

# Modelling linguistic taxonomic dynamics

Søren Wichmann<sup>1,2</sup>, Dietrich Stauffer<sup>3</sup>, F. Wellington S. Lima<sup>4</sup>, and Christian Schulze<sup>3</sup>

<sup>1</sup> Department of Linguistics, Max Planck Institute for Evolutionary Anthropology, Deutscher Platz 6, D-04103 Leipzig, Germany

<sup>2</sup> Languages and Cultures of Indian America (TCIA), Leiden University, P.O. Box 9515, 2300 RA Leiden, The Netherlands

<sup>3</sup> Institute for Theoretical Physics, Cologne University, D-50923 Köln, Euroland

<sup>4</sup> Departamento de Física, Universidade Federal do Piauí, 57072-970 Teresina - PI, Brazil

## Abstract

This paper presents the results of the application of a bit-string model of languages (Schulze and Stauffer 2005) to problems of taxonomic patterns. The questions addressed include the following: (1) Which parameters are minimally needed for the development of a taxonomic dynamics leading to the type of distribution of language family sizes currently attested (as measured in the number of languages per family), which appears to be a power-law? (2) How may such a model be coupled with one of the dynamics of speaker populations leading to the type of language size seen today, which appears to follow a log-normal distribution?

## 1 Introduction

Although computational linguistics is a vast and well-established field, linguists have only recently been concerned with quantitative modeling and simulation of different aspects of language history. Maeda et al. (1997) and Nettle (1999a,b) are notable early exceptions, and more recently much work has been done modelling language change (Niyogi 2002, Baxter et al. 2006) and the evolution of language structure (Cangelosi and Parisi 2002, Nowak et al. 2002, Christiansen and Kirby 2003, Wang et al. 2004, Gong et al. 2005, Niyogi 2006). Within physics an interest has recently developed in simulating language interaction and studying how different situations may lead to various patterns of growth or extermination (see Schulze and Stauffer 2006 for a review). This interest among physicists for modeling language competition was triggered by Abrams and Strogatz (2003), who use differential equations

to describe the vanishing of one language due to the dominance of another. Since then, a series of articles have appeared. For instance, Patriarca and Leppänen (2004) applied the Abrams-Strogatz model to a geographical situation where two languages, X and Y, may dominate in each their region, resulting in the survival of both, rather as in the original model where only one language will survive. Oliveira et al. (2006a,b) have looked at models in which speakers of small languages will tend to switch to geographically more widespread ones, to account for the fact that real geographical areas tend not to show equally-sized languages. Along the lines of the broader cultural model of Axelrod (1997), Teşileanu and Meyer-Ortmanns (2006) looked at the consequences of the possibility that a greater similarity among languages might further language shift. We build on some of this work, but so far the present paper is to our knowledge the first in this recent tradition to address the issue of taxonomic dynamics.

## 2 Why simulate?

Human languages may have existed for at least about  $10^5$  years, possible much longer, possibly a little less. There is no reason to enter into a lengthy discussion about this issue, since in reality not much is known. Johansson (2005: 85) summarizes the state of knowledge as follows: “The evidence concerning the evolution of humanlike speech capabilities is not conclusive; we do not know for sure when it evolved, though it may well have been evolutionarily fairly recent, most likely sometime during the past two million years or so. Some type of speech must have been present in our last common ancestor with Neanderthals, 500,000 years ago or so, though fully human speech with all our articulatory capacity need not be much older than 100,000 years.” Only a few percent of this development is to some extent documented through writing, while another few percent may be inferred by comparative linguistic methods. Thus, we have no clues to aspects of the development of languages for some 90% or more of their history other than what we might infer from abstract extrapolation or from simulations. Like the distant past, the future is also empirically impenetrable. Two aspects of simulations are important. First, we may hope to identify a minimal number of parameters that account for the present state of affairs seen as a result of a long development. Secondly, we may adjust these parameters to test the predictions that different models provide. It should be stressed that

simulations are not necessarily suited to prove any particular model or to make predictions about what is in store for the languages of today; they can only represent tests of different models. Nor can the parameters identified be translated into direct explanatory factors for actual distributions. For instance, a simulation of language competition might restrict its parameters to, say, the relative size of languages and it might stipulate some simple mode of interaction, such as the tendency for speakers of smaller languages to shift to contiguous larger ones. Such a model might lead to a plausible picture of language distributions, perhaps even one resembling the current state of affairs. But this does not mean that this distribution is explained by language sizes and competition alone. The growth or extinction of a given language relates to socioeconomic, historical, geographical, ecological and many other circumstances (Thomason 2001: 38-46, Mufwene 2001: 145-166, McWhorter 2001: 253-262). Since a primary aim of simulation is to reduce the set of parameters, it cannot and should not, however, take into account all relevant factors, but must remain an abstraction.

### 3 The aim of the investigation

Wichmann (2005) made some simple observations about the present-day quantitative distribution of language family sizes, as measured in numbers of languages per family, and about the distribution of language sizes, as measured in numbers of speakers per language (data drawn from *Ethnologue*). It was found that language family sizes approximate a so-called ‘power-law’, that is, a distribution described by the equation  $y = ax^b$ , which corresponds to a straight line on a log-log plot (cf. Zanette 2001 for a similar observation). Such distributions are frequent in both nature and the social world (cf. Newman 2005 for an excellent overview). The slope of the curve on the rank-by-size plot is described by the exponent  $b$ , which was found to be  $-1.905$ .

When testing for the distribution of language sizes, however, no power-law emerged. The absence of a power-law distribution also comes out of studies by Novotny and Drozd (2000) and Sutherland (2003). (Gomes et al. 1999: 493 had earlier plotted the same data on a graph showing the cumulative size distribution,  $n(> S)$ , corresponding to the number of languages with a size greater than  $S$ . Cutting the curve up into different regions and describing each by a separate equation they then made the problematical

claim of the existence of a “composite power-law”). The present paper takes up the challenge of Wichmann (2005: 139) to test, using computer simulations, what the expected past and future distributions of language family sizes and language sizes might look like. The question was raised whether the present distribution of language sizes might be characteristic of a stage of disequilibrium while the expected equilibrium might correspond to a power-law. Stauffer et al. (2006) supported the hypothesis of a disequilibrium. In the present paper we also report on language families.

## 4 The bit-string model

The model used is one eminently suited to computation. It is a variant of that of Schulze and Stauffer (2005), which operates with bit-strings of length  $L$ , where each bit has two values and where the total set of possible dialects has  $2^L$  members. (A precursor to this kind of modelling is Wang and Minett 2005, which used strings of integers to simulate branching by the mutation and transfer of numbers.)

Each bit may be interpreted as the presence or absence of some characteristic grammatical feature. Under this interpretation we might imagine that a number of diagnostic features were identified, the presence or absence of each of which would be sufficient to distinguish among the grammars of the world’s languages. This number corresponds to the length of the bit-string.

Using bits as opposed to features having several values, for instance anything from 2 to 9 values as in Haspelmath et al. (2005), looks to be a move serving more to please our computer rather than to capture reality. We have, however, tested our model with a variety of numbers of possible values and have found no qualitative differences (Schulze and Stauffer 2006, Holman et al. 2006). Moreover, it is impossible to argue that language features are more appropriately described as trinary or quaternary etc. than as binary. The choice of how to describe facts of languages is entirely up to the linguist, who will make whatever choice is appropriate for the purpose at hand. For instance, the feature values in Haspelmath et al. (2005) were chosen largely such as to provide telling maps. Many features could have been broken up into more than 9 values (for instance the first map in the book, showing sizes of consonant inventories, which arbitrarily has 5 values). On the other hand—and this is a good argument for operating with binary features in the present study—all features could be recoded as binary ones (for instance,

presence vs. absence of small vs. moderately small vs. average vs. moderately large vs. large consonant inventories). The fact that a multi-valued feature can always be recoded as binary, i.e. as the presence vs. absence of the given value, speaks in favor of choosing bits for abstract data modelling.

An alternative model of language competition, also allowing for thousands of different languages, is that of de Oliveira et al. (2006a,b). There, however, languages are characterized merely by consecutive numbers 1, 2, 3, ..., which is not suitable for simulating different taxonomic levels. In this model, language families would have to be determined by the history of language dynamics and their genealogical tree (Schulze and Stauffer 2006), and this approach was used later (Stauffer et al. 2007) after the submission of the present work, also for our model. The other recent models of language competition to our knowledge allow only a relatively small number of languages and are, for this reason, also less suitable for taxonomy.

We test two different variants of the model. In one, which we might call the “hierarchical” variant, the bit-string is divided into subsections corresponding to different taxonomic levels. Two languages are defined as belonging to the same family if their “family” parts of the bit-strings agree. In the other, “flat” variant, there is no such partitioning of the string. Instead, taxonomic levels are achieved by defining a certain threshold  $k$  of differences among languages. Differences are measured by comparing two strings and noting the number of positions for which the two strings differ. If the difference is greater than  $k$ , the two languages are said to belong to different taxa. In both versions of the model we only operate with two taxonomic levels, but both could be extended to include more levels. In the following, each variant is described in more detail.

## 4.1 The hierarchical variant

This model achieves two taxonomic levels by partitioning the bit-string. The two levels may be conceptualized as corresponding to language families and languages within one family, respectively, but need not be translated exactly into these concepts (which are themselves not very well defined). The languages of individuals may be classified by comparing the bit-strings representing each individual. In the following we illustrate how the model works if we use a bit-string of length 32. People speaking the same language have to agree in all bits. In our implementation we have chosen to stipulate that people speaking languages belonging to the same family have to agree in the

leftmost 10 bits. For example,

1010100011-01010101101010100111 and  
1010100011-0101110000110101001100

are two different languages belonging to the same family and

1010011010-1001011010101010111101 and  
0101101010-1010001111010101101010

are two different languages belong to two different families. In these examples the dash “-” just indicates the boundary between the two segments of the string, analogously to the convention for phone numbers, which are structured much like our bit-string model.

The choice of lengths of the whole string and its parts is of course arbitrary, and need not be  $10 + 22 = 32$ . Nevertheless, various considerations led to single out certain lengths as more suitable than other. First, the model is computationally most effective for bit-string lengths which are powers of 2. Second, a shorter string is to be preferred to a longer one, all else being equal—again for computational reasons. Third, the string should not be so short that the sheer length imposes artificial constraints on the results. In earlier simulations the effect of different lengths ( $L = 8, 16, 32$ , and  $64$ ) were tested. Since it was found that the results were qualitatively similar for  $L = 16, 32$ , and  $64$ , all values of  $L$  higher than or equal to 16 would be equally suitable. Adding the criterion of minimal computation cost would single out  $L = 16$  as preferable. However, we found that for this length a maximum number of languages was reached before a meaningfully interpretable distribution was found. Unlike the 32 and 64 bit-string models and the real-life present-day distribution, see section 5 below, this did not lead to a power-law distribution, since power-laws require the absence of upper bounds. At a point where either all possible languages or all possible families are filled, the power-law distribution breaks down. Instead, we have chosen a string with the larger length of 32 bits, of which the leading 10 bits define families and the remaining 22 define languages, yielding  $2^{10} = 1024$  possible ‘families’ and  $2^{22} = 4,194,304$  possible ‘languages’. Using an ample  $L$  also ensures that accidental ‘back mutation’, i.e. the phenomenon whereby, by chance, an identical bit-string occurs after some mutations—something which would not happen in real life—will occur so exceedingly rarely that its effects are completely negligible. (Even for  $L = 8$  this situation occurs rarely, cf. Schulze and Stauffer 2006). Simulations using 64 bits, of which 19 bits are

reserved for families and the remaining 45 bits for languages were also made (allowing for  $2^{19} = 524,288$  possible ‘families’ and  $2^{45} = 35,184,372,088,832$  possible ‘languages’). The results were qualitatively similar.

Differentiation is simulated by setting a probability  $p$  for a change in a bit per iteration. An iteration is equivalent to a certain, average time step. Initially each individual randomly selects a language. After some time steps, a bit in either the family sub-string or the language bit-string will change, meaning the creation of a new entity at one of these levels. Given that there are fewer bits in the family bit-string, there is a smaller probability of a change in this part of the string per iteration, and there will therefore be a slower dynamics of families than of languages. In practice, with probability  $pM$  at each iteration, one of the  $L$  bits is selected randomly and then reverted, i.e. changed from 0 to 1 or from 1 to 0. In this process, analogously to biological mutations, all bit positions are equivalent and neither 0 nor 1 is in any way preferred. Testing the effects of different values of  $p$  we have found a transition point where, for small values, the end result of the simulation is a situation where a single language is spoken by most of the population and, for larger values, a situation where all possible languages are spoken by about the same number of speakers. This is equivalent to phase transitions in physical systems, such as when water suddenly begins to turn to vapor as it is being warmed up. The transition point lies around  $p \approx 0.140$ . Thus, to simulate a development leading to something at least roughly comparable to the present-day language distribution a small value of  $p$  is to be preferred; we therefore normally set  $p$  at 0.0001. For the special purpose of studying the phase transition, however, values of  $p$  on either side of the transition point are preferred; thus, towards the end of this paper, where we study the phase transition, we operate with the values  $p = 0.136$  and  $p = 0.144$  (cf. fig. 9 below).

We neglect here for simplicity the diffusion of features from one language to the other used in other simulations involving this model. The simulations may be characterized as agent-based. We start the simulations with  $\sim 10^5$  speakers, each speaking randomly selected languages, and we follow  $10^2$  generations. We assume a shift from small to large populations stipulating that at each iteration with probability  $(1-x)^2$  or  $(1-x^2)$  each individual gives up his/her old language and instead selects the language of one randomly selected individual of the whole population. The variable  $x$  represents the fraction of the total population speaking the old language. Individuals get one child per iteration. This child gets the same bit-string as the

parent except that with some mutation probability one randomly selected bit is changed. Everybody dies with a Verhulst probability proportional to the current population size, something which takes into account factors such as limited food and space. After some time, one language may dominate and be spoken by more than 80 percent of the population. A Fortran program for a similar basic model is published in the appendix of Stauffer et al. (2006a). The histograms of the number of languages spoken by a given number of people are smoothed by random multiplicative noise as in Stauffer et al. (2006b), which may correspond to external perturbations caused by migrations of individuals, intermarriage, changing political circumstances, and other non-systematic factors. Typically, a simulation takes many hours using a fast PC; graphics were produced with gnuplot.

## 4.2 The flat variant

This model is in all but one major respect similar to the hierarchical model. The difference is that taxonomic levels are achieved not by partitioning the string, but by stipulating that two languages which differ in more than one bit belong to different families. The size of each language family (i.e., the numbers of languages in each) is then measured by the number of languages that differ by just one bit from one reference language. We sum over all reference languages, and also over many samples, to get out final statistics. The definition allows one language to belong to different families, just as one person can belong to different friendship groups. Instead, one would get a clear separation into different families without such overlaps if we demand all languages within one family to be separated directly or indirectly by not more than one bit flip. But since we can move from each bit-string of 32 bits to every other possible bit-string through at most 32 such changes of single bits, this definition would mean that all possible languages form one huge family, which is not what we want. (Analogously, on a square lattice we can define a neighbourhood as the set of four nearest neighbours of a given site; then every lattice site belongs to several neighbourhoods. Alternatively, a cluster can be defined as the set of all sites connected directly or indirectly with a given site; then the whole lattice forms one large cluster. Neither definition leads to what we would like to have, which is non-overlapping clusters, corresponding to non-overlapping language families. A further disadvantage of the model is that its equilibrium is either dominance of one language spoken by most people, or fragmentation into numerous languages of about equal size; thus

for dominance there is not much to analyze and for fragmentation nearly all languages could form one cluster, meaning that these more sophisticated definitions might not work better in equilibrium.)

## 5 Results

### 5.1 Results for the application of the hierarchical model

The major results are shown in figs. 1-2 and 4-5. The interest of these are the shapes of the various curves, not the absolute numbers corresponding to each point. The mismatch between large numbers of families and small sizes of languages as compared to the real-world situation is due to the summation over iterations and could be normalized, but this would only serve presentational purposes.

In fig. 1 it is shown how size histograms of families strongly depend on the temporal factor. At the initial stage of the simulation ( $t = 1$ ) we see something close to a normal distribution (the rightmost curve in the diagram). At  $t = 10$  the distribution forms a parabola (curve connecting x's). This distribution is close to what the present-day *language* size distribution looks like (see fig. 6). At  $t = 60$  (stars) a curve resembling the present-day distribution of *language family* sizes (fig. 3) is obtained, but it has a large hump on the right region of the curve. The real-life distribution also has a hump, but it is much smaller. At 300 iterations (squares) there is a discontinuous distribution with a number of small families and a leap up to a number of larger ones, which form a narrow normal distribution. Fig. 2a focuses on the range  $20 \leq t \leq 150$ , where the distribution most closely resembles the present-day one, and varies the population size  $N$  to see the dependency on the graph on that variable. It appears that there is not much influence of  $N$ , provided  $t$  is increased with increasing  $N$ . Moreover, fig. 2a suggests that the closest approximation to the present-day distribution is found around  $10^2$  iterations. Statistically solid results for a long run of the 64-bit model in fig. 2b provide similar results.

We now turn to the results for language sizes. Fig. 4 shows the sizes for the same number of iterations as in fig. 1. Since the simulations start with fragmentation,  $t = 1$  represents a situation with many languages spoken by single speakers (single +). At  $t = 10$  (x symbols) a curve roughly like a parabola and already strongly reminiscent of the present-day situation (fig.

6) has begun to form. At  $t = 60$  (stars) this distribution is beginning to disrupt, as evidenced by the right tail. This situation further develops into one with many large languages and many small ones, with a large gap for language sizes in between, as shown by the curve for  $t = 300$  (squares). Again we narrow in on the range,  $20 \leq t \leq 60$ , where the best approximation of the present situation (fig. 6) is found and vary the population size (fig. 5). For  $t = 40$  and  $N = 50,000$  the distribution closely approximates the present-day one.

By comparing the curves for  $t = 40$  in figs. 2a and 5 an interesting observation is obtained: at identical time steps the curve for language family sizes may approximate a power-law while the curve for language sizes does not, but rather something close to a parabola, as in real life. Wichmann (2005: 128) hypothesized that both curves should approximate a power-law, but the simulations rather suggest that this is only the case for language family sizes, at least given the model and the setting of parameters assumed here.

The overall result, then, suggests that neither the present-day distribution of language family sizes nor that of language sizes are unexpected and that both may have been obtained for a long time and may continue to be obtained. Eventually a dominance of just one large language accompanied by other slightly different languages is possible, but this situation has not yet set in.

## 5.2 Results for the application of the flat model

For investigating the distance among languages the ‘flat’ model is most useful because the distance among two languages belonging to two different families in the hierarchical model cannot easily be measured. (The hierarchical bit-strings representing languages in any two languages belong to two different families are not comparable since the positions no longer mean the same when one moves up one taxonomic level.) Thus we measured differences in simulations implementing the non-hierarchical model, i.e. the standard model of Schulze and Stauffer (2005, 2006) where all bits are equivalent. As in most of our previous studies, only short bit-strings of 8 or 16 bits were used, and the random multiplicative noise was omitted; for these studies we waited until a stationary state after about  $10^3$  iterations was established.

The distance measure used is the so-called ‘Hamming distances’, also investigated by Teşileanu and Meyer-Ortmanns (2006). The Hamming dis-

tance between two bit-strings is the number of bits which are different in a position-by-position comparison of the two strings. For example, the Hamming distance between 01001101 and 11000011 is four.

As explained above, we define a language family in this model as a set of languages differing from a given reference language by not more than  $k$  bits, in this case setting  $k$  to one bit. The results of the simulations of bit-strings of lengths 8 and 16 are shown in fig. 7; as in fig. 2 above, the simulations represent states of non-equilibrium, i.e., they were stopped at some intermediate time and not let run until the distribution no longer changed apart from random fluctuations. These results are not very different from those shown in fig. 2 for the 64 bits string in the hierarchical model.

### 5.3 More on Hamming distances

The above results were obtained by stopping the simulations at a suitable time such that the results are closest to reality. In this section we report on the equilibrium properties for longer times where the distributions no longer change appreciably and where we will have either dominance of one language or fragmentation of the whole population into many different languages.

Fig. 8 nicely shows the phase transition between dominance at low and fragmentation at high mutation rate  $p$  per bit-string when we vary the mutation rate instead of fixing it to only 0.0001 mutations per bit and per iteration. For dominance, nearly everybody speaks one language, and most of the others speak a language differing in only one bit from this dominating language. Fragmentation happens for larger mutation rates; then all possible languages are represented about equally. We see in fig. 8 that dominance is characterized by a small average Hamming distance while for fragmentation the average Hamming distance is about 1/2 (here it is normalized by the length of the bit-string such that two random bit-strings have on average a distance 1/2.) This effect is already seen if one looks only at the two largest languages in the population, as done by Teşileanu and Meyer-Ortmanns (2006).

For fragmentation, the distribution of Hamming distances between two pairs of speakers is roughly Gaussian (normal), shown by a parabola in the semi-logarithmic plot (stars in fig. 8). In the case of dominance, as observed for two lower mutation rates  $p$  in fig. 9, the most probable Hamming distance is zero, and for higher distances the probability to observe them decays very rapidly.

In these simulations we started with one language only and used the prob-

ability  $1 - x^2$  for the shift from small to large languages. We got qualitatively similar results when we started from a population fragmented into many languages, except that then the probability of a shift was set to  $(1 - x)^2$ , to allow for a possible transition from fragmentation to dominance.

## 6 Conclusion

The primary aim of our simulations was to capture, within one and the same model, how two different empirically observed distributions might arise, i.e. a roughly log-normal distribution of language sizes and an approximate power-law for the family sizes. With reasonable lengths of bit-strings, populations and observation times we could, indeed, find the two different behaviours in the same simulation. This suggests, contrary to the hypothesis of Wichmann (2005), that the present-day distribution of language family sizes in combination with that of language sizes may not be unexpected.

In terms of simulation techniques the major contribution of the present paper has been the introduction of new models into the area of linguistic taxonomic dynamics, an area which, to our knowledge, has not previously been investigated by means of computer simulations. The best results were obtained in implementations of the hierarchical bitstring, a model which also has the advantage of being versatile and easy to implement.

The investigations, however, also revealed some problems with the model. If for a fixed length of the bit-strings the population size  $N$  goes to infinity, then in the parameter region of fragmentation all possible languages will be spoken, and all possible families will exist, making taxonomy a mathematical triviality without connection to reality. Thus simulations of large but finite populations, as presented here, may be better than mathematically exact solutions for infinite populations. Moreover, we did find an effective power-law for the family size distribution, but that distribution decayed much faster with increasing number of languages than the real distribution, shown in fig. 3. Thus future research should aim at also applying and testing other models, such as that of de Oliveira et al. (2006a,b), to problems of linguistic taxonomic dynamic. Stauffer et al. (2007) represents a first attempt to apply (a modified version of) the model of de Oliveira et al.

## References

- Abrams, Daniel and Steven H. Strogatz. 2003. Modelling the dynamics of language death. *Nature* 424: 900.
- Axelrod, Robert. 1997. The dissimination of culture: a model with local convergence and global polarization. *The Journal of Conflict Resolution* 41: 203-226.
- Baxter, Gareth J., Richard A. Blythe, William Croft, and Alan J. McKane. 2006. Utterance selection model of language change. *Physical Review E* 73, article no. 046118.
- Cangelosi, Angelo and Domenico Parisi (eds.). 2002. *Simulating the Evolution of Language*. Berlin: Springer-Verlag.
- Christiansen, Morten and Simon Kirby (eds.) 2003. *Language Evolution*. Oxford: Oxford University Press.
- Ethnologue: Languages of the World* (14th edn. edited by Grimes, Barbara F. 2000, 15th edition edited by Raymond, G. Gordon 2005). Dallas, TX: Summer Institute of Linguistics.
- Gomes, Marcelo A. F., G. L. Vasconcelos, I. J. Tsang, and Ing Ren Tsang. 1999. Scaling relations for diversity of languages. *Physica A* 271: 489-495.
- Gong, Tao, James W. Minnett, Jinyun Ke, John H. Holland, and William S.-Y. Wang. 2005. Coevolution of lexicon and syntax from a simulation perspective. *Complexity* 10.6: 50-62.
- Haspelmath, Martin, Matthew Dryer, David Gil, and Bernard Comrie (eds.) (2005). *The World Atlas of Language Structures*. Oxford: Oxford University Press.
- Holman, Eric W., Christian Schulze, Dietrich Stauffer, and Søren Wichmann. 2006. On the relation between structural diversity and geographical distance among languages: observations and computer simulations. Under revision for *Linguistic Typology*, preprint posted at arXiv:physics/0607031.
- Johansson, Sverker. 2005. *Origins of Language. Constraints on Hypotheses*. Amsterdam/Philadelphia: John Benjamins Publishing Company.
- Maeda, Yuichiro, Takahiro Sasaki, Mario Tokoro. 1997. Self re-organizing of language triggered by language contact. In Husbands, Phil and Inman Harvey (eds.). 1997. *Proceedings of the Fourth European Conference on Artificial Life (ECAL'97)*. MIT Press / Bradford Books, Cambridge, MA.

- McWhorter, John H. 2001. *The Power of Babel. A Natural History of Language*. New York: Henry Holt and Company, LLC.
- Mufwene, Salikoko S. *The Ecology of Language Evolution*. Cambridge: Cambridge University Press.
- Nettle, Daniel. 1999a. Linguistic diversity of the Americas can be reconciled with a recent colonization. *Proceedings of the National Academy of Sciences of the U.S.A.* 96: 3325-3329.
- Nettle, Daniel. 1999b. Using social impact theory to simulate language change. *Lingua* 108: 95-117.
- Newman, Mark E. J. 2005. Power laws, Pareto distributions and Zipf's law. *Contemporary Physics* 46: 323-351.
- Niyogi, Partha. 2002. The computational study of diachronic linguistics. In: Lightfoot, David (ed.), *Syntactic Effects of Morphological Change*, 351-365. Cambridge: Cambridge University Press.
- Niyogi, Partha. 2006. *The Computational Nature of Language Learning and Evolution*. Cambridge & London: The MIT Press.
- Novotny, Vojtech and Pavel Drozd. 2000. The size distribution of conspecific populations. *Proceedings of the Royal Society of London B* 267: 947-952.
- Nowak, Martin A., Natalia Komarova, and Partha Niyogi. 2002. Computational and evolutionary aspects of language. *Nature* 417: 611-617.
- Oliveira, Viviane M. de, Marcelo A. F. Gomes, and Ing Ren Tsang. 2006a. Theoretical model for the evolution of the linguistic diversity. *Physica A* 361: 361-370.
- Oliveira, Viviane M. de, Paulo R. A. Campos, Marcelo A. F. Gomes, and Ing Ren Tsang. 2006b. Bounded fitness landscapes and the evolution of the linguistic diversity, e-print physics 0510249 for *Physica A*.
- Patriarca, Marco and Teemu Leppänen. 2004. Modeling language competition. *Physica A* 338: 296-299.
- Schulze, Christian and Dietrich Stauffer. 2005. Monte Carlo simulation of the rise and fall of languages. *International Journal of Modern Physics C* 16: 781-787.
- Schulze, Christian and Dietrich Stauffer. 2006. Computer simulation of language competition by physicists. In: Chakrabarti, B. K., A. Chakraborti

- and A. Chatterjee (eds.), *Econophysics and Sociophysics: Trends and Perspectives*. Weinheim: WILEY-VCH Verlag; and: Recent developments in computer simulations of language competition, *Computing in Science and Engineering* 8 (May/June) 86-93.
- Stauffer, Dietrich, Suzana Moss de Oliveira, Paulo Murilo C. de Oliveira, Jorge S. Sá Martins. 2006a. *Biology, Sociology, Geology by Computational Physicists*. Amsterdam: Elsevier.
- Stauffer, Dietrich, Christian Schulze, F. Wellington S. Lima, Søren Wichmann, and Sorin Solomon. 2006b. Non-equilibrium and irreversible simulation of competition among languages. *Physica A* 371: 719-724.
- Stauffer, Dietrich, Christian Schulze, and Søren Wichmann. 2007. Computer-Simulation des Wettbewerbs zwischen Sprachen. In: Kolling, Stefan (ed.), *Beiträge zur Experimentalphysik, Didaktik und computergestützten Physik*. Berlin: Logos-Verlag. Preprint at arXiv:physics/0611037.
- Sutherland, William J. 2003. Parallel extinction risk and global distribution of languages and species. *Nature* 423: 276-279.
- Teşileanu, Tiberiu and Hildegard Meyer-Ortmanns. 2006. Competition and languages and their Hamming distance. arXiv:physics/0508229, *International Journal of Modern Physics C* 17: 259-278.
- Thomason, Sarah G. 2001. *Language Contact*. Edinburgh: Edinburgh University Press.
- Wang, William S.-Y., Jinyun Ke, and James W. Minett. 2004. Computational studies of language evolution. In: Huang, C. R. and W. Lenders (eds.), *Computational Linguistics and Beyond*, 65-106. Institute of Linguistics, Academia Sinica.
- Wang, William S. Y. and James W. Minett. 2005. The invasion of language: emergence, change and death. *Trends in Ecology and Evolution* 20.5: 263-296.
- Wichmann, Søren. 2005. On the power-law distribution of language family sizes. *Journal of Linguistics* 41: 117-131.
- Zanette, Damian H. 2001. Self-similarity in the taxonomic classification of human languages. *Advances in Complex Systems* 4: 281-286.

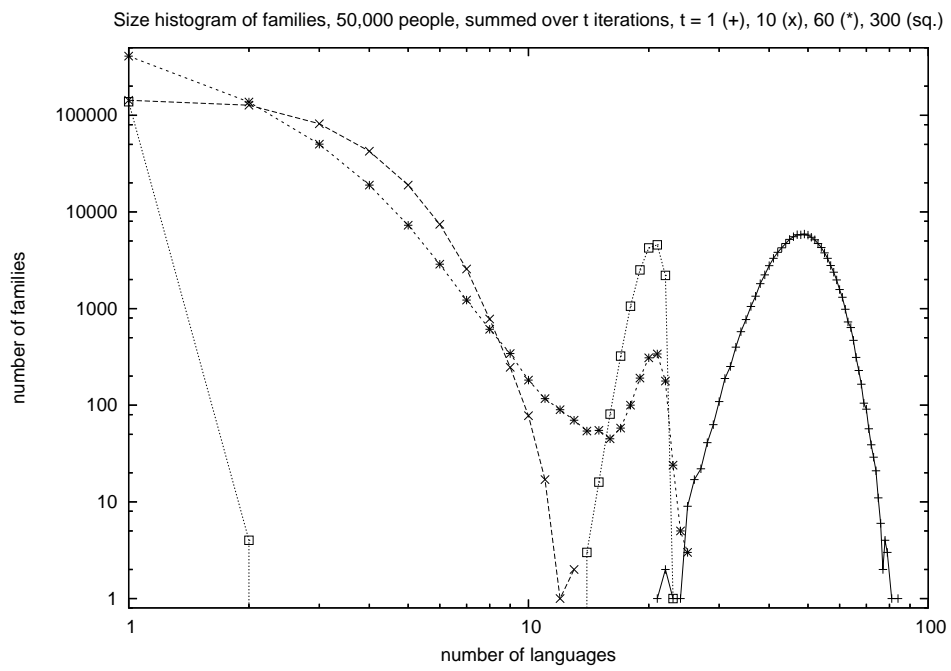


Figure 1: Time dependence of the distribution of family sizes, summed over 100 samples at  $L = 32$ . For long times a narrow peak develops, shown by squares.

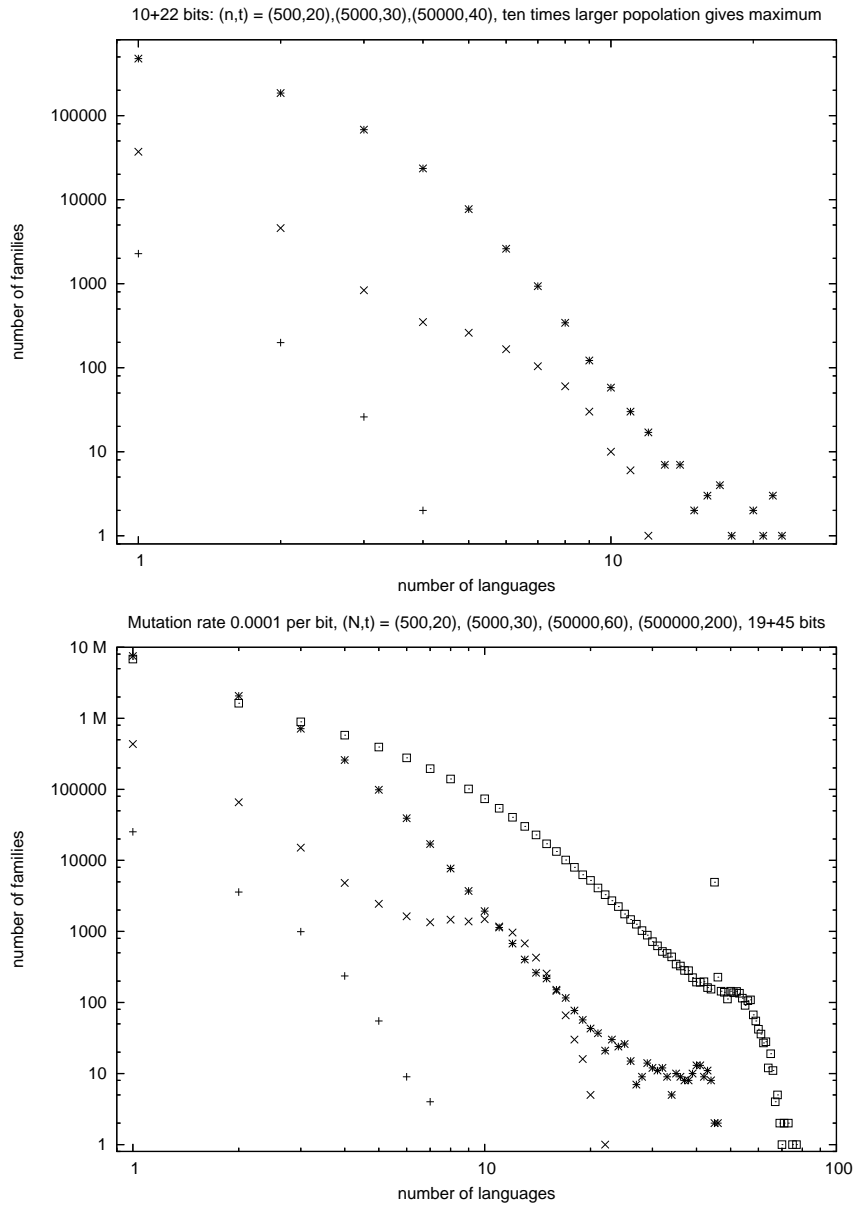


Figure 2: Population-size dependence of the distribution of family sizes, summed over 100 samples, at  $L = 32$  (top) and  $64$  (bottom) and for suitably selected intermediate times. These results are roughly independent of the population.

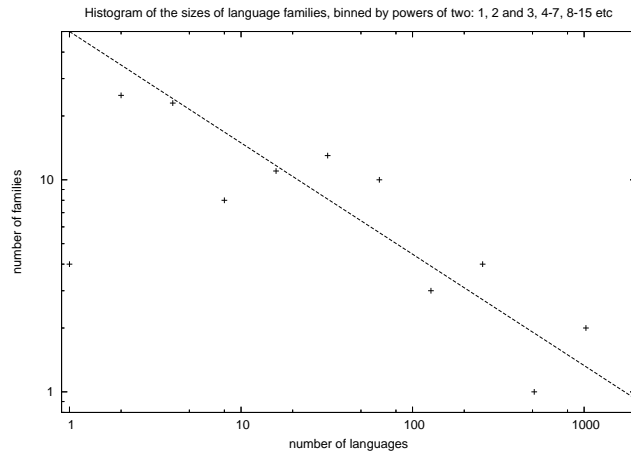


Figure 3: Empirical distribution of family sizes, from *Ethnologue*; see also Wichmann (2005: fig. 2).

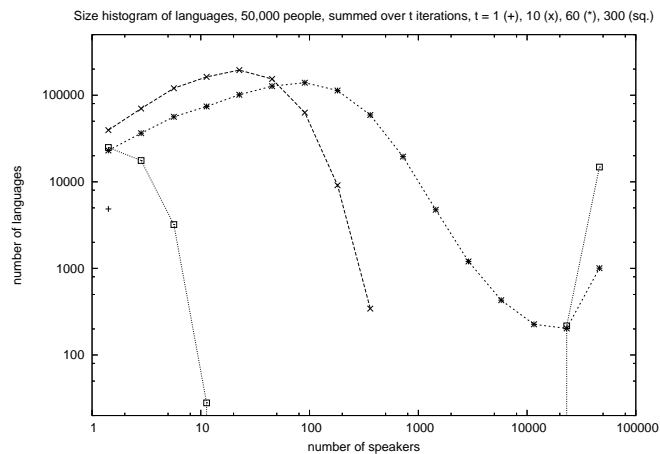


Figure 4: Time dependence of the distribution of language sizes, summed over 100 samples at  $L = 32$ . For long times dominance of one language develops, leading to an isolated peak at language sizes slightly below the total population size. Only intermediate times give the desired roughly parabolic shape. The same simulations were used for these language sizes as for the family sizes in fig. 1.

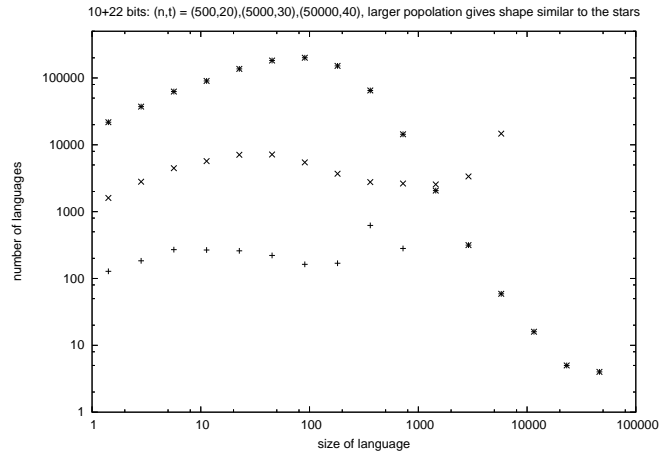


Figure 5: Population-size dependence of the distribution of language sizes. Same simulations as in fig. 2a for family sizes.

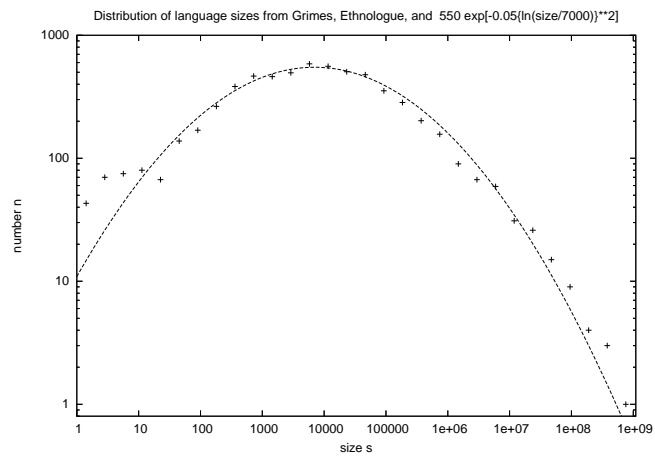


Figure 6: Empirical distribution of language sizes, from *Ethnologue*; see also Sutherland (2003: fig. 1), Wichmann (2005: fig. 6).

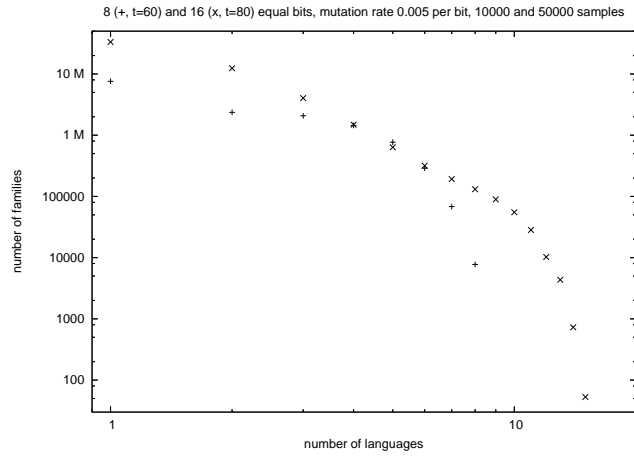


Figure 7: Flat version: distribution of family sizes for 8 and 16 bits.

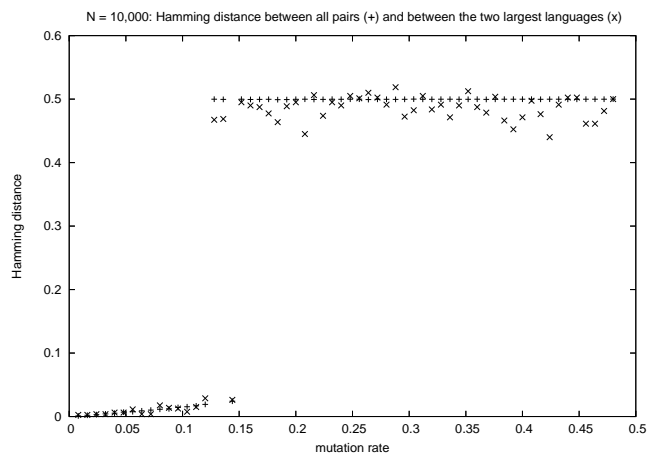


Figure 8: The average normalized Hamming distance for  $L = 8$  jumps from low values (dominance) to nearly  $1/2$  (fragmentation) when the mutation rate increases.

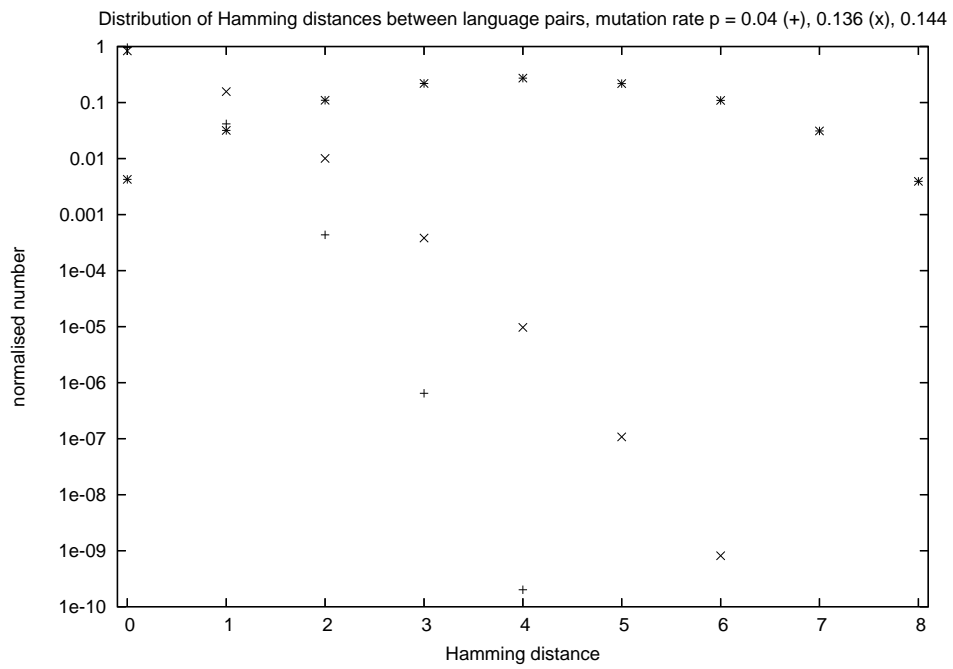


Figure 9: Flat version: distribution of  $k$  values, where  $k$  is the Hamming distance between an arbitrary pair of existing languages, at  $L = 8$ . The parabolic maximum corresponds to fragmentation at a high mutation rate, the two rapidly decaying curves to dominance at lower mutation rates.