

# A Dynamical Systems Model for Language Change

Partha Niyogi\*

Robert C. Berwick†

*Artificial Intelligence Laboratory,  
Center for Biological and Computational Learning,  
Massachusetts Institute of Technology,  
Cambridge, MA 02139*

---

This paper formalizes linguists' intuitions about language change, proposing a dynamical systems model for language change derived from a model for language acquisition. Linguists must explain not only how languages are learned but also how and why they have evolved along certain trajectories and not others. While the language learning problem has focused on the behavior of individuals and how they acquire a particular grammar from a class of grammars  $\mathcal{G}$ , this paper considers a *population* of such learners and investigates the emergent, global population characteristics of linguistic communities over several generations. It is argued that language change follows logically from specific assumptions about grammatical theories and learning paradigms. Roughly, as the end product of two types of learning misconvergence over several generations, *individual* language learner behavior leads to emergent, *population* language community characteristics.

In particular, it is shown that any triple  $\{\mathcal{G}, \mathcal{A}, \mathcal{P}\}$  of grammatical theory, learning algorithm, and initial sentence distributions can be transformed into a dynamical system whose evolution depicts the evolving linguistic composition of a population. It is explicitly shown how this transformation can be carried out for memoryless learning algorithms and parameterized grammatical theories. As the simplest case, the example of two grammars (languages) differing by exactly one binary parameter is formalized, and it is shown that even this situation leads directly to a quadratic (nonlinear) dynamical system, including regions with chaotic behavior. The computational model is applied to some actual data, namely the observed historical loss of "verb second" from old French to modern French. Thus, the formal model allows one to pose new questions about language phenomena that one otherwise could not ask, such as the following.

1. Do languages (grammars) correctly follow observed historical trajectories? This is an evolutionary criteria for the adequacy of grammatical theories.

---

\*Electronic mail address: [niyogi@research.bell-labs.com](mailto:niyogi@research.bell-labs.com).

†Electronic mail address: [berwick@ai.mit.edu](mailto:berwick@ai.mit.edu).

2. What are the logically possible dynamical change envelopes given a posited grammatical theory? These are rates and shapes of linguistic change, including the possibilities for the past and the future.

3. What can be the effect of quantified variation in initial conditions? For example, population differences resulting from socio-political facts.

4. Other intrinsically interesting mathematical questions regarding linguistic dynamical systems.

---

## 1. Introduction: The paradox of language change

---

Much research on language has focused on how children acquire the grammar of their parents from “impoverished” data presented to them during childhood. The logical problem of language acquisition, cast formally, requires the learner to converge (attain) its correct target grammar (i.e., the language of its caretakers, that belongs to a class of possible natural language grammars). However, this learnability problem, if solved perfectly, would lead to a paradox: If generation after generation children successfully attained the grammar of their parents, then languages would never change with time. Yet languages do change.

Language scientists have long been occupied with describing phonological, syntactic, and semantic change, often appealing to an analogy between language change and evolution, but rarely going beyond this. For instance, in [11] language change is talked about in this way:

Some general properties of language change are shared by other dynamic systems in the natural world... In population biology and linguistic change there is constant flux... If one views a language as a totality, as historians often do, one sees a *dynamic system*.

Indeed, entire books have been devoted to the description of language change using the terminology of population biology: genetic drift, clines, and so forth. For a recent example, see Nichols (1992), *Linguistic Diversity in Space and Time*. Other scientists have explicitly made an appeal to dynamical systems in this context; see especially Hawkins and Gell-Mann, 1989. Yet to the best of our knowledge, these intuitions have not been formalized. That is the goal of this paper. A remarkable effort in quantifying cultural change that has many potential unexploited applications to language is developed in [2].

In particular, we show formally that a model of language change emerges as a *logical* consequence of language acquisition, an argument made informally by Lightfoot in [11]. We shall see that Lightfoot’s intuition that languages could behave just as though they were dynamical systems is essentially correct, as is his proposal for turning language ac-

quisition models into language change models. We can provide concrete examples of both “gradual” and “sudden” syntactic changes, occurring over time periods of many generations to just a single generation. In [11] these sudden changes acting over a single generation are referred to as “catastrophic” but this term usually has a different meaning in the dynamical systems literature.

Many interesting points emerge from the formalization, some empirical, some programmatic.

1. Learnability is a well-known criterion for the adequacy of grammatical theories. Our model provides an *evolutionary* criterion: By comparing the trajectories of dynamical linguistic systems to historically observed trajectories, one can determine the adequacy of linguistic theories or learning algorithms.

2. We derive explicit dynamical systems corresponding to parameterized linguistic theories (e.g., the head first/final parameter in head-driven phrase structure grammars or government-binding grammars) and memoryless language learning algorithms (e.g., gradient ascent in parameter space).

3. In the simplest possible case of a two-language (grammar) system differing by exactly one binary parameter, the system reduces to a quadratic map with the usual chaotic properties (dependent on initial conditions). That such complexity can arise even in the simplest case suggests that formally modeling language change may be quite mathematically rich.

4. We illustrate the use of dynamical systems as a research tool by considering the loss of verb second position in old French as compared to modern French. We demonstrate by computer modeling that one grammatical parameterization advanced in the linguistics literature does not seem to permit this historical change, while another does.

5. We can more accurately model the time course of language change. In particular, in contrast to [10] and others, who mimic population biology models by imposing an S-shaped logistic change by *assumption*, we explain the time course of language change, and show that it need not be S-shaped. Rather, language-change envelopes are *derivable* from more fundamental properties of dynamical systems; sometimes they are S-shaped, but they can also be nonmonotonic.

6. We examine by simulation and traditional phase-space plots the form and stability of possible “diachronic envelopes” given varying alternative language distributions, language acquisition algorithms, parameterizations, input noise, and sentence distributions. The results bear on models of language “mixing,” so-called “wave” models for language change, and other proposals in the diachronic literature.

7. As topics for future research, the dynamical system model provides a novel possible source for explaining several linguistic changes including the evolution of modern Greek metrical stress assignment from proto-Indo-European and Bickerton's (1990) "creole hypothesis" concerning the striking fact that all creoles, irrespective of linguistic origin, have exactly the same grammar. In the latter case, the "universality" of creoles could be due to a parameterization corresponding to a common condensation point of a dynamical system, a possibility not considered in Bickerton.

## **2. An acquisition-based model of language change**

How does the combination of a grammatical theory and learning algorithm lead to a model of language change? We first note that just as with language acquisition, there is a seeming paradox in language change: it is generally assumed that children acquire their caretaker (target) grammars without error. However, if this were always true, at first glance grammatical changes within a population could seemingly never occur, since generation after generation children would successfully acquire the grammar of their parents.

Of course, Lightfoot and others have pointed out the obvious solution to this paradox: the possibility of slight misconvergence to target grammars could, over several generations, drive language change, much as speciation occurs in the population biology sense [11]:

As somebody adopts a new parameter setting, say a new verb-object order, the output of that person's grammar often differs from that of other people's. This in turn affects the linguistic environment, which may then be more likely to trigger the new parameter setting in younger people. Thus a chain reaction may be created.

We pursue this point in detail below. Similarly, just as in the biological case, some of the most commonly observed changes in languages seem to occur as the result of the effects of surrounding populations, whose features infiltrate the original language.

We begin our treatment by arguing that the problem of language acquisition at the individual level leads logically to the problem of language change at the group or population level. Consider a population speaking a particular language. In our analysis this implies that all the adult members of this population have internalized the same grammar (corresponding to the language they speak). This is the target language—children are exposed to primary linguistic data (PLD) from this source, typically in the form of sentences uttered by caretakers (adults). The logical problem of language acquisition is how children acquire this target language from their PLD, that is, to come up with

an adequate learning theory. We take a learning theory to be simply a mapping from PLD to the class of grammars, usually effective, and so an algorithm. For example, in a typical inductive inference model, given a stream of sentences, an acquisition algorithm would simply update its grammatical hypothesis with each new sentence according to some pre-programmed procedure. An important criterion for learnability (Gold, 1967) is to require that the algorithm converge to the target as the data goes to infinity (identification in the limit).

Now suppose that we fix an adequate grammatical theory and an adequate acquisition algorithm. There are then essentially two means by which the linguistic composition of the population could change over time. First, if the PLD data presented to the child is altered (due to any number of causes, perhaps to presence of foreign speakers, contact with another population, disfluencies, and the like), the sentences presented to the learner (child) are no longer consistent with a single target grammar. In the face of this input, the learning algorithm might not converge to the target grammar. Indeed, it might converge to some other grammar ( $g_2$ ); or it might converge to  $g_2$  with some probability,  $g_3$  with some other probability, and so forth. In either case, children attempting to solve the acquisition problem using the same learning algorithm could internalize grammars different from the parental (target) grammar. In this way, in one generation the linguistic composition of the population can change. Sociological factors affecting language change affect language acquisition in exactly the same way, yet are abstracted away from the formalization of the logical problem of language acquisition. In this same sense, we similarly abstract away such causes here, though they can be brought into the picture as variation in probability distributions and learning algorithms; we leave this open as a topic for additional research.

Second, even if the PLD comes from a single target grammar, the actual data presented to the learner is truncated, or finite. After a finite sample sequence, children may, with nonzero probability, hypothesize a grammar different from that of their parents. This can again lead to a differing linguistic composition in succeeding generations.

In short, the diachronic model is this: Individual children attempt to attain the target grammar of their caretakers. After a finite number of examples, some are successful, but others may misconverge. The next generation will therefore no longer be linguistically homogeneous. The third generation of children will hear sentences produced by the second—a different distribution—and they, in turn, will attain a different set of grammars. Over successive generations, the linguistic composition evolves as a dynamical system.

In this view, language change is a logical consequence of specific assumptions about the following.

1. The *grammar hypothesis space*—a particular parametrization, in a parametric theory.
2. The *language acquisition device*—the learning algorithm the child uses to develop hypotheses on the basis of data.
3. The *primary linguistic data*—the sentences presented to the children of any one generation.

If we specify 1 through 3 for a particular generation, we should, in principle, be able to compute the linguistic composition for the next generation. In this manner, we can compute the evolving linguistic composition of the population from generation to generation, that is, we arrive at a dynamical system. We now proceed to make this calculation precise. We first review a standard language acquisition framework, and then show how to derive a dynamical system from it.

## ■ 2.1 The language acquisition framework

To formalize the model, we must first state our assumptions about grammatical theories, learning algorithms, and sentence distributions.

1. Denote by  $\mathcal{G}$  a family of possible (target) grammars. Each grammar  $g \in \mathcal{G}$  defines a language  $L(g) \subseteq \Sigma^*$  over some alphabet  $\Sigma$  in the usual way.
2. Denote by  $P$  a distribution on  $\Sigma^*$  according to which sentences are drawn and presented to the learner. Note that if there is a well defined target  $g_t$  and only positive examples from this target are presented to the learner, then  $P$  will have all its measure on  $L(g_t)$ , and zero measure on sentences outside. Suppose  $n$  examples are drawn in this fashion, one can then let  $\mathcal{D}_n = (\Sigma^*)^n$  be the set of all  $n$ -example data sets the learner might be presented with. Thus, if the adult population is linguistically homogeneous (with grammar  $g_1$ ) then  $P = P_1$ . If the adult population speaks 50 percent  $L(g_1)$  and 50 percent  $L(g_2)$  then  $P = 1/2P_1 + 1/2P_2$ .
3. Denote by  $\mathcal{A}$  the acquisition algorithm that children use to hypothesize a grammar on the basis of input data.  $\mathcal{A}$  can be regarded as a mapping from  $\mathcal{D}_n$  to  $\mathcal{G}$ . Acting on a particular presentation sequence  $d_n \in \mathcal{D}_n$ , the learner posits a hypothesis  $\mathcal{A}(d_n) = h_n \in \mathcal{G}$ . Allowing for the possibility of randomization, the learner could, in general, posit  $h_i \in \mathcal{G}$  with probability  $p_i$  for such a presentation sequence  $d_n$ .

The standard (stochastic version) learnability criterion (Gold, 1967) can then be stated as follows.

For every target grammar  $g_t \in \mathcal{G}$  with positive-only examples presented according to  $P$  as above, the learner must converge to the target with probability 1, that is,

$$\text{Prob}[A(d_n) = g_t] \rightarrow_{n \rightarrow \infty} 1.$$

One *particular* way of formulating the learning algorithm is as a local gradient ascent search through a space of target languages (grammars) defined by a one-dimensional  $n$ -length boolean array of *parameters*, with each distinct array fixing a particular grammar (language). With  $n$  parameters, there are  $2^n$  possible grammars (languages). For example, English and Japanese differ in that English is a so-called “verb first” language, while Japanese is “verb final.” Given this framework, we can state the so-called triggering learning algorithm (TLA) from [5] as follows.

- Step 1. *Initialize.* Start at some random point in the (finite) space of possible parameter settings, specifying a single hypothesized grammar with its resulting extension as a language.
- Step 2. *Process input sentence.* Receive a positive example sentence  $s_i$  at time  $t_i$  (examples drawn from the language of a single target grammar  $L(G_t)$ ) from a uniform distribution on unembedded (nonrecursive) sentences of the target language.
- Step 3. *Learnability on error detection.* If the current grammar parses (generates)  $s_i$ , then go to Step 2; otherwise, continue.
- Step 4. *Single-step hill climbing.* Select a single parameter uniformly at random, to flip from its current setting, and change it (0 mapped to 1, 1 to 0) *if and only if that change allows the current sentence to be analyzed*; otherwise, leave the current parameter settings unchanged.
- Step 5. *Iterate.* Go to Step 2.

Of course, this algorithm carries out “identification in the limit” in the standard terminology of learning theory (Gold, 1967); it does not halt in the conventional sense.

It turns out that if the learning algorithm  $\mathcal{A}$  is memoryless (in the sense that previous example sentences are not stored) and  $\mathcal{G}$  can be described by a finite number of parameters, then we can describe the learning system  $\mathcal{A}$ ,  $G_t$ ,  $\mathcal{G}$  as a Markov chain  $M$  with as many states as there are grammars in  $\mathcal{G}$ . More specifically the states in  $M$  are in one-to-one correspondence with grammars  $g \in \mathcal{G}$  and the target grammar  $G_t$  corresponds to a particular target state  $s_t$  of  $M$ . The transition probabilities between states in  $M$  can be computed straightforwardly based on set difference calculations between the languages corresponding to the Markov chain states. We omit the details of this demonstration here; for a simple, explicit calculation in the case of one parameter, see section 3. For a more detailed analysis of learnability issues for memoryless algorithms in finite parameter spaces, consult [14, 18, 19].

## ■ 2.2 From language learning to population dynamics

The framework for language learning has learners attempting to infer grammars on the basis of linguistic data. At any point in time  $n$  (i.e.,

after hearing  $n$  examples) the learner has a current hypothesis  $h$  with probability  $p_n(h)$ . What happens when there is a population of learners? Since an arbitrary learner has a probability  $p_n(h)$  of developing hypothesis  $h$  (for every  $h \in \mathcal{G}$ ), it follows that a fraction  $p_n(h)$  of the population of learners internalize the grammar  $h$  after  $n$  examples. We therefore have a *current state* of the population after  $n$  examples. This state of the population might well be different from the state of the parent population. Assume for now that after  $n$  examples, maturation occurs, that is, after  $n$  examples the learner retains the grammatical hypothesis for the rest of its life. Then one would arrive at the state of the mature population for the next generation. This new generation now produces sentences for the following generation of learners according to the distribution of grammars in its population. Then, the process repeats itself and the linguistic composition of the population evolves from generation to generation.

We can now define a discrete time dynamical system by providing its two necessary components as follows.

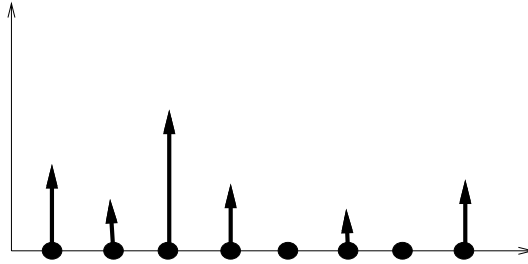
1. *A state space.* A set of system states  $\mathcal{S}$ . Here the state space is the space of possible linguistic compositions of the population. Each state is described by a distribution  $P_{\text{pop}}$  on  $\mathcal{G}$  describing the language spoken by the population. As usual, one needs to be able to define a  $\sigma$ -algebra on the space of grammars, and so on. This is unproblematic for the cases considered here because the set of grammars is finite. At any given point in time  $t$  the system is in exactly one state  $s \in \mathcal{S}$ .
2. *An update rule.* How the system states change from one time step to the next. Typically, this involves specifying a function  $f$  that maps  $s_t \in \mathcal{S}$  to  $s_{t+1}$ . In general, this mapping could be fairly complicated. For example, it could depend on previous states, future states, and so forth; for reasons of space we do not consider all possibilities here. For more information see Strogatz, (1993).

For example, a typical linear dynamical system might consist of state variables  $\mathbf{x}$  (where  $\mathbf{x}$  is a  $k$ -dimensional state vector) and a system of differential equations  $\mathbf{x}' = A\mathbf{x}$  ( $A$  is a matrix operator) which characterize the evolution of the states with time. RC circuits are a simple example of linear dynamical systems. The state (current) evolves as the capacitor discharges through the resistor. Population growth models (e.g., using logistic equations) provide other examples.

As a linguistic example, consider the three-parameter syntactic space described in [5]. This system defines eight possible “natural” grammars, that is,  $\mathcal{G}$  has eight elements. We can picture a distribution on this space as shown in Figure 1. In this particular case, the state space is

$$\mathcal{S} = \left\{ \mathbf{P} \in \mathbb{R}^8 \mid \sum_{i=1}^8 P_i = 1 \right\}.$$





**Figure 1.** A simple illustration of the state space for the three-parameter syntactic case. There are eight grammars. A probability distribution on these eight grammars, as shown above, can be interpreted as the linguistic composition of the population. Thus, a fraction  $P_1$  of the population have internalized grammar  $g_1$  and so on.

Here we interpret the *state* as the linguistic composition of the population. Note that we do not allow for the possibility of a *single* learner having more than one hypothesis at a time; an extension to this case, in which individuals would more closely resemble the “ensembles” of particles in a thermodynamic system, is left for future research. For example, a distribution that puts all its weight on grammar  $g_1$  and 0 everywhere else indicates a homogeneous population that speaks a language corresponding to grammar  $g_1$ . Similarly, a distribution that puts a probability mass of  $1/2$  on  $g_1$  and  $1/2$  on  $g_2$  denotes a population (non-homogeneous) with half its speakers speaking a language corresponding to  $g_1$  and half speaking a language corresponding to  $g_2$ .

To see in detail how the update rule may be computed, consider the acquisition algorithm  $\mathcal{A}$ . For example, given the state at time  $t$ , ( $P_{\text{pop},t}$ ), the distribution of speakers in the parental population, one can obtain the distribution with which sentences from  $\Sigma^*$  will be presented to the learner. To do this, imagine that the  $i$ th linguistic group in the population, speaking language  $L_i$ , produces sentences with distribution  $P_i$ . Then for any  $\omega \in \Sigma^*$ , the probability with which  $\omega$  is presented to the learner is given by

$$P(\omega) = \sum_i P_i(\omega)P_{\text{pop},t}(i).$$

This fixes the distribution with which sentences are presented to the learner. The logical problem of language acquisition also assumes some success criterion for attaining the mature target grammar. For our purposes, we take this as being one of two broad possibilities: either (1) the usual Gold scenario of identification in the limit, what we call the *limiting sample* case; or (2) identification in a fixed, finite time, what we

call the *finite sample* case. Of course, a variety of other success criteria, for example, convergence within some epsilon, or polynomial in the size of the target grammar, are possible; each leads to potentially different language change models. We do not pursue these alternatives here.

Consider case (2) first. Here, one draws  $n$  example sentences according to distribution  $P$ , and the acquisition algorithm develops hypotheses ( $\mathcal{A}(d_n) \in \mathcal{G}$ ). One can, in principle, compute the probability with which the learner will posit hypothesis  $h_i$  after  $n$  examples:

$$\text{Finite Sample: } \text{Prob}[\mathcal{A}(d_n) = h_i] = p_n(h_i). \quad (1)$$

The finite sample situation is always well defined, that is, the probability  $p_n$  always exists. This is easy to see for deterministic algorithms,  $\mathcal{A}_{\text{det}}$ . Such an algorithm would have a precise behavior for every data set of  $n$  examples drawn. In our case, the examples are drawn in i.i.d. fashion according to a distribution  $P$  on  $\Sigma^*$ . It is clear that  $p_n(h_i) = P[\{d_n | \mathcal{A}_{\text{det}}(d_n) = h_i\}]$ . For randomized algorithms, the case is trickier, though tedious, but the probability still exists because all the finite choice paths over all sequences of length  $n$  is enumerable. Previous work [15–19] shows how to compute  $p_n$  for randomized memoryless algorithms.

Now turn to case (1), the limiting case. Here learnability requires  $p_n(g_t)$  to go to 1 for the unique target grammar  $g_t$  if such a grammar exists. However, in general there need not be a unique target grammar since the linguistic population can be nonhomogeneous. Even so, the following limiting behavior might still exist:

$$\text{Limiting Sample: } \lim_{n \rightarrow \infty} \text{Prob}[\mathcal{A}(d_n) = h_i] = p(h_i). \quad (2)$$

Turning from the individual child to the population, since the individual child internalizes grammar  $h_i \in \mathcal{G}$  with probability  $p_n(h_i)$  in the “finite sample” case or with probability  $p(h_i)$  “in the limit,” in a *population* of such individuals one would therefore expect a proportion  $p_n(h_i)$  or  $p(h_i)$  respectively to have internalized grammar  $h_i$ . In other words, the linguistic composition of the *next* generation is given by  $P_{\text{pop},t+1}(h_i) = p_n(h_i)$  for the finite sample case and by  $P_{\text{pop},t+1}(h_i) = p(h_i)$  in the limiting sample case. In this fashion,

$$P_{\text{pop},t} \xrightarrow{\mathcal{A}} P_{\text{pop},t+1}.$$

### Remarks

1. For a Gold-learnable family of languages and a limiting sample assumption, homogeneous populations are always stable. This is simply because each child and therefore the entire population always eventually converges to a single target grammar, generation after generation.

2. However, the finite sample case is different from the limiting sample case. Suppose we have solved the maturation problem, that is, we know

roughly the time, or number of examples  $N$  the learner takes to develop its mature (adult) hypothesis. In that case  $p_N(b)$  is the probability that a child internalizes the grammar  $b$ , and  $p_N(b)$  is the percentage of speakers of  $L_b$  in the next generation. Note that under this finite sample analysis, even for a homogeneous population with all adults speaking a particular language (corresponding to grammar,  $g$ , say),  $p_N(g)$  will not be 1, that is, there will be a small percentage of learners who have misconverged. This percentage could blow up over several generations, and we therefore have potentially unstable languages.

3. The formulation is very general. Any  $\{\mathcal{A}, \mathcal{G}, \mathcal{P}\}$  triple yields a dynamical system. Note that this probability could evolve with generations as well, which would complete all the logical possibilities. However, for simplicity, we assume that this does not happen. In short:

$$(\mathcal{G}, \mathcal{A}, \{P_i\}) \longrightarrow \mathcal{D}(\text{dynamical system}).$$

4. The formulation also does not assume any particular linguistic theory, learning algorithm, or distribution with which sentences are drawn. Of course, we have implicitly assumed a learning model, that is, positive examples are drawn in i.i.d. fashion and presented to the learner. Our dynamical systems formalization follows as a logical consequence of this learning framework. One can conceivably imagine other learning frameworks—these would potentially give rise to other kinds of dynamical systems—but we do not formalize them here.

In previous works [15–19] we investigated the problem of learnability within parametric systems. In particular, we showed that the behavior of any memoryless algorithm can be modeled as a Markov chain. This analysis allows us to solve equations (1) and (2), and thus obtain the update equations for the associated dynamical system. Let us now show how to derive such models in detail. We first provide the particular  $\mathcal{G}, \mathcal{A}, \{P_i\}$  triple, and then give the update rule.

#### The learning system triple

$\mathcal{G}$ : Assume there are  $n$  parameters, this leads to a space  $\mathcal{G}$  with  $2^n$  different grammars.

$\mathcal{A}$ : Let us imagine that the child learner follows some memoryless (incremental) algorithm to set parameters. For the most part, we will assume that the algorithm is the TLA (the single step, gradient-ascent algorithm of [5]) or one of the variants discussed in [18, 19].

$\{P_i\}$ : Let speakers of the  $i$ th language  $L_i$  in the population produce sentences according to the distribution  $P_i$ . For the most part we will assume in our simulations that this distribution is uniform on degree-0 (unembedded) sentences, exactly as in the learnability analysis of [5] or [18, 19].

**The update rule**

We can now compute the update rule associated with this triple. Suppose the state of the parental population is  $P_{\text{pop},n}$  on  $\mathcal{G}$ . Then one can obtain the distribution  $P$  on the sentences of  $\Sigma^*$  according to which sentences will be presented to the learner. Once such a distribution is obtained, then given the Markov equivalence established earlier, we can compute the transition matrix  $T$  according to which the learner updates its hypotheses with each new sentence. From  $T$  one can finally compute the following quantities, one for the finite sample case and one for the limiting sample case:

$$\begin{aligned} & \text{Prob}[\text{Learner's hypothesis} = h_i \in \mathcal{G} \text{ after } m \text{ examples}] \\ &= \left\{ \frac{1}{2^n} (1, \dots, 1)' T^m \right\} [i]. \end{aligned}$$

Similarly, making use of the limiting distributions of Markov chains (Resnick, 1992) one can obtain the following (where ONE is a  $1/2^n \times 1/2^n$  matrix with all ones).

$$\begin{aligned} & \text{Prob}[\text{Learner's hypothesis} = h_i \text{ "in the limit"}] \\ &= (1, \dots, 1)' (I - T + \text{ONE})^{-1}. \end{aligned}$$

These expressions allow us to compute the linguistic composition of the population from one generation to the next according to our analysis of the previous section.

**Remark**

The limiting distribution case is more complex than the finite sample case and requires some careful explanation. There are two possibilities. If there is just a single target grammar, then, by definition, the learners all identify the target correctly in the limit, and there is no further change in the linguistic composition from generation to generation. This case is essentially uninteresting. If there are two or more target grammars, then recalling our analysis of learnability [18, 19], there can be no absorbing states in the Markov chain corresponding to the parametric grammar family. In this situation, a single learner will oscillate between some set of states in the limit. In this sense, learners will not converge to any single, correct target grammar. However, there is a sense in which we can characterize limiting behavior for learners: although a given learner will visit each of these states infinitely often in the limit, it will visit some more often than others. The exact percentage the learner will be in a particular state is given by equation (2). Therefore, since we know the fraction of the time the learner spends in each grammatical state in the limit, we assume that this is the probability with which it internalizes the grammar corresponding to that state in the Markov chain.

The following summarizes the basic computational framework for modeling language change.

1. Let  $\pi_1$  be the initial population mix, that is, the percentage of different language speakers in the community. Assuming that the  $i$ th group of speakers produces sentences with probability  $P_i$ , we can obtain the probability  $P$  with which sentences in  $\Sigma^*$  occur for the next generation of learners.
2. From  $P$  we can obtain the transition matrix  $T$  for the Markov learning model and the limiting distribution of the linguistic composition  $\pi_2$  for the next generation.
3. The second generation now has a population mix of  $\pi_2$ . We repeat step 1 and obtain  $\pi_3$ . Continuing in this fashion, in general we can obtain  $\pi_{i+1}$  from  $\pi_i$ .

This completes the abstract formulation of the dynamical system model. Next, we choose three specific linguistic theories and learning paradigms to model particular kinds of language changes, with the goal of answering the following questions.

- Can we really compute all the relevant quantities to specify the dynamical system?
- Can we evaluate the behavior (phase-space characteristics) of the resulting dynamical system?
- Does the dynamical system model, the formalization, shed light on diachronic models and linguistic theories generally?

In the remainder of this paper we give some concrete answers to these questions within the principles and parameters theory of modern linguistic theory. We turn first to the simplest possible mathematical case, that of two languages (grammars) fixed by a single binary parameter. We then analyze a possibly more relevant, and more complex system, with three binary parameters. Finally, to tackle a more realistic historical problem, we consider a five-parameter system that has actually been used in other contexts to account for language change.

### 3. One-parameter models of language change

Consider the following simple scenario.

- $\mathcal{G}$ : Assume that there are only two possible grammars (parameterized by one boolean valued parameter) associated with two languages in the world,  $L_1$  and  $L_2$ . (This might in fact be true in some limited linguistic contexts.)
- $\mathcal{P}$ : Suppose that speakers who have internalized grammar  $g_1$  produce sentences with a probability distribution  $P_1$  (on the sentences of  $L_1$ ). Similarly, assume that speakers who have internalized grammar  $g_2$  produce

sentences with  $P_2$  (on sentences of  $L_2$ ).

One can now define

$$a = P_1[L_1 \cap L_2]; \quad 1 - a = P_1[L_1 \setminus L_2]$$

and similarly

$$b = P_2[L_1 \cap L_2]; \quad 1 - b = P_2[L_2 \setminus L_1].$$

$\mathcal{A}$ : Assume that the learner uses a one-step, greedy, hill climbing approach to setting target parameters. The TLA described earlier is one such example.

$\mathcal{N}$ : Let the learner receive just two example sentences before maturation occurs, that is, after two example sentences, the current grammatical hypothesis of the learner will be retained for the rest of its life.

Given this framework, the learnability question for this parametric system can be easily formulated and analyzed. Specifically, given a particular target grammar ( $g_i \in \mathcal{G}$ ), and given example sentences drawn according to  $P_i$  and presented to the learner  $\mathcal{A}$ , one can ask whether the hypothesis of the learner will converge to the target.

Now it is possible to characterize the behavior of the individual learner by a Markov chain with two states, one corresponding to each grammar (see section 2 and [18, 19]). With each example the learner moves from state to state according to the transition probabilities of the chain. The transition probabilities can be calculated and depend upon the distribution with which sentences are drawn and the relative overlap between the languages  $L_1$  and  $L_2$ . In particular, if received sentences follow distribution  $P_1$ , the transition matrix is  $T_1$ . This would be the case if the target grammar were  $g_1$ . If the target grammar is  $g_2$  and sentences received according to distribution  $P_2$ , the transition matrix would be  $T_2$  as shown:

$$T_1 = \begin{bmatrix} 1 & 0 \\ 1 - a & a \end{bmatrix}$$

$$T_2 = \begin{bmatrix} b & 1 - b \\ 0 & 1 \end{bmatrix}.$$

Let us examine  $T_1$  in order to understand the behavior of the learner when  $g_1$  is the target grammar. If the learner starts out in state 1 (initial hypothesis  $g_1$ ), then it remains there forever. This is because every sentence that it receives can be analyzed and the learner will never have to entertain an alternative hypothesis. Therefore the transition ( $1 \rightarrow 1$ ) has probability 1 and the transition ( $1 \rightarrow 2$ ) has probability 0. If the learner starts out in state 2, then after one example sentence, the learner will remain there with probability  $a$ —the probability that the learner will receive a sentence that it can analyze, that is, a sentence in

$L_1 \cup L_2$ . Correspondingly, with probability  $1 - a$ , the learner will receive a sentence that it cannot analyze and will have to change its hypothesis. Thus, the transition  $(2 \rightarrow 2)$  has probability  $a$  and the transition  $(2 \rightarrow 1)$  has probability  $1 - a$ .

$T_1$  characterizes the behavior of the learner after *one* example. In general,  $T_1^k$  characterizes the behavior of the learner after  $k$  examples. It may be easily seen that as long as  $a < 1$ , the learner converges to the grammar  $g_1$  in the limit irrespective of its starting state. Thus the grammar  $g_1$  is Gold-learnable.

A similar analysis can be carried out for the case when the target grammar is  $g_2$ . In this case,  $T_2$  describes the corresponding behavior of the learner, and  $g_2$  is Gold-learnable if  $b < 1$ . In short, the entire system is Gold-learnable if  $a, b < 1$ , crucially assuming that maturation occurs and the learner fixes a hypothesis forever after some  $N$  examples, with  $N$  given in advance. Clearly, if  $N$  is very large, then the learner will, with high probability, acquire the unique target grammar, whatever that grammar might be. At the same time, there is a finite probability that the learner will misconverge and this will have consequences for the linguistic composition of the population as discussed in section 2. For the analysis that follows, we will assume that  $N = 2$ .

■ **3.1 One-parameter systems: The linguistic population**

Continuing with our one-parameter model, we next analyze distributions over speakers. At any given point in time, the population consists only of speakers of  $L_1$  and  $L_2$ . Consequently, the linguistic composition can be represented by a single variable,  $p$ : this will denote the fraction of the population speaking  $L_1$ . Clearly  $1 - p$  will speak  $L_2$ . Therefore this community of language composition over time can be explicitly computed as follows.

**Theorem 1.** The linguistic composition in the  $n + 1$ th ( $p_{n+1}$ ) generation is provided by the following transformation on the linguistic composition of  $n$ th generation ( $p_n$ ):

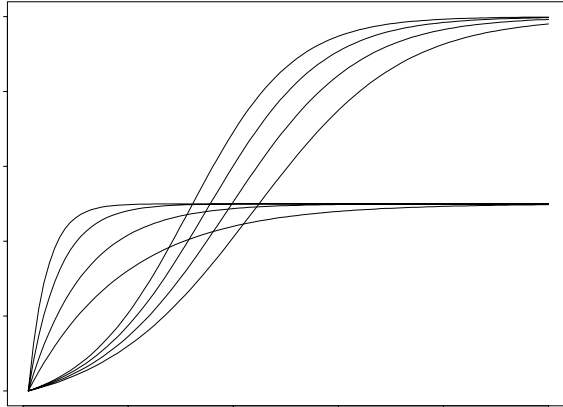
$$p_{n+1} = Ap_n^2 + Bp_n + C$$

where  $A = 1/2((1 - b)^2 - (1 - a)^2)$ ,  $B = b(1 - b) + (1 - a)$ , and  $C = b^2/2$ .

*Proof.* This is a simple specialization of the formula given in section 2. Details are left to the reader. ■

**Remarks**

1. When  $a = b$ , the system has exponential growth. When  $a \neq b$  the dynamical system is a quadratic map (which can be reduced by a transformation of variables to the logistic, and shares the dynamical properties of the logistic). See Figure 2.



**Figure 2.** Evolution of linguistic populations whose speakers differ only in a single, “verb second” parameter value, speaking two languages  $L_1$  and  $L_2$ . This reduces to a one-parameter model as discussed in the text. As we vary the probabilities that speakers of the two languages produce sentences in the intersection of  $L_1$  and  $L_2$ ,  $a$  and  $b$  respectively, we get differently shaped curves. When  $a = b$  the growth is exponential, with different shapes for different values of  $a$  and  $b$  (less than 1.0). When  $a$  is not equal to  $b$  the system has a qualitatively different shape, a logistic growth.

2. The scenario  $a \neq b$ , that the distributions for  $L_1$  and  $L_2$  differ, is much more likely to occur in the real world. Consequently, we are more likely to see logistic growth rather than exponential. Indeed, various parametric shifts observed historically seem to follow an S-shaped curve. Models of these shifts have typically *assumed* the logistic growth [10]. Crucially, in this simple case, the logistic form has now been *derived* as a consequence of specific assumptions about how learning occurs at the individual level, rather than *assumed*, as in all previous models for language change that we are aware of.

3. We get a class of dynamical systems. The quadratic nature of our map comes from the fact that  $N = 2$ . If we choose other values for  $N$  we would get cubic and higher order maps. In other words, there are already an infinite number of maps in the simple one-parameter case. For larger parametric systems the mathematical situation is significantly more complex.

4. Logistic maps are known to be chaotic. In our system, it is possible to show the following.

**Theorem 2.** Due to the fact that  $a, b \leq 1$ , the dynamical system never enters a chaotic regime.



This observation naturally raises the question of whether nonchaotic behavior holds for all grammatical dynamical systems, specifically the linguistically “natural” cases. Or are there linguistic systems where chaos will manifest itself? It would obviously be quite interesting if all the natural grammatical spaces were nonchaotic. We leave these as open questions.

We next turn to more linguistically plausible applications of the dynamical systems model. We begin with a simple three-parameter system as our first example, considering variations on the learning algorithm, sentence distributions, and sample size available for learning. We then consider a different, five-parameter system already presented in the literature [3] as one intended to partially characterize the change from old French to modern French.

#### 4. A three-parameter system

In section 3 we developed the necessary mathematical and computational tools to completely specify the dynamical systems corresponding to memoryless algorithms operating on finite parameter spaces. In this section we investigate the behavior of these dynamical systems. Recall that every choice of  $(\mathcal{G}, \mathcal{A}, \{P_i\})$  gives rise to a unique dynamical system. We start by making specific choices for these three elements as follows.

- $\mathcal{G}$ : This is a three-parameter syntactic subsystem described in [5]. Thus  $\mathcal{G}$  has exactly eight grammars, generating languages from  $L_1$  through  $L_8$ , as shown in the appendix of this paper (taken from [10]).
- $\mathcal{A}$ : The memoryless algorithms we consider are the TLA, and variants by dropping either or both of the single-valued and greediness constraints.
- $\{P_i\}$ : For the most part, we assume sentences are produced according to a uniform distribution on the degree-0 sentences of the relevant language, that is,  $P_i$  is uniform on (degree-0 sentences of)  $L_i$ .

Ideally of course, a complete investigation of diachronic possibilities would involve varying  $\mathcal{G}$ ,  $\mathcal{A}$ , and  $\mathcal{P}$  and characterizing the resulting dynamical systems by their phase-space plots. Rather than explore this entire space, we first consider only systems evolving from homogeneous initial populations, under four basic variants of the learning algorithm  $\mathcal{A}$ . This will give us an initial grasp of how linguistic populations can change. Indeed, linguistic change has been studied before; even the dynamical system metaphor itself has been invoked. Our computational paradigm lets us say much more than these previous descriptions: We can say precisely what the rates of change will be and we can determine what diachronic population curve changes will look like, without stipulating in advance that they must be S-shaped (sigmoid) or not, and without curve fitting to a predefined functional form.

#### 4.1 Homogeneous initial populations

First we consider the case of a homogeneous population, that is, without noise or confounding factors like foreign target languages. How stable are the languages in the three-parameter system in this case? To determine this, we begin with a finite-sample analysis with  $n = 128$  example sentences (recall by the analysis of [15–19] that learners converge to target languages in the three-parameter system with high probability after hearing this many sentences). Some small proportion of the children misconverge; the goal is to see whether this small proportion can drive language change—and if so, in what direction. To give the reader some idea of the possible outcomes, let us consider the four possible variations in the learning algorithm ( $\pm$ Single Step,  $\pm$ Greedy) keeping the sentence distributions and learning sample fixed.

##### 4.1.1 Variation 1: $\mathcal{A} = \text{TLA (+Single Step, +Greedy)}$ ; $P_i = \text{Uniform}$ ; Finite Sample = 128

Suppose the learning algorithm is the TLA. Table 1 shows the language mix after 30 generations. Languages are numbered from 1 to 8. Recall that +V2 refers to a language that has the verb second property, and –V2 one that does not.

##### Observations

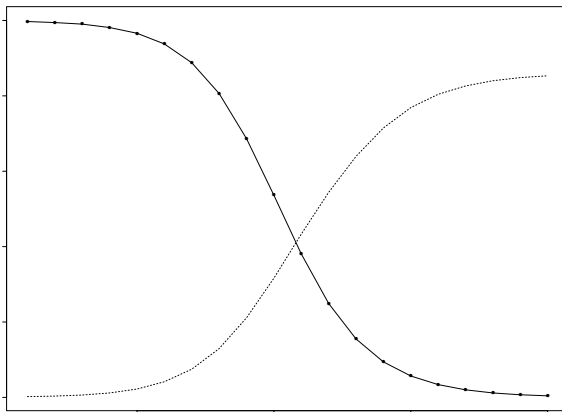
Some striking patterns regarding the resulting population mixes can be noted.

1. All the +V2 languages are relatively stable, that is, the linguistic composition did not vary significantly over 30 generations. This means that every succeeding generation acquired the target parameter settings and no parameter drifts were observed over time.

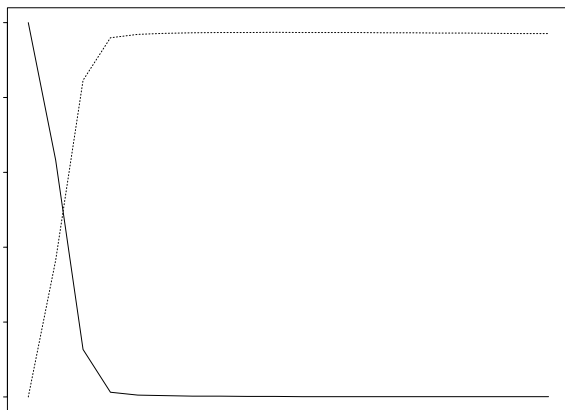
Initial Language	Change to Language?
(–V2) 1	2 (0.85), 6 (0.1)
(+V2) 2	2 (0.98); stable
(–V2) 3	6 (0.48), 8(0.38)
(+V2) 4	4 (0.86); stable
(–V2) 5	2 (0.97)
(+V2) 6	6 (0.92); stable
(–V2) 7	2 (0.54), 4(0.35)
(+V2) 8	8 (0.97); stable

**Table 1.** Language change driven by misconvergence from a homogeneous initial linguistic population. A finite-sample analysis was conducted allowing each child learner 128 examples to internalize its grammar. After 30 generations, initial populations drifted (or not, as shown in the table) to different final linguistic compositions.

2. In contrast, populations speaking  $-V2$  languages all drift to  $+V2$  languages. Thus a population speaking  $L_1$  winds up speaking mostly  $L_2$  (85%). A population speaking language  $L_7$  gradually shifts to a population with 54 percent speaking  $L_2$  and 35 percent speaking  $L_4$  (with a smattering of other speakers) and apparently remains basically stable in this mix thereafter. Note that the relative stability of  $+V2$  languages and the tendency of  $-V2$  languages to drift to  $+V2$  is exactly contrary to evidence in the linguistic literature. For example, in [11] it is claimed that the tendency to lose  $V2$  dominates the reverse tendency in languages of the world. Certainly, both English and French lost the  $V2$  parameter setting, an empirically observed phenomenon that needs to be explained. Immediately then, we see that our dynamical system does not evolve in the expected manner. The reason could be due to any of the assumptions behind the model: the parameter space, the learning algorithm, the initial conditions, or the distributional assumptions about sentences presented to learners. Exactly which is in error remains to be seen, but nonetheless our example shows concretely how assumptions about a grammatical theory and learning theory can make evolutionary, diachronic predictions—in this case, incorrect predictions that falsify the assumptions.



**Figure 3.** Percentage of a population speaking languages  $L_1$  and  $L_2$ , measured on the  $y$ -axis, as the population evolves over some number of generations, measured on the  $x$ -axis. The plot has been shown only up to 20 generations, as the proportions of  $L_1$  and  $L_2$  speakers do not vary significantly thereafter. Note that this curve is S-shaped. In [7] such a shape is *imposed* using models from population biology, while we derive this shape as an emergent property of our dynamical model.  $L_1$  and  $L_2$  differ only in the  $V2$  parameter setting.



**Figure 4.** Percentage of the population speaking languages  $L_5$  and  $L_2$  as the population evolves over a number of generations. Note that a complete shift from  $L_5$  to  $L_2$  occurs after just four generations.

3. The *rates* at which the linguistic composition changes vary significantly from language to language. Consider for example the change of  $L_1$  to  $L_2$ . Figure 3 shows the gradual decrease in speakers of  $L_1$  over successive generations along with the increase in  $L_2$  speakers. We see that over the first six or seven generations very little change occurs, but over the next six or seven generations the population changes at a much faster rate. Note that in this particular case the two languages differ only in the V2 parameter, so the curves essentially plot the gain of V2. In contrast, consider Figure 4 which shows the decrease of  $L_5$  speakers and the shift to  $L_2$ . Here we note a sudden change: over a space of just four generations, the population shifts completely. Analysis of the time course of language change has been given some attention in linguistic analyses of diachronic syntax change, and we return to this issue later.

4. We see that in many cases a homogeneous population splits up into different linguistic groups, and seems to remain stable in that mix. In other words, certain combinations of language speakers seem to asymptote towards equilibrium (at least through 30 generations). For example, a population of  $L_7$  speakers shifts over five or six generations to one with 54 percent speaking  $L_2$  and 35 percent speaking  $L_4$  and remains that way with no shifts in the distribution of speakers. Of course, we do not know for certain whether this is really a stable mixture. It could be that the population mix could suddenly shift after another 100 generations. What we really need to do is characterize the stable points or “limit cycles” of these dynamical systems. Other linguistic mixes can be

inherently unstable; they might drift systematically to stable situations, or might shift dramatically (as with language  $L_1$ ).

5. It seems that the observed instability and drifts are to a large extent an artifact of the learning algorithm. Remember that the TLA suffers from the problem of local maxima. We regard local maxima of a language  $L_i$  to be alternative absorbing states (sinks) in the Markov chain for that target language. This formulation differs slightly from the conception of local maxima in [5], a matter discussed at some length in [15]. Thus, according to our definition,  $L_4$  is not a local maxima for  $L_5$  and consequently no shift is observed. We note that those languages whose acquisition is not impeded by local maxima (the +V2 languages) are stable over time. Languages that have local maxima are unstable; in particular they drift to the local maxima over time. Now consider  $L_7$ . If this is the target language, then there are two local maxima ( $L_2$  and  $L_4$ ) and these are precisely the states to which the system drifts over time. The same is true for languages  $L_5$  and  $L_3$ . In this respect, the behavior of  $L_1$  is quite unusual since it actually does not have any local maxima, yet it tends to flip the V2 parameter over time.

Now let us consider a learning algorithm different from the TLA that does not suffer from local maxima problems, to see whether this changes the dynamical system results.

**4.1.2 Variation 2:  $\mathcal{A} = +\text{Greedy}$ ,  $-\text{Single Value}$ ;  $P_i = \text{Uniform}$ ;  
Finite Sample = 128**

Consider a simple variant of the TLA obtained by dropping the single-valued constraint. This implies that the learner is no longer constrained to change just one parameter at a time: on being presented with a sentence it cannot analyze, it chooses any of the alternative grammars and attempts to analyze the sentence with it. Greediness is retained; thus the learner retains its original hypothesis if the new one is also not able to analyze the sentence. Given this new learning algorithm, and retaining all the other original assumptions, Table 2 shows the distribution of speakers after 30 generations.

**Observations**

In this situation there are no local maxima, and the evolutionary pattern takes on a very different nature. There are two distinct observations to be made.

1. All homogeneous populations eventually drift to a strikingly similar population mix, irrespective of what language they start from. What is unique about this mix? Is it a stable point (or attractor)? Further simulations and theoretical analyses are needed to resolve this question.

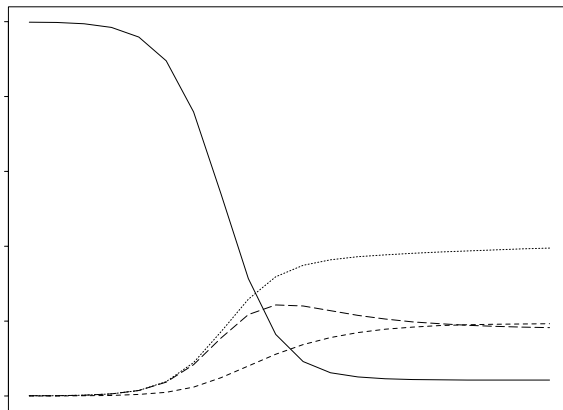
2. All homogeneous populations drift to a population mix of only +V2 languages. Thus, the V2 parameter is gradually set over succeed-

Initial Language	Change to Language?
-V2 1	2 (0.41), 4 (0.19), 6 (0.18), 8 (0.13)
+V2 2	2 (0.42), 4 (0.19), 6 (0.17), 8 (0.12)
-V2 3	2 (0.40), 4 (0.19), 6 (0.18), 8 (0.13)
+V2 4	2 (0.41), 4 (0.19), 6 (0.18), 8 (0.13)
-V2 5	2 (0.40), 4 (0.19), 6 (0.18), 8 (0.13)
+V2 6	2 (0.40), 4 (0.19), 6 (0.18), 8 (0.13)
-V2 7	2 (0.40), 4 (0.19), 6 (0.18), 8 (0.13)
+V2 8	2 (0.40), 4 (0.19), 6 (0.18), 8 (0.13)

**Table 2.** Language change driven by misconvergence. A finite-sample analysis was conducted allowing each child learner (following the TLA with single-value dropped) 128 examples to internalize its grammar. Initial populations were linguistically homogeneous, and they drifted to different linguistic compositions. The major language groups after 30 generations have been listed in this table. Note how all initially homogeneous populations tend to the same composition.

ing generations by all people in the community (irrespective of which language they speak). In other words, as before, there is a tendency to gain V2 rather than lose V2, contrary to the empirical facts.

As an example, Figure 5 shows the changing percentage of the population speaking the different languages starting from a homogeneous population speaking  $L_5$ . As before, learners who have not converged to the target in 128 examples are the driving force for change here. Note again the time evolution of the grammars. For about five generations



**Figure 5.** Time evolution of grammars using a greedy learning algorithm with no single value constraint in place.

there is only a slight decrease in the percentage of speakers of  $L_5$ . Then the linguistic patterns switch rapidly over the next seven generations to a relatively stable mix.

**4.1.3 Variations 3 & 4: –Greedy,  $\pm$ Single Value constraint;  $P_i = \text{Uniform}$ ; Finite Sample = 128**

Having dropped the single-value constraint, we consider the next obvious variation in the learning algorithm: dropping greediness while varying the single-value constraint. Again, our goal is to see whether this makes any difference in the resulting dynamical system. This gives rise to the following two different learning algorithms.

1. Allow the learning algorithm to pick any new grammar at most one parameter value away from its current hypothesis (retaining the single value constraint, but without greediness, i.e., the new grammar does not have to be able to parse the current input sentence).
2. Allow the learning algorithm to pick any new grammar at each step (no matter how far away from its current hypothesis).

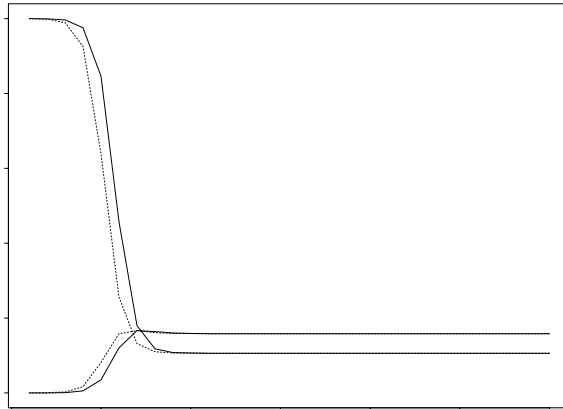
In both cases, the population mix after 30 generations is the same irrespective of the initial language of the homogeneous population. These results are shown in Table 3.

**Observations**

1. Both algorithms yield dynamical systems that arrive at the same population mix after 30 generations. The path by which they arrive at this mix is, however, not the same (see Figure 6).
2. The final population mix contains all languages in significant proportion. This is in distinct contrast to the previous situations, where we saw that –V2 languages were eliminated over time.

Initial Language	Change to Language?
Any Language (Homogeneous)	1 (0.11), 2 (0.16), 3 (0.10), 4 (0.14) 5 (0.12), 6 (0.14), 7 (0.10), 8 (0.13)

**Table 3.** Language change driven by misconvergence, using two different acquisition algorithms that do not obey a local gradient-ascent rule (a greediness constraint). A finite-sample analysis was conducted with the learning algorithm following a random-step algorithm or else a single-step algorithm, along with 128 examples to internalize its grammar. Initial populations were linguistically homogeneous, and they drifted to different linguistic compositions. The major language groups after 30 generations have been listed in this table. Note that all initially homogeneous populations converge to the same final composition.



**Figure 6.** Time evolution of linguistic composition for the situations where the learning algorithm is  $-$ Greedy,  $+$ Single Value constraint (dotted line), and  $-$ Greedy,  $-$ Single Value (solid line). Only the percentage of people speaking  $L_1$  ( $-V2$ ) and  $L_2$  ( $+V2$ ) are shown. The initial population is homogeneous and speaks  $L_1$ . The percentage of  $L_1$  speakers gradually decreases to about 11 percent. The percentage of  $L_2$  speakers rises to about 16 percent from 0 percent. The two dynamical systems converge to the same population mix; however, their trajectories are not the same—the rates of change are different, as shown in this plot.

#### ■ 4.2 Modeling diachronic trajectories

With a basic notion of how diachronic systems can evolve given different learning algorithms, we turn next to the question of population trajectories. While we can already see that some evolutionary trajectories have a “linguistically classical” S-shape, their smoothness can vary. However, our formalization allows us to say much more than this. Unlike the previous work in diachronic linguistics that we are familiar with, we can explore the space of possible trajectories, examining factors that affect their evolutionary time course, without assuming an *a priori* S-shape.

For example, in Bailey (1973) a “wave” model of linguistic change is proposed: linguistic replacements follow an S-shaped curve over time. In Bailey’s own words (from [10]):

A given change begins quite gradually; after reaching a certain point (say, twenty percent), it picks up momentum and proceeds at a much faster rate; and finally tails off slowly before reaching completion. The result is an S-curve: the statistical differences among isolects in the middle relative times of the change will be greater than the statistical differences among the early and late isolects.



The idea that linguistic changes follow an S-curve has also been proposed in [20, 22]. More specific logistic forms have been advanced in [4, 8, 9]. Here, the idea of a logistic functional form is borrowed from population biology where it is demonstrable that the logistic governs the replacement of organisms and of genetic alleles that differ in darwinian fitness. However, it is conceded in [9] that “unlike in the population biology case, no mechanism of change has been proposed from which the logistic form can be deduced.”

Crucially, in our case, we suggest a specific mechanism of change: an acquisition-based model where the combination of grammatical theory, learning algorithms, and distributional assumptions on sentences drive change. The specific form might or might not be S-shaped, and might have varying rates of change. Of course, we do not mean to say that we can simulate *any* possible trajectory—that would make the formalism empty. Rather, we are exploring the initial space of possible trajectories, given some example initial conditions that have been already advanced in the literature. Because the mathematics for dynamical systems is in general quite complex, at present we cannot make general statements of the form, “under these particular initial conditions the trajectory will be sigmoidal, and under these other conditions it will not be.” We have conducted only very preliminary investigations demonstrating that potentially at least, reasonable, distinct initial conditions can lead to demonstrably different trajectories.

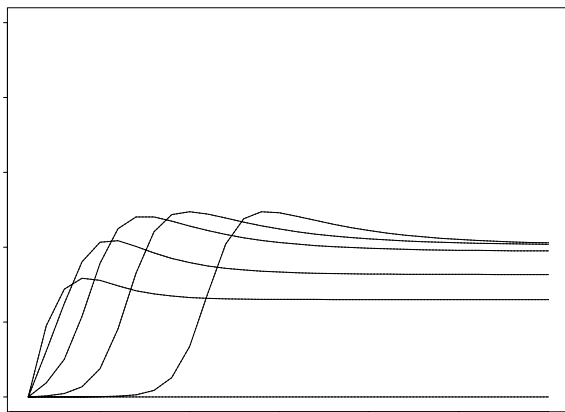
Among the other factors that affect evolutionary trajectories are maturation time, that is, the number of sentences available to the learner before it internalizes its adult grammar, and the distributions with which sentences are presented to the learner. We examine these in turn.

#### 4.2.1 The effect of maturation time or sample size

One obvious factor influencing the evolutionary trajectories is the maturational time, that is, the number  $N$  of sentences the child is allowed to hear before forming its mature hypothesis. This was fixed at 128 in all the systems shown so far (based in part on our explicit computation for the Markov convergence time in this situation). Figure 7 shows the effect of varying  $N$  on the evolutionary trajectories. As usual, we plot only a subspace of the population. In particular, we plot the percentage of  $L_2$  speakers in the population with each succeeding generation. The initial composition of the population was homogeneous (with people speaking  $L_1$ ).

#### Observations

1. The initial rate of change of the population is highest when the maturation time is smallest, that is, the learner is allowed the least amount of time to develop its mature hypothesis. This is not surprising. If the learner were allowed access to a lot of examples to make its mature



**Figure 7.** Time evolution of linguistic composition when varying maturation time (sample size). The learning algorithm used is the +Greedy, –Single Value. Only the percentage of people speaking  $L_2$  (+V2) is shown. The initial population is homogeneous and speaks  $L_1$ . The maturation time was varied through 8, 16, 32, 64, 128, and 256, giving rise to the six curves shown. The curve with the highest initial rate of change corresponds to eight examples for maturation time. The initial rate of change decreases as the maturation time  $N$  increases. The value at which these curves asymptote also seems to vary with the maturation time, and increases monotonically with it.

hypothesis, most learners would reach the target grammar. Very few would misconverge, and the linguistic composition would change little over the next generation. On the other hand, if the learner were allowed very few examples to develop its hypothesis, many would misconverge, possibly causing great change over one generation.

2. The “stable” linguistic compositions seem to depend upon maturation time. For example, if learners are allowed only eight examples, the percentage of  $L_2$  speakers rises quickly to about 0.26. On the other hand, if learners are allowed 128 examples, the percentage of  $L_2$  speakers eventually rises to about 0.41.

3. Note that the trajectories *do not* have an S-shaped curve in contrast to the results in [9].

4. The maturation time is related to the order of the dynamical system.

#### 4.2.2 The effect of sentence distributions ( $P_i$ )

Another important factor influencing evolutionary trajectories is the distribution  $P_i$  with which sentences of the  $i$ th language  $L_i$  are presented to the learner. In a certain sense, the grammatical space and the learning

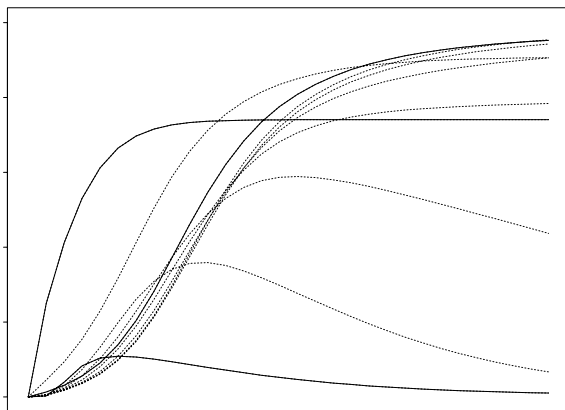
algorithm jointly determine the order of the dynamical system. On the other hand, sentence distributions are much like the parameters of the dynamical system (see section 4.3.2). Clearly the sentence distributions affect rates of convergence within one generation. Further, by putting greater weight on certain word forms rather than others, they might influence systemic evolution in certain directions. While this is again an obvious point, the model lets us consider the alternatives precisely.

To illustrate the idea, consider as an example the interaction between  $L_1$  and  $L_2$  speakers in the community as the sentence distributions with which these speakers produce sentences changes. Recall that so far we have assumed that all speakers produce sentences with uniform distributions on degree-0 sentences of their respective languages. Now we consider alternative distributions, parameterized by a value  $p$  as follows.

1. Let  $L_{1,2} = L_1 \cap L_2$ .
2.  $P_1$ : Speakers of  $L_1$  produce sentences so that all degree-0 sentences of  $L_{1,2}$  are equally likely and their total probability is  $p$ . Further, sentences of  $L_1 \setminus L_{1,2}$  are also equally likely, but their total probability is  $1 - p$ .
3.  $P_2$ : Speakers of  $L_2$  produce sentences so that all degree-0 sentences of  $L_{1,2}$  are equally likely and their total probability is  $p$ . Further, sentences of  $L_2 \setminus L_{1,2}$  are also equally likely, but their total probability is  $1 - p$ .
4. Other  $P_i$  are all uniform over degree-0 sentences.

The parameter  $p$  determines the weight on the sentence patterns in common between the languages  $L_1$  and  $L_2$ . Figure 8 shows the evolution of the  $L_2$  speakers as  $p$  varies. Here the learning algorithm is +Greedy, +Single Value (TLA, or local gradient ascent) and the initial population is homogeneous, 100 percent  $L_1$  and zero percent  $L_2$ . Note that the system moves in different ways as  $p$  varies. When  $p$  is very small (0.05), that is, sentences common to  $L_1$  and  $L_2$  occur infrequently, in the long run the percentage of  $L_2$  speakers *does not* increase; the population stays put with  $L_1$ . However, as  $p$  grows, more strings of  $L_2$  occur, and the dynamical system changes so that the long-term percentage of  $L_1$  speakers decreases and that of  $L_2$  speakers increases. When  $p$  reaches 0.75 the initial population evolves into a completely  $L_2$  speaking community. After this, as  $p$  increases further, we notice that the  $L_2$  speakers increase but can never rise to 100 percent of the population (see  $p = 0.95$ ); there is still a residual  $L_1$  speaking component. This is to be expected, because for such high values of  $p$ , many strings common to  $L_1$  and  $L_2$  occur frequently. This means that a learner could sometimes converge to  $L_1$  just as well as  $L_2$ , and some learners indeed begin to do so, increasing the number of the  $L_1$  speakers.

This example shows us that if we wanted a homogeneous  $L_1$  speaking population to move to a homogeneous  $L_2$  speaking population, by



**Figure 8.** The evolution of  $L_2$  speakers in the community for various values of  $p$  (a parameter related to the sentence distributions  $P_i$ , see text). The algorithm used was the TLA, the initial population was homogeneous, speaking only  $L_1$ . The curves for  $p = 0.05, 0.75$ , and  $0.95$  have been plotted as solid lines.

choosing our distributions appropriately, we could drive the grammatical dynamical system in the appropriate direction. It suggests another important application of the dynamical system approach: one can work backwards, and examine the conditions needed to generate a change of a certain kind. By checking whether such conditions could have possibly existed historically, we can falsify a grammatical theory or a learning paradigm. Note that this example showed the effect of sentence distributions, and how to alter them to obtain desired evolutionary envelopes. One could, in principle, alter the grammatical theory or the learning algorithm in the same fashion—leading to a tool to aid the search for an adequate linguistic theory. Again, we stress that we obviously do not want so weak a theory that we can arrive at any possible initial conditions simply by carrying out reasonable changes to the sentence distributions. This may, of course, be possible; we have not yet examined the general case.

#### ■ 4.3 Nonhomogeneous populations: Phase-space plots

For our three-parameter system we have been able to characterize the update rules for the dynamical systems corresponding to a variety of learning algorithms. Each dynamical system has a specific update procedure according to which the states evolve from some *homogeneous* initial population. A more complete characterization of the dynamical system would be achieved by obtaining *phase-space* plots of this system. Such phase-space plots are pictures of the state-space  $\mathcal{S}$  filled with

trajectories obtained by letting the system evolve from various initial points (states) in the state space.

#### 4.3.1 Phase-space plots: Grammatical trajectories

We described earlier the relationship between the state of the population in one generation and the next. In our case, let  $\Pi$  denote an eight-dimensional vector variable (state variable). Specifically,  $\Pi = (\pi_1, \dots, \pi_8)'$  (with  $\sum_{i=1}^8 \pi_i$ ) as discussed before. The following schema reiterates the chain of dependencies involved in the update rule governing system evolution. The state of the population at time  $t$  (in generations), allows us to compute the transition matrix  $T$  for the Markov chain associated with the memoryless learner. Now, depending upon whether we want (1) an asymptotic analysis or (2) a finite sample analysis, we compute (1) the limiting behavior of  $T^m$  as  $m$  (the number of examples) goes to infinity (for an asymptotic analysis), or (2) the value of  $T^N$  (where  $N$  is the number of examples after which maturation occurs). This allows us to compute the next state of the population. Thus  $\Pi(t+1) = g(\Pi(t))$  where  $g$  is a complex nonlinear relation:

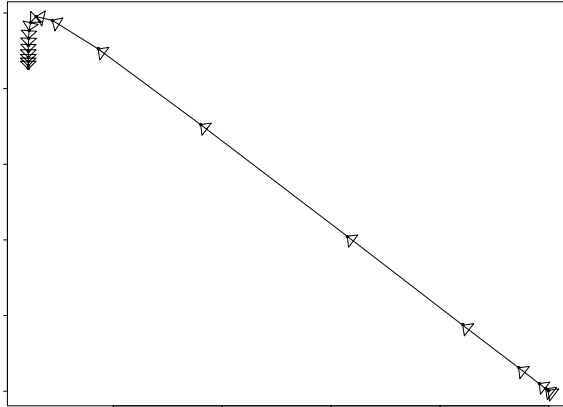
$$\Pi(t) \implies P \text{ on } \Sigma^* \implies T \implies T^m \implies \Pi(t+1).$$

If we choose a certain initial condition  $\Pi_1$ , the system will evolve according to the above relation and one can obtain a trajectory of  $\Pi$  in the eight-dimensional space over time. Each initial condition yields a unique trajectory and one can then plot these trajectories obtaining a phase-space plot. Each such trajectory corresponds to a line in the eight-dimensional plane given by  $\sum_{i=1}^8 \pi_i = 1$ . One cannot directly display such a high dimensional object, but we plot in Figure 9 the projection of a particular trajectory onto a two-dimensional subspace given by  $(\pi_1(t), \pi_2(t))$  (the proportion of speakers of  $L_1$  and  $L_2$ ) at different points in time.

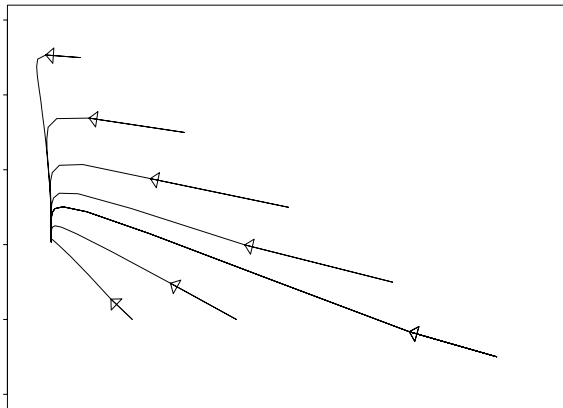
As mentioned earlier, with a different initial condition we get a different grammatical trajectory. The complete state-space picture is thus filled with all the different trajectories corresponding to different initial conditions. Figure 10 shows this.

#### 4.3.2 Stability issues

The phase-space plots show that many initial conditions yield trajectories that seem to converge to a single point in the state space. In the dynamical systems terminology, this corresponds to a fixed point of the system—a population mix that stays at the same composition. Many natural questions arise at this stage. What are the conditions for stability? How many fixed points are there in a given system? How can we solve for them? These are interesting questions but detailed answers are not within the scope of the current paper. In lieu of a more complete analysis we merely state here the equations that allow one to characterize the stable population mixes.



**Figure 9.** Subspace of a phase-space plot. The plot shows  $(\pi_1(t), \pi_2(t))$  as  $t$  varies, that is, the proportion of speakers speaking languages  $L_1$  and  $L_2$  in the population. The initial state of the population was homogeneous (speaking language  $L_1$ ). The algorithm used was +Greedy, -Single Value.



**Figure 10.** Subspace of a phase-space plot. The plot shows  $(\pi_1(t), \pi_2(t))$  as  $t$  varies for different nonhomogeneous initial population conditions. The algorithm used was +Greedy, -Single Value.

First, some notational preliminaries. As before, let  $P_i$  be the distribution on the sentences of the  $i$ th language  $L_i$ . From  $P_i$ , we can construct  $T_i$ , the transition matrix whose elements are given by the explicit procedure documented in [15–19]. The matrix  $T_i$  models a +Greedy, -Single Value learner if the target language is  $L_i$  (with sentences from the target produced with  $P_i$ ). Similarly, one can obtain the matrices for other

learning variants. Note that fixing the  $P_i$  fixes the  $T_i$  and in doing so the  $P_i$  are a different sort of “parameter” that characterize how the dynamical system evolves. There are thus two *distinct* kinds of parameters in our model: first, parameters that define the  $2^n$  languages and define the state-space of the system; and second, the  $P_i$  that characterize the way in which the system evolves and are therefore the parameters of the complete grammatical dynamical system. If the state of the parent population at time  $t$  is  $\Pi(t)$ , then it is possible to show that the (true) transition matrix for  $\pm$ Greedy,  $\pm$ Single Value learners is  $T = \sum_{i=1}^8 \pi_i(t)T_i$ . For the finite case analysis, the following holds.

**Statement 1 (finite case)**

A fixed point of the grammatical dynamical system (obtained by a  $\pm$ Greedy,  $\pm$ Single Value learner operating on the eight-parameter space with  $k$  examples to choose its final hypothesis) is a solution of the following equation:

$$\Pi' = (\pi_1, \dots, \pi_8) = (1, \dots, 1)' \left( \sum_{i=1}^8 \pi_i T_i \right)^k .$$

This equation is obtained simply by setting  $\Pi(t + 1) = \Pi(t)$ . Note however, that this is an example of a nonlinear multidimensional iterated function map. The analysis of such dynamical systems is nontrivial and beyond the scope of the current paper.

Similarly, for the limiting (asymptotic) case, the following holds.

**Statement 2 (limiting or asymptotic analysis)**

A fixed point of the grammatical dynamical system (obtained by a  $\pm$ Greedy,  $\pm$ Single Value learner operating on the eight-parameter space (given infinite examples to choose its mature hypothesis) is a solution of the following equation:

$$\Pi' = (\pi_1, \dots, \pi_8) = (1, \dots, 1)' \left( I - \sum_{i=1}^8 \pi_i T_i + \text{ONE} \right)^{-1} ,$$

where ONE is the  $8 \times 8$  matrix with all its entries equal to 1.

Again this is trivially obtained by setting  $\Pi(t + 1) = \Pi(t)$ . The expression on the right provides an analytical expression for the update equation in the asymptotic case. See [21] for details. All the caveats mentioned before in the finite case statement apply here as well.

**Remark**

We have just touched the surface as far as the theoretical characterization of these grammatical dynamical systems are concerned. The main purpose of this paper is to show that these dynamical systems exist as

a logical consequence of assumptions about the grammatical space and an acquisition theory. We have exhibited only some preliminary simulations with these systems. From a theoretical perspective, it would be much more valuable to have complete characterizations of such systems. In Strogatz (1993) it is suggested that nonlinear multidimensional mappings with greater than three dimensions are likely to be chaotic. It is also interesting to note that iterated function maps define fractal sets. Such investigations are beyond the scope of this paper, and might well be a fruitful area for further research.

## 5. From old French to modern French: An analysis revisited

So far, our examples have been based on a three-parameter linguistic theory for which we derived several different dynamical systems. Our goal was to concretely instantiate our philosophical arguments, sketching the factors that influence evolutionary trajectories. In this section, we briefly consider a different parametric linguistic system studied in [3]. The historical context in which Clark and Roberts advanced their linguistic proposal is the evolution of modern French from old French. Their parameters are intended to capture some, but of course not all, of this change. They too use a learning algorithm—in their case, a genetic algorithm—to account for historical change but do not analyze their model from the dynamical systems viewpoint. Here we adopt their parameterization, with all its strengths and weaknesses, but consider an alternative learning paradigm and the dynamical systems approach.

Extensive simulations in section 4 reveal that while the learnability problem of the three-parameter space can be solved by stochastic hill climbing algorithms, the long-term evolution of these algorithms have a behavior that is at variance with the diachronic change actually observed in historical linguistics. In particular, we saw how there was a tendency to gain rather than lose the V2 parameter setting. While this could well be an artifact of the class of learning algorithms considered, a more likely explanation is that loss of V2 (observed in many languages of the world such as French, English, and so forth) is due to an interaction of parameters and triggers other than those considered in section 4. We investigate this possibility and begin by first reviewing the alternative parametric theory in [2].

### 5.1 The parametric subspace and data

We now consider a syntactic space involving five (boolean-valued) parameters. We do not attempt to describe these parameters. The interested reader should consult [3, 6] for details.

$p_1$ : Case assignment under agreement ( $p_1 = 1$ ) or not ( $p_1 = 0$ ).

$p_2$ : Case assignment under government ( $p_2 = 1$ ) or not ( $p_2 = 0$ ). Relevant triggers for this parameter include “Adv V S” and “S V O.”



$p_3$ : Nominative clitics.

$p_4$ : Null Subject. Here relevant triggers would include “wh V S O.”

$p_5$ : Verb-second V2. Triggers include “Adv V S” and “S V O.”

These five parameters define a 32-grammar space. Each grammar in this parametrized system can be represented by a string of five bits depending upon the values of  $p_1, \dots, p_5$ , for instance, the first bit position corresponds to case assignment under agreement. We can now look at the surface strings (sentences) generated by each such grammar. For the purpose of explaining how old French changed to modern French, the following key sentences are considered in [2]. The parameter settings required to generate each sentence are provided in brackets; an asterisk is a “does not matter” value and an “X” means any phrase.

The relevant data

adv V S	[*1**1]
SVO	[*1**1] or [1***0]
wh V S O	[*1***]
wh V S O	[**1**]
X (pro) V O	[*1*11] or [1**10]
X V s	[**1*1]
X s V	[**1*0]
X S V	[1***0]
(S) V Y	[*1*11]

The parameter settings provided in brackets set the grammars which generate the sentence. For example, the sentence form “adv V S” (corresponding to *quickly ran John*, an incorrect word order in English), is generated by all grammars that have case assignment under government (the second element of the array set to 1,  $p_2 = 1$ ) and verb second movement ( $p_5 = 1$ ). The other parameters can be set to any value. Clearly there are eight different grammars that can generate (alternatively parse) this sentence. Similarly there are 16 grammars that generate the form S V O (eight corresponding to parameter settings of [\*1\*\*1] and eight corresponding to parameter settings of [1\*\*\*0]) and four grammars that generate (S) V Y.

#### Remark

Note that the sentence set considered in [2] is only a subset of the total number of degree-0 sentences generated by the 32 grammars in question. In order to directly compare their model with ours, we have not attempted to expand the data set or fill out the space any further. As a result, all the grammars do not have unique extensional properties, that is, some generate the same set of sentences.

## ■ 5.2 The case of diachronic syntax change in French

Continuing with the analysis from [2], within this parameter space, it is historically observed that the language spoken in France underwent a parametric change from the twelfth century to modern times. In particular, they point out that both V2 and prodrop are lost, illustrated by examples similar to the following.

### *Loss of null subjects: pro-drop*

1. (old French; +pro drop)  
Si firent (pro) grant joie la nuit  
'thus (they) made great joy the night'
2. (modern French; -pro drop)  
\* Ainsi s'amusaient bien cette nuit  
'thus (they) had fun that night'

### *Loss of V2*

3. (old French; +V2)  
Lors oient ils venir un escoiz de tonnoire  
'then they heard come a clap of thunder'
4. (modern French; -V2)  
\* Puis entendirent-ils un coup de tonnerre  
'then they heard a clap of thunder'

In [2] it is observed that it has been argued this transition was brought about by the introduction of new word orders during the fifteenth and sixteenth centuries resulting in generations of children acquiring slightly different grammars and eventually culminating in the grammar of modern French. A brief reconstruction of the historical process (after [3]) runs as follows.

### **Old French; setting [11011]**

The language spoken in the twelfth and thirteenth centuries had verb-second movement and null subjects, both of which were dropped by the twentieth century. The sentences generated by the parameter settings corresponding to old French are the following.

Old French	
adv V S –	[* 1 ** 1]
S V O –	[* 1 ** 1] or [1 ** * 0]
wh V S O	[* 1 ***]
X (pro) V O	[* 1 * 11] or [1 ** 10]

Note that from this data set it appears that both the case agreement and nominative clitics parameters remain ambiguous. In particular, old

French is in a subset-superset relation with another language (generated by the parameter settings of 11111). In this case, possibly some kind of subset principle [1] could be used by the learner; otherwise it is not clear how the data would allow the learner to converge to the old French grammar in the first place. None of the  $\pm$ Greedy,  $\pm$ Single Value algorithms would converge uniquely to the grammar of old French.

The string (X)VS occurs with a frequency of 58 percent and SV(X) occurs with 34 percent in old French texts. It is argued that this frequency of (X)VS is high enough to cause the V2 parameter to trigger to +V2.

### Middle French

In middle French, the data is not consistent with any of the 32 target grammars (equivalent to a heterogenous population). Analysis of texts from that period reveal that some old forms (such as adv V S) decreased in frequency and new forms (like adv S V) increased. It is argued in [3] that such a frequency shift causes “erosion” of V2, brings about parameter instability, and ultimately convergence to the grammar of modern French. In this transition period (i.e., when middle French was spoken and written) the data is of the following form:

adv V S [ $*1^{**}1$ ]; SVO [ $*1^{**}1$ ] or [ $1^{***}0$ ]; wh V S O [ $*1^{***}$ ];  
wh V s O [ $**1^{**}$ ]; X (pro)V O [ $*1^{*}11$ ] or [ $1^{**}10$ ]; X V s [ $**1^{*}1$ ];  
X s V [ $**1^{*}0$ ]; X S V [ $1^{***}0$ ]; (s)VY [ $*1^{*}11$ ]

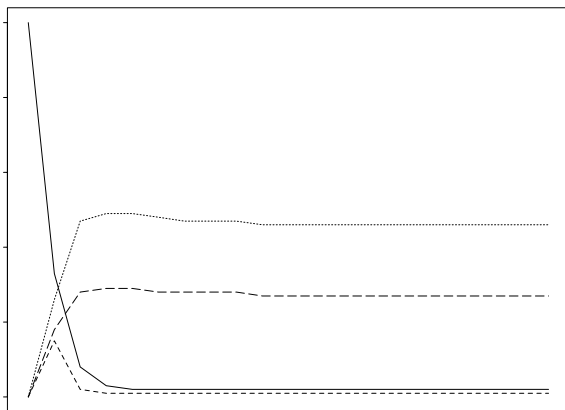
Thus, we have old sentence patterns like adv V S (though it decreases in frequency and becomes only 10%), SVO, X (pro)V O, and wh V S O. The new sentence patterns which emerge at this stage are adv S V (increases in frequency to become 60%), X subjclitic V, V subjclitic (pro)V Y (null subjects), and wh V subjclitic O.

### Modern French [10100]

By the eighteenth century, French had lost both the V2 parameter setting as well as the null subject parameter setting. The sentence patterns consistent with modern French parameter settings are SVO [ $*1^{**}1$ ] or [ $1^{***}0$ ], X S V [ $1^{***}0$ ], and V s O [ $**1^{**}$ ]. Note that this data, though consistent with modern French, will not trigger all the parameter settings. In this sense, modern French (just like old French) is not uniquely learnable from data. However, as before, we shall not concern ourselves overly with this, for the relevant parameters (V2 and null subject) are uniquely set by the data here.

## ■ 5.3 Some dynamical system simulations

We can obtain dynamical systems for this parametric space, for a TLA (or TLA-like) algorithm in a straightforward fashion. We show the results of two simulations conducted with such dynamical systems.



**Figure 11.** Evolution of speakers of different languages in a population starting off with speakers only of old French.

### 5.3.1 Homogeneous populations [initial–old French]

We conducted a simulation on this new parameter space using the TLA. Recall that the relevant Markov chain in this case has 32 states. We start the simulation with a homogeneous population speaking old French (parameter setting = 11011). Our goal is to see if misconvergence alone could drive old French to modern French.

Just as before, we can observe the linguistic composition of the population over several generations. It is observed that in one generation, 15 percent of the children converge to grammar 01011; 18 percent to grammar 01111; 33 percent to grammar 11011 (target) and 26 percent to grammar 11111 with very few having converged to other grammars. Thereafter, the population consists mostly of speakers of these four languages, with one important difference: 15 percent of the speakers eventually *lose* V2. In particular, they have acquired the grammar 11110. In Figure 11 the percentage of the population speaking the four languages mentioned above as they evolve over 20 generations is shown. Notice that in the space of a few generations the speakers of 11011 and 01011 have dropped out altogether. Most of the population now speaks language 1111 (46%) and 01111 (27%). Fifteen percent of the population speaks 11110 and there is a smattering of other speakers. The population remains roughly stable in this configuration thereafter.

#### Observations

1. On examining the four languages to which the system converges after one generation, we notice that they share the same settings for the principles: [case assignment under government], [pro drop], and [V2]. These correspond to the three parameters which are uniquely set

by data from old French. The other two parameters can take on any value. Consequently, four languages are generated, all of which satisfy the data from old French.

2. Recall our earlier remark that, due to insufficient data, there were equivalent grammars in the parameter system. It turns out that in this particular case, the grammars (01011) and (11011) are identical as far as their extensional properties are concerned; as are the grammars (11111) and (01111).

3. There is subset relation between the two sets described in 2. The grammar (11011) is in a subset relation with (11111). This explains why after a few generations most of the population switches to either (11111) or (01111) (the superset grammars).

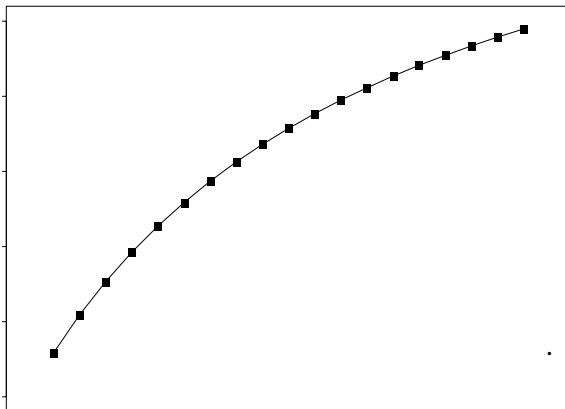
4. An interesting feature of the simulation is that 15 percent of the population eventually acquires the grammar (11110), that is, they have lost the V2 parameter setting. This is the first sign of instability of V2 that we have seen in our simulations so far (for greedy algorithms which are psychologically preferred). Recall that for such algorithms, the V2 parameter was very stable in our previous example.

### 5.3.2 Heterogenous populations (mixtures)

Section 5.3.1 showed that with no new (foreign) sentence patterns the grammatical system starting out with only old French speakers showed some tendency to lose V2. However, the grammatical trajectory did not terminate in modern French. In order to more closely duplicate this historically observed trajectory, we examine alternative initial conditions. We start our simulations with an initial condition which is a mixture of two sources; data from old French and data from new French (reproducing in this sense, data similar to that obtained from the middle French period). Thus, children in the next generation observe new surface forms. Most of the surface forms observed in middle French are covered by this mixture.

#### Observations

1. On performing the simulations using the TLA as a learning algorithm on this parameter space, an interesting pattern is observed. Suppose the learner is exposed to sentences with 90 percent generated by old French grammar (11011) and 10 percent by modern French grammar (10100), within one generation 22 percent of the learners have converged to the grammar (11110) and 78 percent to the grammar (11111). Thus, the learners set each of the parameter values to 1 except the V2 parameter setting. Now, modern French is a nonV2 language; and 10 percent of data from modern French is sufficient to cause 22 percent of the speakers to lose V2. This is the behavior over one generation.



**Figure 12.** Tendency to lose V2 as a result of new word orders introduced by modern French source in our Markov model.

The new population (consisting of 78% speaking grammar (11111) and 22% speaking grammar (11110)) remains stable forever.

2. Figure 12 shows the proportion of speakers who have lost V2 after one generation as a function of the proportion of sentences from the modern French source. The shape of the curve is interesting. For small values of the proportion of the modern French source the slope of the curve is greater than 1. Thus, there is a greater tendency of speakers to lose V2 than to retain it. This results in 10 percent of novel sentences from the modern French source causing 20 percent of the population to lose V2; similarly 20 percent of novel sentences from the modern French source causes 40 percent of the speakers to lose V2. This effect wears off later. This seems to capture computationally the intuitive notion of many linguists that a small change in inputs provided to children could drive the system towards larger change.

3. Unfortunately, there are several shortcomings in this particular simulation. First, we notice that mixing old and modern French sources does not cause the desired (historically observed) grammatical trajectory from old to modern French (corresponding in our system to movement from state (11011) to state (10100) in our Markov chain). Although we find that a small injection of sentences from modern French causes a larger percentage of the population to lose V2 and gain subject clitics (which are historically observed phenomena), nevertheless, the entire population retains the null subject setting and case assignment under government. It should be mentioned that it is argued in [3] that the change in case assignment under government is the driving force which

allows alternate parse-trees to be formed and causes the parametric loss of V2 and null subject. In this sense, it is a more fundamental change.

4. If the dynamical system is allowed to evolve, it ends up in either of the two states (11111) or (11110). This is essentially due to the subset relations these states (languages) have with other languages in the system. Another complication in the system is the equivalence of several different grammars (with respect to their surface extensions), for example, given the data we are considering, the grammars (01011) and (11011) (old French) generate the same sentences. This leads to multiplicity of paths, convergence to more than one target grammar, and general inelegance of the state-space description.

### Future Directions

There are several possibilities to consider here.

1. Using more data and filling out the state-space might yield greater insight. Note that studies in the development of other languages such as Italian or Spanish within this framework might also be useful.

2. TLA-like hill climbing algorithms do not pay attention to the subset principle explicitly. It would be interesting to explicitly program this into the learning algorithm and observe the evolution thereafter.

3. There are often cases when several different grammars generate the same sentences or at least fit the data equally well. Algorithms that look only at surface strings are unable then to distinguish between them resulting in convergence to all of them with different probabilities in our stochastic setting. We showed an example of this for convergence to four states earlier. An elegance criterion is suggested in [3] that looks at the parse-trees to decide between these grammars. This difference between *strong* generative capacity and *weak* generative capacity can easily be incorporated into the Markov model as well. The transition probabilities, now, will not depend upon the surface properties of the grammars alone, but also upon the elegance of derivation for each surface string.

4. Rather than the evolution of the population, one could look at the evolution of the distribution of words. One can also obtain bounds on frequencies with which the new data in the middle French period must occur so that the correct drift is observed.

## 6. Conclusions and directions for future research

In this paper, we have argued that any combination of a grammatical theory and a learning paradigm leads to a model of grammatical evolution and diachronic change. A learning theory (paradigm) attempts

to account for how children (the individual child) solve the problem of language acquisition. By considering a population of such “child learners,” we have arrived at a model of the *emergent*, global, population behavior. The key point is that such a model is a logical consequence of grammatical and learning theories. Consequently, whenever a linguist suggests a new grammatical, or learning theory, they are also suggesting a particular evolutionary theory—and the consequences of this need to be examined.

### ■ 6.1 Historical linguistics and diachronic criteria

From a programmatic perspective, this paper has two important consequences. First, it allows one to take a formal, analytic view of historical linguistics. Most accounts of language change have tended to be descriptive in nature. In contrast, we place the study of historical or diachronic linguistics in a formal framework. In this sense, our conception of historical linguistics is closest in spirit to evolutionary theory and population biology. Indeed, most previous attempts to model language change, such as [3, 10] have been influenced by the evolutionary models.

Second, this approach allows us to formally pose a *diachronic* criterion for the adequacy of grammatical theories. A significant body of work in learning theory has already sharpened the *learnability* criterion for grammatical theories, in other words, the class of grammars  $\mathcal{G}$  must be learnable by some psychologically plausible algorithm from primary linguistic data. Now we can go one step further. The class of grammars  $\mathcal{G}$  (along with a proposed learning algorithm  $\mathcal{A}$ ) can be reduced to a dynamical system whose evolution must match that of the true evolution of human languages (as reconstructed from historical data).

We have attempted to lay the framework for the development of research tools to study historical phenomena. To concretely demonstrate that the grammatical dynamical systems need not be impossibly difficult to compute (or simulate), we explicitly showed how to transform parametrized theories, and memoryless learning algorithms to dynamical systems. The specific simulations of this paper are far too incomplete to have any long term linguistic implications, though, we hope, it certainly forms a starting point for research in this direction. Nevertheless, certain interesting results were obtained.

1. It was shown that the V2 parameter was more stable in the three-parameter case than it was in the five-parameter case. This suggests that the loss of V2 (actually observed in history) might have more to do with the choice of parametrizations than learning algorithms, or primary linguistic data (though we must suggest great caution before drawing strong conclusions on the basis of this study).

2. Some light was shed on the time course of evolution. In particular, it was shown how this was a derivative of more fundamental assump-



tions about initial population conditions, sentence distributions, and learning algorithms.

3. Notions of system stability were formally developed. Thus, certain parameters could change with time, others might remain stable. This can now be measured, and the conditions for stability or change can be investigated.

4. It was demonstrated that one could manipulate the system (by changing the algorithm, sentence distributions, or maturational time) to allow evolution in certain directions. These logical possibilities suggest the kinds of changes needed in linguistics for greater explanatory adequacy.

## ■ 6.2 Further research

This has been our first attempt to define the boundaries of the language change problem as a dynamical system, there are several directions for further research.

1. From a linguistic perspective, one could examine alternative parametrized theories, and track the change of certain languages in the context of these theories (much like our attempt to track the change of French in this paper). Some worthwhile attempts could include: (a) The study of parametric stress systems [7]—and in particular, the evolution of modern Greek stress patterns from proto-Indo European. (b) The investigation of the possibility that creoles correspond to fixed points in parametric dynamical systems, a possibility which might explain the striking fact that all creoles (irrespective of the linguistic origin, i.e., initial linguistic composition of the population) have the same grammar. (c) The evolution of modern Urdu, with Hindi syntax, and Persian vocabulary; a cross-comparison of so-called “phylogenetic” descriptive methods currently used to trace back the development of an early, common, proto-language.

2. From a mathematical perspective, one could take this research in many directions including: (a) The formalization of the update rule for other grammatical theories and learning algorithms, and the characterization of the dynamical systems implied. (b) A better characterization of stability issues and phase-space plots. (c) The possibility of true chaotic behavior—recall that our dynamical systems are multidimensional non-linear iterated function mappings.

It is our hope that research in this line will mature to make useful contributions, both to linguistics, and in view of the unusual nature of the dynamical systems involved, to the study of such systems from a mathematical perspective.

## Acknowledgments

---

This report describes research done at the Center for Biological and Computational Learning and the Artificial Intelligence Laboratory of the Massachusetts Institute of Technology. Support for the Center is provided in part by a grant from the National Science Foundation under contract ASC-9217041. R. C. Berwick was supported by a Vinton-Hayes Fellowship. P. Niyogi was supported by a grant from the NEC Corporate Foundation.

## Appendix

### A. The three-parameter system of Gibson and Wexler

---

The three-parameter system discussed in [5] includes two parameters from  $X$ -bar theory. Specifically, they relate to specifier-head relations, and head-complement relations in phrase structure. The following parametrized production rules denote this:

$$\begin{aligned} XP &\rightarrow \text{Spec}X'(p_1 = 0) \text{ or } X'\text{Spec}(p_1 = 1) \\ X' &\rightarrow \text{Comp}X'(p_2 = 0) \text{ or } X'\text{Comp}(p_2 = 1) \\ X' &\rightarrow X. \end{aligned}$$

A third parameter is related to verb movement. In German, and Dutch root declarative clauses, it is observed that the verb occupies exactly the second position. This verb-second phenomenon might or might not be present in languages of the world, and this variation is captured by means of the V2 parameter.

Table 4 provides the unembedded (degree-0) sentences from each of the eight grammars (languages) obtained by setting the three parameters of section 4 to different values. The languages are referred to as  $L_1$  through  $L_8$ .

## References

---

- [1] R. C. Berwick, *The Acquisition of Syntactic Knowledge* (MIT Press, 1985).
- [2] L. Cavalli-Sforza and M. Feldman, *Cultural Transmission and Evolution: A Quantitative Approach* (Princeton University Press, 1981).
- [3] R. Clark and I. Roberts, "A Computational Model of Language Learnability and Language Change," *Linguistic Inquiry*, 24(2) (1993) 299-345.
- [4] G. Altmann *et al.*, "A Law of Change in Language," in *Historical Linguistics*, edited by B. Brainard (Studienverlag Dr. N. Brockmeyer., Bochum, FRG, 1982).

Language	Spec	Comp	V2	Degree-0 unembedded sentences
L <sub>1</sub>	1	1	0	"V S" "V O S" "V O1 O2 S" "AUX V S" "AUX V O S" "AUX V O1 O2 S" "ADV V S" "ADV V O S" "ADV V O1 O2 S" "ADV AUX V S" "ADV AUX V O S" "ADV AUX V O1 O2 S"
L <sub>2</sub>	1	1	1	"S V" "S V O" "O V S" "S V O1 O2" "O1 V O2 S" "O2 V O1 S" "S AUX V" "S AUX V O" "O AUX V S" "S AUX V O1 O2" "O1 AUX V O2 S" "O2 AUX V O1 S" "ADV V S" "ADV V O S" "ADV V O1 O2 S" "ADV AUX V S" "ADV AUX V O S" "ADV AUX V O1 O2 S"
L <sub>3</sub>	1	0	0	"V S" "O V S" "O2 O1 V S" "V AUX S" "O V AUX S" "O2 O1 V AUX S" "ADV V S" "ADV O V S" "ADV O2 O1 V S" "ADV V AUX S" "ADV O V AUX S" "ADV O2 O1 V AUX S"
L <sub>4</sub>	1	0	1	"S V" "O V S" "S V O" "S V O2 O1" "O1 V O2 S" "O2 V O1 S" "S AUX V" "S AUX O V" "O AUX V S" "S AUX O2 O1 V" "O1 AUX O2 V S" "O2 AUX O1 V S" "ADV V S" "ADV V O S" "ADV V O2 O1 S" "ADV AUX V S" "ADV AUX O V S" "ADV AUX O2 O1 V S"
L <sub>5</sub> (English, French)	0	1	0	"S V" "S V O" "S V O1 O2" "S AUX V" "S AUX V O" "S AUX V O1 O2" "ADV S V" "ADV S V O" "ADV S V O1 O2" "ADV S AUX V" "ADV S AUX V O" "ADV S AUX V O1 O2"
L <sub>6</sub>	0	1	1	"S V" "S V O" "O V S" "S V O1 O2" "O1 V S O2" "O2 V S O1" "S AUX V" "S AUX V O" "O AUX S V" "S AUX V O1 O2" "O1 AUX S V O2" "O2 AUX S V O1" "ADV V S" "ADV V S O" "ADV V S O1 O2" "ADV AUX S V" "ADV AUX S V O" "ADV AUX S V O1 O2"
L <sub>7</sub> (Bengali, Hindi)	0	0	0	"S V" "S O V" "S O2 O1 V" "S V AUX" "S O2 O1 V AUX" "ADV S V" "ADV S O V" "ADV S O2 O1 V" "ADV S V AUX" "ADV S O V AUX" "ADV S O2 O1 V AUX"
L <sub>8</sub> (German, Dutch)	0	0	1	"S V" "S V O" "O V S" "S V O2 O1" "O1 V S O2" "O2 V S O1" "S AUX V" "S AUX O V" "O AUX S V" "O1 AUX S O2 V" "O2 AUX S O1 V" "ADV V S" "ADV V S O" "ADV V S O2 O1" "ADV AUX S V" "ADV AUX S O V" "S AUX O2 O1 V" "S O V AUX" "ADV AUX S O2 O1 V"

**Table 4.** Unembedded sentences from each of the eight grammars in the three-parameter system.

[5] E. Gibson and K. Wexler, "Triggers," *Linguistic Inquiry*, **25** (1994).

[6] L. Haegeman, *Introduction to Government and Binding Theory* (Blackwell, Cambridge, USA, 1991).

[7] M. Halle and W. Idsardi, "General Properties of Stress and Metrical Structure," in *DIMA CS Workshop on Human Language*, Princeton, NJ, 1991.

[8] A. S. Kroch, "Grammatical Theory and the Quantitative Study of Syntactic Change," from a paper presented at NWAVE 11, Georgetown University, 1982.

[9] A. S. Kroch, "Function and Grammar in the History of English: Periphrastic 'do,'" in *Language Change and Variation*, edited by Ralph Fasold (Benjamins, Amsterdam, 1989).

[10] Anthony S. Kroch, "Reflexes of Grammar in Patterns of Language Change," *Language Variation and Change*, (1990) 199–243.

[11] D. Lightfoot, *How to Set Parameters* (MIT Press, Cambridge, MA, 1991).

[12] B. Mandelbrot, *The Fractal Geometry of Nature* (W. H. Freeman and Co., New York, 1982).

- [13] P. Niyogi, *The Informational Complexity of Learning from Examples*, PhD thesis, Massachusetts Institute of Technology, Cambridge, MA, 1994.
- [14] P. Niyogi, *The Informational Complexity of Learning: Perspectives on Neural Networks and Generative Grammar* (Kluwer Academic Publishers, Boston, 1997).
- [15] P. Niyogi and R. C. Berwick, "Formalizing Triggers: A Learning Model for Finite Parameter Spaces," Technical Report 1449, AI Laboratory, MIT, 1993.
- [16] P. Niyogi and R. C. Berwick, "A Markov Language Learning Model for Finite Parameter Spaces," in *Proceedings of 32nd Meeting of Association for Computational Linguistics*, 1994.
- [17] P. Niyogi and R. C. Berwick, "Formal Models for Learning Finite Parameter Spaces," Workshop on Cognitive Models of Language, Tilburg University, Ne., 1994. Also in *Models of Language Learning: Inductive and Deductive Approaches*, edited by P. Broeder and J. Murre (MIT Press, Cambridge, MA, to appear).
- [18] P. Niyogi and R. C. Berwick, "Learning from Triggers," *Linguistic Inquiry*, 27(4) (1996) 605–622.
- [19] P. Niyogi and R. C. Berwick, "A Language Learning Model for Finite Parameter Spaces," *Cognition*, 61 (1996) 161–193.
- [20] C. Osgood and T. Sebeok, "Psycholinguistics: A Survey of Theory and Research Problems," *Journal of Abnormal and Social Psychology*, 49(4) (1954) 1–203.
- [21] S. Resnick, *Adventures in Stochastic Processes* (Birkhauser, 1992).
- [22] U. Weinreich, W. Labov, and M. I. Herzog, "Empirical Foundations for a Theory of Language Change," in *Directions for Historical Linguistics: A Symposium*, edited by W. P. Lehmann (University of Texas Press, Austin, 1968).