

# Language acquisition with communication between learners

Rasmus Ibsen-Jensen<sup>1,†,\*</sup>, Josef Tkadlec<sup>1,†</sup>, Krishnendu Chatterjee<sup>1</sup>, and Martin A. Nowak<sup>2</sup>

<sup>1</sup>IST Austria, Klosterneuburg, A-3400, Austria

<sup>2</sup>Program for Evolutionary Dynamics, Department of Organismic and Evolutionary Biology, Department of Mathematics, Harvard University, Cambridge, MA 02138, USA

<sup>†</sup>These authors contributed equally to this work.

\*Corresponding author: ribsens@ist.ac.at

**We consider a class of students learning a language from a teacher. The situation can be interpreted as a group of child learners receiving input from the linguistic environment. The teacher provides sample sentences. The students try to learn the grammar of the teacher. In addition to just listening to the teacher, the students can also communicate with each other. The students hold hypotheses about the grammar and change them if they receive counter evidence. The process stops when all students have converged to the correct grammar. We study how the time to convergence depends on the structure of the class room by introducing and evaluating various complexity measures. We find that structured communication between students, although potentially introducing confusion, can greatly reduce some of the complexity measures. Our theory can also be interpreted as applying to the scientific process, where nature is the teacher and the scientists are the students.**

Keywords: Language learning; Inductive Inference; Population structures in learning

## 23 1 Introduction

24 In traditional language learning theory [1, 2, 3], there is a teacher and a learner [4, 5, 6]. The  
25 teacher uses a particular grammar and provides sample sentences from the corresponding language.  
26 A language is a set of finitely or infinitely many sentences. A grammar is a finite list of rules  
27 that specifies the language. The learner has a search space of candidate grammars. The task for  
28 the learner is to converge to the grammar of the teacher after having heard a sufficient number  
29 of sentences. This setting for learning is called “inductive inference” [7, 8]. The goal is to infer  
30 the underlying rules from examples. The teacher cannot directly communicate the rules of the  
31 grammar, (s)he only provides sample sentences consistent with it.

32 Learning by inductive inference is more general than natural language acquisition. It arises  
33 whenever generative rules are supposed to be inferred from examples. It is the basis for mutual  
34 understanding in human communication. It is also the activity of scientists searching for the laws  
35 of nature [9]. The scientists conduct experiments and the nature gives the answers. Then the  
36 scientists seek to formulate the underlying rules, the grammar of nature. In the present work, we  
37 focus on language learning as a particular case of cultural transmission.

38 Learning theory is often concerned with positive or negative results about the learnability of  
39 sets of grammars [10, 11, 12, 13, 14, 15, 16, 17]. It is the basis for a mathematical formalization of  
40 what Chomsky calls “universal grammar” [8, 18, 19]. Several works also considered the computa-  
41 tional problems related to learning [20, 21, 22]. In the evolutionary dynamics of human language  
42 acquisition, the question is extended to asking under which conditions a population of speakers  
43 learning from each other can converge to a coherent language [23, 24, 25, 26, 27].

44 In this paper we explore a new setting. There is a teacher (either a person, or a body of  
45 knowledge, or the linguistic environment or nature) and a population of learners. In addition to  
46 just listening to the teacher, the learners can also communicate with each other. At each moment,  
47 each learner holds a hypothesis as for what is the teacher’s grammar and can update this hypothesis  
48 upon hearing a single sentence from the teacher or some other learner. The learners and the teacher  
49 speak and listen to one another until, eventually, all learners successfully learn the grammar of the  
50 teacher. In the next section we introduce a model in which the communication among learners and  
51 the teacher proceeds in an organised way. We study which communication structures improve – or  
52 obstruct – the efficiency of this learning process.

53 The efficiency of the learning process also depends on the power of individual learners. Here  
54 we consider learners of two different types: weak memoryless learners and powerful batch learners.  
55 As far as memory is concerned, these two types of learners serve as a lower and upper bound for  
56 human learning capacity [5, Section 13.3.3]. Memoryless learners hold, at any moment, a candidate  
57 grammar. Whenever they receive a counterexample (a sentence that doesn't belong to the language  
58 corresponding with their current grammar) they randomly choose another grammar from their  
59 search space. They are called "memoryless" because they could pick a grammar which they have  
60 already rejected. In contrast, batch learners keep track of all the inputs they have received so far  
61 and for their hypothesis they always select grammar that is most consistent with the sentences  
62 they have observed so far. When learning from a single teacher without other inputs, both types  
63 of learners have the property of consistency: once they find the right grammar they do not change  
64 it anymore.

65 The underlying dynamical system can be seen as a new kind of evolutionary process. Candi-  
66 date grammars spread in the population of learners. The teacher, or the environment, selects for  
67 particular grammars. The process stops when all learners have adopted the correct grammar. The  
68 basic question is: How is the time to linguistic coherence affected by the population structure?

## 69 2 Model

70 In this section, we first introduce a general model for language learning with structured commu-  
71 nication between learners. Next we present two types of learners (memoryless  $(p, q)$ -learners and  
72 powerful batch learners) that we later analyze in detail. Finally, we introduce a complexity mea-  
73 sure called *rounds complexity* that we use to evaluate the efficiency of the learning process for  
74 different communication structures and types of learners. Our main scientific finding is as follows:  
75 while communication between learners can potentially cause confusion and certain communication  
76 structures between learners indeed do slow down the learning process, we present communication  
77 structures that can significantly expedite the learning process.

78 The process of learning a language can be modelled in a variety of ways [28, 29, 30, 31, 32,  
79 33]. In the traditional setting there is a single teacher and a single learner, and only the teacher  
80 communicates with the learner. Here we extend the traditional setting as follows:

- 81 1. We consider a single teacher and a population of learners.

82 2. The population of learners can communicate among each other.

83 3. We consider structured communication between the learners and study whether such com-  
84 munication can improve the efficiency of the process.

85 For clarity of presentation, we identify a grammar (a list of rules) with the language (a set of  
86 sentences) it generates. The hypothesis of each individual at each time is thus a language. (Recall  
87 that the units passed at each communication event are sentences.)

88 *Single learner.* In the traditional “single teacher – single learner” scenario, the teacher speaks some  
89 language  $L_1$  unknown to the learner and repeatedly generates sentences from  $L_1$ . The learner has a  
90 search space of possible languages  $L_1, L_2, \dots$  and initially holds an arbitrary hypothesis as for what  
91 the teacher’s language is. Upon hearing each sentence from the teacher, the learner can update  
92 this hypothesis. The process ends when the learner’s hypothesis becomes  $L_1$ .

93 *Structured learning for multiple learners.* In our case, there is a group of  $n + 1$  individuals (one  
94 teacher and  $n$  learners). There is a set  $L$  of  $\ell$  languages  $L_1, \dots, L_\ell$ . Each language consists of  
95 sentences (one sentence can belong to multiple languages).

96 The communication structure among learners is represented by a directed graph (network)  
97 where nodes correspond to individuals (including the teacher) and an edge (arrow) from individual  
98  $A$  to  $B$  means that  $A$  listens to  $B$ . At each moment, each learner holds a hypothesis  $L_i \in L$   
99 regarding what the teacher’s language is. Initially, teacher holds  $L_1$  and the hypotheses of the  
100 learners are arbitrary. In every round of the learning process we pick all the edges of the graph  
101 one by one, in random order. Every time an edge is picked, the speaker of that edge generates a  
102 sentence from the language she is currently hypothesizing and the listener of the edge can update  
103 his hypothesis. The process stops when all the learners learned the teacher’s language  $L_1$ .

104 *Example.* As a toy example, consider a single teacher  $T$  and two learners  $A$  (Alice) and  $B$  (Bob)  
105 such that both  $A$  and  $B$  listen to  $T$  and moreover  $B$  listens to  $A$ . Suppose that there are two  
106 languages  $L_1, L_2$  that don’t overlap at all. Suppose that  $A$ ’s initial hypothesis is  $L_2$  while  $B$  starts  
107 with  $L_1$  ( $T$  starts with  $L_1$  too). Finally, suppose that both learners follow the same simple update  
108 rule: whenever they hear a sentence they can not parse, they switch their hypothesis to the other  
109 possible language with probability 80 % (and keep it otherwise).

110 In this example, a single round can play out as follows (see Figure 1(b)): First we pick the edge  
111 between  $B$  and  $T$ .  $B$  receives a sentence he understands and keeps his hypothesis  $L_1$ . Next we

112 pick the edge between  $B$  and  $A$ .  $B$  receives a sentence from  $A$ 's language  $L_2$ . He can't parse it  
113 and (with probability 80 %) he switches his hypothesis to  $L_2$ . Finally, we pick the edge between  $A$   
114 and  $T$ .  $A$  receives a sentence she can't parse, still (with probability 20 %) she sticks to her current  
115 hypothesis  $L_2$ . As an outcome of the round, both  $A$  and  $B$  now hold the wrong hypothesis  $L_2$ .

116 Note that had we first picked the edge between  $A$  and  $T$ ,  $A$  could have switched to  $L_1$  with  
117 probability 80 % and the whole process would have finished in a single round. Allowing learners to  
118 speak among themselves can create confusion and can result in less efficient learning.

119 *Memoryless learners:  $(p, q)$ -learning.* Here we describe a type of a memoryless learner that we call a  
120  $(p, q)$ -learner. There are two positive numbers  $p, q \in [0, 1]$  with  $p + q \leq 1$ . Upon hearing a sentence,  
121 a  $(p, q)$ -learner updates her hypothesis as follows: (a) if the learner holds the same language as the  
122 speaker, then nothing changes; (b) if the learner holds a different language from the speaker, then:

- 123 1. with probability  $p$  the learner's hypothesis changes to the language of the speaker;
- 124 2. with probability  $q$  the learner's hypothesis does not change;
- 125 3. with probability  $(1 - p - q)/(\ell - 2)$  the learner switches to one of the remaining languages (i.e.,  
126 with the remaining probability one of the other languages is chosen uniformly at random).

127 An illustration is presented in Figure 1(a).

128 The parameters  $p, q$  can model various features of language learning. (a) The parameter  $q$  can  
129 represent the overlap between different languages, such that even if the languages of the speaker  
130 and the listener are different, the sentence from the speaker can be parsed by the listener and hence  
131 the listener does not switch. (b) The parameter  $p$  represents the bias to switch to the language of  
132 the speaker by listening to a single sentence. Note that since the switch happens by listening to a  
133 single sentence we consider that  $p$  is proportional to  $1/\ell$ .

134 *Discussion of  $(p, q)$ -learners.* We explain how our model of a  $(p, q)$ -learner generalises several  
135 classical language learning scenarios considered in the literature.

- 136 • *RWA:* A model of random walk (without greediness and single-value constraints) (RWA) on  
137 languages has been considered in [6, Section 4.2.1] where if the speaker and the listener have  
138 different languages, then the switch is uniformly at random among all languages. In the above  
139 setting we achieve this with  $p = q = 1/\ell$ .

- 140 • *SS*: A model of language learning with symmetric language overlap (SS) was considered in [5,  
141 Section 13.3.2]. The overlap was characterised by parameter  $a$  in [5, Eqn. (13.26)], which  
142 precisely corresponds to the parameter  $q$  in our model.
- 143 • A speaker can speak sentences that are either helpful or hindering to learning. For example,  
144 with helpful sentences, the switching probability  $p$  can increase to  $c/\ell$ , where  $c > 1$ . In  
145 contrast, with hindering sentences, it can decrease to  $c/\ell$ , where  $c < 1$ .
- 146 • Another aspect in communication that has been considered in [6, Section 3.3] is the presence  
147 of noise. Due to the presence of noise, the sentence from a speaker might not be received by a  
148 listener, and hence the listener does not switch. The parameter  $q$  in our model can represent  
149 such noise in the communication.

150 The symmetry (SS) generalises (RWA) with overlap between languages. RWA and SS represent  
151 the simplest examples of language learning. Extension to the case of non-symmetrically overlapping  
152 languages is discussed in Supplementary Information (SI) Section 3.7.

153 *Batch learners.* The other type of the learner we consider is a powerful batch learner. A batch  
154 learner remembers all the inputs she received so far and for her hypothesis, she always selects the  
155 language that is most consistent with all her observations (initially, her memory is empty). More  
156 formally, having observed sentences  $s_1, s_2, \dots, s_n$ , the batch learner updates her hypothesis to a  
157 language  $L_i$  from her search space for which the size of the set  $L_i \cap \{s_1, s_2, \dots, s_n\}$  is maximised.  
158 We consider batch learning in the case of symmetric language overlap  $q < 1$ . That is, the size of the  
159 overlap of any  $k$  languages is equal to  $q^{k-1}$  times the size of any of the languages (see SI Section 2.2  
160 for details).

161 *The main scientific question: Rounds complexity.* While a basic question in learning theory is  
162 about identification of the correct language in the limit, an equally important question is about  
163 the efficiency of the learning process, which has been described in details in [21, Chapter 2]. The  
164 efficiency of the learning process is determined by the speed of convergence to the correct language  
165 by the whole population. The main scientific question we investigate in this work is the effect of  
166 communication structures in the learning process. More precisely, we are interested in communi-  
167 cation structures that speed up the learning process. In order to assess the efficiency of the process,  
168 we compute the expected (average) number of rounds until the process has converged (that is, all

169 learners learned the teacher’s language). We refer to this as the *rounds complexity* of the process.  
170 We discuss other relevant measures later.

171 *Illustration of the scientific question.* We illustrate our scientific question on a small example with  
172 four learners for RWA learning model of [6, Section 4.2.1]. As baseline we consider that there is  
173 no communication between the learners (denoted as the empty graph). We illustrate four possible  
174 communication structures in Figure 2. We observe that with respect to the expected number of  
175 rounds the communication structures Graph B and Graph C are worse than the empty graph,  
176 whereas the communication structure Graph D is better than the empty graph. The main take  
177 away message is: while some communication structures are worse for the learning process, others  
178 can lead to more efficient learning.

### 179 3 Results

180 Remember that  $n$  is the number of learners. We present both theoretical results and simulation re-  
181 sults. In theoretical results we introduce several communication structures (empty graph, complete  
182 graph, tree graph, Layered Hierarchy graphs). For each communication structure we analyze the  
183 rounds complexity (i.e. the expected number of rounds until all individuals have learned teacher’s  
184 language). Then we compare the rounds complexities in the limit of large  $n$ . Later we show  
185 matching numerical simulations for small  $n$ .

186 Our theoretical results are presented in terms of  $n$  and  $T$ , where  $T$  denotes the expected number  
187 of rounds in the single teacher and single learner case ( $T$  also corresponds to the sample complexity  
188 of [22]). For example, in case of single learner and RWA or SS with  $\ell$  languages we have  $T \approx c \cdot \ell$   
189 for some constant  $c > 0$ . First we consider  $(p, q)$ -learners.

190 *Remark on asymptotic complexity.* When comparing the rounds complexity of two processes  $A$  and  
191  $B$  in the limit of large population size  $n$ , the improvement can be either a *constant-factor* if the  
192 dependency on  $n$  is the same (e.g.  $A = 10 \cdot n$  vs.  $B = 5 \cdot n$ ), or *asymptotic* if the dependency on  
193  $n$  is different (e.g.  $A = 10 \cdot n$  vs.  $B = 10 \cdot \sqrt{n}$ ). In the former case we say that the asymptotic  
194 complexities match. In the latter case we say that  $B$  has better asymptotic complexity than  $A$   
195 (expression  $\sqrt{n}$  is much smaller than  $n$  for large  $n$ ). For detailed treatment see [34, Section 1.3]

196 *Classroom teaching: empty graph (Figure 3(a)).* For the baseline comparison we consider the most

197 natural extension of the single learner scenario: The empty graph consists of multiple learners who  
198 all listen to the same teacher and don't communicate among each other at all.

199 The rounds complexity is at most  $c_1 \cdot T \cdot \log n$ , where  $c_1 > 0$  is a constant (see SI Section 3.2).  
200 Hence the rounds complexity is linear in  $T$  and logarithmic in  $n$ . In particular, for RWA and SS,  
201 the upper bound is  $c_1 \cdot \ell \cdot \log n$ . Moreover, for RWA and SS, we provide matching lower bounds to  
202 show that the upper bound is optimal, and hence the upper bound cannot be improved in general.

203 *Complete graph (Figure 3(b)).* The opposite extreme is the complete graph where all learners  
204 speak to each other. Even in the simplest RWA and SS models, the complete graph has rounds  
205 complexity that is exponential in  $n$  (see SI Section 3.4). Hence it is extremely inefficient for the  
206 learning process and we will not discuss complete graphs further.

207 *Tree graph (Figure 3(c)).* Speaking to many other individuals is more demanding for the speaker. If  
208 we insist that every individual speaks to only a constant number of other individuals, we naturally  
209 obtain a tree graph. In terms of rounds complexity, the tree graph is worse than the empty graph  
210 but only by a constant factor (not asymptotically).

211 For simplicity we consider the binary tree (every individual speaks to at most two others). The  
212 vertices are organised in levels, and the teacher has level 0. Every vertex at level  $i$  has at most two  
213 incoming edges from vertices of level  $i + 1$ , and each vertex (other than the teacher) has exactly  
214 one outgoing edge. Vertices without incoming edges are called leaves. For every  $n$ , we construct  
215 a binary tree which has at most  $\log n$  levels. We show that the rounds complexity is at most  
216  $c_2 \cdot T \cdot \log n$ , where  $c_2 > 0$  is a constant (see SI Section 3.5). Hence, as for the empty graph, the  
217 dependency is linear in  $T$  and logarithmic in  $n$ . The constant  $c_2$  is greater than  $c_1$ , and thus the  
218 tree is worse than the empty graph by a constant factor, although asymptotic complexities are the  
219 same. Moreover, for RWA and SS, we establish similar lower bounds as in the case of empty graph.

220 *Layered Hierarchies.* Our most interesting results are related to certain hierarchical structures  
221 that we call *Layered Hierarchies*. We show that certain Layered Hierarchies might improve the  
222 rounds complexity, but do not improve the asymptotic complexity, whereas Layered Hierarchies  
223 with quickly growing group sizes improve even the asymptotic complexity.

224 *Description of Layered Hierarchies (Figure 3(d),(e)).* We start with a general description of Layered  
225 Hierarchies. In a  $k$ -Layered Hierarchy graph the learners are partitioned into groups (or layers)  
226  $S_1, S_2, \dots, S_k$ . The edges go from each group  $S_i$  to the previous group  $S_{i-1}$ , for  $2 \leq i \leq k$ , and

227 the edges from the first group  $S_1$  go to the teacher. An illustration of 2-Layered Hierarchy and  
 228  $k$ -Layered Hierarchy graphs are shown in Figure 3(d),(e), respectively. Incidentally, the empty  
 229 graph can be called the 1-Hierarchy. We have described the principle of Layered Hierarchy graphs  
 230 without specifying the sizes of the groups which we discuss below.

231 “*Slowly growing*” *Layered Hierarchies*. The group sizes can be of various types, and we discuss  
 232 the simple ones below: (a) *Constant size*. All group sizes are the same. (b) *Additive growth*. The  
 233 next group size is a constant more than the current group size. (c) *Multiplicative growth*. The next  
 234 group size is a constant times larger than the current group size. Let us consider the above group  
 235 sizes for three layers ( $k = 3$ ).

- 236 • *Constant size*. In this case, each group has  $n/3$  learners. In particular the first group has  
 237  $n/3$  learners, and even just considering the time to convergence for the first group, in general  
 238 the rounds complexity is at least  $c_1 \cdot T \cdot \log(n/3)$ . Thus the asymptotic complexity does not  
 239 change with respect to the empty graph.

- 240 • *Additive growth*. Let the group sizes be  $x$ ,  $2 \cdot x$ , and  $3 \cdot x$ . Since the sum of the group sizes is  
 241  $n$ , the first group size is  $n/6$ . Similarly, to the above item, in general the rounds complexity  
 242 is at least  $c_1 \cdot T \cdot \log(n/6)$ . Again the asymptotic complexity does not change with respect to  
 243 the empty graph.

- 244 • *Multiplicative growth*. Let the group sizes be  $x$ ,  $x^2$ ,  $x^3$ . Since the sum of the group sizes is  
 245  $n$ , the first group size is  $x \approx n^{1/3}$ , and similarly to the previous items, in general the rounds  
 246 complexity is at least  $c_1 \cdot T \cdot \log n^{1/3} = \frac{1}{3} \cdot c_1 \cdot T \cdot \log n$ . We observe even in this case the  
 247 asymptotic complexity does not change as compared to the empty graph.

248 We remark that even though the asymptotic complexity doesn’t change, the rounds complexity  
 249 of Layered Hierarchies is in practice often smaller than that of an empty graph by a constant factor.  
 250 The corresponding simulation results are presented in SI Section 5.2 (Figure SI.3).

251 *Exponentially growing Layered Hierarchy*. We now consider Layered Hierarchy graphs where the  
 252 group sizes grow exponentially, and show that they provide a significant asymptotic improvement  
 253 over the empty graph among learners. We start with the simpler case of Exponential 2-Layered  
 254 Hierarchy (for brevity 2-Hierarchy in the sequel), then describe the general case of Exponential  
 255 Layered Hierarchy (for brevity, Hierarchy). In the 2-Hierarchy, intuitively, the teacher quickly

256 teaches a small group of learners and then uses them as additional teachers to speed up the teaching  
 257 of the rest of the population. The Hierarchy iterates this construction. The precise descriptions  
 258 are as follows:

- 259 • *2-Hierarchy*. We split the learners into two groups  $S_1, S_2$ , where the size of  $S_1$  is proportional  
 260 to  $\log n$ , which is written as  $|S_1| \propto \log n$ . The graph then consists of all the edges from  $S_1$  to  
 261 the teacher and all the edges from  $S_2$  to  $S_1$ ; see Figure 3(d) with  $|S_1| \propto \log n$  and  $|S_2| \propto n$ .  
 262 For example, a 2-Hierarchy of 1 000 learners has  $|S_1| = 10$  and  $|S_2| = 990$ .
- 263 • *Hierarchy*. Hierarchy is obtained by iterating the construction of the 2-Hierarchy. We split  
 264 the learners into groups  $S_1, \dots, S_k$  such that the first group consists of 2 learners and that  
 265 each following group is exponentially larger than the previous group:  $|S_{i+1}| \propto 2^{|S_i|}$ . The  
 266 edges go from each group to the previous group and from the first group to the teacher; see  
 267 Figure 3(e) with  $|S_1| = 2$  and  $|S_{i+1}| \propto 2^{|S_i|}$  for  $i = 1, \dots, k - 1$ . A Hierarchy of 1 000 learners  
 268 would include 2, 4, 16, and 978 learners in the respective groups.

269 We establish the following results (see SI Section 3.6).

- 270 • For the 2-Hierarchy the expected number of rounds is at most  $c_3 \cdot T \cdot \log \log n$ , where  $c_3 > 0$  is  
 271 a constant. While the rounds complexity dependency is linear in  $T$ , the dependency is double  
 272 logarithmic in  $n$ , which is significantly better than logarithmic. Moreover, even if we interpret  
 273 dependency in  $T$ , for large  $n$ , we have  $c_1 \cdot \log n > c_3 \cdot \log \log n$ . Thus, for a reasonably large  
 274 population the 2-Hierarchy is better than the empty graph.
- 275 • For Hierarchy we show the expected number of rounds is at most  $c_4 \cdot T \cdot \log^* n$ , where  $c_4 > 0$   
 276 is a constant and  $\log^*$  (“log star”) is the iterated logarithm, which is a *very* slowly increasing  
 277 function that appears in many computer science applications. Formally,  $\log^* n$  is the number  
 278 of times the logarithm function must be iteratively applied to number  $n$  before the result is  
 279 less than or equal to 1. For any  $1 \leq n \leq 2^{256} \sim 10^{77}$  we have  $1 \leq \log^* n \leq 4$ , and thus  $\log^*(n)$   
 280 is effectively constant for all practical purposes. The Hierarchy therefore provides dramatic  
 281 improvements over the empty graph.

282 For 2-Hierarchy we again provide matching lower bounds for RWA and SS to show that the  
 283 upper bound cannot be improved in general.

284 *Remark on rounds complexity.* If we compare the empty graph and the 2-Hierarchy for RWA or SS,  
285 where the number of languages is finite and equal to  $\ell$ , for memoryless learners we obtain that the  
286 rounds complexity is proportional to  $\log n \cdot \ell$  for empty graph, and proportional to  $\log \log n \cdot \ell$  for 2-  
287 Hierarchy. Note that our results establish how the population structure influences the dependency  
288 on  $n$ . The improvement of  $\log n$  to  $\log \log n$  can be significant when  $\ell$  is large. For example, if  
289  $n = 16$ , then  $\log n$  is 4 whereas  $\log \log n$  is 2. Hence the rounds complexity decreases from  $4\ell$  to  $2\ell$ ,  
290 which can be significant speedup in practice.

291 *Other complexity measures.* The expected number of rounds (i.e. rounds complexity) is the most  
292 natural measure for the efficiency of the learning process. However, there are other relevant mea-  
293 sures which we discuss now.

294 1. The *communication complexity* is the expected number of communication events until the  
295 process converges. Each communication event represents one usage of one edge in the graph.  
296 The measure represents the total amount of sentences that need to be exchanged in the whole  
297 population.

298 2. The *bottleneck complexity* is the expected maximum number of communication events that  
299 need to be done by a single individual, which could be the teacher or one of the learners,  
300 until the process converges. If the bottleneck is the teacher then this measure relates to the  
301 amount of sentences that need to be extracted from the environment.

302 *Relevance of the complexity measures.* In distributed computing and network computation, rounds  
303 complexity is a very relevant notion, and communication complexity (or message complexity) is also  
304 well-studied [35, 36]. Typically, in distributed computing the communication structures are sym-  
305 metric and bottleneck is not widely studied, however in hierarchical network structures, bottleneck  
306 is an important complexity measure [37]. This work shows that these complexity measures from  
307 network theory become relevant for language learning in population structures, and in particular,  
308 the population structure can affect the complexity measures.

309 *Results for other complexity measures.* We now present our results for the other complexity mea-  
310 sures for the graphs we consider. We first note the following:

311 1. *Communication complexity.* The communication complexity is always the rounds complexity  
312 times the number of edges in the graph (including the edges to the teacher).

313 2. *Bottleneck complexity.* The bottleneck complexity is always the rounds complexity times the  
314 max-degree of the graph.

315 We show that the empty graph is optimal with respect to communication complexity (see SI  
316 Section 3.3). There is no graph that can be better than the empty graph for the communication  
317 complexity. The bounds for communication and bottleneck complexity for all the graphs are ob-  
318 tained from our results on rounds complexity. Note that the asymptotic communication complexity  
319 has the same dependency on  $T$  and  $n$  in all cases except for the complete graph. However, the  
320 associated constants are different, with the empty graph having the least constant among them.  
321 All the results are presented in Table 1.

322 *Discussion of the results for  $(p, q)$ -learners.* As mentioned above, the empty graph is optimal with  
323 respect to the communication complexity. The complete graph is worse in terms of all complexity  
324 measures. The tree graph matches the asymptotic complexity of the empty graph with respect  
325 to communication and rounds complexity, and improves the bottleneck complexity from  $n \log n$  to  
326  $\log n$ . The 2-Hierarchy matches the asymptotic complexity of the empty graph with respect to  
327 communication complexity, significantly improves the round complexity dependency from  $\log n$  to  
328  $\log \log n$  and improves the bottleneck complexity from  $n \log n$  to  $n \log \log n$ . The Hierarchy matches  
329 the asymptotic communication complexity of the empty graph and significantly improves the round  
330 complexity from  $\log n$  to  $\log^* n$  and the bottleneck complexity from  $n \log n$  to  $n \log^* n$ .

331 *Results for batch learners.* For batch learners under the assumption of symmetrically overlapping  
332 languages we obtain results that are similar in spirit to those for  $(p, q)$ -learners. The complete graph  
333 is much worse than the empty graph in terms of all complexity measures. The tree graph improves  
334 the bottleneck complexity as compared to the empty graph. The 2-Hierarchy graph improves both  
335 the rounds complexity and the bottleneck complexity as compared to the empty graph. The results  
336 are summarised in Table 1 (see SI Section 4 for details).

337 *Numerical simulations (Figure 4).* Our theoretical results establish asymptotic complexity bounds  
338 that apply in the limit of large population sizes. To complement them, we present numerical simula-  
339 tions for small population sizes. Since for the complete graph, the complexities grow exponentially,  
340 it is not possible to simulate the process even for small population sizes. Moreover, for small pop-  
341 ulation sizes the 2-Hierarchy and the Hierarchy coincide. Hence we present simulation results for  
342 the empty graph, the binary tree, and the 2-Hierarchy.

343 1. *Fixed  $\ell$  and varying  $n$ .* We consider  $\ell = 10$ , and vary population sizes from 10 to 1000.  
344 For each population size and graph, we run 10000 trials, and then take the average of the  
345 complexity measures. Our results are shown in Figure 4(a,d). We observe that 2-Hierarchy  
346 significantly improves over the empty graph in terms of rounds complexity.

347 2. *Fixed  $n$  and varying  $\ell$ .* In Figure 4(b,c,e,f), we present the rounds complexity for fixed  $n$  and  
348 varying  $\ell$  from 2 to 100. We use two different values of  $n$ : 30 and 100. We observe that even  
349 for  $n = 30$  the 2-Hierarchy is better than the empty graph. Thus, even for small population  
350 the 2-Hierarchy graph is better than the empty graph.

351 Furthermore, in SI Section 5.3 we present simulation results for randomly generated population  
352 structures. Random graphs do not improve the complexity measures compared to the empty graph.  
353 In SI Section 5.4 we show the full distribution of the number of rounds to fixation, comparing empty  
354 graph, the 2-Hierarchy, and the Tree graph. Therein we also present analogous simulations for the  
355 case of non-symmetric overlaps among languages.

## 356 4 Further Directions

357 There are many possible directions for further research. Here we list those related to other types  
358 of learners and models of learning (see SI Section 6 for more suggestions):

359 One direction is to consider other types of learners, presumably with intermediate capabilities  
360 as compared to memoryless  $(p, q)$ -learners and powerful batch learners. Another direction is to  
361 consider populations comprising learners of different types.

362 Yet another direction is to extend the model by defining a notion of similarity among the  
363 languages in the search space of the learners. The potential implications of such a generalization  
364 are two-fold: First, one could consider learners who, when updating their hypothesis, preferably  
365 update to a language similar either to their current language or to the language of the speaker [38].  
366 Second, instead of insisting that the learners converge to (exactly) the teacher's language, one could  
367 ask for the time to convergence to a language sufficiently similar to that of the teacher.

## 368 5 Discussion

369 A group of individuals, learning language from a teacher or from their linguistic environment,  
370 instantiate a novel evolutionary process. The learners formulate hypotheses, which get dismissed  
371 (or modified) if sentences are received that cannot be parsed. In a sense, the linguistic environment  
372 selects the correct grammar in an iterated, population based process over time. While the wrong  
373 grammars become extinct eventually, the correct grammar proliferates by eliciting copies of itself  
374 in other learners.

375 In the classical setting, the theory of learning by inductive inference considers a teacher and  
376 a learner. But here we have considered a group of learners. A new twist arises naturally: the  
377 learners not only listen to the teacher (or the environment) but also to each other. Communication  
378 between learners can be problematic, because a learner already holding the correct hypothesis can  
379 be thrown off by listening to another learner who entertains an incorrect hypothesis. We show  
380 that certain population structures increase the complexity of the overall learning task, while others  
381 reduce it. Hierarchical structures, which consist of layers of learners where each layer listens to the  
382 layer above, can be extremely efficient. Such structures might help in other types of structured  
383 cultural transmission.

384 In evolutionary graph theory, a population structure is represented by a graph, where each node  
385 is a type of an individual (such as either wild type or mutant), and the underlying evolutionary  
386 stochastic process in essence picks edges to update the type of individuals (for example in Moran  
387 process, an individual reproduces and then an edge is chosen for replacing one of its neighbours).  
388 In our scenario, each language hypothesis defines a type of the node of the graph and a stochastic  
389 process updates the language hypotheses. In evolutionary graph theory, fixation time represents  
390 the time till the population is homogeneous, which is precisely what we study as rounds complexity.

391 The process of learning language is akin to the endeavour of the scientific progress. Here nature is  
392 the teacher, natural laws are the grammatical rules, and scientists are the learners. Scientists listen  
393 to evidence from nature and also listen to each other. Sometimes scientists hold wrong hypothesis  
394 and thereby confuse others. The communication of scientific knowledge has some hierarchical  
395 structures: from scientists to science teachers to students. Our results suggest that communication  
396 between individuals, although potentially confounding, can increase the overall efficiency of the  
397 process.

## 398 **6 Methods**

399 In this section we briefly describe our key methods to establish both upper and lower bounds for  
400 the various complexity measures.

401 *Construction of graphs.* The first key step in achieving our results is the construction of the graphs.  
402 Intuitively the tree graph presents an approach of learning in different levels with distributed  
403 responsibility for teaching. The 2-Hierarchy graph is based on the intuition that we first make a  
404 small group of people learn, and then they become teachers as well. The Hierarchy extends the  
405 idea of 2-Hierarchy iteratively.

406 *Bounds for measures.* Our upper bound for  $(p, q)$ -learners on the tree graph is based on an analysis  
407 of the process and uses Chernoff bound [39]. For the 2-Hierarchy and Hierarchy graphs, the principle  
408 is that once a group learns the language of the teacher, it teaches the next group. For every group of  
409 learners, we define its *phase* as lasting from the moment everyone in all the previous groups speaks  
410 the right language until everyone in that group also speaks the right language. We establish the  
411 number of rounds each phase takes and obtain the desired result by summing over all the groups.  
412 For batch learners, we proceed similarly. See SI for details.

413 *Lower bound.* The most interesting lower bound we establish is on the communication complexity,  
414 as we derive all other lower bounds from it. We actually show that for  $(p, q)$ -learners, no graph  
415 can achieve a communication complexity better than  $c \cdot n \log n$ , for some constant  $c > 0$ . For the  
416 result we use a coupling argument [40] to compare an arbitrary graph with the empty graph, and  
417 use Markov's inequality [39].

## 418 **Author contributions**

419 R.I.-J., J.T., K.C., and M.A.N. designed research, performed research, and wrote the paper.

## 420 **Data Accessibility**

421 Data and scripts for plotting figures have been uploaded as part of the electronic supplementary  
422 material.

## 423 **Funding Statement**

424 A.P., J.T. and K.C. acknowledge support from ERC Start grant no. (279307: Graph Games),  
425 Austrian Science Fund (FWF) grant no. P23499-N23 and S11407-N23 (RiSE). M.A.N. acknowl-  
426 edges support from Office of Naval Research grant N00014-16-1-2914 and from the John Templeton  
427 Foundation. The Program for Evolutionary Dynamics is supported in part by a gift from B. Wu  
428 and E. Larson.

## 429 **Competing interests**

430 We have no competing interests.

## 431 **References**

- 432 [1] Lightfoot D. The development of language: Acquisition, change, and evolution. Wiley-  
433 Blackwell; 1999.
- 434 [2] Wexler K, Culicover P. Formal principles of language acquisition. MIT Press; 1980.
- 435 [3] Smith JM. Evolution and the Theory of Games. Cambridge university press; 1982.
- 436 [4] Komarova NL, Niyogi P, Nowak MA. The evolutionary dynamics of grammar acquisition.  
437 Journal of theoretical biology. 2001;209(1):43–59.
- 438 [5] Nowak MA. Evolutionary dynamics. Harvard University Press; 2006.
- 439 [6] Niyogi P. The computational nature of language learning and evolution. MIT press Cambridge,  
440 MA.; 2006.
- 441 [7] Nowak MA, Komarova NL, Niyogi P. Computational and evolutionary aspects of language.  
442 Nature. 2002;417(6889):611–617.
- 443 [8] Chomsky N, DiNozzi R. Language and mind. Harcourt Brace Jovanovich New York; 1972.
- 444 [9] Jain S, Osherson D, Royer JS, Sharma A. Systems That Learn: An Introduction to Learning  
445 Theory (Learning, Development and Conceptual Change). 2nd ed. The MIT Press; 1999.

- 446 [10] Vapnik VN, Vapnik V. Statistical learning theory. vol. 1. Wiley New York; 1998.
- 447 [11] Gold EM. Language identification in the limit. *Information and control*. 1967;10(5):447–474.
- 448 [12] Osherson DN, Stob M, Weinstein S. *Systems that learn: An introduction to learning theory*  
449 *for cognitive and computer scientists*. The MIT Press; 1986.
- 450 [13] Pinker S. Formal models of language learning. *Cognition*. 1979;7(3):217–283.
- 451 [14] Niyogi P, Berwick RC. A language learning model for finite parameter spaces. *Cognition*.  
452 1996;61(1):161–193.
- 453 [15] Osherson DN, Stob M, Weinstein S. Learning theory and natural language. *Cognition*.  
454 1984;17(1):1–28.
- 455 [16] Case J, Moelius Iii SE. Optimal language learning. In: *Algorithmic Learning Theory*. Springer;  
456 2008. p. 419–433.
- 457 [17] Heinz J, Kasprzik A, Kötzing T. Learning in the limit with lattice-structured hypothesis  
458 spaces. *Theoretical Computer Science*. 2012;457:111–127.
- 459 [18] Chomsky N. Principles and parameters in syntactic theory. *Explanation in linguistics: The*  
460 *logical problem of language acquisition*. 1981;32:75.
- 461 [19] Yang CD. *Knowledge and learning in natural language*. Oxford University Press on Demand;  
462 2002.
- 463 [20] De la Higuera C. *Grammatical inference: learning automata and grammars*. Cambridge  
464 University Press; 2010.
- 465 [21] Heinz J, Sempere JM. *Topics in grammatical inference*. Springer; 2016.
- 466 [22] Zeugmann T. From learning in the limit to stochastic finite learning. *Theoretical Computer*  
467 *Science*. 2006;364(1):77–97.
- 468 [23] Nowak MA, Komarova NL, Niyogi P. Evolution of universal grammar. *Science*.  
469 2001;291(5501):114–118.
- 470 [24] Komarova N, Rivin I. *Mathematics of learning*. arXiv preprint math/0105235. 2001;.

- 471 [25] Christiansen MH, Dale RA, Ellefson MR, Conway CM. The role of sequential learning in  
472 language evolution: Computational and experimental studies. In: *Simulating the evolution of*  
473 *language*. Springer; 2002. p. 165–187.
- 474 [26] Komarova NL, Nowak MA. Language dynamics in finite populations. *Journal of Theoretical*  
475 *Biology*. 2003;221(3):445–457.
- 476 [27] Lee Y, Stabler TCCEP, Taylor CE. The role of population structure in language evolution.  
477 *language*. 2005;22(23):24–25.
- 478 [28] Nowak MA, Krakauer DC. The evolution of language. *Proceedings of the National Academy*  
479 *of Sciences*. 1999;96(14):8028–8033.
- 480 [29] Stabler EP. Mathematics of language learning. *Histoire Épistémologie Langage*.  
481 2009;31(1):127–145.
- 482 [30] Niyogi P, Berwick RC. Evolutionary consequences of language learning. *Linguistics and*  
483 *Philosophy*. 1997;20(6):697–719.
- 484 [31] Niyogi P, Berwick RC. The proper treatment of language acquisition and change in a popula-  
485 tion setting. *Proceedings of the National Academy of Sciences*. 2009;106(25):10124–10129.
- 486 [32] Kirby S. Spontaneous evolution of linguistic structure-an iterated learning model of the  
487 emergence of regularity and irregularity. *Evolutionary Computation, IEEE Transactions on*.  
488 2001;5(2):102–110.
- 489 [33] Clark R, Roberts I. A computational model of language learnability and language change.  
490 *Linguistic Inquiry*. 1993;24(2):299–345.
- 491 [34] Cormen TH. *Introduction to algorithms*. MIT press; 2009.
- 492 [35] Attiya H, Welch J. *Distributed computing: fundamentals, simulations, and advanced topics*.  
493 vol. 19. John Wiley & Sons; 2004.
- 494 [36] Lynch NA. *Distributed algorithms*. Morgan Kaufmann; 1996.
- 495 [37] Tanenbaum AS, Wetherall D. *Computer networks*. Prentice hall; 1996.

- 496 [38] Bryden J, Wright SP, Jansen VA. How humans transmit language: horizontal transmission  
 497 matches word frequencies among peers on Twitter. *Journal of The Royal Society Interface*.  
 498 2018;15(139):20170738.
- 499 [39] Mitzenmacher M, Upfal E. *Probability and computing: Randomized algorithms and proba-*  
 500 *bilistic analysis*. Cambridge University Press; 2005.
- 501 [40] Lindvall T. *Lectures on the coupling method*. Courier Corporation; 2002.

## 502 Figure and table captions

503 **Figure 1. A teacher and a group of learners.** The teacher is represented as a square and  
 504 learners as circles. Individuals whose hypothesis is the teacher’s language  $L_1$  are shown in red,  
 505 others in blue (teacher is always red). Possible communications are indicated by edges. When an  
 506 edge is selected for the communication event, it is shown in bold. **(a)** An illustration of  $(p, q)$ -  
 507 learning. In one step of the learning process, we select an edge (indicated in bold) and then the  
 508 listener of that edge updates their language hypothesis. (i) Learner  $X$  listens to the teacher and  
 509 switches to the teacher’s language with probability  $p$ . (ii) Learner  $Y$  already has the same language  
 510 as the teacher, but due to listening to a learner  $X$  who speaks a ‘wrong’ language,  $Y$  switches with  
 511 probability  $1 - q$  to a (possibly different) wrong language. **(b)** An illustration of one possible run of  
 512 a single round as described in the paragraph Example. Population structure consists of a teacher,  
 513 Alice, and Bob. There are two non-overlapping languages  $L_1, L_2$ . When a learner hears a sentence  
 514 they don’t understand, they switch their hypothesis to the other language with probability 80 %  
 515 (and keep it otherwise). We picked the edges in order  $B \rightarrow T, B \rightarrow A, A \rightarrow T$ . In the second step,  
 516  $B$  switched from correct  $L_1$  to incorrect  $L_2$ .

517 **Figure 2. Simulations for small graphs.** **(a)** Four distinct structures of the class room, each  
 518 with one teacher and four learners. Note that Graph  $A$  is the ‘empty graph’ because there are no  
 519 communications between the learners. **(b)** Simulation results for these four graphs showing the  
 520 average number of rounds that are needed for all learners to converge to the correct language versus  
 521 the number of languages  $\ell$  in the search space. Here we consider  $(p, q)$ -learners with  $p = q = 1/\ell$ .  
 522 Each point is an average over 100 000 trials. In each round, the communication happens along each  
 523 edge once, in random order. Graphs  $B$  and  $C$  are much worse than the empty graph,  $A$ , but graph

524  $D$  is faster. This simple example shows that communication between learners can both accelerate  
525 and decelerate the process.

526 **Figure 3. Different population structures of language learning.** The teacher is shown in  
527 red and the learners in blue. **(a)** The empty graph represents the case where learners only listen to  
528 the teacher and do not communicate with each other. **(b)** The opposite extreme is the complete  
529 graph where all possible communications between learners are realized. **(c)** In the tree graph with  
530 branching factor  $k = 2$ , the teacher speaks to two learners, who each speak to two learners and so  
531 on. **(d, e)** The 2-Layered Hierarchy and the  $k$ -Layered Hierarchy consist of layers such that each  
532 learner from a given layer listens to all individuals from the previous layer. In the special case of  
533 Exponentially growing Layered Hierarchies (2-Hierarchy and Hierarchy), each layer is exponentially  
534 bigger than the previous one.

535 **Figure 4. Numerical simulation results.** The colours represent different graph families: Blue:  
536 Empty graph; Orange: 2-Hierarchy; Green: Tree graph. The empty graphs is shown in bold since  
537 it is the baseline comparison. First, we consider memoryless learners with helpful teacher, that  
538 is  $p = 2/\ell$ ,  $q = 1/\ell$  **(a)** Rounds complexity against the population size  $n$ , for fixed number of  
539 languages  $\ell = 10$ . For empty graph the dependency on  $n$  is logarithmic, for tree graph it is also  
540 logarithmic but worse by a constant factor, and for the 2-Hierarchy graph it is asymptotically better  
541 (namely doubly logarithmic). **(b), (c)**, Rounds complexity against the number of languages  $\ell$ , for  
542 fixed population size  $n = 30$  and  $n = 100$ . The 2-Hierarchy beats the empty graph in both cases.  
543 Since the dependency on  $\ell$  in all cases is linear, any value of  $\ell$  would yield analogous outcome in  
544 **(a)**. **(d), (e), (f)** Similar plots for batch learners under symmetric language overlap  $q = 0.1$ . **(d)**  
545 Rounds complexity against the population size  $n$ , for fixed number of languages  $\ell = 10$ . As in **(a)**,  
546 for the empty graph the dependency is logarithmic whereas for the 2-Hierarchy it is asymptotically  
547 better. However, for tree graph the dependency is linear in  $n$ . **(e), (f)** This time the dependency  
548 on  $\ell$  is logarithmic in all cases (batch learners are more powerful than memoryless learners). All  
549 the values shown are averages over 10 000 trials.

550 **Table 1. Complexity bounds for language learning.** The tables show the various complexity  
551 measures for different graphs as function of population size,  $n$ , and expected time to teach one  
552 learner in a single teacher single learner model,  $T$ . The first table refers to  $(p, q)$ -learners, the second  
553 table refers to batch learners under symmetric language overlap. Rounds complexity denotes the

554 average number of rounds until all learners hold the correct grammar. Communication complexity  
555 denotes the average number of communications until this state is reached and bottleneck complexity  
556 denotes the average maximum number of communications produced from a single person. There  
557 exist constants  $c_1, c_2, c_3, c_4$  such that the complexity measures are lower bounded by the expressions.  
558 Except for batch learners on tree graphs, all bounds are tight up to a constant, which means there  
559 exist positive constants for which the corresponding expressions are upper bounds. The expression  
560  $\log^* n$  denotes the iterated logarithm of  $n$  (see text).