AUTOMATIC SIMULATION OF HISTORICAL CHANGE

Raoul N. Smith Northwestern University

0.0 <u>Purpose</u>. One of the principal reasons for studying the history of a language has been to explain the system of its modern reflex, the contemporary language. This has been especially true in attempting to deal with certain anomalies in the modern language. But the role, if appropriate, of utilizing information concerning diachronic processes in a synchronic description is not at all clear. Recent studies describing contemporary languages, based purely on synchronically motivated grounds, suggest a much more intimate relation between a synchronic grammar and what has been 2 previously posited as a diachronic description of that language.

The two major problems involved in historical studies have been the statement of the sound change (or, as this has been reinterpreted, the grammar change) and the relation of this change to other diachronic changes, that is, its relative chronology. A great deal of attention has been paid to the former but very little to the latter, whose significance has greatly increased

T

See, for example, Lightner, 1965 and Schachter and Fromkin, 1968.

Since a naive native speaker of a language can not be expecte t know the history of his language, the reason for this reletion may lie in the manner in which the rules were added to the grammars of his ancestors. It is hoped that the future results of this study will help to shed light on this relation by comparing them and their associated grammars with synchronic descriptions.

due to recent results in generative grammar. One of the reasons for the lack of rigor in stating the relative chronology has probably been the large amounts of data required for input to the set of rules and the very large number of stages/rules which must be accounted for and the many permutations of these rules which should be tested. This lack of rigor, in turn, has made it very difficult to discuss coherently the historical development of a language.

The purpose of this paper is to discuss certain limited aspects of historical language change and suggest the possible use of the computer in approaching their solution. The types of problems discussed are only phonological and include only those changes conditioned by phonetic environment and which do not require syntactic information (for example, the change of the Old Russian unstressed infinitive ending /t'i/ to /t'/).

1.0 Language change. The discovery or postulation of the history of a language has been approached by two rather well-known methods. These are the reconstruction of the parent language of a set of genetically related languages and the reconstruction of an older stage of a language given a later stage. Both assume that languages evolve, that is, change (either gradually or abruptly) and that the relation of one language stage to the other is that one preceded the other. By comparing these stages in the development of the same or related sister languages one can

therefore reconstruct or recover the parent or proto language from which it or its sister languages evolved.

The two problems or reconstruction -- the reconstruction of a parent language of related sister languages and the reconstruction of an earlier form of a single language -- have been approached by separate methods.

The problem of reconstructing the parent of a set of related sister languages has been formalized by the so-called comparative method. The comparative method assumes, among other things and in a simplified version, that by comparing sets of sounds occurring in the same positions of the same words in the sister languages one can reconstruct the sound from whi h these sister sounds evolved. ("Same position" and "same word" may be difficult to define in a particular case.) For example, the word for "three" in various Indo-European languages is

Sanskrit	trayah
Greek	TPETS
Latin	tres
Gothic	preis
Lithuanian	trỹs
Old Church Slavonic	tr⊧je

Since five of the six forms have some kind of voiceless, dental stop in word initial position and that the fricative in Gothic can be accounted for by a change particular to Gothic, the assumption is made that the parent language of these, Proto-Indo-European, had

a voiceless, dental stop symbolized by *t, and this fact is highlighted by arranging these correspondences in tabular form:

PIE	Skr	GR	Lat	Go	Lith	ocs
*t	t	ĩ	t	Þ	t .	t

Once the proto language has been reconstructed, the correspondences used in the reconstruction can be reinterpreted as the results of phonological changes, for example, Proto-Indo-European *t became t in Old Church Slavonic, or, as it is usually expressed

PIE *t > OCS t

The problem of discovering an older stage of a language given only a later stage of that language is approached by examining certain alternations in the language at the later stage and from these postulate a proto form from which it could have evolved. The alternations in the later stage which are usually chosen are in the form of morphophonemic alternations, that is, phonemically differently shaped forms of the same morpheme. The assumption is made that these irregularities in the shape of one and the same morpheme must have been conditioned regularly so that by postulating one proto form and accounting for the change by a general rule, we have successfully reconstructed that form of the earlier language.

1

The forms of the proto language reconstructed by the comparative method can be interpreted as a statement of the state of the art in reconstruction for that language family.

For example, in Modern Russian the first person singular present of the verb "to be able" is /mogu/ and the second person singular is /možiš/. The first singular of "to read" is /čitaju/ and the second singular is /čitajiš/. From these and other sets of verbs we would conclude that the first singular ending is /u/ and that of the second singular is /iš/. Therefore, the stem morpheme alternants must be $\{mog \sim mož\}$ and $\{čitaj \sim čitaj\}$. The first of these two sets exhibits a morphophonemic alternation of g/ž. This same alternation in the first person singular occurs in other sets of verbs when the stem ends in a velar. From this (and other corroborative forms) we postulate that an earlier stage of Russian had one form for this verb stem, namely /mog/, and that before the vowel /i/ /g/ later became /ž/.

The proposal in this paper is to reverse the bottom-to-top model of the comparative method and that of internal reconstruction into a top-to-bottom generative model where the input forms are reconstructed lexical items and the rules are the set of postulated sound changes for the language. But there are two major difficulties in reversing the older models. One, the documented changes have often been incorrectly or incompletely stated and, two, the relative chronology of various rules has not been adequately . described. We hope to show how the computer can be used at least to test the accuracy of the rules and secondly to test or, hopefully, to help discover the relative chronology of the rules exhibited by their ordering.

The model proposed here is one where the proto language is a set of reconstructed forms (chosen, for example, from a standard reference work). The rules describing the phonological changes in that language are then described and ordered. As the program operates on these forms, the output from each rule represents a 2 synchronic stage in the development of that language. As final output one hopes to get the modern language. If any of the output is incorrect, then it is assumed to be from one of four possible sources: an incorrectly formulated rule (including analogical formation), a non-existent rule, an incorrectly ordered rule, or an incorrectly reconstructed form. Being able to differentiate which of these is the actual cause for the incorrect output is simplest only in the case where all of the output was the result of the application of only one rule.

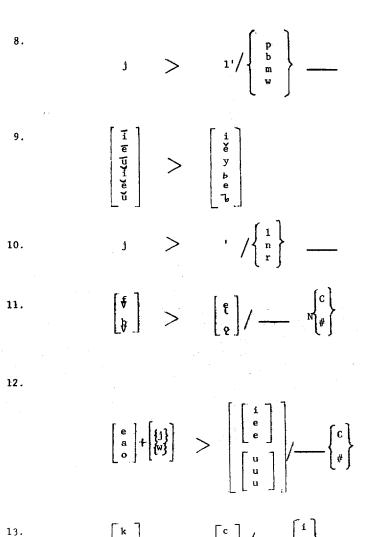
2.0 <u>A sketch of the phonological history of Russian</u>. The rules which were tested were an abridged version of a set presented by 1 the author in a recent paper. The rules attempt to account for certain aspects of the development of the phonological system of Céntemporary Standard Russian from a late form of Proto-Indo-European. These rules were:

1

Kantor, Marvin and R.N. Smith, "A sketch of the major developments in Russian historical phonology" (to appear). The original formulation was in terms of distinctive features; however, for this programmatic study a segmental notation has been used for ease of statement, etc.

1. $\begin{bmatrix} k^{h} \\ g^{h} \\ k^{\mu} \\ g^{\mu} \\ g^{\mu} \\ g^{\mu} \end{bmatrix} > \begin{bmatrix} k \\ g \\ k^{\mu} \\ g^{\mu} \\ g^{\mu} \end{bmatrix}$ $\begin{bmatrix} k y \\ g^{u} \end{bmatrix} > \begin{bmatrix} k \\ g \end{bmatrix}$ 2. s > x $\left| \begin{cases} i \\ u \\ r \\ k \end{cases} \right|$ non-obstruer 3. $\begin{bmatrix} \hat{k} \\ \hat{g} \end{bmatrix} > \begin{bmatrix} s \\ z \end{bmatrix}$ 4. $\begin{bmatrix} \left\{ \breve{a} \right\} \\ \left\{ \frac{\breve{a}}{a} \right\} \end{bmatrix} > \begin{bmatrix} \circ \\ a \end{bmatrix}$ 5. $\begin{bmatrix} k \\ g \\ x \end{bmatrix} > \begin{bmatrix} \check{c} \\ \check{z} \\ \check{s} \end{bmatrix} / - \begin{bmatrix} f \\ V \\ j \end{bmatrix}$ 6. $\begin{bmatrix} t \\ d \\ s \\ z \end{bmatrix} > \begin{bmatrix} \check{c} \\ \check{z} \\ \check{z} \\ \check{z} \end{bmatrix} / _ j$ 7.

ç



 $\begin{bmatrix} k \\ g \\ x \end{bmatrix} > \begin{bmatrix} c \\ z \\ s \end{bmatrix} / - \begin{bmatrix} i \\ \tilde{e} \end{bmatrix}$

· · · ·

k g x $\begin{bmatrix} c \\ z \\ s \end{bmatrix} / \notin __$ 14. > $\left. \right| \left\{ \begin{array}{c} \check{s} \\ \check{z} \\ \check{z} \\ \check{c} \\ j \end{array} \right\}$ > ě a 15. oro olo ere olo or ol er el 16. / >С _ c d t 17. >1 18. W v 1 С C' ţ 19. [e e a u 20. >ě 21. >e

9

Some of these rules appear in a less formal and explicit way in all standard texts of the history of the Russian language but none of them has all of these rules and in most cases there is little or no discussion of ordering. 3.0 <u>Description of test</u>. Approximately five hundred reconstructed Proto-Indo-European forms were chosen from Walde and Pokorny (1932). These were punched onto cards along with their English glosses. A separate Russian gloss was typed onto a print-out of the PIE forms. The transcription of Walde and Pokorny for the PIE lexical items was maintained as closely as possible, including such notations as subscript e and o. The only criterion imposed on choice of words was that they be as long as possible, so as to have a variety of environments.

The program was written in SNOBOL4 for the CDC 6400. Each rule set was numbered so as to coincide with the set of rules listed in section 2, with a zero appended to each rule number so as to allow for later insertions. Changing a rule consists at the moment of simple removal and replacement of cards. The history of a word or set of words can be gotten by allowing it to be processed with accompanying output generated by each rule set. Similarly, the lexicon for a particular stage can be generated by allowing the input to be processed up through the rule covering that stage and, if wanted, suppressing output from intermediate stages. With the availability of larger storage capacity the output from each stage can be generated once and stored in such a way that it can be referenced simply and thereby eliminate regeneration of input forms when the need for a rule change arises.

Frequency counters will be added in ofder to measure the functional load of a rule, at least in terms of dictionary

frequency. How this can be incorporated meaningfully into a theory of language change is not clear at this time.

The effect of borrowing can be simulated by introduction of lexical items just prior to a specific stage. There are too many variables involved in this case and the predictions have been poor. The effects of loss of original PIE are even more obvious but will require much further study.

4.0 <u>Discussion of results</u>. The program is obviously language dependent but the basic conception is of general applicability. The set of rules described in section 2.0 has been programmed and has successfully predicted the Modern Russian form from the PIE input in many cases including the following:

PIE	<u>Mod. Russian</u>
*b ^h e1ĝ ^h -	boloz-
*med ^h i-	mež-
*sloyg-	slug-
*ang ^{yh} i-	už
*korm-	sram
*gomb ^h -	zub
*g ^u rūg ^h ~	griz-

The number of forms of PIE which were not related to Modern Russian was much greater. The main reason is probably due to loss (again assuming a uniform parent PIE as the sole source of the lexicon). The differences between generated and actual Modern Russian could be accounted for in a few instances, for example *apsa did not become Mod. R 'osina' in part because there are no explicit rules in the program for the simplification of consonant clusters. But other incorrect outputs can not be accounted for in many instances by any easily accessible, documented rule. For example, there is no rule to handle the two disparate outputs from similar input forms with respect to the initial cluster 'sp-' in the forms

<u>PIE</u><u>Mod. Russian</u> *sperg-pr'ad-*spleng^h-selez-

Similarly, there is a rule eu>u, for example, to account for *leut- becoming l'ud-, and others, but there are at least two exceptions to this rule:

*b ^h reuk-	bros-
*b ^h leu-	bles-

Whether the original rule has too general an environmental condition or whether u was generated and later underwent some other changes, heretofore unpostulated, is unknown. It is possible that these forms should not have been considered as correspondences.

No case of an error in rule ordering has been found yet for this sample.

Examining the output is at times a forbidding task. It may help to simplify the discovery of causes of error by generating the output of each rule in the form of a KWIC index where the occurrences of each phone would be grouped together and the environments then made quite clear. When input and output are then compared one might find more easily why a rule was not applied or why it should be generalized, etc.

5.0 <u>Conclusions</u>. This paper is necessarily meant to be only a preliminary progress report and as such has raised many other questions in addition to its concrete results. Some of these questions are very basic, in particular, since many of the output forms could not be accounted for, should one really attempt to generate forms of a modern language from reconstructed lexical items if the rules used are not those postulated during the process of reconstruction since the former should be a record of the latter. Given, therefore, a set of reconstructed forms and a separate set of rules, it becomes very difficult to account for the source of errors in the output. Also, the assumption of a uniform, single-stage proto language may require many restrictions.

The computer can under ideal conditions be successfully used in testing hypothesized changes in the history of a language, given certain simplifying assumptions. It can be expected to operate best when the rules and reconstructed proto forms are established by the same investigator, working within the bounds established by a general theory of historical change.

Bibliography

- Kantor, Marvin and R. N. Smith, "A sketch of the major developments in Russian historical phonology" (to appear).
- Lightner, Theodore M., <u>Segmental</u> phonology of modern standard Russian. MIT dissertation, 1965.

Schachter, Paul and V. Fromkin, <u>A phonology of Akan: Akuapem</u>, <u>Asante and Fante</u>. (Working Papers in Phonetics No. 9) UCLA, 1968.

Walde, A. and J. Pokorny, <u>Vergleichendes</u> <u>Worterbuch</u> <u>der indo-</u> <u>germanischen</u> <u>Sprachen</u>. Berlin, 1932.