

1     **A STOCHASTIC MODEL OF LANGUAGE CHANGE THROUGH**  
2     **SOCIAL STRUCTURE AND PREDICTION-DRIVEN INSTABILITY**

3                     W. GARRETT MITCHENER\*

4     **Abstract.** Children reliably learn their community’s language; consequently human languages  
5 are relatively stable on short time scales. However, languages can change dramatically over the course  
6 of centuries, and once begun, such changes generally run monotonically to completion. We consider  
7 a stochastic model that reproduces this pattern of fluctuations via large deviations. We begin with  
8 a Markov chain that represents an age-structured population in which children learn which of two  
9 grammars their community prefers, but are aware of age-correlated usage patterns and will use the  
10 dispreferred grammar more often if they infer that its use is spreading. The Markov chain is shown  
11 to converge in the limit of an infinite population to a stochastic differential equation that generalizes  
12 the Wright-Fisher SDE for population genetics. This proof is not routine because the dynamics are  
13 only defined in a Cartesian product of simplexes, and it must be verified that trajectories of the  
14 SDE cannot escape. Results are proved by changing variables in a way that expands each simplex  
15 to an entire plane, yielding reasonable constraints on the dynamics that ensure that a standard but  
16 sophisticated theorem for well-posedness of SDEs can be applied. The SDE yields a phase portrait  
17 that reveals the mechanism that causes these models to produce sporadic, monotone, population-  
18 wide transitions between grammars. A further simplification results in a stochastic functional-delay  
19 differential equation that shows how population-level memory effects and the attempt by learners to  
20 avoid sounding outdated results in prediction-driven instability.

21     **Key words.** language change, prediction-driven instability, population dynamics, stochastic  
22 differential equation, noise-activated transitions

23     **AMS subject classifications.** 37H10, 60H10, 60J20, 91F20

24     **1. The paradox of language change.** A primary tool in the field of linguistics  
25 is the *idealized grammar*, that is, a formalism that distinguishes correctly formed  
26 utterances from ill-formed utterances [8, 17]. Historically, much of the research on  
27 how children acquire their native language has focused on how they might choose one  
28 idealized grammar from many innate possibilities on the basis of example sentences  
29 from the surrounding society [1, 57, 60]. From the perspective of idealized grammar,  
30 language change is paradoxical: Children acquire their native language accurately and  
31 communicate with adults from preceding generations, yet over time, the language can  
32 change drastically. Some changes may be attributed to an external event, such as  
33 political upheaval, but not every instance of language change seems to have an external  
34 cause. Despite their variability, languages do maintain considerable short-term  
35 stability, consistently accepting and rejecting large classes of sentences for centuries.  
36 The primary challenge addressed by the model discussed in this article is to capture  
37 this meta-stability.

38     Many existing models of language learning in a population focus on characterizing  
39 stable properties of languages. For example, naming games and other lexical models focus on the  
40 process by which a population forms a permanent consensus on a vocabulary, and how effective  
41 that vocabulary is at representing meanings [9, 23, 49, 50, 51, 56, 61]. Related models focus on  
42 the structure of lexeme or phoneme inventories once a stable equilibrium is reached [22, 21, 66].  
43 Several algorithms have been proposed as models for the acquisition of idealized grammars  
44 [5, 6, 7, 14, 21, 22, 44, 43, 60]. These focus on details of the acquisition process and follow  
45 the *probably almost correct* (PAC) learning framework [15], in which the learner’s input is  
46 a list of grammatically correct utterances called the *primary lin-*

---

\*Department of Mathematics, College of Charleston, SC (MitchenerG@cofc.edu)

48 *guistic data* (PLD), and the learner is required to choose a single idealized grammar  
 49 from a limited set that is somehow maximally consistent with that input. Learners  
 50 are given no data on social structure, or any negative evidence, that is, information  
 51 stating that a possible utterance is ungrammatical. The input may be from a single  
 52 individual [21, 22] or a population, perhaps consisting of adults that collectively use  
 53 several idealized grammars [43]. Other proposed algorithms address the sensitivity  
 54 of the PAC framework to noise in the PLD by including means of ignoring rarely  
 55 occurring constructions [28, 29, 54, 65].

56 There is room to improve on these models. In contrast to actual human language,  
 57 these models typically have stable equilibrium states from which the learner  
 58 or population cannot escape. Furthermore, PAC learning algorithms typically make  
 59 use of the *subset principle*: Out of all the available idealized grammars, the “correct”  
 60 choice is the one that generates the smallest set of utterances including the input. The  
 61 subset principle is frequently included in language learning models because children  
 62 typically ignore assertions by adults that a particular utterance is ungrammatical.  
 63 However, there is evidence that the subset principle does not accurately reproduce  
 64 certain features of child language acquisition, and that children make use of statistical  
 65 patterns in adult speech to determine that utterances they previously accepted are  
 66 actually ungrammatical [39, 4].

67 Many language models for populations are adapted from deterministic, continu-  
 68 ous, biological population models and represent language by communication games.  
 69 These focus on stable behavior in an infinite homogeneous population, although some  
 70 exhibit ongoing fluctuations [40, 33, 41, 34, 35, 37, 48, 46, 47, 45, 53]. Some are  
 71 designed to represent a single change [24]. In these models, children learn from an  
 72 average of speech patterns, and except for [41], these do not model the origins of  
 73 language changes directly. Instead, an external event must disturb the system and  
 74 push it from one stable state to another.

75 As we will see in [section 2](#), a general mean-field model in which children learn  
 76 from the entire population equally does not lead to spontaneous change, even in the  
 77 presence of random variation. It appears that spontaneous changes can only arise  
 78 from random fluctuations in combination with some sort of momentum driven by  
 79 social structure.

80 Based on extensive field studies, Labov [26] proposes a model in which phonetic  
 81 change is driven by females who naturally change their individual speech over time,  
 82 a force called *incrementation*. A semi-structured approach as in [36] assumes a fully  
 83 interconnected finite population but agents vary in their influence on learners. These  
 84 models approximate the time course of a single change, in qualitative agreement with  
 85 data, but neither addresses the origin of the change.

86 Some models use network dynamics rather than a mean-field assumption and  
 87 allow learners to collect input disproportionately from nearby members of the popu-  
 88 lation [12, 59]. These models incorporate observations made by Labov and others that  
 89 certain individuals tend to lead the population in adopting a new language variant,  
 90 and the change spreads along the friendship network [25, 26, 27].

91 In contrast, the model analyzed in this article is built from an alternative perspec-  
 92 tive in an attempt to resolve the language change paradox. Utterances may be drawn  
 93 from multiple idealized grammars and classified as more or less archaic or innovative.  
 94 Such an approach can consider the variation present in natural speech and model it  
 95 as a *stochastic grammar*, that is, a collection of similar idealized grammars, each of  
 96 which is used randomly at a particular rate [24, 25, 26, 62]. From this continuous  
 97 perspective, language change is no longer a paradox, but acquisition requires more

98 than selecting a single idealized grammar as in the PAC framework. Instead, children  
 99 must learn multiple idealized grammars, plus the usage rates and whatever conditions  
 100 affect them.

101 Crucially, instead of limiting learners' input to example sentences, we will assume  
 102 that children also know something about the ages of speakers and prefer not to sound  
 103 outdated. They bias their speech against archaic forms by incorporating a prediction  
 104 step into their acquisition of a stochastic grammar, which introduces incrementation  
 105 without directly imposing it as in [26]. The age structure and bias against archaic  
 106 forms introduce momentum into the dynamics, which generates the desired meta-  
 107 stability. The population tends to hover near a state where one idealized grammar  
 108 is highly preferred. However, children occasionally detect accidental correlations be-  
 109 tween age and speech, predict that the population is undergoing a language change,  
 110 and accelerate the change. This feature will be called *prediction-driven instability*.

111 The majority of the language modeling literature does not focus on the formal  
 112 aspects of mathematical models, such as confirming that the dynamics are well-posed  
 113 or deriving a continuous model as the limit of a discrete model, even though such  
 114 details are known to be generally important [11]. Numerical simulations of the discrete  
 115 form of the age-structured stochastic model developed in this article confirm that it  
 116 has the desired behavior [38] but its continuous form has yet to be placed on a sound  
 117 theoretical foundation. So, in section 2 we formulate a discrete mean-field model as a  
 118 Markov chain and discuss its weaknesses. Then in section 3, we extend it to include  
 119 age-structure, then rigorously consider the limit of an infinitely large population and  
 120 reformulate the Markov chain as a continuous-time martingale problem.

121 We rewrite this martingale problem as a system of stochastic differential equations  
 122 (SDEs), show that it has a unique solution for all initial values, and show that paths  
 123 of the Markov chain converge weakly to solutions of the SDEs. The proofs make use of  
 124 theorems in [10] for the existence and uniqueness of solutions to SDEs and convergence  
 125 of discrete Markov chains to such solutions. However, the SDEs of interest take  
 126 values in a phase space consisting of Cartesian products of simplexes, and changes-  
 127 of-variables are required to derive SDEs taking values in a plane as required by the  
 128 standard theorems. Furthermore, the drift and volatility terms in the resulting SDEs  
 129 grow too quickly in magnitude at infinity for the most commonly used theorems to  
 130 be directly applied. Instead, asymptotic estimates must be used to verify that the  
 131 drift terms push solutions back toward the origin, in which case a more general result  
 132 presented in [10] guarantees the existence of unique solutions for all time. These  
 133 results confirm that solutions to the SDEs are at no risk of straying into unrealistic  
 134 territory where the usage rate of some grammar has escaped from  $[0, 1]$ . Furthermore,  
 135 they make minimal assumptions about the vector field and are applicable to other  
 136 dynamical systems on simplexes.

137 In the two dimensional case, in which agents use one grammar or the other exclu-  
 138 sively, it is possible to see in the phase portrait that proximity of stable equilibria to  
 139 the boundaries of their basins of attraction is what facilitates spontaneous language  
 140 change. A final modification to the two-dimensional SDEs allows them to be reformu-  
 141 lated as a one-dimensional functional-delay SDE. In this form, it becomes clear that  
 142 the population switches from one meta-stable state to another when children detect a  
 143 chance fluctuation in the usage rate of the dominant grammar away from the running  
 144 average, and amplify it.

145 **2. First stage: An unstructured mean-field model.** Let us suppose initially  
 146 that individuals have a choice between two similar idealized grammars  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .

147 Each simulated agent uses  $\mathcal{G}_2$  in forming an individual-specific fraction of spoken  
 148 sentences, and  $\mathcal{G}_1$  in forming the rest. Assume that children are always able to acquire  
 149 both idealized grammars and the only challenge is learning the usage rates. Assume  
 150 that the population consists of  $N$  adult agents, each of which is one of  $K + 1$  types,  
 151 numbered 0 to  $K$ , where type  $k$  means that the individual uses  $\mathcal{G}_2$  at a rate  $k/K$  and  
 152  $\mathcal{G}_1$  at a rate  $1 - k/K$ . The state of the chain at time step  $j$  is a vector  $T$  where  $T_n(j)$   
 153 is the type of the  $n$ -th agent. Define the count vector  $C$  where  $C_k(j)$  is the number  
 154 of agents of type  $k$ ,

$$155 \quad C_k(j) = \sum_n \mathbf{1}(T_n(j) = k).$$

156 Dividing the count vector by the population size yields the speech distribution vector  
 157  $X = C/N$  such that an agent selected at random from the population uniformly at  
 158 time  $j$  is of type  $k$  with probability  $X_k(j)$ .

159 The mean usage rate of  $\mathcal{G}_2$  at step  $j$  is therefore

$$160 \quad (1) \quad M(j) = \sum_{k=0}^K \left( \frac{k}{K} \right) X_k(j)$$

161 Children are assumed to learn the usage rates of the two grammars based only on  
 162  $M(j)$ , the mean usage rate of  $\mathcal{G}_2$  in the adult population at time  $j$ . Children are  
 163 assumed to be exposed to enough sample utterances from across the entire population  
 164 to accurately estimate  $M(j)$ . The model requires a *mean learning function*  $q(m)$  that  
 165 gives the mean usage rate of children learning from a population with a mean rate  $m$ .

166 The transition process from step  $j$  to  $j+1$  is as follows. Two additional parameters  
 167 are required, a birth-and-death rate  $r_D$  and a resampling rate  $r_R$ . At each time step,  
 168 each individual agent is examined and one of these three operations is randomly  
 169 applied to it:

- 170 • With probability  $p_D = r_D/N$  it dies and is replaced.
- 171 • With probability  $r_R$  it is resampled.
- 172 • With probability  $1 - p_D - r_R$  it is unchanged.

173 Details are given in the following subsections.

174 **2.1. Time, learning, and the birth-death operation.** Each time step is  
 175 interpreted as  $1/N$  years. The lifespan of an individual in time steps has a geometric  
 176 distribution with parameter  $p_D$ . The average life span is therefore  $1/p_D$  time steps  
 177 or  $1/r_D$  years.

178 When an agent dies, a replacement agent is created and its type is selected at  
 179 random based on a discrete distribution vector  $Q(M(j))$ . That is,  $Q_k(m)$  is the  
 180 probability that a child learning from a population with mean usage rate  $m$  is of type  
 181  $k$ , and therefore uses  $\mathcal{G}_2$  at rate  $k/K$ . As a specific example,  $Q(m)$  could be the mass  
 182 function for a binomial distribution with parameters  $q(m)$  and  $K$ ,

$$183 \quad (2) \quad Q_k(m) = \binom{K}{k} q(m)^k (1 - q(m))^{K-k}.$$

184 Since the mean of such a distribution is  $q(m)K$ , it follows that  $q$  and  $Q$  satisfy the  
 185 identity

$$186 \quad (3) \quad q(m) = \sum_{k=0}^K \left( \frac{k}{K} \right) Q_k(m)$$

187 which confirms that  $q(m)$  is indeed the mean usage rate of  $\mathcal{G}_2$  by children learning  
 188 from adults with mean usage rate  $m$ .

189 The mean learning function must be S-shaped to ensure that there are two equi-  
 190 librium states, representing populations dominated by one grammar or the other. In  
 191 general,  $q$  is assumed to be smooth, strictly increasing, with one inflection point, and

$$192 \quad (4) \quad \begin{aligned} &0 < q(0) < 1/2 \\ &1/2 < q(1) < 1 \end{aligned}$$

193 In practice,  $q(0)$  will be close to 0 and  $q(1)$  will be close to 1. A curved mean learning  
 194 function means that the more commonly used idealized grammar becomes even more  
 195 commonly used, until the other grammar all but disappears. This tendency is in  
 196 agreement with the observation that children regularize language: A growing body  
 197 of evidence [19] indicates that for the task of learning a language with multiple ways  
 198 to say something, adults tend to use all the options and match the usage rates in the  
 199 given data, but children prefer to pick one option and stick with it. Beyond these  
 200 general properties, this learning model makes no attempt to directly represent the  
 201 neurological details of language acquisition, although researchers are exploring this  
 202 area [2, 20, 39, 55, 65].

203 **2.2. Resampling of adults.** When an agent is resampled, its new state is copied  
 204 from another agent picked uniformly at random. The average time an agent spends  
 205 between resamplings is  $1/r_R$  time steps. This feature of the transition process incor-  
 206 porates the fact that as an adult, an individual’s language is not entirely fixed [26, 25].  
 207 Furthermore, as will be explained in section 4, without this resampling feature, the  
 208 random fluctuations of this Markov chain diminish to 0 in the limit as  $N \rightarrow \infty$ , which  
 209 would defeat the purpose of developing a stochastic model. This consideration leads  
 210 to the peculiar fact that in formulating the Markov chain,  $p_D$  must scale as  $1/N$  but  
 211 the probability  $r_R$  of an agent being resampled must remain constant. The Wright-  
 212 Fisher model [10] includes a similar feature: In the discrete formulation, each time  
 213 step is considered a single generation and each agent is always resampled, akin to  
 214 setting  $r_R = 1$ , but when passing to the limit  $N \rightarrow \infty$ , the generation time is taken  
 215 to scale as  $1/N$  without scaling the resampling process.

216 It is possible that in contrast to standard practice in the population genetics  
 217 literature,  $r_D$  should also scale as  $1/N$ . That would cause fluctuations in grammar  
 218 use to shrink as the population size grows, in agreement with anecdotal reports that  
 219 languages spoken by only a small number of native speakers change rapidly compared  
 220 to those with larger populations, but in disagreement with other studies [63, 64].  
 221 Resolution of this issue is beyond the scope of this article.

222 **2.3. Behavior of the model.** This Markov chain is regular. Although it spends  
 223 most of its time hovering near a state dominated by one idealized grammar or the  
 224 other, it must eventually exhibit spontaneous language change by switching to the  
 225 other. However, computer experiments confirm that under this model, a population  
 226 takes an enormous amount of time to switch dominant grammars. This model is  
 227 therefore unsuitable for understanding language change on historical time scales. A  
 228 further undesirable property is that when a population does manage to shift to an  
 229 intermediate state, it is just as likely to return to the original grammar as to complete  
 230 the shift to the other grammar. Historical studies [24, 65] show that language changes  
 231 typically run to completion monotonically and do not reverse themselves partway  
 232 through (but see [62] for some evidence to the contrary), so again this model is  
 233 unsatisfactory.

234 **3. Second stage: An age-structured model.** One way to remedy the weak-  
 235 nesses of these mean-field models is to introduce social structure into the population.  
 236 According to sociolinguistics, ongoing language change is reflected in variation, so  
 237 there is reason to believe children are aware of socially correlated speech variation  
 238 and use it during acquisition [25].

239 There are many ways to formulate a socially structured population, and not all  
 240 formulations apply to all societies. For this article, let us assume that there are  
 241 two age groups, roughly representing youth and their parents, and that children can  
 242 detect systematic differences in their speech. We also assume that there are social  
 243 forces leading children to avoid sounding out-dated.

244 Let us adapt the Markov chain from section 2 to include age structure. To rep-  
 245 resent the population at time  $j$ , fix the total number of youth and the total number  
 246 of parents at  $N$ , so there are  $2N$  agents total. To make the notation systematic,  
 247 superscript labels  $Y$  and  $A$  will be used, referring to the youth and adult generations,  
 248 respectively. Let  $T_n^Y(j)$  be the type of the  $n$ -th youth and  $T_n^A(j)$  be the type of the  
 249  $n$ -th adult, all between 0 and  $K$ . Define  $C_k^Y(j)$  to be the number of youth of type  $k$ ,  
 250 and define  $C_k^A(j)$  to be the number of adults of type  $k$ . Let

$$251 \quad (5) \quad X^Y = \frac{1}{N}C^Y \quad \text{and} \quad X^A = \frac{1}{N}C^A$$

252 be the probability distribution vectors of the two generations. Assume that apart from  
 253 age, children make no distinction among individuals. Thus, they learn essentially from  
 254 the mean usage rates of the two generations,

$$255 \quad (6) \quad \begin{aligned} M^Y(j) &= \sum_{k=0}^K \left(\frac{k}{K}\right) X_k^Y(j) \\ M^A(j) &= \sum_{k=0}^K \left(\frac{k}{K}\right) X_k^A(j) \end{aligned}$$

256 The modified transition process from time  $j$  to  $j + 1$  is as follows. Each adult is  
 257 examined:

- 258 • With probability  $p_D = r_D/N$  it is replaced to simulate death and aging.
- 259 • With probability  $r_R$  it is resampled from the adult population.
- 260 • With probability  $1 - p_D - r_R$  it is unchanged.

261 Each youth is examined:

- 262 • With probability  $p_D = r_D/N$  it is replaced to simulate birth and aging.
- 263 • With probability  $r_R$  it is resampled from the youth population.
- 264 • With probability  $1 - p_D - r_R$  it is unchanged.

265 Each time step is interpreted as  $1/N$  years. The number of time steps spent by an  
 266 individual in each age group has a geometric distribution with parameter  $p_D$ . The  
 267 average time spent as an adult and as a youth is therefore  $1/p_D$  time steps or  $1/r_D$   
 268 years, so the average life span is now  $2/r_D$ .

269 When an agent is resampled, its new state is copied from another agent from the  
 270 same generation selected uniformly at random. As before, resampling leaves the mean  
 271 behavior unchanged while introducing volatility.

272 It is certainly possible to incorporate birth, aging, and death into the model by  
 273 deleting an adult, directly moving someone from the youth generation to the adult  
 274 generation, and creating a new youth. However, the calculations are simplified if birth

275 and death are handled separately, resulting in mathematically trivial differences to  
276 the Markov chain.

277 When an adult dies, rather than moving a youth, a replacement is created by  
278 sampling from an aging distribution  $V(X^Y)$ , that is very close to  $X^Y$  but gives  
279 at least a minimal probability to every type. This feature allows for innovation in  
280 adults, and avoids a technical problem that would cause the model to fall outside the  
281 hypotheses of [Lemma 4.6](#). The examples in this article use

$$282 \quad (7) \quad V_k(X) = X_k(1 - (K + 1)\eta) + \eta$$

283 with  $\eta = 1/1000$ .

284 For birth and aging, a randomly selected youth is deleted, and a replacement  
285 youth is created based on the discrete probability vector  $R(M^Y(j), M^A(j))$ . Here,  
286  $R(x, y)$  represents the acquisition process, together with prediction: Children hear  
287 that the younger generation uses  $\mathcal{G}_2$  at a rate  $x$ , and the older generation uses a rate  
288  $y$ . Based on  $x$  and  $y$  and any trend those numbers indicate, they predict a rate that  
289 their generation should use, and learn based on that predicted target value. Let the  
290 predicted mean usage rate be given by a smooth function  $r(x, y)$  that is increasing  
291 with respect to  $x$ , decreasing with respect to  $y$ , and satisfies

$$292 \quad \forall x, y: \quad y < x \implies x < r(x, y)$$

293 and

$$294 \quad \forall x, y: \quad y > x \implies x > r(x, y).$$

295 That is, any trend from the past  $y$  compared to the present  $x$  should continue to  
296 the future  $r(x, y)$ . Then, our assumptions on learning based on prediction can be  
297 incorporated into the mathematics by setting  $R(x, y) = Q(r(x, y))$ .

298 For a specific example, let us consider a population of 1000 agents, 500 in each  
299 age group, with a birth-death rate of  $r_D = 1/20$ . Therefore, the mean lifespan of  
300 an agent is 40 years. The resampling rate is  $r_R = 0.0001$ . There are 6 types of  
301 agents, representing speech patterns that use  $\mathcal{G}_2$  for a fraction  $0, 1/5, \dots, 1$  of spoken  
302 sentences.

303 The learning distribution  $Q(m)$  is a binomial distribution with parameters  $q(m)$   
304 and 5. The example  $q$  in this article is

$$305 \quad (8) \quad q(m) = \frac{1}{32} + \frac{3600}{751} \left( \frac{33m}{1280} + \frac{161m^2}{320} - \frac{m^3}{3} \right)$$

306 This polynomial was constructed to be slightly asymmetric and strictly increasing on  
307  $[0, 1]$ . Its range is  $[1/32, 31/32]$ , so it satisfies [\(4\)](#) and conditions that will be needed  
308 to apply [Proposition 4.1](#).

309 The example prediction function  $r(x, y)$  is based on an exponential sigmoid. Given  
310  $s(t) = 1/(1 + \exp(-t))$ , define  $t_1 = s^{-1}(x)$  and  $t_2 = s^{-1}(y)$ . Then  $h = t_1 - t_2$  is a  
311 measure of the trend between the generations. A scale factor  $\alpha$  is applied to  $h$ , and  
312 the scaled trend is added to  $t_1$ . After some simplification,

$$313 \quad (9) \quad r_0(x, y) = s(t_1 + \alpha h) = \frac{1}{1 + \left(\frac{1-x}{x}\right)^{\alpha+1} \left(\frac{y}{1-y}\right)^\alpha}$$

314 For the example calculations in this paper,  $\alpha = 3$ . Observe that  $r_0$  is a rational  
315 function, defined and continuous everywhere in  $[0, 1] \times [0, 1]$  except at the corners.

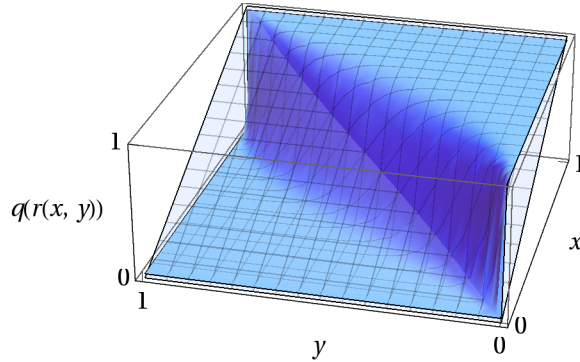


FIG. 1. The learning-prediction function  $q(r(x, y))$  and the plane given by the graph of  $(x, y) \mapsto x$ .

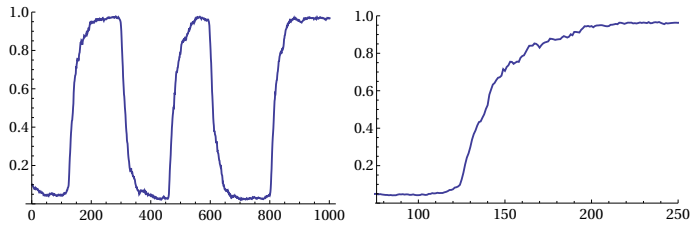


FIG. 2. Trajectory of the mean usage rate  $M^Y(t)$  of  $\mathcal{G}_2$  in the young generation from a sample path of the age-structured Markov chain. Left: The path from time 0 to 1000 years, showing several changes between  $\mathcal{G}_1$  (low) and  $\mathcal{G}_2$  (high). Right: The path from time 75 to 250 years, showing a single grammar change.

316 This definition may be smoothly extended to include  $r_0(1, 0) = 1$  and  $r_0(0, 1) = 0$ ,  
 317 but no extension is possible to  $(0, 0)$  and  $(1, 1)$ . To remedy this, we will assume that  
 318 agents have slightly imperfect perception and introduce

$$319 \quad (10) \quad w(x) = \frac{1}{2} + (1 - \delta) \left( x - \frac{1}{2} \right)$$

320 which maps  $[0, 1]$  to  $[\delta/2, 1 - \delta/2]$ . Pictures in this article will use  $\delta = 1/1000$ . Thus  
 321 a suitable prediction function that is defined and smooth on all of  $[0, 1] \times [0, 1]$  is

$$322 \quad (11) \quad r(x, y) = r_0(w(x), w(y)).$$

323 The combined mean learning-prediction function  $q(r(x, y))$  is plotted in [Figure 1](#).  
 324 An important feature is that since  $q(0) > 0$  and  $q(1) < 1$ , the graph is slightly above  
 325 the plane given by  $(x, y) \mapsto x$  along the edge where  $x = 0$ , and is slightly below that  
 326 plane along the edge where  $x = 1$ . This means that given an initial condition where  
 327 one of the idealized grammars is not used at all, there is a non-zero probability that  
 328 it will appear spontaneously.

329 This model turns out to exhibit the desired properties. The population can spon-  
 330 taneously change from one language to the other and back within a reasonable amount  
 331 of time, and once initiated the change runs to completion without turning back.  
 332 See [Figure 2](#) for a graph of the mean usage rate of  $\mathcal{G}_2$  among the younger age group  
 333 as a function of time for a typical run of this Markov chain.



334 **4. Diffusion limit.** To better understand why spontaneous change happens in  
335 this model, we develop a continuous limit for the Markov chain governing the speech  
336 distributions  $X^Y$  and  $X^A$  of the younger and older generations, respectively, which  
337 are points in the open simplex,

$$338 \quad \mathcal{S}^K = \left\{ (x_0, \dots, x_K) \mid x_k \in (0, 1), \sum_{k=0}^K x_k = 1 \right\}.$$

339 In the limit as the population size  $N$  increases without bound, the Markov chain  
340  $(X^Y(j), X^A(j)) : \mathbf{N} \rightarrow \mathcal{S}^K \times \mathcal{S}^K$  ought to converge to the solution  $(\xi^Y(t), \xi^A(t)) :$   
341  $[0, \infty) \rightarrow \mathcal{S}^K \times \mathcal{S}^K$  of a martingale problem. To formulate it, we must calculate the  
342 infinitesimal drift and covariance functions.

343 **4.1. Notation.** To reduce notational clutter in this subsection, all time-dependent  
344 quantities at time step  $j$  will be written without a time index, as in  $T_n^Y$ ,  $C_k^Y$ , and  
345  $X_k^Y$ . The learning distribution  $Q(r(M^Y, M^A))$  will be written as just  $Q$ , and the  
346 aging distribution  $V(X^Y)$  will be written as just  $V$ . Time-dependent quantities at  
347 time step  $j+1$  will be written with a bar, as in  $\bar{T}_n^Y$ ,  $\bar{C}_k^Y$ , and  $\bar{X}_k^Y$ . Expectations and  
348 variances with a  $j$  subscript are conditioned on the information available at time step  
349  $j$ .

350 **4.2. Infinitesimal mean and variance.** Conditioning on time step  $j$ ,  $\mathbf{1}(\bar{T}_n^Y = k)$   
351 is a Bernoulli random variable that takes on the value 1 with probability

$$352 \quad g(n, k) = (1 - p_D - r_R) \mathbf{1}(T_n^Y = k) + p_D Q_k + r_R X_k^Y$$

353 that is, either  $T_n^Y = k$  and it didn't change, or it died and was replaced by a child of  
354 type  $k$ , or it was resampled and became type  $k$ . With this observation, the mean and  
355 variance of  $\bar{C}_k^Y$  conditioned on information known at time step  $j$  can be calculated as  
356 follows.

$$357 \quad (12) \quad \mathbf{E}_j(\bar{C}_k^Y) = \sum_n g(n, k) = (1 - p_D) C_k^Y + p_D N Q_k$$

358

$$359 \quad (13) \quad \begin{aligned} \text{Var}_j(\bar{C}_k^Y) &= \sum_n g(n, k) - g(n, k)^2 \\ &= (1 - p_D - r_R) C_k^Y + p_D N Q_k + r_R N X_k^Y \\ &\quad - (1 - p_D - r_R)^2 C_k^Y \\ &\quad - 2(1 - p_D - r_R) C_k^Y (p_D Q_k + r_R X_k^Y) \\ &\quad - N (p_D Q_k + r_R X_k^Y)^2 \end{aligned}$$

360 If  $m \neq n$  then  $\bar{T}_m^Y$  and  $\bar{T}_n^Y$  are conditionally independent given the information avail-  
 361 able at time  $j + 1$ . If  $h \neq k$  then  $\mathbf{1}(\bar{T}_n^Y = k) \mathbf{1}(\bar{T}_n^Y = h) = 0$ . Therefore,

$$\begin{aligned}
 \text{Cov}_j(\bar{C}_k^Y, \bar{C}_h^Y) &= \sum_n \text{Cov}_j(\mathbf{1}(\bar{T}_n^Y = k), \mathbf{1}(\bar{T}_n^Y = h)) \\
 &= - \sum_n g(n, k)g(n, h) \\
 362 \quad (14) \quad &= - \left( (1 - p_D - r_R) C_k^Y (p_D Q_h + r_R X_h^Y) \right. \\
 &\quad + (1 - p_D - r_R) C_h^Y (p_D Q_k + r_R X_k^Y) \\
 &\quad \left. + N (p_D Q_k + r_R X_k^Y) (p_D Q_h + r_R X_h^Y) \right)
 \end{aligned}$$

363 It follows that

$$364 \quad (15) \quad \mathbf{E}_j \left( \frac{\bar{X}_k^Y - X_k^Y}{1/N} \right) = r_D (Q_k - X_k^Y),$$

365 which gives the infinitesimal drift component for a martingale problem. We also need  
 366 an estimate of the covariance matrix for  $X^Y$ :

$$367 \quad (16) \quad \text{Var}_j(\bar{X}_k^Y) = \frac{1}{N} (2r_R - r_R^2) (X_k^Y - (X_k^Y)^2) + O\left(\frac{1}{N^2}\right)$$

$$368 \quad (17) \quad \text{Cov}_j(\bar{X}_k^Y, \bar{X}_h^Y) = -\frac{1}{N} (2r_R - r_R^2) X_k^Y X_h^Y + O\left(\frac{1}{N^2}\right)$$

370 Similar drift and covariance formulas can be derived for  $X^A$ ,

$$371 \quad (18) \quad \mathbf{E}_j \left( \frac{\bar{X}_k^A - X_k^A}{1/N} \right) = r_D (V_k - X_k^A)$$

$$372 \quad (19) \quad \text{Var}_j(\bar{X}_k^A) = \frac{1}{N} (2r_R - r_R^2) (X_k^A - (X_k^A)^2) + O\left(\frac{1}{N^2}\right)$$

$$373 \quad (20) \quad \text{Cov}_j(\bar{X}_k^A, \bar{X}_h^A) = -\frac{1}{N} (2r_R - r_R^2) X_k^A X_h^A + O\left(\frac{1}{N^2}\right)$$

375 As a further simplification, we can rescale time by a factor of  $r_D$ . This finally  
 376 yields the infinitesimal drift function

$$377 \quad b : \mathcal{S}^K \times \mathcal{S}^K \rightarrow \mathbf{R}^{2K}$$

378

$$379 \quad (21) \quad b \begin{pmatrix} \xi^Y \\ \xi^A \end{pmatrix} = \begin{pmatrix} b^Y \\ b^A \end{pmatrix} = \begin{pmatrix} Q - \xi^Y \\ V - \xi^A \end{pmatrix}$$

380 and the infinitesimal covariance function  $\varepsilon^2 A$

$$381 \quad A : \mathcal{S}^K \times \mathcal{S}^K \rightarrow \mathbf{M}(\mathbf{R}, 2K \times 2K)$$

382

(22)

$$383 \quad A \begin{pmatrix} \xi^Y \\ \xi^A \end{pmatrix} = \begin{pmatrix} A^Y & 0 \\ 0 & A^A \end{pmatrix} = \begin{pmatrix} \xi_1^Y - (\xi_1^Y)^2 & -\xi_1^Y \xi_2^Y & \dots & & & \\ -\xi_1^Y \xi_2^Y & \ddots & & & & 0 \\ \vdots & & & & & \\ & & & 0 & \xi_1^A - (\xi_1^A)^2 & -\xi_1^A \xi_2^A & \dots \\ & & & & -\xi_1^A \xi_2^A & \ddots & \\ & & & & \vdots & & \end{pmatrix}$$

384 and

$$385 \quad (23) \quad \varepsilon = \sqrt{\frac{2r_R - r_R^2}{r_D}} = \sqrt{\frac{1 - (1 - r_R)^2}{r_D}}$$

386 It can be verified by direct calculation that  $A$  is positive definite. The dimensions  
387 given here use the convention that  $\xi_0^Y$  and  $\xi_0^A$  are omitted from the dynamics. They  
388 will not be considered independent variables because of the population size constraints

$$389 \quad (24) \quad \begin{aligned} \xi_0^Y &= 1 - (\xi_1^Y + \dots + \xi_K^Y) \\ \xi_0^A &= 1 - (\xi_1^A + \dots + \xi_K^A) \end{aligned}$$

390 The drift function can be augmented by defining

$$391 \quad b_0^Y = -\sum_{j=1}^K b_j^Y \quad \text{and} \quad b_0^A = -\sum_{j=1}^K b_j^A$$

392 so that deterministic dynamics under the vector field on  $\mathbf{R}^{K+1} \times \mathbf{R}^{K+1}$  defined by  
393 the augmented  $b$  preserve (24).

394 If the resampling feature is removed by setting  $r_R = 0$ , then  $\varepsilon = 0$  and the  
395 dynamics become deterministic. The resampling feature can also be removed from  
396 just the older generation by zeroing out  $A^A$ , or from just the younger generation by  
397 zeroing out  $A^Y$ .

398 **4.3. Convergence to system of SDEs.** The discrete time Markov chain de-  
399 fined in section 3 converges to a system of stochastic differential equations (SDEs) in  
400 the limit as the population size  $N \rightarrow \infty$  and the physical time of a transition step  
401 goes to 0. The time associated with step  $j$  of the Markov chain is  $t = j/N$ , so to  
402 properly express the convergence of the Markov chain to a process in continuous time  
403 and space, we need the auxiliary processes  $\hat{X}^Y$  and  $\hat{X}^A$  that map continuous time to  
404 discrete steps,

$$405 \quad (25) \quad \begin{aligned} \hat{X}^Y(t) &= X^Y(\lfloor Nt \rfloor) \\ \hat{X}^A(t) &= X^A(\lfloor Nt \rfloor) \end{aligned}$$

406 The limiting initial value problem for  $(\xi^Y, \xi^A) \in \mathcal{S}^K \times \mathcal{S}^K$  is built from the infinitesimal  
407 vector field (21) and covariance matrix (22):

$$\begin{aligned}
d\xi_k^Y(t) &= b^Y(\xi^Y, \xi^A) dt + \varepsilon \sigma^Y(t) dB^Y(t) \\
\xi_0^Y &= 1 - \sum_{k=1}^K \xi_k^Y \\
d\xi_k^A(t) &= b^A(\xi^Y, \xi^A) dt + \varepsilon \sigma^A(t) dB^Y(t) \\
\xi_0^A &= 1 - \sum_{k=1}^K \xi_k^A \\
\xi^Y(0) &= \xi_{\text{init}}^Y \\
\xi^A(0) &= \xi_{\text{init}}^A
\end{aligned}
\tag{26}$$

409 Here  $B^Y$  and  $B^A$  are independent  $K$ -dimensional Brownian motions, and  $\sigma^Y$  and  $\sigma^A$   
410 are the unique positive-definite, symmetric square-roots of  $A^Y$  and  $A^A$ . There is no  
411 general closed form for  $\sigma^Y$  and  $\sigma^A$ , but the theory turns out to only require  $A^Y$  and  
412  $A^A$ .

413 **PROPOSITION 4.1.** *Suppose  $(\hat{X}^Y(0), \hat{X}^A(0))$  converges to  $(\xi_{\text{init}}^Y, \xi_{\text{init}}^A)$  as  $N \rightarrow$   
414  $\infty$ . Suppose  $(b^Y, b^A)$  satisfies the hypotheses of [Proposition 4.8](#). Then for each  
415  $\varepsilon_0 > \varepsilon > 0$ , the process  $(\hat{X}^Y(t), \hat{X}^A(t))$  converges weakly as  $N \rightarrow \infty$  to the solution  
416 to (26).*

417 *Proof.* We apply theorem 7.1 from Chapter 8 of [10] as follows. The calcula-  
418 tions (15), (16), and (17) in [section 4](#) verify that the step-to-step drift, variances,  
419 and covariances of the Markov chain converge to the corresponding functions in the  
420 SDE (26) as the time step size  $1/N$  goes to zero. The remaining condition to check  
421 is Durrett's hypothesis (A), which is that the martingale problem associated to the  
422 SDE is well posed. The SDE has pathwise-unique strong solutions, as we will prove in  
423 [Proposition 4.8](#). That implies uniqueness in distribution [10, §5.4 theorem 4.1] which  
424 implies that the martingale problem is well posed [10, §5.4 theorem 4.5] which implies  
425 the desired convergence.  $\square$

426 The commonly referenced theorem for existence and uniqueness of solutions to  
427 initial value problems for SDEs (see [52, theorem 5.2.1], for example) is not sufficient  
428 for (26). It applies to dynamics on Euclidean space, but the dynamics of interest  
429 here are restricted to  $\mathcal{S}^K \times \mathcal{S}^K$ . We can change variables to expand the simplices to  
430 whole spaces, but then the global Lipschitz property and global growth constraints  
431 required by that theorem are not met. We must therefore apply more general theorems  
432 from [10] instead.

433 **4.4. Change of variables.** First, we deal with phase space, as (26) only makes  
434 sense for  $(\xi^Y, \xi^A) \in \mathcal{S}^K \times \mathcal{S}^K$ . We change variables so as to push the boundary of the  
435 phase space off to infinity. Since the formulas are exactly parallel for each generation,  
436 the generation label superscripts will be omitted where possible. To further conserve  
437 space, let  $\gamma = 1/(K + 1)$ . Each vector  $\xi \in \mathcal{S}^K$  is mapped to a vector  $\lambda$ ,

$$\lambda_k = \tilde{\xi}(\xi_k - \gamma)
\tag{27}$$

439 where

$$440 \quad (28) \quad \tilde{\xi} = \left( \prod_{k=0}^K \xi_k \right)^{-\gamma}$$

441 The interior of the simplex expands to the entire plane

$$442 \quad \left\{ \lambda \in \mathbf{R}^{K+1} \mid \sum_{k=0}^K \lambda_k = 0 \right\}.$$

443 Let us also define

$$444 \quad (29) \quad \xi_{\min} = \min_k \xi_k \quad \xi_{\max} = \max_k \xi_k \quad \lambda_{\min} = \min_k \lambda_k \quad \lambda_{\max} = \max_k \lambda_k$$

445 Note that the extrema for  $\xi_k$  and  $\lambda_k$  occur at the same value of the index  $k$ . Since  
446  $\sum_{k=0}^K \xi_k = 1$ , it follows immediately that

$$447 \quad (30) \quad \xi_{\min} \leq \gamma \leq \xi_{\max} \quad \lambda_{\min} \leq 0 \leq \lambda_{\max}$$

448 Furthermore, since  $\lambda_{\min} = \tilde{\xi}(\xi_{\min} - \gamma)$ ,

$$449 \quad (31) \quad \tilde{\xi} = \frac{-\lambda_{\min}}{\gamma - \xi_{\min}} > -\lambda_{\min}(K+1)$$

450 LEMMA 4.2. *The change of variables is smooth and smoothly invertible provided*  
451 *none of the  $\xi_k$ 's are zero, although the inverse does not have a closed form.*

452 *Proof.* To prove the existence of the inverse, note that if there is a solution for  
453 the  $\xi_k$ 's in terms of  $\lambda_k$ 's, it must hold that

$$454 \quad \tilde{\xi}^{-(K+1)} = \prod_{k=0}^K \xi_k = \prod_{k=0}^K (\tilde{\xi}^{-1} \lambda_k + \gamma) = \tilde{\xi}^{-(K+1)} \prod_{k=0}^K (\lambda_k + \gamma \tilde{\xi})$$

455 Thus  $f(\tilde{\xi}) = 1$  where  $f$  is the polynomial

$$456 \quad (32) \quad f(x) = \prod_{k=0}^K (\lambda_k + \gamma x)$$

457 Assuming that the  $\lambda_k$ 's are known, note that  $f(-\lambda_{\min}(K+1)) = 0$ , and for  $x >$   
458  $-\lambda_{\min}(K+1)$ ,  $f(x)$  is product of strictly positive terms, all of which are strictly  
459 increasing in  $x$ , and it is unbounded as  $x \rightarrow \infty$ . There is therefore a unique solution  
460 to  $f(x) = 1$  with  $x > -\lambda_{\min}(K+1)$ . Let  $\tilde{\xi}$  be this solution, and recover  $\xi_k = \tilde{\xi}^{-1} \lambda_k + \gamma$ .  
461 This change of variables is smooth and locally Lipschitz, but not globally Lipschitz  
462 because each partial derivative (40) is unbounded as  $\xi_j \rightarrow 0$ .  $\square$

463 Several additional inequalities relating  $\xi$  and  $\lambda$  will be required. First, to avoid  
464 confusion about whether the 0th element of a vector is included in a dot product or  
465 magnitude, let us define

$$466 \quad (33) \quad \|v\|^2 = \sum_{k=1}^K v_k^2 \quad \|v\|_0^2 = \sum_{k=0}^K v_k^2 = \|v\|^2 + (v_0)^2$$

$$467 \quad (34) \quad u \cdot v = \sum_{k=1}^K u_k v_k \quad u \odot v = \sum_{k=0}^K u_k v_k$$

468

469 For a general vector  $v = (v_0, \dots, v_K)^\top$ , with extreme elements  $v_{\min}$  and  $v_{\max}$ , it is  
 470 elementary to verify that

$$471 \quad (35) \quad \|v\|_0^2 - v_{\max}^2 \leq \|v\|^2 \leq \|v\|^2 + v_{\min}^2 \leq \|v\|_0^2 \leq \|v\|^2 + v_{\max}^2 \leq 2\|v\|^2$$

LEMMA 4.3.

$$472 \quad (36) \quad \tilde{\xi} \leq \frac{1 - \lambda_{\min}}{\gamma} \leq \frac{1 + \|\lambda\|_0}{\gamma} \leq \frac{1 + \sqrt{2}\|\lambda\|}{\gamma}$$

473 *Proof.* From the definition of  $\tilde{\xi}$ , it is clear that

$$474 \quad (37) \quad 1 < \xi_{\max}^{-1} \leq \tilde{\xi} \leq \xi_{\min}^{-1}$$

475 Building from (37),

$$476 \quad \tilde{\xi} \leq \xi_{\min}^{-1} = \left( \gamma + \frac{\lambda_{\min}}{\tilde{\xi}} \right)^{-1}$$

477 It follows that

$$478 \quad \tilde{\xi}\gamma + \lambda_{\min} = \tilde{\xi} \left( \gamma + \frac{\lambda_{\min}}{\tilde{\xi}} \right) \leq 1$$

479 which, in conjunction with (35), yields the bounds (37).  $\square$

480 LEMMA 4.4. *There is a constant  $\rho > 0$  such that for all  $\xi$*

$$481 \quad (38) \quad (K+1)\tilde{\xi} \leq \sum_{k=0}^K \xi_k^{-1} \leq \rho \tilde{\xi}^{K+1} \leq \rho \left( \frac{1 + \sqrt{2}\|\lambda\|}{\gamma} \right)^{K+1}$$

482 *Proof.* The lower bound on  $\sum \xi_k^{-1}$  comes from the standard harmonic-geometric  
 483 mean inequality. For the upper bound, note that

$$484 \quad f(\xi) = \left( \sum_{k=0}^K \frac{1}{\xi_k} \right) \left( \prod_{k=0}^K \xi_k \right)$$

485 is a polynomial, so it has an absolute maximum  $\rho$  on the closure of  $\mathcal{S}^K$ .  $\square$

486 It is important to note that the power  $K+1$  of  $\|\lambda\|$  in the upper bound (38) is the  
 487 best possible. Consider the case of  $\xi_0 = \delta$ ,  $\xi_k = (1 - \delta)/K$  for  $k > 0$  and small  $\delta > 0$ .  
 488 Then  $\sum \xi_k^{-1} \approx \delta^{-1} + K$ ,  $\tilde{\xi} \approx K^{K\gamma} \delta^{-\gamma}$ , and  $\lambda_k \approx K^{K\gamma} (\xi_k - \gamma) \delta^{-\gamma}$ . In this case,  $\sum \xi_k^{-1}$   
 489 is on the order of  $\|\lambda\|^{1/\gamma}$ . This power is why so much care must be taken to establish  
 490 the well-posedness of (42).

491 **4.5. Itô's formula.** The following partial derivative formulas are needed in the  
 492 application of Itô's formula, and are written here assuming  $i \geq 1$ ,  $j \geq 1$ ,  $k \geq 1$ . Recall

493 that  $\xi_0$  is not considered a separate independent variables because of (24).

$$494 \quad (39) \quad \partial_{\xi_j} \tilde{\xi} = \gamma \tilde{\xi} (\xi_0^{-1} - \xi_j^{-1})$$

$$495 \quad (40) \quad \begin{aligned} \partial_{\xi_j} \lambda_k &= \gamma \tilde{\xi} (\xi_0^{-1} - \xi_j^{-1}) (\xi_k - \gamma) + \mathbf{1}(j = k) \tilde{\xi} \\ &= \gamma \lambda_k (\xi_0^{-1} - \xi_j^{-1}) + \mathbf{1}(j = k) \tilde{\xi} \end{aligned}$$

$$496 \quad (41) \quad \begin{aligned} \partial_{\xi_i \xi_j} \lambda_k &= \gamma^2 \tilde{\xi} (\xi_0^{-1} - \xi_i^{-1}) (\xi_0^{-1} - \xi_j^{-1}) (\xi_k - \gamma) + \gamma \tilde{\xi} \xi_0^{-2} (\xi_k - \gamma) \\ &\quad + \mathbf{1}(i = j) \left( \gamma \tilde{\xi} \xi_j^{-2} (\xi_k - \gamma) \right) \\ &\quad + \mathbf{1}(i = k) \left( \gamma \tilde{\xi} (\xi_0^{-1} - \xi_j^{-1}) \right) \\ &\quad + \mathbf{1}(j = k) \left( \gamma \tilde{\xi} (\xi_0^{-1} - \xi_i^{-1}) \right) \end{aligned}$$

502 Applying Itô's formula to change variables to  $\lambda$  yields, for  $k \geq 1$ ,

$$503 \quad (42) \quad d\lambda_k = \left( D_\xi \lambda_k \cdot b + \frac{\varepsilon^2}{2} \text{tr} (\sigma^\top (D_\xi^2 \lambda_k) \sigma) \right) dt + (D_\xi \lambda_k)^\top \sigma dB$$

504 where  $D_\xi$  is the gradient with respect to  $\xi$  and  $D_\xi^2$  is the Hessian matrix with respect  
505 to  $\xi$ . No particular form of  $b$  is assumed.

506 Since  $\sigma^Y$  is symmetric and the trace has the general property that  $\text{tr}(PQR) =$   
507  $\text{tr}(QRP)$ , the trace term may be evaluated as follows despite the fact that no explicit  
508 form is possible for  $\sigma$ :

$$509 \quad \text{tr} (\sigma^\top (D_\xi^2 \lambda_k) \sigma) = \text{tr} ((D_\xi^2 \lambda_k) \sigma \sigma^\top) = \text{tr} ((D_\xi^2 \lambda_k) A)$$

510 After a laborious simplification,

$$511 \quad (43) \quad \text{tr} (\sigma^\top (D_\xi^2 \lambda_k) \sigma) = \gamma(\gamma + 1) \lambda_k \sum_{j=0}^K \xi_j^{-1}$$

512 **4.6. Well-posedness of the SDEs.** The drift and volatility terms of (42) are  
513 continuously differentiable, so they automatically satisfy a local Lipschitz inequality,  
514 as required by the general theorem concerning the existence and uniqueness of  
515 solutions in [10, §5.3].

516 The theorem also requires a growth constraint formulated as follows. Let us adapt  
517 the usual big-O notation, using

$$518 \quad f(\lambda^Y, \lambda^A) = g(\lambda^Y, \lambda^A) + \mathcal{O}^2$$

519 to mean that there exists a constant  $H > 0$  such that for all  $\lambda^Y$  and  $\lambda^A$ ,

$$520 \quad f(\lambda^Y, \lambda^A) - g(\lambda^Y, \lambda^A) < H \left( 1 + \|\lambda^Y\|^2 + \|\lambda^A\|^2 \right)$$

521 The growth constraint required in [10, §5.3] is  $\beta^Y + \beta^A = \mathcal{O}^2$  where

$$\begin{aligned}
\beta^Y &= \sum_{k=1}^K \lambda_k^Y \left( D_{\xi^Y} \lambda^Y \cdot b^Y + \frac{\varepsilon^2}{2} \operatorname{tr} \left( (\sigma^Y)^\top \left( D_{\xi^Y}^2 \lambda_k^Y \right) \sigma^Y \right) \right) \\
&\quad + \varepsilon^2 \operatorname{tr} \left( (\sigma^Y)^\top (D_{\xi^Y} \lambda^Y)^\top (D_{\xi^Y} \lambda^Y) \sigma^Y \right) \\
\beta^A &= \sum_{k=1}^K \lambda_k^A \left( D_{\xi^A} \lambda^A \cdot b^A + \frac{\varepsilon^2}{2} \operatorname{tr} \left( (\sigma^A)^\top \left( D_{\xi^A}^2 \lambda_k^A \right) \sigma^A \right) \right) \\
&\quad + \varepsilon^2 \operatorname{tr} \left( (\sigma^A)^\top (D_{\xi^A} \lambda^A)^\top (D_{\xi^A} \lambda^A) \sigma^A \right)
\end{aligned}
\tag{44}$$

523  $D_{\xi^Y} \lambda^Y$  and  $D_{\xi^A} \lambda^A$  are Jacobian matrices, and  $D_{\xi^Y}^2 \lambda_k^Y$  and  $D_{\xi^A}^2 \lambda_k^A$  are Hessian ma-  
524 trices. The difficulty here is that  $\beta^Y + \beta^A$  turns out to contain terms of degree greater  
525 than 2, so we must confirm that these are negative for large  $\lambda$ . The following estimates  
526 are derived omitting the generation label where possible, as parallel logic applies to  
527  $\beta^Y$  and  $\beta^A$ .

528 Incorporating (43), the generic  $\beta$  term is

$$\beta = \sum_{k=1}^K \lambda_k \left( D_\xi \lambda \cdot b + \frac{\varepsilon^2}{2} \gamma(\gamma+1) \lambda_k \sum_{j=0}^K \xi_j^{-1} \right) + \varepsilon^2 \operatorname{tr} \left( \sigma (D_\xi \lambda) (D_\xi \lambda)^\top \sigma \right)
\tag{45}$$

530 The remaining trace term can be evaluated by cyclically reordering the matrices

$$\operatorname{tr} \left( \sigma^\top (D_\xi \lambda)^\top (D_\xi \lambda) \sigma \right) = \operatorname{tr} \left( (D_\xi \lambda)^\top (D_\xi \lambda) \sigma \sigma^\top \right) = \operatorname{tr} \left( (D_\xi \lambda)^\top (D_\xi \lambda) A \right)$$

532 After a massive amount of simplification,

$$\begin{aligned}
\beta &= -\gamma \|\lambda\|^2 \sum_{j=0}^K b_j \xi_j^{-1} + \varepsilon^2 \left( \frac{\gamma(\gamma+1)}{2} \lambda_0 + \gamma^2 \|\lambda\|^2 \right) \sum_{j=0}^K \xi_j^{-1} + \tilde{\xi} (\lambda \cdot b + 2\varepsilon^2 \lambda \cdot \xi) \\
&\quad + \varepsilon^2 \left( -\|\lambda\|^2 + \tilde{\xi}^2 (1 - \xi_0 - \|\xi\|^2) + 2\gamma \tilde{\xi} \lambda_0 \right)
\end{aligned}
\tag{46}$$

534 The largest magnitude terms are those that include  $\xi_j^{-1}$ , and those must be handled  
535 carefully. The others are  $\mathcal{O}^2$  in light of inequalities proved in subsection 4.4, and the  
536 assumption that the  $b_j$ 's are bounded.

$$\beta = -\gamma \|\lambda\|^2 \sum_{k=0}^K \frac{b_k}{\xi_k} + \varepsilon^2 \gamma \|\lambda\|^2 \left( \frac{\gamma+1}{2 \|\lambda\|^2} + \gamma \right) \sum_{k=0}^K \frac{1}{\xi_k} + \mathcal{O}^2
\tag{47}$$

538 To express the constraints on  $b^Y$  and  $b^A$  that are necessary to guarantee that the  
539 remaining large magnitude terms in  $\beta^Y$  and  $\beta^A$  are negative overall, the following  
540 definitions are required. Given  $\mu > 0$ , define the  $\mu$ -border of  $\mathcal{S}^K$  to be

$$\mathcal{S}_\mu^K = \{ x \in \mathcal{S}^K \mid \exists k : x_k < \mu \}
\tag{48}$$

542 The *border class* of  $x \in \mathcal{S}_\mu^K$  is  $\operatorname{BC}(x; \mu) = \sum_k \mathbf{1}(x_k < \mu)$ . The parameter  $\mu$  will be  
543 omitted when it is clear from context.



544 LEMMA 4.5. If  $\xi \in \mathcal{S}_\mu^K$  then

545 (49) 
$$\tilde{\xi} \geq \mu^{-\gamma \text{BC}(\xi)} \geq \mu^{-\gamma}$$

546 *Proof.* Let  $c = \text{BC}(\xi)$ . Then there are  $c$  indices  $k$  for which  $1/\mu < \xi_k^{-1}$  and  $K - c$   
547 indices for which  $1 < \xi_k^{-1}$ . Taking the  $\gamma$  power of the product yields (49).  $\square$

548 LEMMA 4.6. Suppose  $(b^Y, b^A)$  is bounded. Suppose there exist numbers  $G > 0$ ,  
549  $F$ , and  $\gamma > \mu > 0$  such that

550 (50) 
$$\begin{aligned} \text{if } \xi^Y \in \mathcal{S}_\mu^K \text{ then } \sum_{k=0}^K \frac{b_k^Y}{\xi_k^Y} &\geq G \sum_{k=0}^K \frac{1}{\xi_k^Y} + F \\ \text{and if } \xi^A \in \mathcal{S}_\mu^K \text{ then } \sum_{k=0}^K \frac{b_k^A}{\xi_k^A} &\geq G \sum_{k=0}^K \frac{1}{\xi_k^A} + F \end{aligned}$$

551 Then there exists an  $\varepsilon_0 > 0$  such that for each  $\varepsilon_0 > \varepsilon > 0$ ,  $\beta^Y + \beta^A = \mathcal{O}^2$ .

552 *Proof.* If  $\xi \in \mathcal{S}^K \setminus \mathcal{S}_\mu^K$ , then  $\lambda$  is bounded and each  $\xi_k$  satisfies  $1/\xi_k < 1/\mu$ . Since  
553  $b$  is assumed to be bounded, it is straightforward to confirm that  $\beta = \mathcal{O}^2$  in this case.

554 Suppose  $\xi \in \mathcal{S}_\mu^K$ . Then from (49),  $\tilde{\xi} \geq \mu^{-\gamma}$ . Consequently, (36) implies  
555  $\|\lambda\| \geq (\gamma\mu^{-\gamma} - 1)/\sqrt{2}$ . Using the lower bound  $G$  to replace the  $b_k$  terms and pushing  
556 degree 2 terms into  $\mathcal{O}^2$ ,

557 
$$\beta = \gamma \|\lambda\|^2 \left( \sum_{k=0}^K \frac{1}{\xi_k} \right) \left[ -G + \varepsilon^2 \left( \frac{\gamma + 1}{2\|\lambda\|} + \gamma \right) \right] + \mathcal{O}^2$$

558 If  $\varepsilon$  is small enough,

559 
$$\varepsilon \leq \sqrt{\frac{G}{\frac{\gamma+1}{\sqrt{2}(\gamma\mu^{-\gamma}-1)} + \gamma}} = \varepsilon_0$$

560 then the factor in square brackets is negative and  $\beta = \mathcal{O}^2$ .

561 Since the above arguments apply to both  $\beta^Y$  and  $\beta^A$ , the sum satisfies  $\beta^Y + \beta^A =$   
562  $\mathcal{O}^2$ .  $\square$

563 LEMMA 4.7. If the vector field has the form

564 
$$\begin{aligned} b^Y(\xi^Y, \xi^A) &= U^Y(\xi^Y, \xi^A) - \xi^Y \\ b^A(\xi^Y, \xi^A) &= U^A(\xi^Y, \xi^A) - \xi^A \end{aligned}$$

565 where  $U^Y$  and  $U^A$  are probability vectors with uniform positive lower bounds

566 
$$\begin{aligned} \forall \xi^Y, \xi^A : U^Y(\xi^Y, \xi^A) &\geq U_{\min}^Y > 0 \\ \forall \xi^Y, \xi^A : U^A(\xi^Y, \xi^A) &\geq U_{\min}^A > 0 \end{aligned}$$

567 then it satisfies (50).

568 *Proof.* For either generation,  $\square$

569 
$$\sum_{k=0}^K \frac{U_k - \xi_k}{\xi_k} \geq U_{\min} \sum_{k=0}^K \frac{1}{\xi_k} - (K + 1)$$

570 The example vector field (21) satisfies (50). Since the example  $Q$  from (2) is the  
 571 probability vector for a binomial distribution, its least element is either  $Q_0$  or  $Q_K$ .  
 572 Therefore, each element of  $Q$  satisfies

$$573 \quad Q_k \geq Q_{\min} = \min \{q(0)^K, (1 - q(0))^K, q(1)^K, (1 - q(1))^K\}.$$

574 Each element of the distribution vector  $V$  as in (7) satisfies  $V_k \geq \eta$ . If we try to  
 575 set  $V(\xi^Y, \xi^A) = \xi^Y$ , then there is no way to choose  $G$ , hence the need for  $\eta > 0$ .

576 **PROPOSITION 4.8.** *If  $b$  satisfies the hypotheses of Lemma 4.6, then for each  $\varepsilon_0 >$   
 577  $\varepsilon > 0$ , the SDEs (42) and (26) have pathwise-unique strong solutions for all positive  
 578 time starting from each suitable initial value.*

579 *Proof.* The theorem from [10, §5.3] in conjunction with Lemma 4.6 confirms the  
 580 result for  $(\lambda^Y, \lambda^A)$ , and the change of variables from subsection 4.4 maps those solu-  
 581 tions to solutions of (26).  $\square$

582 **4.7. Generalizations.** The results in this section generalize to many other sit-  
 583 uations, since many of the proofs make no assumptions on the specific form of  $b$ ,  
 584 although they were developed to apply to (21). For example, if there are more than  
 585 two grammars of interest, the indices 0 through  $K$  can be remapped to any mixtures  
 586 of grammars and the learning function  $Q$  can be adjusted accordingly, resulting in a  
 587 discrete time model that converges to continuous time process with the same form as  
 588 (26). There's also no need to restrict  $Q$  to be the mass function for any particular  
 589 distribution.

590 In formulating (26), it was assumed that both generations were subdivided into  
 591 the same types, that is, everyone of type  $k$  uses  $\mathcal{G}_1$  with probability  $k/K$ . The results  
 592 in this section do not depend on requiring all sub-populations to have states with  
 593 same interpretation, or even to lie in simplexes of the same dimension.

594 These results also generalize immediately to a population divided into any number  
 595 of sub-populations, such as multiple age groups, geographic regions, or social classes.  
 596 The key theorem 7.1 from Chapter 8 of [10] requires that the time step size be  $\frac{1}{N}$ .  
 597 It would continue to apply if the sub-populations were of different sizes but all were  
 598 proportional to  $N$ .

599 **5. Dynamics in a 2-dimensional case.** We will continue by restricting our  
 600 attention to the case of  $K = 1$ . That is, simulated individuals use  $\mathcal{G}_2$  exclusively or  
 601 not at all, and in discrete time,  $X_0^Y$  is the fraction of the young generation that never  
 602 uses  $\mathcal{G}_2$  and  $X_1^Y$  is the fraction that always uses  $\mathcal{G}_2$ . Since  $X_0^Y + X_1^Y = 1$ , it is only  
 603 necessary to deal with  $X_1^Y$ . Likewise we may focus on  $\xi_1^Y$ ,  $\xi_1^A$ , and  $Q_1 = q(r(X, Y))$ .

604 The covariance function (22) reduces to a 2-by-2 diagonal matrix so it has a very  
 605 simple square-root:

$$606 \quad (51) \quad \sigma \begin{pmatrix} \xi^Y \\ \xi^A \end{pmatrix} = \begin{pmatrix} \sqrt{\xi_1^Y - (\xi_1^Y)^2} & 0 \\ 0 & \sqrt{\xi_1^A - (\xi_1^A)^2} \end{pmatrix}$$

607 As  $N \rightarrow \infty$ , the discrete time process converges weakly to the solution  $(\xi^Y, \xi^A) :$   
 608  $[0, \infty) \rightarrow (0, 1) \times (0, 1)$  of

$$609 \quad (52) \quad \begin{aligned} d\xi_1^Y &= (q(r(\xi_1^Y, \xi_1^A)) - \xi_1^Y) dt + \varepsilon \sqrt{\xi_1^Y(1 - \xi_1^Y)} dB^Y \\ d\xi_1^A &= (\xi_1^Y(1 - \eta/2) + \eta - \xi_1^A) dt + \varepsilon \sqrt{\xi_1^A(1 - \xi_1^A)} dB^A \\ \xi^Y(0) &= \xi_{\text{init}}^Y \text{ and } \xi_1^A(0) = \xi_{\text{init}}^A \end{aligned}$$

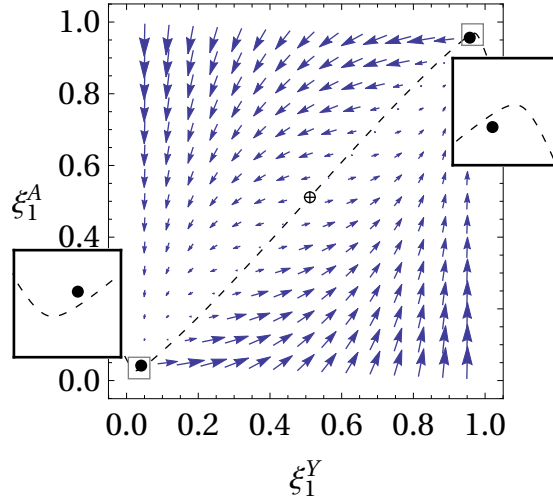


FIG. 3. Phase portrait for (52) with  $\varepsilon = 0$ . The crossed dot  $\oplus$  is a saddle point, the two dots  $\bullet$  are sinks, and the dashed curve is the separatrix between their basins of attraction. The arrows indicate the direction of the vector field. Inset boxes show magnified pictures of the areas around the sinks.

610 where  $B^Y$  and  $B^A$  are independent one-dimensional Brownian motions.

611 **5.1. Comparison to the deterministic limit.** In the deterministic limit  $\varepsilon = 0$   
612 and returning to the specific learning process described in section 2 and section 3, the  
613 dynamical system (52) has two stable equilibria representing populations where both  
614 generations are dominated by one grammar or the other. The separatrix forming  
615 the boundary between the two basins of attraction passes very close to the stable  
616 equilibria. See Figure 3. Under the stochastic dynamics, the population will hover  
617 near an equilibrium until random fluctuations cause it to stray across the separatrix,  
618 where it will be blown toward the other. It will continue to oscillate irregularly  
619 between the two equilibria for all time. These separatrix-crossing events generate  
620 spontaneous monotonic language changes separated by reasonably long intervals of  
621 temporary stability.

622 **5.2. Memory kernel form.** Another way to understand this form of instability  
623 is to express  $\xi_1^A$  as an average of  $\xi_1^Y$  over its past, with an exponential kernel giving  
624 greater weight to the recent past. This is accomplished by making two simplifications.  
625 First, the resampling step from the Markov chain will be applied only to the younger  
626 generation, which removes the random term from  $d\xi_1^A$  in (52) but not from  $d\xi_1^Y$ .  
627 Second,  $\eta$  in the aging distribution  $V$  will be set to 0. This yields a linear ordinary  
628 differential equation for  $\xi_1^A$  with  $\xi_1^Y$  acting as an inhomogeneity

$$629 \quad \frac{d\xi_1^A}{dt} = \xi_1^Y - \xi_1^A, \text{ with solution } \xi_1^A(t) = e^{-t}\xi_{\text{init}}^A + \int_0^t e^{-(t-s)}\xi_1^Y(s)ds.$$

630 With this simplification, the dynamics for  $\xi_1^Y$  take the form of a stochastic functional-  
631 delay differential equation

$$632 \quad (53) \quad d\xi_1^Y(t) = \left( q \left( r \left( \xi_1^Y(t), K_t \xi_1^Y \right) \right) - \xi_1^Y(t) \right) dt + \varepsilon \sqrt{\xi_1^Y(t)(1 - \xi_1^Y(t))} dB$$

633 where the delay appears through convolution with a memory kernel

$$634 \quad K_t f = e^{-t} \xi_{\text{init}}^A + \int_0^t e^{-(t-s)} f(s) ds.$$

635 The age structure serves to give the population a memory, so that the speech pattern  
 636  $\xi_1^Y$  of the young generation changes depending on how the current young generation  
 637 deviates from its recent past average. Chance deviations of sufficient size are am-  
 638 plified when children detect them and predict that the trend will continue, yielding  
 639 prediction-driven instability.

## 640 **6. Discussion.**

641 **6.1. Comparison to other models.** The discrete and continuous models as  
 642 described in [sections 2](#) and [3](#) are based on the Wright-Fisher model of population  
 643 genetics as described in [\[10\]](#), which is formulated as a Markov chain and its limit as a  
 644 stochastic differential equation for an infinite population. The original Wright-Fisher  
 645 model takes values on an interval, which makes the theoretical analysis much simpler  
 646 than for [\(26\)](#). A similar derivation to that of [section 4](#) resulting in a Fokker-Planck  
 647 equation is given in [\[3\]](#), without the theoretical treatment given here. The model in  
 648 [\[58\]](#) derives a similar model, grounding the learning process in Bayesian inference.  
 649 Neither these nor the Wright-Fisher model incorporate age structure or forces such  
 650 as learning and prediction that are not present in biological birth-death processes.

651 A related dynamical system is the FitzHugh-Nagumo model for a spiking neuron  
 652 [\[30, 42\]](#), which is a general family of two-variable dynamical systems. Its structure  
 653 is similar to [Figure 3](#) except that it has only the lower left stable equilibrium, which  
 654 represents a resting neuron. A disturbance causes the neuron's state to stray away  
 655 from that rest state and go on a long excursion known as an action potential or spike.

656 The language change model examined here differs from the stochastic FitzHugh-  
 657 Nagumo model in several ways. It is derived as a continuous limit of a Markov chain  
 658 rather than from adding noise to an existing dynamical system. It has two stable  
 659 equilibria rather than one as long as  $\varepsilon$  is sufficiently small (although it is conceivable  
 660 that some linguistic phenomenon might exhibit the single stable equilibrium). It is  
 661 naturally confined to  $\mathcal{S}^K \times \mathcal{S}^K$ , where FitzHugh-Nagumo models occupy an entire  
 662 plane. The random term added to a FitzHugh-Nagumo model is normally Brownian  
 663 motion multiplied by a small constant. The change of variables  $\theta = \arcsin(2\xi - 1)$ ,  
 664  $\phi = \arcsin(2\zeta - 1)$  transforms the low dimensional case [\(52\)](#) to that form but the  
 665 system remains confined to a square, and the change of variables to [\(42\)](#) on the whole  
 666 plane has a non-constant coefficient on the Brownian motion. Thus, the theory of  
 667 FitzHugh-Nagumo models must be adapted before it can be applied to this language  
 668 model.

669 Population-level memory has been used to model other social trends that exhibit  
 670 momentum. For example, the authors of [\[16\]](#) develop a model in which parents use  
 671 a discrete-time memory kernel analogous to [\(53\)](#) to compute running averages of the  
 672 popularity of given names, and use this information when naming babies. The case of  
 673 language change is different because children seem to be capable of contributing with-  
 674 out decades of accumulated experience. They must get historical information from  
 675 some other source, and age-correlated differences in speech is a reasonable hypothesis.

676 **7. Conclusion.** The main goal of this article is to begin with a discrete time, fi-  
 677 nite population model that can represent spontaneous language change in a population  
 678 between meta-stable states, each dominated by one idealized grammar, and connect

679 it via solid theory to a continuous time, infinite population model. Language is repre-  
 680 sented as a mixture of the idealized grammars to reflect the variability of speech seen  
 681 in manuscripts and social data. A Markov chain that includes age structure has all the  
 682 desired properties for the first model. The population can switch spontaneously from  
 683 one language to the other and the transition is monotonic. Intuitively, the mechanism  
 684 of these spontaneous changes is that every so often, children pick up on an accidental  
 685 correlation between age and speech, creating the beginning of a trend. The prediction  
 686 step in the learning process amplifies the trend, and moves the population away from  
 687 equilibrium, which suggests the term *prediction-driven instability* for this effect.

688 Fundamental results were proved. Specifically, in the limit as the number of agents  
 689 goes to infinity, sample paths of the Markov chain converge weakly to solutions to a  
 690 system of well-posed SDEs, which have the form of drift terms plus a small stochastic  
 691 perturbation. The derivation of the correct SDEs and the proof that the convergence  
 692 happens as intended require a change of variables specifically tailored to the geometry  
 693 of the simplex, together with theoretical tools more sophisticated than those typically  
 694 needed for population dynamics models. The proof that the system of SDEs is well-  
 695 posed relies only on general properties of the drift vector field and the specific form  
 696 of the infinitesimal covariance matrix.

697 Looking at a low dimensional case, in the limit of zero noise, the prediction-driven  
 698 instability comes from the proximity of stable sinks to the separatrix of their basins  
 699 of attraction. The instability comes from the general geometry of the phase space as  
 700 in [Figure 3](#). Alternatively, the prediction process may be understood as comparing  
 701 the current state of the population to an average emphasizing its recent past, and  
 702 chance deviations trigger the instability. Concrete formulas were given for  $q$ ,  $r$ , and  
 703  $Q$ , but the interesting behavior is not limited to these examples.

704 Future studies of this model could include adapting and applying techniques for  
 705 studying noise-activated transitions among meta-stable states, including exit time  
 706 problems [[13](#), [31](#), [32](#)]. For example, it is possible to numerically estimate the time  
 707 between transitions using a partial differential equation or a variational technique.  
 708 The change of variables and associated theory may be of use to other dynamical  
 709 systems whose phase space is a simplex, such as replicator dynamics [[18](#)].

710

## REFERENCES

- 711 [1] D. ADGER, *Core Syntax: A Minimalist Approach*, Oxford University Press, Oxford, May 2003.  
 712 [2] A. ALISHAHI AND S. STEVENSON, *A Computational Model of Early Argument Structure Acqui-*  
 713 *sition*, *Cognitive Science*, 32 (2008), pp. 789–834.  
 714 [3] G. J. BAXTER, R. A. BLYTHE, W. CROFT, AND A. J. MCKANE, *Utterance selection model of*  
 715 *language change*, *Physical Review E*, 73 (2006), p. 046118.  
 716 [4] M. BECKER, *There Began to Be a Learnability Puzzle*, *Linguistic Inquiry*, 37 (2006),  
 717 pp. 441–456.  
 718 [5] R. C. BERWICK AND P. NIYOGI, *Learning from Triggers*, *Linguistic Inquiry*, 27 (1996),  
 719 pp. 605–622.  
 720 [6] T. BRISCOE, *Grammatical acquisition: Inductive bias and coevolution of language and the*  
 721 *language acquisition device*, *Language*, 76 (2000), pp. 245–296.  
 722 [7] T. BRISCOE, *Grammatical acquisition and linguistic selection*, in *Linguistic Evolution through*  
 723 *Language Acquisition: Formal and Computational Models*, T. Briscoe, ed., Cambridge  
 724 University Press, Cambridge, UK, 2002, pp. 255–300.  
 725 [8] N. CHOMSKY, *Aspects of the Theory of Syntax*, MIT Press, Cambridge, MA, 1965.  
 726 [9] F. CUCKER, S. SMALE, AND D.-X. ZHOU, *Modeling Language Evolution*, *Foundations of Com-*  
 727 *putational Mathematics*, 4 (2004), pp. 315–343.  
 728 [10] R. DURRETT, *Stochastic Calculus: A Practical Introduction*, Probability and stochastics series,  
 729 CRC Press, Boca Raton, 1996.

- 730 [11] R. DURRETT AND S. LEVIN, *The Importance of Being Discrete (and Spatial)*, Theoretical  
731 Population Biology, 46 (1994), pp. 363–394.
- 732 [12] Z. FAGYAL, S. SWARUP, A. M. ESCOBAR, L. GASSER, AND K. LAKKARAJU, *Centers and periph-*  
733 *eries: Network roles in language change*, *Lingua*, 120 (2010), pp. 2061–2079.
- 734 [13] M. I. FREIDLIN AND A. D. WENTZELL, *Random perturbations of dynamical systems*, Springer,  
735 1998.
- 736 [14] E. GIBSON AND K. WEXLER, *Triggers*, *Linguistic Inquiry*, 25 (1994), pp. 407–454.
- 737 [15] E. M. GOLD, *Language Identification in the Limit*, *Information and Control*, 10 (1967),  
738 pp. 447–474.
- 739 [16] T. M. GURECKIS AND R. L. GOLDSTONE, *How You Named Your Child: Understanding the*  
740 *Relationship Between Individual Decision Making and Collective Outcomes*, *Topics in*  
741 *Cognitive Science*, 1 (2009), pp. 651–674.
- 742 [17] M. D. HAUSER, N. CHOMSKY, AND W. T. FITCH, *The Faculty of language: What is it, who*  
743 *has it, and how did it evolve?*, *Science*, 298 (2002), pp. 1569–1579.
- 744 [18] J. HOFBAUER AND K. SIGMUND, *Evolutionary Games and Population Dynamics*, Cambridge  
745 University Press, 1998.
- 746 [19] C. L. H. KAM AND E. L. NEWPORT, *Regularizing Unpredictable Variation: The Roles of*  
747 *Adult and Child Learners in Language Formation and Change.*, *Language Learning &*  
748 *Development*, 1 (2005), pp. 151–195.
- 749 [20] C. KEMP, A. PERFORIS, AND J. B. TENENBAUM, *Learning overhypotheses with hierarchical*  
750 *Bayesian models*, *Developmental Science*, 10 (2007), pp. 307–321.
- 751 [21] S. KIRBY, *Spontaneous evolution of linguistic structure: An iterated learning model of the*  
752 *emergence of regularity and irregularity*, *IEEE Transactions on Evolutionary Computation*,  
753 5 (2001), pp. 102–110.
- 754 [22] S. KIRBY AND J. R. HURFORD, *The Emergence of Structure: An overview of the iterated*  
755 *learning model*, in *Simulating the Evolution of Language*, 2002, pp. 121–148.
- 756 [23] N. L. KOMAROVA, P. NIYOGI, AND M. A. NOWAK, *The Evolutionary dynamics of grammar*  
757 *acquisition*, *Journal of Theoretical Biology*, 209 (2001), pp. 43–59.
- 758 [24] A. KROCH, *Reflexes of Grammar in Patterns of Language Change*, *Language Variation and*  
759 *Change*, 1 (1989), pp. 199–244.
- 760 [25] W. LABOV, *Principles of Linguistic Change: Internal Factors*, vol. 1, Blackwell, Malden, MA,  
761 1994.
- 762 [26] W. LABOV, *Principles of Linguistic Change: Social Factors*, vol. 2, Blackwell, Malden, MA,  
763 2001.
- 764 [27] W. LABOV, *Transmission and Diffusion*, *Language*, 83 (2007), pp. 344–387.
- 765 [28] D. LIGHTFOOT, *How to Set Parameters: Arguments from Language Change*, MIT Press, Cam-  
766 bridge, MA, 1991.
- 767 [29] D. LIGHTFOOT, *The Development of Language: Acquisition, Changes and Evolution*, Blackwell,  
768 Malden, MA, 1999.
- 769 [30] B. LINDNER AND L. SCHIMANSKY-GEIER, *Analytical approach to the stochastic FitzHugh-*  
770 *Nagumo system and coherence resonance*, *Physical Review E*, 60 (1999), p. 7270.
- 771 [31] R. S. MAIER AND D. L. STEIN, *The Escape Problem for Irreversible Systems*, *chao-dyn/9303017*,  
772 (1993).
- 773 [32] R. S. MAIER AND D. L. STEIN, *A Scaling Theory of Bifurcations in the Symmetric Weak-Noise*  
774 *Escape Problem*, *cond-mat/9506097*, (1995).
- 775 [33] W. G. MITCHENER, *Bifurcation analysis of the fully symmetric language dynamical equation*,  
776 *Journal of Mathematical Biology*, 46 (2003), pp. 265–285.
- 777 [34] W. G. MITCHENER, *A Mathematical model of the loss of verb-second in Middle English*, in  
778 *Medieval English and its Heritage*, no. 16 in *Studies in English Medieval Language and*  
779 *Literature*, Peter Lang Pub Inc, 2006, pp. 189–202.
- 780 [35] W. G. MITCHENER, *Game dynamics with learning and evolution of universal grammar*, *Bulletin*  
781 *of Mathematical Biology*, 69 (2007), pp. 1093–1118.
- 782 [36] W. G. MITCHENER, *Inferring leadership structure from data on a syntax change in English*, in  
783 *Scientific applications of language methods*, no. 2 in *Frontiers in Mathematical Linguistics*  
784 *and Language Theory*, Imperial College Press, London, 2010, pp. 633–662.
- 785 [37] W. G. MITCHENER, *Mean-field and measure-valued differential equation models for language*  
786 *variation and change in a spatially distributed population*, *SIAM Journal on Mathematical*  
787 *Analysis*, 42 (2010), pp. 1899–1933.
- 788 [38] W. G. MITCHENER, *A Mathematical model of prediction-driven instability: How social struc-*  
789 *ture can drive language change*, *Journal of Logic, Language and Information*, 20 (2011),  
790 pp. 385–396.
- 791 [39] W. G. MITCHENER AND M. BECKER, *Computational models of learning the raising-control*

- 792 *distinction*, Research on Language and Computation, 8 (2011), pp. 169–207.
- 793 [40] W. G. MITCHENER AND M. A. NOWAK, *Competitive exclusion and coexistence of universal*  
794 *grammars*, Bulletin of Mathematical Biology, 65 (2003), pp. 67–93.
- 795 [41] W. G. MITCHENER AND M. A. NOWAK, *Chaos and language*, Proceedings of the Royal Society  
796 of London B: Biological Sciences, 271 (2004), pp. 701–704.
- 797 [42] J. D. MURRAY, *Mathematical Biology*, vol. 1, Springer-Verlag, New York, 2002.
- 798 [43] P. NIYOGI, *The Computational Nature of Language Learning and Evolution*, MIT Press, Cam-  
799 bridge, MA, 2006.
- 800 [44] P. NIYOGI AND R. C. BERWICK, *A Language learning model for finite parameter spaces*, Cog-  
801 nition, 61 (1996), pp. 161–193.
- 802 [45] M. A. NOWAK, *Evolutionary Dynamics: Exploring the equations of life*, Harvard University  
803 Press, Cambridge, MA, 2006.
- 804 [46] M. A. NOWAK AND N. L. KOMAROVA, *Towards an evolutionary theory of language*, Trends in  
805 Cognitive Sciences, 5 (2001), pp. 288–295.
- 806 [47] M. A. NOWAK, N. L. KOMAROVA, AND P. NIYOGI, *Evolution of universal grammar*, Science,  
807 291 (2001), pp. 114–118.
- 808 [48] M. A. NOWAK, N. L. KOMAROVA, AND P. NIYOGI, *Computational and evolutionary aspects of*  
809 *language*, Nature, 417 (2002), pp. 611–617.
- 810 [49] M. A. NOWAK, D. C. KRAKAUER, AND A. DRESS, *An error limit for the evolution of language*,  
811 Proceedings of the Royal Society of London. Series B: Biological Sciences, 266 (1999),  
812 pp. 2131–2136.
- 813 [50] M. A. NOWAK, J. PLOTKIN, AND D. C. KRAKAUER, *The Evolutionary Language Game*, Journal  
814 of Theoretical Biology, 200 (1999), pp. 147–162.
- 815 [51] M. A. NOWAK, J. B. PLOTKIN, AND V. A. A. JANSEN, *The evolution of syntactic communica-*  
816 *tion*, Nature, 404 (2000), pp. 495–498.
- 817 [52] B. ØKSENDAL, *Stochastic Differential Equations*, Universitext, Springer Berlin Heidelberg,  
818 Berlin, Heidelberg, 1985.
- 819 [53] K. M. PAGE AND M. A. NOWAK, *Unifying Evolutionary Dynamics*, Journal of Theoretical  
820 Biology, 219 (2002), pp. 93–98.
- 821 [54] L. PEARL AND A. WEINBERG, *Input Filtering in Syntactic Acquisition: Answers From Language*  
822 *Change Modeling*, Language Learning and Development, 3 (2007), pp. 43–72.
- 823 [55] A. PERFORIS, J. B. TENNENBAUM, AND E. WONNACOTT, *Variability, Negative Evidence, and*  
824 *the Acquisition of Verb Argument Constructions*, Journal of Child Language, 37 (2010),  
825 pp. 607–642.
- 826 [56] J. B. PLOTKIN AND M. A. NOWAK, *Language evolution and information theory*, Journal of  
827 Theoretical Biology, 205 (2000), pp. 147–159.
- 828 [57] A. RADFORD, *Minimalist Syntax: Exploring the Structure of English*, Cambridge University  
829 Press, Cambridge, UK, June 2004.
- 830 [58] F. REALI AND T. L. GRIFFITHS, *Words as alleles: connecting language evolution with Bayesian*  
831 *learners to models of genetic drift*, Proceedings of the Royal Society of London B: Biological  
832 Sciences, 277 (2010), pp. 429–436.
- 833 [59] S. SWARUP AND L. GASSER, *Unifying evolutionary and network dynamics*, Physical Review E,  
834 75 (2007), p. 066114.
- 835 [60] B. TESAR AND P. SMOLENSKY, *Learnability in Optimality Theory*, MIT Press, Cambridge, MA,  
836 May 2000.
- 837 [61] P. E. TRAPA AND M. A. NOWAK, *Nash Equilibria for an evolutionary language game*, Journal  
838 of Mathematical Biology, 41 (2000), pp. 172–1888.
- 839 [62] A. WARNER, *Why DO Dove: Evidence for Register Variation in Early Modern English Nega-*  
840 *tives*, Language Variation and Change, 17 (2005), pp. 257–280.
- 841 [63] S. WICHMANN AND E. W. HOLMAN, *Population Size and Rates of Language Change*, Human  
842 Biology, 81 (2009), pp. 259–274.
- 843 [64] S. WICHMANN, D. STAUFFER, C. SCHULZE, AND E. W. HOLMAN, *Do language change rates*  
844 *depend on population size?*, Advances in Complex Systems, 11 (2008), pp. 357–369.
- 845 [65] C. D. YANG, *Knowledge and Learning in Natural Language*, Oxford University Press, Oxford,  
846 2002.
- 847 [66] W. ZUIDEMA AND B. DE BOER, *The evolution of combinatorial phonology*, Journal of Phonetics,  
848 37 (2009), pp. 125–144.