carried out. Despite the variety of modeling solutions proposed, no attempt has been made yet to evaluate the research impact and emerging priorities in this field by conducting a bibliometric analysis of scientific publications. Such a kind of analysis, indeed, is a precious aid to reveal the developments in the field and to plan future research directions. Therefore, this study can be valuable to many language evolution researchers that have to determine previous and current research highlights and possible topics for future researches.

## 3    Previous surveys on language evolution models

Many researchers of language evolution, mainly linguists and computer scientists, have paid considerable attention to understanding how the evolution of language can be computationally represented through a formal model. In the past years, several models of language evolution have been produced. These models have been differently classified by several authors [Vogt, 2009] [Jaeger et al., 2009] [Grifoni et al., 2016] according to various features of the evolution process. These three classifications involve a partial overlap of models mainly because the surveyed period is different for the three classifications (i.e. 1999-2008 for Vogt, 1995-2009 for Jaeger et al., and 2003-2012 for Grifoni et al.). Moreover, these classifications have different objectives (i.e. the complexity of interactions for Vogt, the linguistic representation and the modeling paradigm for Jaeger et al., the grammatical formalism and the computational method for Grifoni et al.) and, therefore, only the models facing the specific objective have been considered by the authors.

For the sake of completeness, in this section, we discuss in more detail these three main classifications.

Vogt (2009) provided a classification of language evolution models based on the complexity of interactions that they can handle. Specifically, he classifies these models into two classes: *macro-evolutionary analytical* models and *micro-evolutionary agent-based* models; the second class is further subdivided into *agent-based analytical* models and *agent-based cognitive* models. Table 1 shows the summary of the models reviewed by Vogt (2009) along with a set of references to papers analyzed in the original review.

| **COMPUTATIONAL MODELING** | | |
|---|---|---|
| Macro-evolutionary analytical models | Micro-evolutionary agent-based models | |
| | Agent-based analytical models | Agent-based cognitive models |
| Abrams and Strogatz (2003) Kandler and Steele (2008) Minett and Wang (2008) Nowak et al. (2002) Patriarca and Leppänen (2004) | Nettle (1999) Baxter et al. (2009) Minett and Wang (2008) | Kaplan (2005) Baronchelli et al. (2006) Vogt (2006) Briscoe (2000) Parisi et al. (2008) |

*Table 1: Surveyed models of language evolution by Vogt (2009)*

Jaeger et al. (2009) proposed a taxonomy of language evolution models, which consists of the following two dimensions: modeling paradigms and linguistic representations. According to the modeling paradigms, language evolution models are classified as *macroscopic* and *agent-based*. The latter is further classified in *iterated learning* models, *language games*, and *genetic evolution* models. According to the linguistic representations, language evolution models are classified into *symbolic grammars*, *simple recurrent networks* or *emergent grammars*. Table 2 summarizes the models reviewed by Jaeger et al. (2009).

| MODELING PARADIGM | | | | LINGUISTIC REPRESENTATION | | |
|---|---|---|---|---|---|---|
| Agent-based | | | Macro-scopic | Symbolic grammars | Simple recurrent network | Emergent grammars |
| Iterated learning | Naming game | Genetic evolution | Game theoretic | | | |
| Kirby et al. (2009) | Steels (1995) Wellens et al. (2008) VanTrijp (2008) | Cangelosi and Parisi (1998) | Jäger (2007) | Jäger (2004) Briscoe (2000) Zuidema (2002) | Christiansen and Chater (1999 ) | Steels and de Beule (2006) |

*Table 2: Surveyed models of language evolution in [Jaeger et al., 2009]*

Grifoni et al. (2016) proposed a taxonomy based on computational methods and grammatical representations. According to the computational methods (see the columns of Table 3), language evolution models are classified as *agent-based*, *evolutionary computation-based*, and *game-theoretic*. Agent-based models are further classified in *iterated learning* and *naming games*, while evolutionary computation-based models are further classified in *genetic algorithm* and *grammatical evolution*. When considering the grammatical representations (see the rows of Table 3), they have been classified in *Context-free grammar-based*, *attribute grammar-based*, *Christiansen grammar-based*, *fluid construction grammar-based*, and *universal grammar-based* models. Although several models that do not rely on grammars exist in the literature, Grifoni et al. focused only on language evolution models that have a grammatical representation. Their study identified 52 articles published between 2003 and 2012.

Previous surveys are devoted to analyze and classify language evolution models, but none of them provides a bibliometric analysis to evaluate the research impact and emerging priorities in this field. Compared to previous surveys, therefore, the research contribution of this article is shifted from a classification toward a bibliometric work. Specifically, in this paper, we perform a bibliometric study of language evolution models by analyzing 102 papers selected from 607 articles gathered from Web of Science (WoS) and Scopus databases and published from 2000 and 2015 in order to extract some conclusions about the research impact, emerging priorities, and possible future perspectives, also considering some elements arising from the answers of some authors of the language evolution models to a short interview. To gain both quantitative and qualitative evidence on the productivity and research impact of language evolution models, indeed, in this study we have applied a hybrid methodology that is composed of the following two steps:

| COMPUTATIONAL METHODS | | | | |
|---|---|---|---|---|
| Agent-based | | Evolutionary computation | | Game theoretic |
| Iterated learning | Naming game | Genetic algorithm | Grammatical evolution | |
| **Context-free grammar** Cultural grammar system (CGS) [Jimenez-Lopez, 2012] | GRAmmar EvoLution (GRAEL) [De Pauw, 2003] | GRAmmar EvoLution (GRAEL) [De Pauw, 2003] | Grammatical evolution by grammatical evolution (GE)² [O'Neill and Ryan, 2004] | |
| **Attribute grammar** | | | Language Evolver (LEVER) [Juergens and Pizka, 2006] Attribute Grammar Evolution (AGE) [de la Cruz et al., 2005] | |
| **Christiansen grammar** | | | Christiansen Grammar Evolution (CGE) [Ortega et al. 2007] | |
| **Fluid construction grammar** | FCGlight [Saveluc and Ciortuz, 2010] | | | |
| **Universal grammar** Iterated Learning Model (ILM) [Smith et al., 2003] | | | | Game dynamics (GD) [Mitchener, 2007] Evolution-ary game theory (EGT) [Jäger, 2007] |

*The left side of the table is labeled vertically: GRAMMATICAL REPRESENTATION*

*Table 3: Surveyed models of language evolution in [Grifoni et al., 2016]*

1) a bibliometric analysis of the scientific production of language evolution models in order to determine: (i) current tendencies and future research trends within the research topic of language evolution modelling; (ii) to what extent the models fulfill the research challenges related to the need of empirical evidence, semantically-enriched linguistic representations, and adaptable grammatical formalisms, as introduced in Section 2. This step is detailed in Section 4;

2) a qualitative analysis based on a short interview with the authors of some analyzed models, along with an analysis of the future work of the surveyed papers, in order to identify future research directions, as described in Section 5.

The applied methodology allows gathering more objective evidence compared to the previous surveys because it integrates bibliometric and social research techniques in order to provide quantitative data as evidence to support the results of qualitative analysis provided by the interviews.

In the following two sections we detail the implementation of each step of the methodology and we provide a discussion of the obtained results.

## 4    The bibliometric analysis: research impact of language evolution models

The analysis of scientific production, based on bibliometric data, is one of the most widely used methods for obtaining indicators about temporal evolution, variations, and trends in a specific field of research. Several works [Xie et al., 2008] [Chen et al., 2015] have applied this kind of analysis in the study of research trends.

Consistent with the approach applied in these works, we conducted a systematic search for scientific papers published from 2000 to 2015 using two relevant search engines (Web of Science and Scopus) by including the following keywords in the search: "language evolution" OR "evolution of language" AND "computational model*". Only peer-reviewed articles written in English were included in the analysis.

A total of 607 articles were returned from Web of Science and Scopus by using these search keywords, respectively 400 from Web of Science and 207 from Scopus. 10 articles from Scopus were excluded because they were duplicate publications or whole conference proceedings. The resulting set of 197 articles from Scopus and the 420 from Web of Science had an overlap of 47 articles. Therefore, a total of 550 papers were examined. Reading in detail the content of the papers, we further reduced this set to 332 papers by including only those that actually discuss language evolution theories or models. For each article we examined if it proposes a new computational model, if it uses an existing computational model, if it deals with human or animal language evolution, if it is a review or contributed paper, and if it discusses language evolution from a neurological point of view, if it relies on a linguistic representation. From this extracted information we computed the following bibliometric indicators that are used to evaluate the quantity and the quality of the scientific production, as well as whether and how the papers solve the research challenges introduced in Section 2:

- temporal distribution of publications: it gives a measure of the scientific production on language evolution models over the years;
- temporal distribution of publications by computational methods: it gives a measure of the most applied computational methods over the years (to answer to RC1);
- temporal distribution of publications by linguistic representations: it gives a measure of the most applied linguistic representations over the years (to answer to RC2 and RC3);
- number of citations (retrieved from WoS at the end of October 2016) of published papers: it gives a measure of the scientific impact of the models.

The process of gathering and selection of papers is shown in Figure 1. The constraints used to select analyzed papers are shown in the labels near the arrows. The bibliometric indicators applied to the analyzed papers are shown in the blue rectangles.
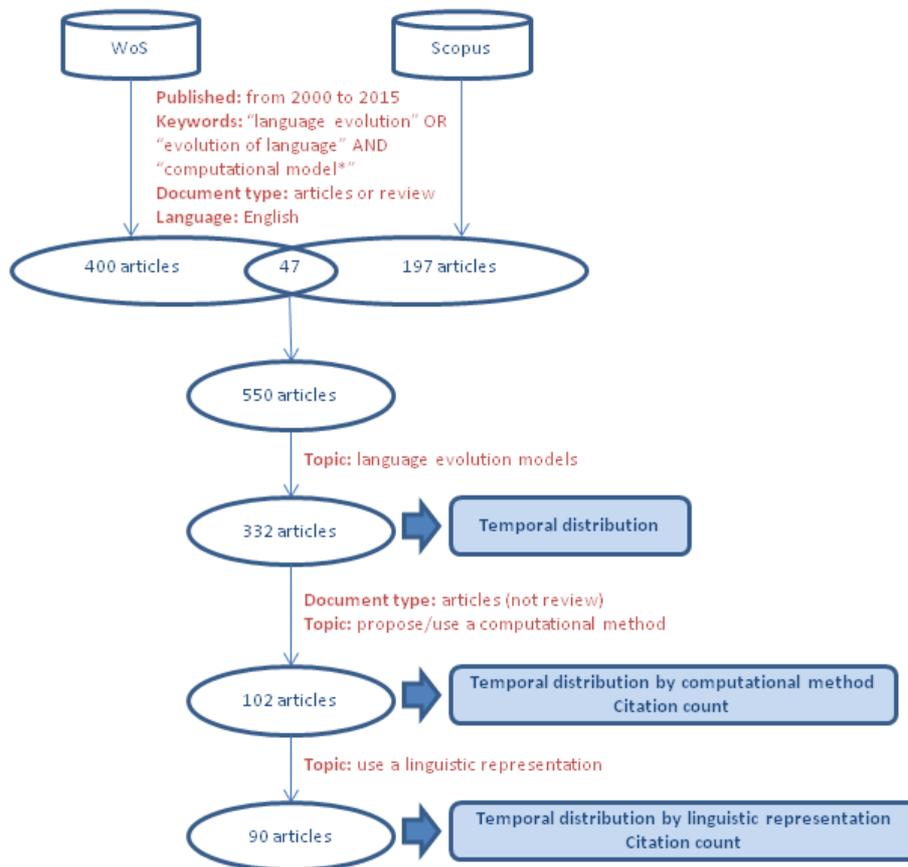


*Figure 1: The process of gathering and selection of papers for the bibliometric analysis.*

In the following sections, the results of the bibliometric analysis of language evolution models are presented. The results have been grouped according to the aforementioned bibliometric indicators.

## 4.1      Temporal distribution of publications

The temporal distribution of these 332 papers in the period 2000–2015 is shown in Figure 2. From 2000 to 2004 the trend indicates a slight increase of the scientific production that stayed between 3 (in 2001) and 12 (in 2003) papers with an average of around 7 papers per year. The number of publications increased substantially from 2005 to 2015 ranging from 22 published papers in 2005 and 33 papers in 2015 and reaching a peak of 46 papers in 2014, a minimum of 16 papers in 2008, and an average of around 27 papers per year. Not surprisingly, over 88% of the analyzed articles were published in this last period.
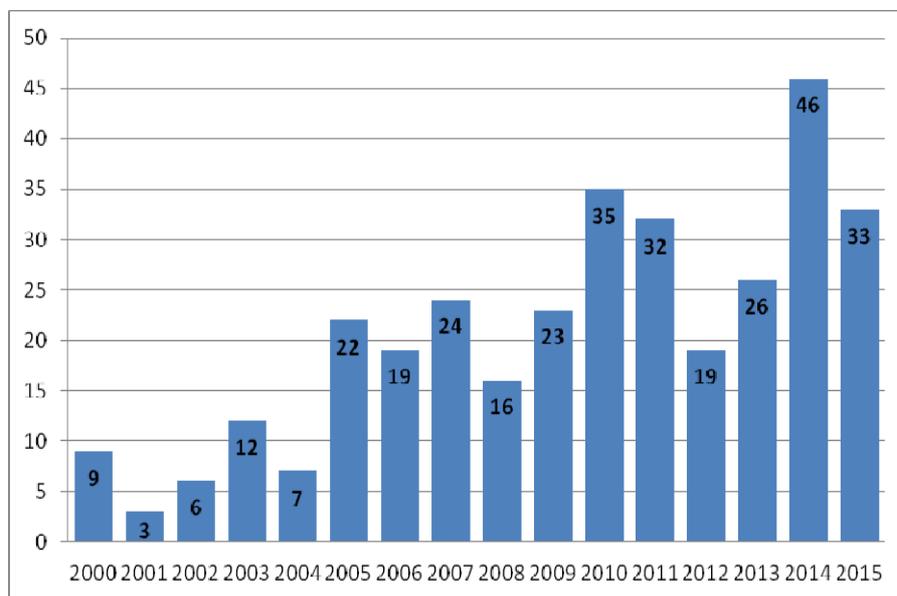


*Figure 2: Total number of papers on language evolution modeling gathered from WoS and Scopus and published from 2000 to 2015*

The choice of the timespan from 2000 to 2015 is justified by the fact that computational modeling has started to be extensively applied to language evolution mainly from 2000. Indeed, by searching for scientific papers published in journals from 1985 to 1999 by using the same two search engines (WoS and Scopus) and the same keywords, described at the beginning of Section 4, a total of 27 papers were returned, 21 of which deal with language evolution models. The temporal distribution of these 21 papers in the period 1985–1999 is shown in Figure 3.
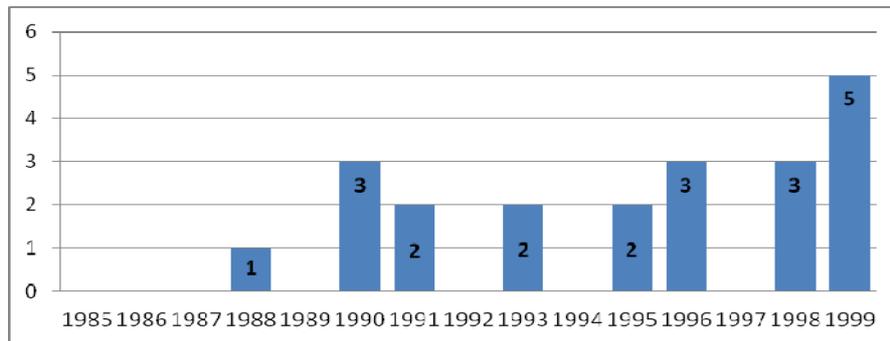
*Figure 3: Total number of papers on language evolution modeling gathered from WoS and Scopus and published from 1985 to 1999*

The chart shows that the first paper retrieved from the two search engines was published in 1988, followed by a low number of publications until 1998 (around 2-3 papers per year with four years (1989, 1992, 1994, 1997) in which no articles were published). Therefore, although the first language evolution models were developed in the early 90's, a significant increase of the publications on these models is observed only after 2000. Therefore, omitting the period before 2000 does not influence the quality of the results obtained by the performed analysis since only a minor fragment of research has not been considered. These considerations lead us to set the timespan for the subsequent analyses to January 2000 - December 2015.

The chart in Figure 2 further reveals that researches on language evolution modeling are experiencing now a fruitful period and we can expect that a greater number of scientific articles would be published in the coming years. This expectation is supported also by the high number of articles (20, to be precise) on this topic published in the first five months of 2016 that we have extracted from WoS and Scopus.

## 4.2    Temporal distribution of publications by computational methods

To evaluate the first research challenge related to the need of empirical evidence, we have investigated which computational methods are more/less used by recent language evolution models and to what extent these methods fulfill the need of empirical evidence. To achieve that, we have first analyzed whether and which computational methods are applied in the extracted publications. Specifically, we started from the set of 332 papers, extracted as described in Section 4, and we selected from them only the papers with the following characteristics: (i) original research article (not review), (ii) dealing with language evolution, and (iii) proposing/using a computational method. This selection yields 102 papers as final set for evaluating the temporal distribution of publications by computational methods.

The computational methods applied in these 102 papers are shown in Figure 4. We can observe that the majority of the papers (about 58%) were based on agents, followed by machine learning methods with 22 papers (about 22%) and game theoretic methods with 14 papers (about 14%). The remaining computational methods

have a scientific production that ranges from 1 to 5 articles (with an average of about 2 articles per method) published in the period 2000-2015. Note that several papers apply more than one computational method for modeling language evolution. A brief description of the most applied classes of computational methods is given in Appendix.A.

Considering the three most applied computational methods, we can observe that the scientific production of agent-based models is distributed across 14 years from 2002 to 2015 (see Figure 5), with a peak of 7 papers in 2010 and 2015. The scientific production of machine learning methods has been concentrated in three periods: in 2003-2007 it grows from 1 paper in 2003 and 2004 to 3 papers in 2005-2007; in 2009-2010, 2 papers per year were published; in 2013-2015 it grows from 2 papers in 2013 to 4 papers in 2015. Finally, game theoretic methods had the less continuous scientific production, with the first publications in 2000 and the last ones in 2014, reaching a peak of 3 published papers in 2007 and no papers in the periods 2001-2003, 2006, 2009-2010, 2013, and 2015.

This analysis shows that the agent-based models are the most prolific in terms of published papers and they have the most continuous bibliographic production throughout the period 2002-2015.
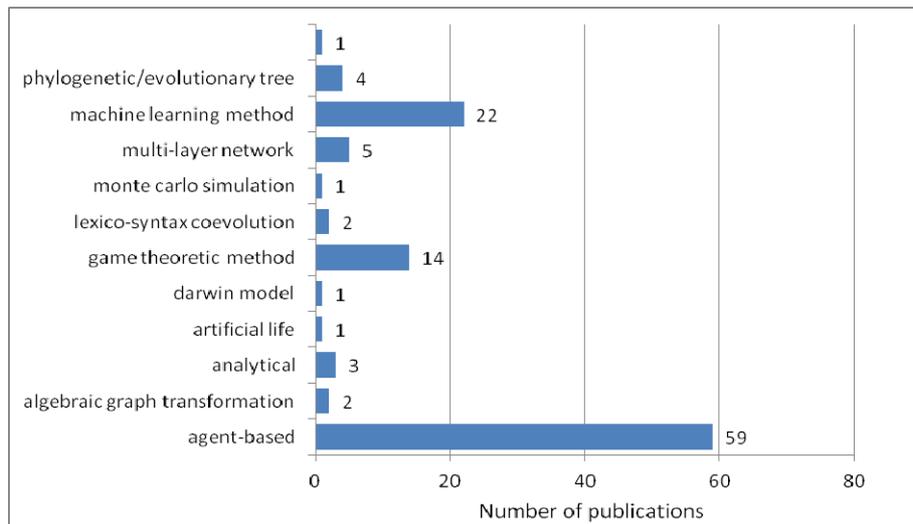


*Figure 4: Total number of papers (published from 2000 to 2015) applying a computational method for modeling language evolution*
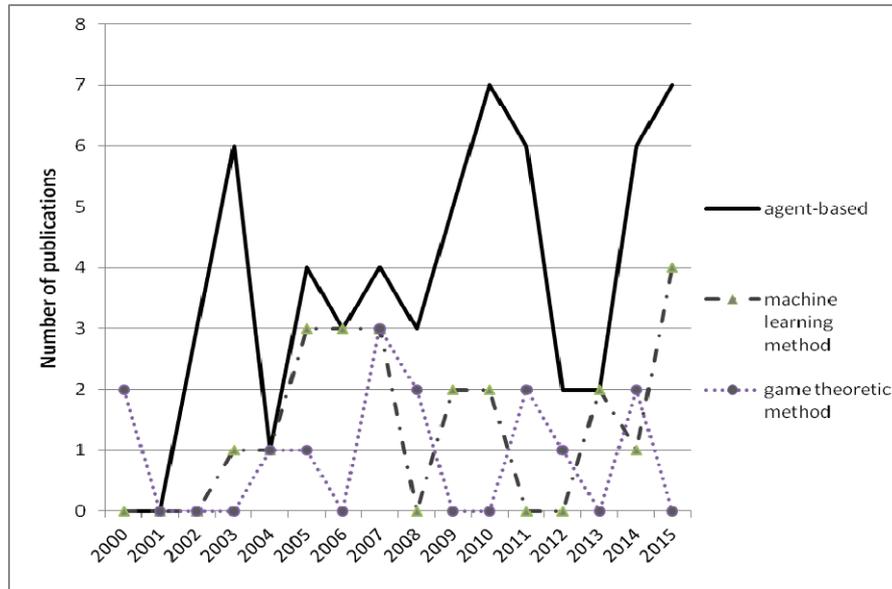
*Figure 5: Temporal distribution of the bibliographic production of the most applied computational methods. Papers applying more than one computational method are counted in all the corresponding classes of methods*

## 4.3 Temporal distribution of publications by linguistic representations

To evaluate the second research challenge RC2 related to the need of representing semantic aspects of the language we have investigated which linguistic representations are more/less used by language evolution models and to what extent these representations are semantically-enriched. By "linguistic representation" we mean a formalism to represent the linguistic knowledge within the language evolution model (e.g. grammars, signal-meaning matrices, cognitive representations, etc).

To achieve that, we started from the 102 publications extracted in the previous phase (see Section 4.2) and we analyzed whether and which linguistic representations are used in these papers. The results are shown in Figure 6.

We can observe that 12 papers (about 11.7%) do not specify the linguistic representation. The lack of this information is mainly because these models analyze the evolution of macroscopic features of language without deepening on the linguistic formalism used to represent it. For instance, they model the population dynamics of species speaking different languages, or the changes occurring in bilingual communities due to language competition.

Therefore, the total number of gathered papers using a linguistic representation is 90. The majority of these papers (about 31.4%) were based on abstract associations between meanings and forms (without a reference grammatical formalism), followed by grammatical representations with 31 papers (about 30.4%) and graph-based representations with 9 papers (about 8.8%). A lower scientific production has been

achieved by cognitive representations and linguistic traces, with 6 papers (about 5.88%) for each one. The remaining papers use bit strings (3 papers – 2,9%) and triples composed of meaning, forms and association weights (3 papers – 2.9%). A brief description of the most applied classes of linguistic representations is given in Appendix.B.
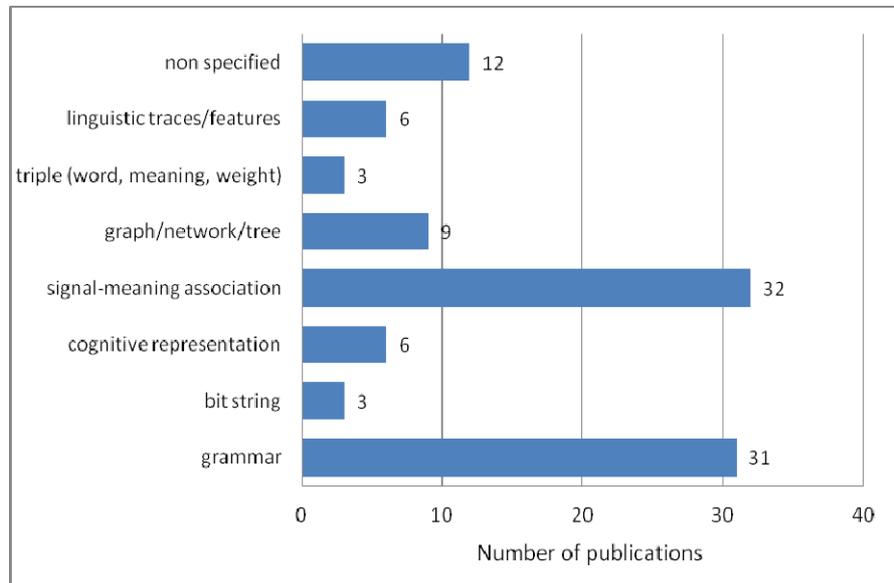


*Figure 6: Total number of papers (published from 2000 to 2015) using a linguistic representation*

Considering the three most applied linguistic representations, we can observe that the scientific production of signal-meaning association-based models is distributed across all 16 years (see Figure 7), with a peak of 4 papers in 2008, 2010, and 2011. The scientific production of grammar-based methods has been concentrated in two periods: in 2002-2007 it grows from 1 paper in 2002 to 4 papers in 2005, while in 2009-2015 it grows from 2 papers in 2009 to 4 papers in 2014. Finally, graph-based methods had the less continuous scientific production, with the first publications in 2002 and the last one in 2014, reaching a peak of 2 published papers in 2010 and 2013 and no papers in the years 2003-2004, 2006, 2008-2009, and 2011.
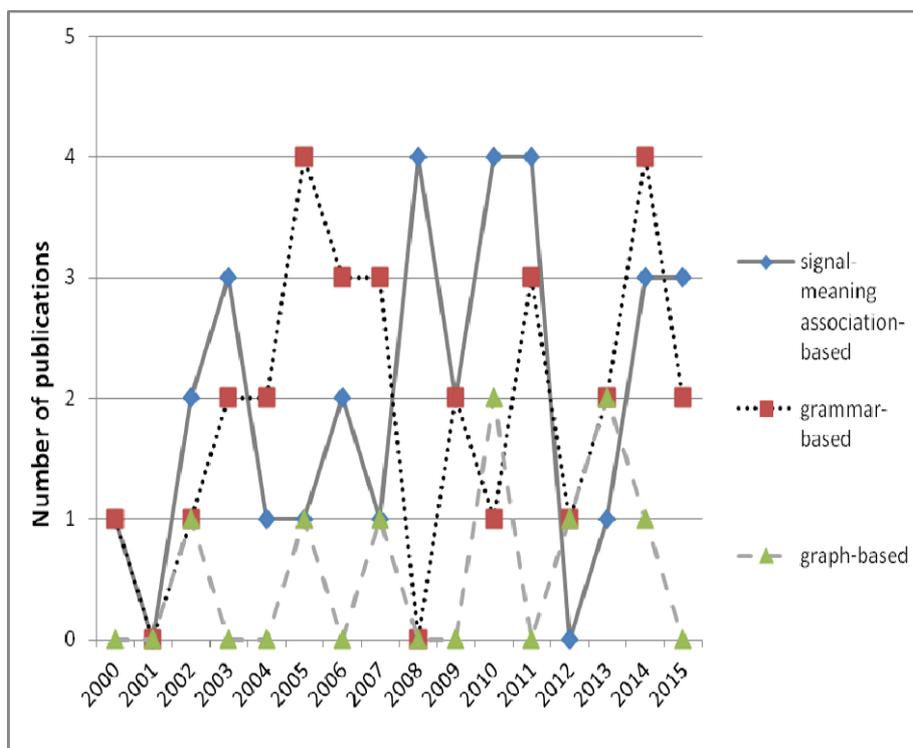
*Figure 7: Temporal distribution of the bibliographic production of the most applied linguistic representations*

## 4.4 Number of citations of published papers

To have a measure of the scientific impact of language evolution models we have analyzed the number of citations of both the set of 102 papers applying a computational method, extracted as described in Section 4.2, and the set of 90 papers using a linguistic representation, described in Section 4.3. The number of citations has been retrieved from WoS at the end of October 2016.

Considering the citation count of the 102 publications applying a computational method, they received a total of 1255 citations. Analyzing the citations of the papers applying the three most applied computational methods, Figure 8 shows that agent-based models had the highest scientific impact with 707 citations (about 56%), followed by game-theoretic models with 308 citations (about 24.5%), and finally machine learning-based models with 240 citations (about 19%). The remaining computational methods reach a total of 116 citations (9.2%). However, since the papers belonging to the game theoretic class are only 14, compared with 59 papers of the agent-based group (see Figure 4), we consider the average citations per paper. According to that, game theoretic methods are the one with the highest citations (22

citations/paper), followed by agent-based models (12 citations/paper), and machine learning-based models (10.9 citations/paper).
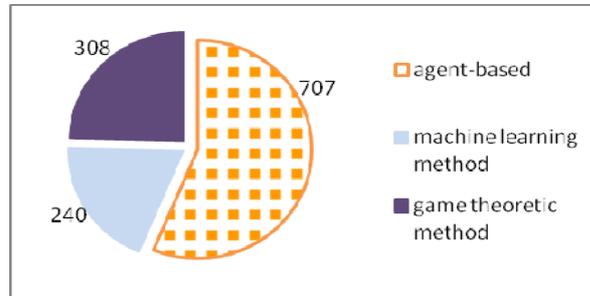


*Figure 8: Number of citations of papers published in 2000-2015 applying a computational modeling approach. Papers applying more than one computational method are counted in all the corresponding classes of methods*

Analyzing the citation count of all the 90 publications applying a linguistic representation, they received a total of 1102 citations. Analyzing the citations of the papers applying the three most applied linguistic representations, Figure 9 shows that signal-meaning association-based models had the highest scientific impact with 634 citations (about 57.5%), followed by grammar-based models with 255 citations (about 23.1%), and finally graph-based models with 103 citations (about 9.3%). The remaining linguistic representations reach a total of 110 citations (about 10%). Considering the average citations per model, the signal-meaning association-based models remain the class with the highest citations (19.8 citations/paper), followed by graph-based models (11.4 citations/paper), and grammar-based models (8.2 citations/paper).
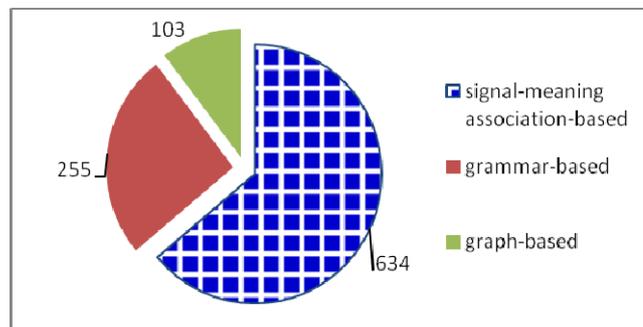


*Figure 9: Number of citations of papers published in 2000-2015 using a linguistic representation*

### 4.5 Discussion on the results of the bibliometric analysis

In this section, we provide a discussion on the results obtained from the bibliometric analysis that allows answering to each research challenge identified in Section 2.

#### 4.5.1 RC1: Which computational methods are more/less used by recent language evolution models and to what extent these methods fulfill the need of empirical evidence?

One of the aspects that the bibliometric analysis reveals is that agent-based models have the greatest research impact, both for the highest number of published papers and the high citation count (in second place after game theoretic methods). The main reason stays in the numerous advantages of agent-based methods. First of all, they allow compensating for the lack of empirical evidence present in many language evolution theories developed during the 80s and 90s. Moreover, agent-based methods are best suited to simulate communicative interactions among a population of two or more individuals and study under which conditions the human linguistic behavior may be reproduced. They allow modeling the influence of the individual actions and interactions on the evolution of language (at the individual level). Further advantages of agent-based language evolution models rely on their ability both to provide reproducible and testable dynamics, and to lead to mechanisms and algorithms robust against noise and perception. Spranger and Steels (2012) listed the following four benefits of agent-based modeling: (1) to make implicit assumptions explicit, (2) to test theories for coherence and consistency, (3) to allow for manipulation of model conditions which are difficult to manipulate with humans and (4) to generate new hypotheses. Moreover, Vogt (2009) argues that agent-based models provide a significant step towards more realism. However, they have the main shortcoming that they can become extremely complicated, resulting in behavior that is difficult to describe and interpret [de Boer, 2006]. Steels (2011) and Lipowska (2011) distinguish agent-based techniques in two groups, according to the way in which the agents interact to arrive at a shared linguistic knowledge (if using an iterated transmission or using self-organization): iterated learning and naming game. These two classes of methods turn out to be the most applied also in the 59 agent-based papers analyzed in this survey, as shown in Figure 10. However, Figure 10 also shows that further computational methods have been integrated with the agent-based approach, such as machine learning methods (mainly Bayesian methods), lexico-syntactic co-evolution, etc. Therefore, the partition of agent-based techniques in two groups does not match with the extracted data. It is correct to say that agent-based language evolution models can be combined with further techniques that range from iterated learning and language game to machine learning in order to explore how language evolves through learning and use.
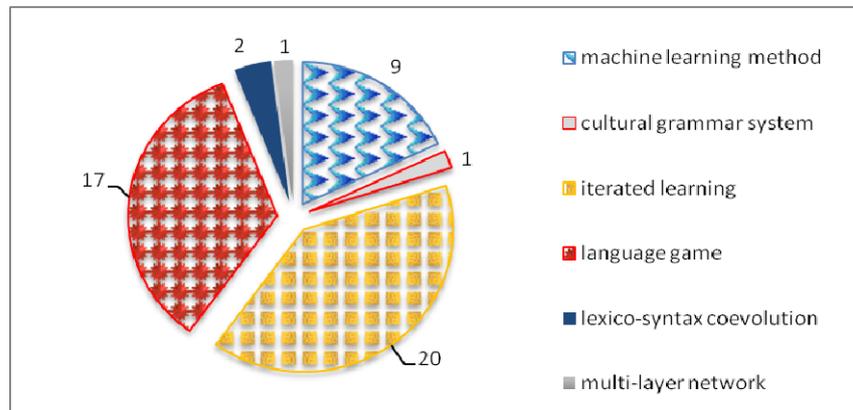
*Figure 10: Number of agent-based papers published in 2000-2015 applying a further computational method*

A great research impact has been obtained also by game-theoretic models, as they have the highest citation count. The game theoretic perspective, first developed by Smith (1982) in the early 80's for modeling the evolution of behavior, has been applied in game-theoretic models with the aim of aggregating the linguistic behavior of a population and defining general mathematical equations that model the evolution of this behavior. They have been successfully applied to describe the costs and benefits of varying strategies and the dynamics for establishing equilibria in language evolution [Watumull and Hauser, 2014]. The greatest advantage of game-theoretic methods relies on the possibility to reuse the rich body of results established by game theorists.

A lower research impact has been obtained by machine learning-based models. They have been used mainly to simplify complex agent-based models by applying algorithms and optimization techniques from evolutionary computing, Bayesian networks, and artificial neural networks.

These considerations allow answering to RC1 by arguing that current language evolution models tend to bridge the gap of empirical evidence. Indeed, they rely mainly on computational methods (firstly agent-based methods) that go beyond the idealization and approximation of some language evolution theories not very well supported empirically, by transforming theoretical models into simulation models. As noted by Spranger and Steels (2012), indeed, "*every instantiated theory that is shown to work using agent-based modeling can immediately be accepted as a coherent and stringent proposal that (at least in principle) reveals all underlying assumptions and provides reproducible and testable dynamics*".

### 4.5.2 RC2: Which linguistic representations are more/less used by recent language evolution models and to what extent these formalisms fulfill the need of representing semantic aspects of the language?

On the side of linguistic representation, the bibliometric analysis reveals that signal-meaning association-based models have the greatest research impact, both for the highest number of published papers (32) and the highest citation count (19.8 citations/paper). The linguistic representation based on signal-meaning associations is strictly related to the symbol grounding problem [Harnad, 1990] that investigates how to ground a set of symbols (e.g. forms, signals, utterances, words, etc.) in a set of possible meanings. An additional finding of the analysis is that 78% of the 32 papers based on signal-meaning associations apply agent-based modeling by investigating how a population of agents evolves the communication system through interactions with the environment and among individuals and how it arrives at shared symbolic conventions by constructing a set of relations between symbols and meanings.

A high scientific production (31 papers) has been obtained also by grammar-based methods that move from signal-meaning associations to more complex structures of language using grammatical constructions to link syntactic and semantic categories intervening between symbols and meanings.

Considering graph-based representation, we can observe that this is used mainly by phylogenetic approaches that model longer timescale of the history of populations by investigating how language originates and develops during the evolutionary history of the species. In these papers, indeed, the linguistic knowledge is mainly represented through graphs, trees or networks that model the relationships between groups of related languages.

In order to assess if linguistic representations are equipped with structures and constructions able to represent semantic features of the language, we have analyzed the presence of the word "semantics" "syntactic" "lexical" and "phonological" in the text of the 90 gathered articles to identify the level(s) of language [Hickey and Puppel, 1997] addressed by the linguistic representations. The results are summarized in Table 4. The majority of the papers (about 18%) address lexical representations as they use predominantly signal-meaning associations without syntax in order to map uniquely the set of signals (e.g. objects, forms, utterances, words, etc.) to the set of possible meanings and to reach a shared lexicon. Another 18% of the 90 analyzed papers adopt a more sophisticated linguistic representation that implements lexical knowledge, syntactic constructions, and semantic categorizations. The majority of these papers apply a grammatical approach to realize that. Finally, about 15.5% of the papers apply a lexico-syntactic representation that allows implementing both lexical knowledge and syntactic constructions for linking different meanings/signals. In total, 53 papers (about 59%) represent the lexical level of language, 38 papers (about 42%) the syntactic level, and 32 papers (about 35.5%) the semantic level.

| Level of Language | Counts |
|---|---|
| phonological | 6 |
| lexical | 16 |
| syntactic | 4 |
| semantic | 5 |
| lexical, syntactic | 14 |
| lexical, semantic | 7 |
| syntactic, semantic | 4 |
| lexical, syntactic, semantic | 16 |
| not specified | 18 |

*Table 4: Levels of language addressed by the gathered articles*

Therefore, the analysis reveals that most of the papers do not apply linguistic representations able to deal with semantic features of the language. The main motivation that leads researchers to choose this strategy instead of syntactic and semantic representations is based on the greater effort required to implement linguistic representations dealing with all the levels of language. Modeling the evolution of language considering syntactic and semantic features implies a more expressive power of the language used but also a more sophisticated linguistic framework that requires an extra effort to be implemented. Therefore, the majority of the analyzed researches investigate the evolution of lexical communication systems that are easily implementable but with a low expressivity of the language. Despite that, several researches, although less in number, address syntactic and semantic representations of the language (mainly based on grammatical formalisms) that increase the expressivity of the language and the complexity of the model.

### 4.5.3    RC3: Which grammatical formalisms are more/less used by recent language evolution models and to what extent these formalisms fulfill the need of adaptability?

To evaluate the third research challenge we have investigated which grammatical formalisms are used in language evolution models and whether they are adaptable (see RC3 defined in Section 2). Therefore, we have analyzed the 31 papers applying a grammar identified in Section 4.3, and we have evaluated whether they provide explicit constructions and operators able to expand the grammar in order to cope with new, incoming needs [Steels, 2010].

The results of our analysis are shown in Figure 11. The majority of surveyed models (about 48.4%) use context-free grammars (CFGs), followed by fluid construction grammars (FCGs) (about 12.9%), universal grammar (UG) (about 9.7%), phrase-structure grammars, graph grammars, attribute grammars (about 6.5% each one) and, finally, definite-clause grammar, fuzzy grammar and Christiansen grammar (about 3.2% each one).

Most of the papers (about 83.9%) use non-adaptable grammatical representations (CFGs, UGs, phrase-structure grammars, graph grammars, attribute grammars, definite-clause grammar, and fuzzy grammar). Only five papers (about 16.1%) apply adaptable grammatical formalisms (FCGs and Christiansen grammar) having constructions that allow extending the grammar on the fly. The main motivations that lead researchers to choose non-adaptable grammatical representations are twofold. The former is the lower complexity of the model since it does not require further grammatical constructions and mechanisms that allow extending the grammar on the fly. The latter is the minor effort necessary for implementing efficiently this kind of linguistic representations since a huge amount of linguistic tools (e.g. parsers, taggers, disambiguators, etc.) are available for these more traditional formalisms (mainly CFGs).
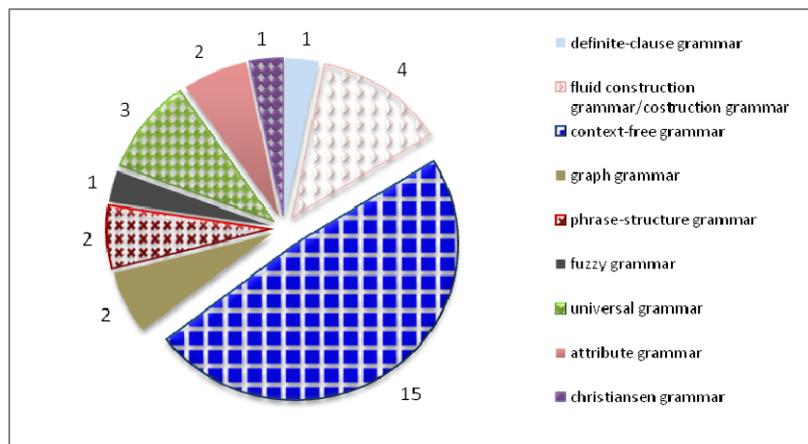


*Figure 11: Number of papers published in 2000-2015 applying a linguistic representation based on grammars*

## 5   Future research directions

In the previous section, we have analyzed current trends of language evolution research linked with the fulfillment of three main research challenges that are the need of empirical evidence, the need for semantically-enriched linguistic representations and adaptable grammatical formalisms. However, many open challenges still remain to be solved. In this section, we briefly discuss the future research directions resulting from the analysis of three kinds of data sources: the answers to a brief interview with the authors of some of the surveyed models, the future work section of some of the surveyed papers, and the literature on language evolution.

First of all, we have administrated a brief interview to the authors of some language evolution models asking whether and how their research on language evolution models has evolved in recent years. Specifically, during the interview we preliminary asked the authors for a list of their published papers on the language

evolution model and then, we asked if their research evolved toward some emerging topics. In case of an affirmative answer, we asked to specify which the emerging research directions are and to provide a list of their published papers on that. The interview was administered to the authors of 9 papers selected from the 90 papers proposing/using a computational method and relying on a linguistic representation (see Section 4.3). The selection has been performed by preferring the papers proposing a new language evolution model instead of using an existing one. Table 5 provides a summary of the answers received from the authors of these 9 papers about the evolution of their research. From these answers, the following considerations emerge. Most of the authors did not continue this research after the development of the model due to various reasons, mainly the end of project funding and a different research agenda (GRAEL, FCGLight, Iterated Learning Model (ILM), and Evolutionary game theory (EGT)). Some authors (GE$^2$) are focusing their research on alternative grammatical formalisms by experimenting with how the language evolution model performs with different kinds of context-free and context-sensitive grammars. Some authors (LEVER) are working on a further abstraction of the language evolution process by using metamodels and representing the language evolution as a transformation between metamodels of language. Finally, the authors of Game Dynamics (GD) are focusing their research on the cognitive aspects of language evolution studying the phenomenon at the neural synaptic level and trying to simulate through neural networks the evolution that happens in human language.

Moreover, to have further indications about future trends, we have analyzed the future work described in the 9 selected papers. From this analysis, three kinds of future challenges emerge (summarized in Table 6):

- application challenges: many language evolution models (LEVER and CGE) will investigate how to apply the model to new evolutionary problems;
- modeling challenges: several models (LEVER, GD, EGT, (GE)2, and FCGlight) will address new solutions for improving the representational ability of the model;

semantic challenges: some models (AGE and CGE) will investigate new solutions to specify semantics in the model.

| RESEARCH EVOLUTION | |
|---|---|
| GRAEL [De Pauw, 2003] | The authors did not continue research on language evolution, due to reasons of time and a different research agenda. |
| LEVER [Juergens and Pizka, 2006] | The authors are evolving their research on language evolution towards modeling languages tailored to a specific domain and defined by a metamodel. They are facing the problem of migrating existing models to a new version of their metamodel and proposed an approach, named COPE, which specifies the coupled evolution of metamodels and models to reduce migration effort. |
| $(GE)^2$ [O'Neill and Ryan 2004] | After the publication of the $(GE)^2$ model, the authors are exploring the use of different types of grammars ranging from context-free to context-sensitive including Tree-adjunct Grammars [Murphy et al., 2010], Attribute Grammars [O'Neill et al., 2004] and Shape Grammars [O'Neill et al., 2010] [McDermott et al., 2012]. |
| Attribute Grammar Evolution (AGE) [de la Cruz et al. 2005] | The authors are exploring the use of grammatical evolution to obtain an ecology of artificial beings associated with mathematical functions [Alfonseca and Soler Gil, 2013] in order to generate artificial ecologies that exhibit some of the features of natural evolution. |
| Christiansen Grammar Evolution (CGE) [Ortega et al. 2007] | The authors did not have new projects in this direction. |
| FCGlight [Saveluc and Ciortuz 2010] | The authors did not continue to work on the FCGlight topic, due to the end of project funding and involvement in other research areas. |
| Game dynamics (GD) [Mitchener 2007] | The authors are exploring how to set up stochastic dynamics to represent a population of language learners that can spontaneously move from one equilibrium point to another in a way that resembles documented language change. They are working also on a simulation of the evolution of neural synaptic coding, in order to evolve neural networks and manipulate information in something vaguely like what happens in language. |
| Iterated Learning Model (ILM) [Smith et al. 2003] | The authors did not continue to work on the ILM model for language evolution. |
| Evolutionary game theory (EGT) [Jäger, 2007] | The authors did not continue to work on the EGT model for language evolution. |

*Table 5: Answers received from some authors about the evolution of their research on language evolution models*

| | Future challenges from the surveyed papers |
|---|---|
| GRAEL | ---- |
| LEVER | • Application of LEVER to the development of real world DSLs to increase its expressiveness.<br>• Automatic adaptation of path expressions for those commands that merely refactor a language. |
| (GE)$^2$ | • Improvement of the scalability and flexibility of the model by evolving the number of functions along with their respective parameters, outputs and data types. |
| AGE | • Determination of the adequate semantics needed to get the optimal performance of the evolution model. |
| CGE | • Application of the CGE to new evolutionary problems.<br>• Investigation of other ways to specify semantics in language evolution models. |
| FCGlight | • Development of new strategies for language games (e.g. learning of clitic pronouns, etc.). |
| GD | • Improvement of the level of detail of the linguistic environment, including features such as noisy linguistic data and social and spatial structure. |
| ILM | ---- |
| EGT | • The interaction between first language acquisition and adult learning, and its modeling in terms of EGT. |

*Table 6: Future challenges from the 9 selected papers*

Finally, we have analyzed the language evolution literature, to identify further general challenges (not related to the surveyed models). Specifically, we perform a search using Google Scholar, Web of Science and Scopus for published research papers that deal with future directions/perspectives of language evolution modeling. A total of 7 research papers were selected and included in the analysis as they actually discuss future directions/perspectives of language evolution modeling. From this analysis, three main open research directions emerged.

First, standard approaches for validation of language evolution models are missing in the literature. As resulted from an analysis of the literature, indeed, there is not a standard validation strategy and standard metrics that can be applied to evaluate the performance of language evolution models. This lack of standard experimental approaches for model evaluation is also highlighted by Gong et al. [2014] that suggest exploring experimental semiotics and artificial language learning as a basis for possible experimental validation strategies. Therefore, future research could define standard strategies and metrics that help to evaluate and to compare language evolution models.

Another open challenge in language evolution research consists of a close collaboration among linguists, modelers, and neuro-scientists in order to define advanced language evolution models that integrate linguistics, modeling, and neuro-

scientific knowledge. Interdisciplinary competencies have to work together in order to obtain "*a biologically plausible, computationally feasible, and behaviorally adequate understanding of language evolution*" process [Gong et al., 2014]. Specifically, new findings of brain structure and neural processing are necessary that can lead to the generation of new hypotheses about language evolution [Niederhut, 2014]. As claimed by Niederhut (2014), "*by incorporating knowledge and methods from neuroscience, language evolution research can gain a strong foothold in the possible and the testable*". This challenge matches also with the research undertaken by the authors of GD (see Table 5).

As a further future perspective, language evolution models should take into account multimodal aspects of language [Grifoni et al., 2016]. Current language evolution models, indeed, rely on the assumption that language is identified by three main components: speech (signal), syntax (structure) and semantics (meaning) [Fitch, 2005]. Therefore, they address the language evolution problem using a linguistic approach focused on speech and/or text. However, human beings communicate using the five senses and there is an emergent tendency to embed different modalities (such as gestures, facial expressions, etc) into the language [D'Ulizia and Ferri, 2006, D'Ulizia et al., 2007; D'Ulizia et al., 2008; Paulmann et al., 2009; Regenbogen et al., 2012, Ferri et al., 2012; Kanero, 2014; D'Andrea et al., 2017]. This tendency is confirmed also by several recent researches [Vigliocco et al., 2014; Levinson and Holler J, 2014; Waller et al., 2013] that highlight the multimodal nature of language and the relevance of multimodality for language learning and evolution. In fact, by the historical point of view "*language research has focused predominantly on speech and/or text, thus ignoring the wealth of additional information available in face-to-face communication*"; however, the analysis of literature carried out in [Vigliocco et al., 2014] makes evident "*that speech and gesture are part and parcel of the same system and together constitute a tightly integrated processing unit, thus underscoring the need for a multimodal approach to the study of language*". Caschera et al. [2012, 2018] have highlighted the necessity of tools for modeling the evolution of multimodal human interaction in long-term changing situations. As a future research perspective, therefore, we envision language evolution models that embed multimodal aspects as a crucial part of the language evolution.

# 6 Conclusion

In this paper, we have investigated how language evolution modelers have faced three main research challenges in language evolution research, that are the need for empirical evidence, the need of semantically-enriched linguistic representations and the need of adaptable grammatical formalisms. To achieve that, a hybrid methodology has been used, which integrates bibliometric analysis of several scientific papers proposing language evolution models, and interviews with some authors of the papers. The surveyed models have been analyzed considering their scientific production, scientific impact, and to what extent they fulfill the challenges.

From this analysis, and considering answers to the interview, we have extracted some conclusions about the research impact of current language evolution models and future research directions.

About the research impact, agent-based models turned out to be the most published class of language evolution models. Moreover, the results showed that about 70% of the agent-based language evolution models are combined with further techniques that range from iterated learning and language game to machine learning. Another result of the bibliometric analysis showed that signal-meaning association-based models were the class of language evolution models most published and most cited. Moreover, grammar-based models also turned out to be a highly published class of language evolution models. Moreover, the analysis of the level(s) of language (lexical, syntactic, and semantic) addressed by the linguistic representations showed that language evolution models were oriented mainly towards either lexical communication systems or lexical-syntactic-semantic communication systems. Finally, the results of the analysis showed that CFGs are the most used grammatical formalism due to their simplicity and intuitively appealing formalism.

About future research directions, the following challenges are emerging: (i) standard validation approaches, (ii) a multidisciplinary collaboration among linguists, modelers and neuro-scientists, and (iii) the embedding of multimodal aspects.

The results of this study can be valuable to many language evolution researchers that have to determine previous and current research highlights and possible topics for future researches.

# References

[Abrams and Strogatz, 2003] Abrams, D.M., & Strogatz, S.H. (2003). Modeling the dynamics of language death. Nature 424:900.

[Alfonseca and Soler Gil, 2013] Alfonseca, M., & Soler Gil, F. J. (2013). Evolving an ecology of mathematical expressions with grammatical evolution. Biosystems, 111(2), pp 111-119.

[Baronchelli et al., 2006] Baronchelli, A., Felici, M., Loreto, V., Caglioti, E., & Steels, L. (2006). Sharp transition towards shared vocabularies in multiagent systems. Journal of Statistical Mechanics 6:6–14.

[Baxter et al., 2009] Baxter, G.J., Blythe, R.A., Croft, W., & McKane, A. J. (2009). Modeling language change: An evaluation of Trudgill's theory of the emergence of New Zealand English. Lang Var Change 21:257–296.

[Benz et al., 2011] Benz, A., Ebert, C., Jäger, G., & van Rooij, R. (2011) Language, games and evolution, LNCS 6207, Springer Berlin Heidelberg.

[Briscoe, 2000] Briscoe, T. (2000). Grammatical acquisition: Inductive bias and coevolution of language and the language acquisition device. Language 76(2):245–296.

[Cangelosi and Parisi, 1998] Cangelosi, A., & Parisi, D. (1998). The emergence of a "language" in an evolving population of neural networks. Connection Science 10(2):83–89.

[Cangelosi, 2001] Cangelosi, A. (2001). Evolution of communication and language using signals, symbols, and words. IEEE Transactions on Evolutionary Computation, 5(2), 93-101.

[Caschera et al., 2012] Caschera, M. C., D'Ulizia, A., Ferri, F., & Grifoni, P. (2012). Towards Evolutionary Multimodal Interaction. In On the Move to Meaningful Internet Systems: OTM 2012 Workshops, Springer Berlin Heidelberg, pp. 608-616.

[Caschera et al., 2018] Caschera, M. C., D'Ulizia, A., Ferri, F., & Grifoni, P. (2018). MONDE: a method for predicting social network dynamics and evolution. Evolving Systems, 1-17.

[Chen et al., 2015] Chen, H., Jiang, W., Yang, Y., Yang, Y., & Man, X. (2015). Global trends of municipal solid waste research from 1997 to 2014 using bibliometric analysis, Journal of the Air & Waste Management Association, 65:10, 1161-1170.

[Chomsky, 1986] Chomsky N (1986). Knowledge of language. New York: Praeger.

[Christiansen, 1990] Christiansen, M.H. (1990). A survey of adaptable grammars. ACM SIGPLAN Notices 25(11):35–44.

[Christiansen and Chater, 1999] Christiansen, M.H., & Chater, N. (1999). Toward a connectionist model of recursion in human linguistic performance. Cognitive Science 23:157–205.

[Christiansen and Kirby, 2003] Christiansen, M.H., & Kirby, S. (2003). Language evolution: Consensus and controversies. Trends in Cognitive Sciences 7(7):300–307.

[D'Andrea et al., 2017] D'Andrea A, D'Ulizia A, Ferri F, Grifoni P. (2017). EMAG: An Extended Multimodal Attribute Grammar for Behavioural Features. Digital Scholarship in the Humanities, 32(2), 251-275.

[de Boer, 2006] de Boer, B. (2006). Computer modelling as a tool for understanding language evolution. In Evolutionary Epistemology, Language and Culture – A non-adaptationist, systems theoretical approach. Springer, Dordrecht, pp 381-406.

[Dediu and de Boer, 2016] Dediu, D., & de Boer, B. (2016). Language evolution needs its own journal. Journal of Language Evolution, 1(1), 1-6.

[de la Cruz et al., 2005] de la Cruz, M., de la Puente, A.O., & Alfonseca, M. (2005). Attribute grammar evolution, artificial intelligence and knowledge engineering applications: A bioinspired approach. First International Work-Conference on the Interplay Between Natural and Artificial Computation, IWINAC 2005, Las Palmas, Canary Islands, Spain, June 2005, pp 182–191.

[De Pauw, 2003] De Pauw, G. (2003). GRAEL: an agent-based evolutionary computing approach for natural language grammar development, Proceedings of the 18th international joint conference on Artificial intelligence, August 09-15, 2003, Acapulco, Mexico, pp 823-828.

[D'Ulizia and Ferri, 2006] D'Ulizia, A., & Ferri, F. (2006). Formalization of multimodal languages in pervasive computing paradigm, Advanced Internet Based Systems and Applications, Second International Conference on Signal-Image Technology and Internet-Based Systems (SITIS 2006), Revised Selected Papers, Springer, Lecture Notes in Computer Science 4879, pp 126-136.

[D'Ulizia et al., 2007] D'Ulizia, A., Ferri, F., & Grifoni, P. (2007). A Hybrid Grammar-Based Approach to Multimodal Languages Specification, OTM 2007 Workshop Proceedings, 25-30 November 2007, Vilamoura, Portugal, Springer-Verlag, Lecture Notes in Computer Science 4805, pp 367-376.

[D'Ulizia et al., 2008] D'Ulizia, A., Ferri, F., & Grifoni, P. (2008). Toward the Development of an Integrative Framework for Multimodal Dialogue Processing. In On the Move to Meaningful Internet Systems: OTM 2008 Workshops Springer Berlin Heidelberg, pp. 509-518.

[Ferri et al., 2012] Ferri, F., D'Ulizia, A., & Grifoni, P. (2012). Multimodal Language Specification for Human Adaptive Mechatronics. Journal of Next Generation Information Technology, 3(1), pp 47-57.

[Fitch, 2005] Fitch, W.T. (2005). The evolution of language: A comparative review. Biology and Philosophy 20(2–3):193–203.

[Garcia, 2014] Garcia J. (2014). Machine Learning and Cognitive Systems: The Next Evolution of Enterprise Intelligence (Part I). *Wired Innovation Insights*.

[Gong et al., 2006] Gong, T., Minett, J.W., & Wang, W.S. (2006). Computational simulation on the coevolution of compositionality and regularity. In: Proceedings of the 6th international conference on the evolution of language, 12–15 Apr 2006, Rome, Italy, pp 99–106.

[Gong et al., 2014] Gong, T., Shuai, L., & Zhang, M. (2014). Modelling language evolution: Examples and predictions. Physics of life reviews, 11(2), 280-302.

[Grifoni et al., 2016] Grifoni, P., D'Ulizia, A., & Ferri, F. (2016). Computational methods and grammars in language evolution: a survey, Artificial Intelligence Review. Volume 45, Issue 3, pp 369-403.

[Harnad, 1990] Harnad, S. (1990) The Symbol Grounding Problem. Physica D 42: 335-346.

[Hauser et al., 2014] Hauser, M. D., Yang, C., Berwick, R. C., Tattersall, I., Ryan, M. J., Watumull, J., ... & Lewontin, R. C. (2014). The mystery of language evolution. Frontiers in psychology, 5, 401.

[Hickey and Puppel, 1997] Hickey, R., & Puppel S. (eds) 1997. Language History and Linguistic Modelling. A Festschrift for Jacek Fisiak on his 60th Birthday. Berlin: Mouton de Gruyter, 2 vols., 2121 pages.

[Jackendoff, 1999] Jackendoff, Ray. 1999. Possible stages in the evolution of the language capacity. Trends in Cognitive Sciences 3, 272–279.

[Jäger, 2004] Jäger, G. (2004). Learning constraing subhierarchies: The bidirectional gradual learning algorithm. In: R. Blutner & H. Zeevat (eds.), Optimality theory and pragmatics. Basingstoke: Palgrave MacMillan, pp 251–287.

[Jäger, 2007] Jäger, G. (2007). Evolutionary Game Theory and Typology: A Case Study. Language 83(1), pp 74-109.

[Jaeger et al., 2009] Jaeger, H., Baronchelli, A., Briscoe, E., Christiansen, M. H., Griffiths, T., Jäger, G., Kirby, S., Komarova, N., Richerson, P. J., Steels, L., Triesch, J. (2009). What can mathematical, computational and robotic models tell us about the origins of syntax?, book chapter in "Biological Foundations and Origin of Syntax", eds. D. Bickerton and E. Szathmáry. Strüngmann Forum Reports, vol. 3. Cambridge, MA: MIT Press, pp 385-410.

[Jimenez-Lopez, 2012] Jimenez-Lopez, M. D. (2012). A Grammar-Based Multi-Agent System for Language Evolution, Highlights on PAAMS, AISC 156, pp 45–52.

[Juola, 2003] Juola, P. (2003). The time course of language change. Computers and the Humanities, 37(1), 77-96.

[Juergens and Pizka, 2006] Juergens, E., & Pizka, M. (2006). The Language Evolver Lever - Tool Demonstration, Electronic Notes in Theoretical Computer Science, Volume 164, Issue 2, pp 55-60.

[Kandler and Steele, 2008]  Kandler, A., & Steele, J. (2008). Ecological models of language competition. Biol Theor 3:164–173.

[Kanero, 2014] Kanero, J. (2014). The gesture theory of language origins: Current issues and beyond, In: McCrohon L, Thompson B, Verhoef T and Yamauchi H (eds) The past, present and future of language evolution research, Tokyo: EvoLang9 Organising Committee, pp 1–7.

[Kaplan, 2005] Kaplan, F. (2005). Simple models of distributed coordination. Connect Sci 17(3–4):249–270.

[Kirby et al., 2009] Kirby, S., Christiansen, M., & Chater, N. (2009). Syntax as an adaptation to the learner. In Bickerton D and Szathmáry E (eds) Biological foundations and origin of syntax, Strüngmann Forum Reports, vol. 3. Cambridge, MA: MIT Press.

[Lipowska D, 2011] Lipowska D (2011) Naming game and computational modelling of language evolution.

[Levinson and Holler J, 2014] Levinson, S.C., & Holler, J. (2014). The origin of human multi-modal communication. Phil. Trans. R. Soc. B 369.

[Luuk and Luuk, 2014] Luuk, E., & Luuk, H. (2014). The evolution of syntax: Signs, concatenation and embedding. Cognitive Systems Research, 27, 1-10.

[McDermott et al., 2012] McDermott, J., Swafford, J. M., Hemberg, M., Byrne, J., Hemberg, E., Fenton, M., McNally, C., Shotton, E., & O'Neill, M. (2012). An Assessment of String-Rewriting Grammars for Evolutionary Architectural Design. Environment and Planning B, 39(4), pp 713-731.

[Minett and Wang, 2008] Minett, JW., & Wang, W.S. (2008). Modeling endangered languages: The effects of bilingualism and social structure. Lingua 118(1):19–45.

[Mitchener, 2007] Mitchener, W. G. (2007). Game dynamics with learning and evolution of universal grammar. Bulletin of Mathematical Biology, 69(3), pp 1093-1118.

[Murphy et al., 2010] Murphy, E., O'Neill, M., Galvan-Lopez, E., & Brabazon, A. (2010). Tree-Adjunct Grammatical Evolution. IEEE Congress on Evolutionary Computation 2010 IEEE Press Barcelona, Spain.

[Nettle, 1999] Nettle, D. (1999). Is the rate of linguistic change constant? Lingua 108:119–136.

[Niederhut, 2014] Niederhut, D. (2014). Beyond "neuroevidence". In: McCrohon L, Thompson B, Verhoef T, Yamauchi B (eds) The past, present and future of language evolution research, Tokyo: EvoLang9 Organising Committee, pp. 102–109.

[Nowak et al., 2002] Nowak, M.A., Komarova, N.L., & Niyogi, P. (2002). Computational and evolutionary aspects of language. Nature 417(6889):611–617.

[O'Neill et al., 2010] O'Neill, M., McDermott, J., Swafford, J. M., Byrne, J., Hemberg, E., Shotton, E., McNally, C., Brabazon, A., & Hemberg, M. (2010). Evolutionary Design using Grammatical Evolution and Shape Grammars: Designing a Shelter. International Journal of Design Engineering, 3(1), pp 4-24.

[O'Neill and Ryan, 2004] O'Neill, M., & Ryan, C. (2004). Grammatical evolution by grammatical evolution: The evolution of grammar and Genetic Code. LNCS 3003, pp 138-149.

[Ortega et al., 2007] Ortega, A., De La Cruz, M., & Alfonseca, M. (2007). Christiansen grammar evolution: Grammatical evolution with semantics. Evolutionary Computation, IEEE Transactions on, 11(1):77–90.

[Parisi et al., 2008] Parisi, D., Antinucci, F., Natale, F. et al. (2008). Simulating the expansion of farming and the differentiation of European languages. In Origin and Evolution of Languages: Approaches, Models, Paradigms, B. Laks (ed.), *Equinox* Publishing, pp 234–258.

[Patriarca and Leppänen, 2004] Patriarca, M., & Leppänen, T. (2004). Modeling language competition. Physica A 338(1–2):296–299.

[Paulmann et al., 2009] Paulmann, S., Jessen, S., & Kotz, S. A. (2009). Investigating the multimodal nature of human communication: Insights from ERPs. Journal of Psychophysiology, 23(2), 63-76.

[Regenbogen et al., 2012] Regenbogen, C., Schneider, D. A., Gur, R. E., Schneider, F., Habel, U., & Kellermann, T. (2012). Multimodal human communication — Targeting facial expressions, speech content and prosody. NeuroImage 60(4), 2346-2356.

[Saveluc and Ciortuz, 2010] Saveluc, V., & Ciortuz, L. (2010). FCGlight: A system for studying the evolution of natural language. 12th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing SYNASC 2010, Timisoara, Romania, 23-26 September 2010, IEEE, pp. 188–193.

[Singh, 2005] Singh, Y. N. (2005). Computational Modelling of Evolution of Language. Retrieved                                                                    from http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.97.7997&rep=rep1&type=pdf on 10 February 2015.

[Smith, 1982] Smith, J. M. (1982). Evolution and the Theory of Games. Darwin (Vol. 13, pp 224). Cambridge University Press.

[Smith et al., 2003] Smith, K., Kirby, S., & Brighton, H. (2003). Iterated learning: a framework for the emergence of language. Artificial Life, 9(4), pp 371–86.

[Soules, 2002] Soules, M. (2002). Animating the language machine: Computers and performance. Computers and the Humanities, 36(3), 319-344.

[Spranger and Steels, 2012] Spranger, M., & Steels, L. (2012). Synthetic modeling of cultural language evolution. In: Five approaches to language evolution, Tokyo, Evolang Organization Committee, pp. 130–139.

[Steels, 1995] Steels, L. (1995). A self-organizing spatial vocabulary. Artificial Life Journal 2(3):319–332.

[Steels, 1997] Steels, L. (1997). The synthetic modelling of language origins. Evolution of Communication 1:1–34.

[Steels, 2010] Steels, L. (2010). Modeling the formation of language in embodied agents: Methods and Open Challenges. In: Nolfi S, Mirolli M (eds.), Evolution of communication and language in embodied agents, Springer-Verlag Berlin Heidelberg, pp 223–233.

[Steels, 2011] Steels, L. (2011). Introducing Fluid Construction Grammar. In: Steels L (Ed.), Design patterns in fluid construction grammar, 3–30. Amsterdam: John Benjamins.

[Steels and de Beule, 2006] Steels, L., & De Beule, J. (2006). A (very) brief introduction to fluid construction grammar. In: Proceedings of the 3rd Workshop on Scalable Natural Language Understanding, New York City, June 2006, pp 73–80.

[Vigliocco et al., 2014] Vigliocco, G., Perniss, P., & Vinson, D. (2014). Language as a multimodal phenomenon: implications for language learning, processing and evolution. Phil.Trans. R. Soc. B, 369 (1651).

[Vogt, 2006] Vogt, P. (2006). Language evolution and robotics: Issues in symbol grounding and language acquisition. In: Artificial Cognition Systems. A. Loula, R. Gudwin, & J. Queiroz (Eds.) Idea Group, pp176–209.

[Vogt, 2009] Vogt, P. (2009). Modeling interactions between language evolution and demography. Human Biology 81(2): 237–258.

[Vogt and de Boer, 2010] Vogt, P., & de Boer, B. (2010Vogt, P., & de Boer, B. (2010). Language evolution: computer models for empirical data. Adaptive Behavior, 18.

[Von Neumann and Morgenstern, 1944] Von Neumann, J., Morgenstern, O. (1944). Theory of games and economic behavior. Princeton University Press 2:625.

[Waller et al., 2013] Waller, B., Liebal, K., Burrows, A., & Slocombe, K. (2013). How can a multimodal approach to primate communication help us understand the evolution of communication?. Evolutionary Psychology, 11(3), 538-549.

[Wang et al., 1978] Wang, W.S., Liao, C.C., Gaskins, R., & Wang, M.S. (1978). QUINCE system: State-of-the-art review. California Univ Berkeley Dept of Linguistics.

[Watumull and Hauser, 2014] Watumull, J., & Hauser, M. D. (2014). Conceptual and empirical problems with game theoretic approaches to language evolution. Frontiers in psychology, 5.

[Wellens et al., 2008] Wellens, P., Loetzsch, M., & Steels, L. (2008). Flexible word meaning in embodied agents. Connection Science 20(2):173–191.

[Xie et al., 2008] Xie, S., Zhang, J., & Ho, Y. S. (2008). Assessment of world aerosol research trends by bibliometric analysis. Scientometrics 77(1):113–30.

[Zuidema, 2002] Zuidema, W. (2002). Language adaptation helps language acquisition – A computational model study. In: Proceedings of the Seventh International Conference on Simulation of Adaptive Behavior. En B. Hallam, D. Floreano, J. Hallam, G. Hayes y J.Meyer (Eds.), Cambridge, MA: MIT Press.

## Appendix

### I. Computational methods in language evolution models

In this section we give some details on the most applied computational methods, resulting from the performed bibliometric analysis.

**Agent-based methods** allow the emergence of a language to be represented in a bottom-up fashion as a result of the interaction of a group of agents. These agents represent humans with different kinds of cognitive and social behavior that interact with each other in a population. Each agent is equipped with linguistic abilities for conceptualization, production, parsing, and interpretation. Moreover, they are endowed with learning mechanisms that allow expanding their basic linguistic knowledge. Agent-based methods are naturally adopted by those who want to develop real-world applications, such as software agents or robots evolving shared communication systems.

**Machine learning-based methods** are used to automate the individual learning of adaptive systems (mainly agents) involved in language evolution by emulating the human linguistic behavior. The major machine learning techniques used in language evolution include neural networks, genetic algorithms, Bayesian networks, and rule induction. While in the past they were applied independently, in recent times these techniques are being used in a hybrid fashion, closing the boundaries between them and enabling the development of more effective models [Garcia, 2014].

**Game-theoretic methods** address language evolution from the most general perspective by aggregating the behavior of a population and defining general mathematical equations that model the evolution of this behavior. Specifically, they

study language evolution from the viewpoint of mathematical game theory, which was developed by Von Neumann and Morgenstern (1944) to describe human and economic behavior. Game-theoretic research in language evolution [Benz et al., 2011] has led to the definition of a framework, whose functioning can be summarized as follows. There are two players, the sender (or teacher) and the receiver (or learner), and a finite alphabet (a set of symbols). A language is a probability distribution defined on a set of strings composed of the symbols of the alphabet. One player has to choose between several strings compatible with the teacher's grammar such that its preferences over its own strings depend on which strings the other player chooses. Preferences are real numbers that are attached to each choice of each player and represent the "utility". The aim is to predict how the players behave in order to maximize their expected utility.

## II. Linguistic representations in language evolution models

In this section we give some details on the most applied linguistic representations used in language evolution models, resulting from the performed bibliometric analysis. As argued by Cangelosi (2001), "some rely on the use of simple signals, while others use symbolic communication systems or complex syntactical structures".
**Signal-meaning association-based methods** are strictly related to the symbol grounding problem [Harnad, 1990] that investigates how to ground a set of symbols (e.g. forms, signals, utterances, words, etc.) in a set of possible meanings. In this method, communication relies on simple associations between symbols and meanings. This linguistic representation is mainly applied by language evolution models that only focus on lexicon emergence, without making any explicit reference to the role of syntax. It is best suited to model the early stages of the human language evolution (e.g. primates communication).
**Grammar-based methods** move from signal-meaning associations to more complex structures of language using grammatical constructions to link syntactic and semantic categories intervening between symbols and meanings. In these methods, grammar emerges in the attempt to convey more information by combining words into phrases or sentences [Nowak et al., 2002]. Depending on the kind of grammatical formalism applied, this linguistic representation allows representing syntactic and/or semantic relationships between symbols. It is best suited to simulate complex languages in which it is possible to identify "words", i.e. symbols that belong to specific grammatical classes, such as verbs, nouns, prepositions, etc. [Cangelosi, 2001].
**Graph-based methods** allow representing linguistic units (e.g. words, sentences, etc.) as nodes in a graph and relations between them as edges. This linguistic representation is used mainly by phylogenetic approaches that model longer timescale of the history of populations by investigating how language originates and develops during the evolutionary history of the species. In these papers, indeed, the linguistic knowledge is mainly represented through graphs, trees or networks modelling the relationships between a group of related languages.

## III. Grammatical formalisms in language evolution models

In this section we give some details on the most applied grammatical formalisms used in language evolution models, resulting from the performed bibliometric analysis.

**Context-free grammars** (CFGs) are the class of Chomsky grammars most used in language evolution models due to their simplicity and intuitively appealing formalism. CFGs have the main benefit that they are able to model all frequent linguistic phenomena of human (natural) language assuring, at the same time, a low parsing complexity. They are easy for human beings to understand and for computers to manipulate [Wang et al., 1978]. CFGs are, indeed, computationally tractable, i.e. it is possible to write a computer program which determines whether sentences are grammatical or not.

**Fluid Construction Grammars** (FCGs) were developed with the aim of incorporating the cognitive and interactional foundations of language into the grammatical model. FCGs [Steels and de Beule, 2006] are a kind of construction grammar that offers a way of structuring and representing semantic meaning into patterns and named constructions. One of the main benefits of FCGs is the flexibility of the language processing in the sense that sentences can be understood even if they are partially ungrammatical or incomplete. Next, FCGs adopt the reversibility principle which means that the same constructional definition must be usable without change both in parsing and production and without compromising efficiency or generating unnecessary searches [Steels, 2011]. The main limitation of FCGs was highlighted by Gong et al. (2006) and relies on the lack of matching between syntactic structures and semantic categories.

**Universal Grammar** (UG) is a theory, developed by Chomsky, which explains how humans acquire languages. It arises from the necessity of justifying the "poverty of stimulus" to which children's brains are exposed. Hearing sentences coming from the environment does not allow children to uniquely and correctly specify the underlying grammatical rules. Therefore, children must have some innate capacity to learn the correct grammar from a set of candidate grammars. UG is the theory of this restricted set. Formally, UG is not a grammar, but a theory of a collection of grammars [Nowak et al., 2002]. According to Chomsky (1986, pp. 150–151), UG consists of "a system of principles with parameters to be fixed, along with a periphery of marked exceptions". The "core grammar" entails a set of universal principles, which apply to all languages, and a set of parameters which may vary from language to language. By contrast, the "peripheral grammar" is made up of quirks and irregularities of language. The main advantage of UG stays in the fact that it helps the learner to generalize rules and allows creative use of the language. On the negative side, this theory does not take into account social and environmental factors for differentiating learners, but only ideal learners with ideal grammars are considered.