

DEFAULT FINITE STATE MACHINES AND FINITE STATE PHONOLOGY

Gerald Penn
Computational Linguistics Program
Carnegie Mellon University
Pittsburgh, PA 15213
Internet: penn@lcl.cmu.edu

Richmond Thomason
Intelligent Systems Program
University of Pittsburgh
Pittsburgh, PA 15260
Internet: thomason@isp.pitt.edu

Abstract

We propose DFSM's as an extension of finite state machines, explore some of their properties, and indicate how they can be used to formalize naturally occurring linguistic systems. We feel that this implementation of two-level rules may be more linguistically natural and easier to work with computationally. We provide complexity results that shed light on the computational situation.

INTRODUCTION

Two-level phonology combines the computational advantages of finite state technology with a formalism that permits phenomena to be described with familiar-looking rules. The problem with such a scheme is that, in practice, the finite state machines (FSM's) can grow too large to be manageable; one wants to describe them and to run them without having to deal with them directly. The KIMMO approach¹ seeks to achieve this by (1) decomposing the computational process into a battery of parallel finite state machines and (2) compiling rules (which notationally resemble familiar phonological rules, but which are interpreted declaratively) into these parallel finite state implementations. But the KIMMO formalism unfortunately gains no tractability in the process of compilation. Moreover, the compiler is complex enough to create software engineering problems, and this has led to practical difficulties, which in turn have made the KIMMO technology less generally available than one might wish. Here, we describe a different finite-state foundation for two-level rules, involving generalizations of FSM's which we call *Default Finite State Machines* (DFSM's). Whether or not this approach remains intractable after compilation is an open question; but even without compilation, we believe that it has some conceptual advantages as well.

¹See the discussion and references in (Sproat, 1992).

DFSM's extend FSM's (specifically, finite-state transducers) so that transitions can be context-sensitive, and enforce a preference for the maximally specific transitions. The first change allows phonological rules to appear as labels of transition arcs in transducers; the second change incorporates the elsewhere condition into the computational model.² DFSM's can be implemented directly, although there may be a method to compile them into a more efficient machine. We believe that either approach will be feasible for realistic linguistic applications (though, of course, not in the theoretically worst case). In particular, the direct implementation of DFSM's is very straightforward; no rule compiler is needed, since rules are labels on the arcs of the machines themselves. This implementation may not provide an optimal allocation of space and time usage at run time, but we believe that it will be adequate for testing and research purposes.

This presentation of DFSM's is confined to defining the basic ideas, presenting some examples of linguistic description, and providing a partial complexity analysis. In later work, we hope to explore descriptive and implementational issues further.

NOTATIONAL PRELIMINARIES

We assume an alphabet \mathcal{L} , with a reserved symbol $0 \notin \mathcal{L}$ for insertions and deletions. A *replacement* over \mathcal{L} is a pair of the form $l = \langle l, l' \rangle$ where (1) $l \in \mathcal{L}$ and (2) $l' \in \mathcal{L}$ or $l' = 0$; $\text{Replacements}_{\mathcal{L}}$ is the set of replacements over \mathcal{L} . $\text{US-strings}_{\mathcal{L}}$ is the set of strings over the set $\mathcal{L}^2 \cup [\mathcal{L} \times \{0\}]$ of replacements.

²The elsewhere condition is built into an implementation (due to Karttunen) of the TWOL rule compiler; see (Dalrymple et al., 1993), pp. 28-32. But on this approach, default reasoning and the elsewhere condition are not employed at a level of computation that is theoretically modeled; this reasoning is simply a convenient feature of the code that translates rules into finite state automata.

We use roman letters to denote themselves: for instance, 'l' denotes the letter l. Boldface letters denote constant replacements: for instance, **l** is the pair $\langle l, l \rangle$. Moreover, ϵ is the empty string over \mathcal{L} , and ϵ is the empty string over the \mathcal{L} replacements. When the name of a subset of \mathcal{L} (e.g. C) is written in boldface, (e.g. C), the set of identity pairings is intended (e.g., $C = \{l : l / l \in C\}$).

We use ordinary italics as variables over letters, and boldface italics as variables over replacements and strings of replacements. Ordinarily, we will use l for replacements and x, y for strings of replacements. Finally, we use ' $l:l$ ' for the pair $\langle l, l \rangle$.

Where $x \in \text{US-strings}_{\mathcal{L}}$, $U\text{-String}(x)$ is the underlying projection of x , and $S\text{-String}(x)$ is its surface projection. That is, if $x = \langle x, x' \rangle$, then $U\text{-String}(x) = x$ and $S\text{-String}(x) = x'$.

RULE NOTATION AND EXAMPLES

The rules with which we are concerned are like the rewrite rules of generative phonology; they are general, context-conditioned replacements. That is, a rule allows a replacement if (1) the replacement belongs to a certain type, and (2) the surrounding context meets certain constraints.

If we represent the contextual constraints extensionally, as sets of strings, a rule will consist of three things: a replacement type, and two sets of US-Strings. Thus, we can think of a rule as a triple $\langle X, Y, F \rangle$, where X and Y are sets of US-strings. Imagine that we are given a replacement instance l in a context $\langle x, y \rangle$, where x and y are US-strings. This contextualized replacement $\langle x, l, y \rangle$ satisfies the rule if $x \in X$, $y \in Y$, and $l \in F$.

For linguistic and computational purposes, the sets that figure in rules must somehow be finitely represented. The KIMMO tradition uses regular sets, which can of course be represented by regular expressions, for this purpose. We have not been able to convince ourselves that regular sets are needed in phonological applications.³ In-

³The issue here is whether there are any linguistically plausible or well-motivated applications of the Kleene star in stating phonological rules. For instance, take the English rule that replaces "e by 0 after a morpheme boundary preceded by one or more consonants preceded by a vowel." You could represent the context in question with the regular expression VC^*C ; but you could equally well use $VC|VCC|VCCC|VCCCC$. The only way to distinguish the two rule formulations is by considering strings that violate the phonotactic constraints of English; but as far as we can see, there are no intuitions about the results of applying English rules to underlying strings like $typppppe+ed$. We do not question the usefulness

of regular expressions in many computational applications, but are not convinced that they are needed in a linguistic setting. We would be interested to see a well motivated case in which the Kleene star is linguistically indispensable in formulating a two-level phonological rule. Such a case would create problems for the approach that we adopt here.

stead, we make the stronger assumption that contexts can be encoded by finite sets of strings. A string satisfies such a context when its left (or its right) matches one of the strings in this set. (Note that satisfaction is not the same as membership; infinitely many strings can satisfy a finite set of strings.) Assuming a finite alphabet, all replacement types will be finite sets. With these assumptions, a rule can be finitely encoded as a pair $\langle \langle X, Y \rangle, F \rangle$, where the sets X and Y are finite, and F is a replacement type.

Language:

Let $\mathcal{L} = \{a, b, \dots, z, +, \#, ', '\}$

Declare the following subsets of \mathcal{L} :

$C = \{b, c, d, f, g, h, j, k, l, m, n, p, q, r, s, t, v, w, x, y, z\}$
 $Csib = \{s, x, z\}$

Example rules:

Example 1

Rule encoding: $\langle \langle \{\epsilon\}, \{\epsilon\} \rangle, \langle \langle +, 0 \rangle \rangle$

Rule notation: $+ \rightarrow 0 / _$

Rule description: Delete +.

Example 2

Rule encoding: $\langle \langle C, \langle \langle +, 0 \rangle \rangle \rangle, \langle \langle y, i \rangle \rangle$

Rule notation: $y \rightarrow i / C_+ : 0$

Rule description: Replace y by i before a morpheme boundary and after a constant US-consonant, i.e. after $\langle l, l \rangle$, where $l \in C$.

Example 3

Rule encoding: $\langle \langle \langle \{sh\}, \{l \wedge (\# , 0) / l \in Csib\} \rangle, \langle \langle +, e \rangle \rangle \rangle$

Rule notation: $+ \rightarrow e / sh_Csib \# : 0$

Rule description: Replace + with e after sh and before a suffix in Csib.

Example rule applications:

1. The rule encoded in Example 1 is satisfied by $\langle +, 0 \rangle$ in the context $\langle \text{cat}, s \rangle$ because (1) for some x , $\text{cat} = x \wedge \epsilon$, (2) for some y , $s = \epsilon \wedge y$, and (3) $\langle +, 0 \rangle \in \langle \langle +, 0 \rangle \rangle$.

of regular expressions in many computational applications, but are not convinced that they are needed in a linguistic setting. We would be interested to see a well motivated case in which the Kleene star is linguistically indispensable in formulating a two-level phonological rule. Such a case would create problems for the approach that we adopt here.

2. The rule encoded in Example 2 is not satisfied by $\langle y, i \rangle$ in the context $\langle \text{spot} + :t, +:0\text{ness} \rangle$ because there is no α such that $\text{spot} + :t = \alpha \hat{\ } l$, where $l \in C$.
3. The rule encoded in Example 3 is not satisfied by $\langle +, 0 \rangle$ in the context $\langle \text{ash}, s \# :0 \rangle$. In fact, the context is satisfied: (1) $\text{sh} = \alpha \hat{\ } \text{sh}$ for some α and (2) $s \# :0 \in C \text{sib} \hat{\ } y$ for some y . (3.1) Moreover, the underlying symbol of the replacement (namely, $+$) matches the argument of the rule's replacement function. Under these circumstances, we will say that the rule is *applicable*. But the rule is not satisfied, because (3.2) the surface symbol of the replacement (namely, 0) does not match the value of the rule's replacement function (namely, e): thus, $\langle +, 0 \rangle \notin \{ \langle +, e \rangle \}$.

INDEXED STRINGS AND RULES

We now restate the above ideas in the form of formal definitions.

Definition 1. Context type.

A context type is a pair $C = \langle X, Y \rangle$, where X and Y are sets of US-Strings.

Definition 2. Indexed US-strings.

An indexed US-String over \mathcal{L} is a triple $\langle \alpha, l, y \rangle$, where $\alpha, y \in \text{US-strings}_{\mathcal{L}}$ and $l \in \text{Replacements}_{\mathcal{L}}$.

An indexed US-string is a presentation of a nonempty US-string that divides the string into three components: (1) a replacement occurring in the string, (2) the material to the left of that replacement, and (3) the material to the right of it. Where $\langle \alpha, l, y \rangle$ is an indexed string, we call α the *left context* of the string, y the *right context* of the string, and l the *designated replacement* of the string.

A rule licenses certain sorts of replacements in designated sorts of environments, or *context types*. For instance, we may be interested in the environment after a consonant and before a morpheme boundary. Here, the phrase "after a consonant" amounts to saying that the string before the replacement must end in a consonant, and the phrase "before a morpheme boundary" says that the string after the replacement must begin in a morpheme boundary. Thus, we can think of a context type as a pair of constraints, one on the US-string to the left of the replacement, and the other on the US-string to its right. If we identify such constraints with the set of strings that satisfy them, a context type is then a pair of sets of US-strings; and an indexed string satisfies a context type in case its left and right context belong to the corresponding types.

Definition 3. Replacement types.

A replacement type over \mathcal{L} is a partial function F from $\mathcal{L} \cup \{0\}$ to $\mathcal{L} \cup \{0\}$. (Thus, a replacement type is a certain set of replacements.) $\text{Dom}(F)$ is the domain of F .

Definition 4. Rules.

A rule is a pair $\mathcal{R} = \langle C, F \rangle$, where C is a context type and F is a replacement type.

Definition 5. Rule applicability.

A rule $\langle \langle X, Y \rangle, F \rangle$ is applicable to an indexed string $\langle \alpha, l, y \rangle$ if and only if $\alpha \in X$, $y \in Y$, and $F(l)$ is defined, i.e., $l \in \text{Dom}(F)$.

Definition 6. Rule satisfaction.

An indexed string $\langle \alpha, l, y \rangle$ satisfies a rule $\langle C, F \rangle$ if and only if $\alpha \in X$, $y \in Y$, and $F(l) = l'$.

The above definitions do not assume that the contexts are finitely encodable. But, as we said, we are assuming as a working hypotheses that phonological contexts are finitely encodable; this idea was incorporated in the method of rule encoding that we presented above. We now make this idea explicit by defining the notion of a finitely encodable rule.

Definition 7. LeftExp(X), RightExp(X)

$$\text{LeftExp}(X) = \{ \alpha \hat{\ } \alpha / \alpha \in X \}$$

$$\text{RightExp}(X) = \{ \alpha \hat{\ } \alpha / \alpha \in X \}$$

Definition 8. Finite encodability

A subset X of $\text{US-strings}_{\mathcal{L}}$ is *left-encoded* by a set U in case $X = \text{LeftExp}(U)$, and is *right-encoded* by V in case $X = \text{RightExp}(V)$. (It is easy to get confused about the usage of "left" and "right" here; in left encoding, the left of the encoded string is arbitrary, and the right must match the encoding set. We have chosen our terminology so that a left context type will be left-encoded and a right context type will be right-encoded.)

A context type $C = \langle X, Y \rangle$ is encoded by a pair $\langle U, V \rangle$ of sets in case U left-encodes X and V right-encodes Y .

A rule $\mathcal{R} = \langle C, F \rangle$ is finitely encoded by a rule encoding structure $\langle \langle U, V \rangle, g \rangle$ in case $\langle U, V \rangle$ encodes C , $g = F$, and U and V are finite.

In the following material, we will not only confine our attention to finitely encodable rules, but will refer to rules by their encodings; when the notation $\langle \langle X, Y \rangle, F \rangle$ appears below, it should be read as a rule encoding, not as a rule. Thus, for instance, the indexed string $\langle \text{cat}, +:0, s \rangle$ satisfies the rule (encoded by) $\langle \langle \{ \epsilon \}, \{ \epsilon \} \rangle, \{ \langle +, 0 \rangle \} \rangle$, even though $\text{cat} \notin \{ \epsilon \}$.

SPECIFICITY OF CONTEXT TYPES AND RULES

We have a good intuitive grasp of when one context type is more specific than another. For instance, the context type *preceded by a back vowel* is more specific than the type *preceded by a vowel*; the context type *followed by an obstruent* is neither more nor less specific than the type *followed by a voiced consonant*; the context type *preceded by a vowel* is neither more nor less specific than the type *followed by a vowel*.

Since we have identified context types with pairs of sets of strings, we have a very natural way of defining specificity relations such as "more specific than", "equivalent", and "more specific than or equivalent": we simply use the subset relation.

Definition 9. $C \leq C'$.

Let $C = \langle X_1, Y_1 \rangle$ and $C' = \langle X_2, Y_2 \rangle$ be context types. $C \leq C'$ if and only if $X_1 \subseteq X_2$ and $Y_1 \subseteq Y_2$.

Definition 10. $C \equiv C'$.

$C \equiv C'$ if and only if $C \leq C'$ and $C' \leq C$.

Definition 11. $C < C'$.

$C < C'$ if and only if $C \leq C'$ and $C' \not\leq C$.

It is not in general true that if $LeftExp(X) \subseteq LeftExp(Y)$, then $X \subseteq Y$; for instance, $LeftExp(\{aa, ba\}) \subseteq LeftExp(\{a\})$, but $\{aa, ba\} \not\subseteq \{a\}$. However, we can easily determine the specificity relations of two contexts from their finite encodings:

Lemma 1. $LeftExp(X) \subseteq LeftExp(Y)$ iff for all $x \in X$ there is a $y \in Y$ such that for some z , $x = z \hat{\ } y$. Similarly, $RightExp(X) \subseteq RightExp(Y)$ iff for all $x \in X$ there is a $y \in Y$ such that for some z , $x = y \hat{\ } z$.

Proof of the lemma is immediate from the definitions. It follows from the lemma that there is a tractable algorithm for testing specificity relations on finitely encodable contexts:

Lemma 2. Let C be finitely encoded by $\langle X_1, X_2 \rangle$ and C' be finitely encoded by $\langle Y_1, Y_2 \rangle$. Then there is an algorithm for testing whether $C \leq C'$ that is no more complex than $O(m \times n \times k)$, where $m = \max(|X_1|, |X_2|)$, $n = \max(|Y_1|, |Y_2|)$, and k is the length of the longest string in $Y_1 \cup Y_2$.

Proof. Test whether for each $x_1 \in X_1$ there is a $y_1 \in Y_1$ that matches the end of x_1 . Then perform a similar test on X_2 and Y_2 .

DFSM'S

A DFSM's transitions are labelled with finitely encodable rules rather than with pairs of symbols. Moreover, nondeterminism is restricted so that in case of conflicting transitions, a maximally specific transition must be selected. The critical definition is that of *minimal satisfaction of an arc by an indexed path*, where an indexed path represents a DFSM derivation, by recording the state transitions and replacements that are traversed in processing a US-String.

Definition 12. *Arcs.*

An arc over a set S of states and alphabet \mathcal{L} is a triple $A = \langle s, s', \mathcal{R} \rangle$, where $s, s' \in S$ and \mathcal{R} is a rule over \mathcal{L} .

Definition 13. *DFSMs.*

A DFSM on \mathcal{L} is a structure $\mathcal{M} = \langle S, i, T, \mathcal{A} \rangle$, where S is a finite set of states, $i \in S$ is the initial state, $T \subseteq S$ is the set of terminal states, and \mathcal{A} is a set of arcs over S on \mathcal{L} .

Definition 14. *Paths.*

A path π or $\pi(s_0, s_n)$ over \mathcal{M} from state s_0 to state s_n is a string $s_0 l_1 s_1 l_1 \dots l_n s_n$, where for all m , $0 \leq m \leq n$, s_m is a state of \mathcal{M} and $l_m \in US\text{-strings}_{\mathcal{L}}$.

Remark 1: $n \geq 0$, so that the simplest possible path has the form s , where s is a state. *Remark 2:* we use the notations π and $\pi(s, s')$ alternatively for the same path; the second notation provides a way of referring to the beginning and end states of the path.

Definition 15. *Recovery of strings from paths.*

Let $\pi = s_0 l_1 s_1 l_1 \dots l_n s_n$. Then $String(\pi) = l_1 \dots l_n$.

Definition 16. *Indexed paths.*

An indexed path over \mathcal{M} is a triple $\langle \pi, l, \pi' \rangle$ where π, π' are paths, and $l_m \in US\text{-strings}_{\mathcal{L}}$. $\langle \pi, l, \pi' \rangle$ is an indexing of path σ if and only if $\sigma = \pi \hat{\ } l \hat{\ } \pi'$.

Definition 17. *Applicability of an arc to an indexed path.*

An arc $\langle u, u', \mathcal{R} \rangle$ is applicable to an indexed path $\langle \pi(s, t), l, \pi'(s', t') \rangle$ if and only if $t = u$ and the rule \mathcal{R} is applicable to the indexed string $\langle String(\pi), l, String(\pi') \rangle$.

Definition 18. *Satisfaction of an arc by an indexed path.*

$\langle \pi(s, t), l, \pi'(s', t') \rangle$ satisfies an arc $\langle u, u', \mathcal{R} \rangle$ if and only if $t = u$, $s' = u'$, and the indexed string $\langle String(\pi), l, String(\pi') \rangle$ satisfies the rule \mathcal{R} .

Definition 19. *Minimal satisfaction of an arc by an indexed path.*

$\langle \pi, l, \pi' \rangle$ minimally satisfies an arc $A = \langle s, s', \mathcal{R} \rangle$ of \mathcal{M} if and only if $\langle \pi, l, \pi' \rangle$ satisfies A and there is no state s'' and arc $A' = \langle s, s'', \mathcal{R}' \rangle$ of \mathcal{M} such that $A' = \langle s, s'', \mathcal{R}' \rangle$ is applicable to $\langle \pi, l, \pi' \rangle$ and $\mathcal{R}' < \mathcal{R}$.

As we said, the above definition is the crucial component of the definition of DFSM's. According to this definition, to see whether a DFSM derivation is correct, you must check that each state transition represents a maximally specific rule application. This means that at each stage the DFSM does not provide another arc with a competing replacement and a more specific context. ("Competing" means that the underlying symbols of the replacement match; a replacement competes even if the surface symbols does not match the letter in the US-String being tested.)⁴

Definition 20. *Indexed path acceptance by a DFSM.*

$\mathcal{M} = \langle \mathcal{S}, i, T, \mathcal{A} \rangle$ accepts an indexed path $\langle \pi, l, \pi' \rangle$ if and only if there is an arc $A' = \langle s, s', \mathcal{R}' \rangle$ of \mathcal{M} that is minimally satisfied by $\langle \pi, l, \pi' \rangle$.

Definition 21. *Path acceptance by a DFSM.*

$\mathcal{M} = \langle \mathcal{S}, i, T, \mathcal{A} \rangle$ accepts a path $\pi(s, s')$ if and only if \mathcal{M} accepts every indexing of π , $s = i$, and $s' \in T$.

Definition 22. *US-String acceptance by a DFSM.*

\mathcal{M} accepts $\mathfrak{s} \in US\text{-strings}_{\mathcal{L}}$ if and only if there is a path π such that \mathcal{M} accepts π , where $\mathfrak{s} = String(\pi)$.

Definition 23. *Generation of SF from UF by a DFSM.*

\mathcal{M} generates a surface form x' from an underlying form x (where x and x' are strings over \mathcal{L}) if and only if there is a $\mathfrak{s} \in US\text{-strings}_{\mathcal{L}}$ such that \mathcal{M} accepts \mathfrak{s} , where $U\text{-String}(\mathfrak{s}) = x$ and $S\text{-String}(\mathfrak{s}) = x'$.

EXAMPLE: SPELLING RULES FOR ENGLISH STEM+SUFFIX COMBINATIONS

The following is an adaptation of the treatment in Antworth (1990) of English spelling rules, which

⁴This use of competition builds some directional bias into the definition of DFSM's, i.e., some preference for their use in generation. Even if we are using DFSM's for recognition, we will need to verify that the recognized string is generated from an underlying form by a derivation that does not allow more specific competing derivations.

in turn is taken from Karttunen and Wittenburg (1983).

- $\mathcal{M} = \langle \mathcal{S}, i, T, \mathcal{A} \rangle$, where $\mathcal{S} = \{i, s, t\}$. $T = \{t\}$.
 - *Task of i:* Begin and process left word boundary.
 - *Task of s:* Process stem and suffixes.
 - *Task of t:* Quit, having processed right word boundary.

• *Remark:* the small number of states is deceptive, since contexts are allowed on the arcs. An equivalent finite-state transducer would have many hundreds of states at least.

• *Remark:* the relatively small number of arcs enumerated below is also deceptive, since two of these "arcs," arc 3 and arc 13, are actually schemes. In the following discussion we will speak loosely and refer to these schemes as arcs; this will simplify the discussion and should create no confusion.

• Declare the following subsets of \mathcal{L} :

$Ltr = \{a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z\}$

$C = \{b, c, d, f, g, h, j, k, l, m, n, p, q, r, s, t, v, w, x, y, z\}$

$Csib = \{s, x, z\}$

$Cpal = \{c, g\}$

$V = \{a, e, i, o, u\}$

$Vbk = \{a, o, u\}$;

• Where $s, s' \in \mathcal{S}$, let $\mathcal{A}_{s,s'} = \{A / A \in \mathcal{A} \text{ and for some } \mathcal{R}, A = \langle s, s', \mathcal{R} \rangle\}$. We present arcs by listing the rules associated with the arcs, for each appropriate pair $\langle s, s' \rangle$ of states. We will give each arc a numerical label, and give a brief explanation of the purpose of the arc.

• Arcs in $\mathcal{A}_{i,s}$:

1. $\# \rightarrow 0 / _$
Delete left word boundary.

• Arcs in $\mathcal{A}_{s,s}$:

2. $+ \rightarrow 0 / _$
Delete morpheme boundary.
3. $l \rightarrow l / _ : l \in Ltr$
Any underlying letter is normally unchanged.
4. $' \rightarrow ' / _$
Apostrophe is normally unchanged.
5. $' \rightarrow ' / _$
Stress is normally unchanged.
6. $+ \rightarrow e / [Csib | ch | sh | y:i] _ s [+:0 | \#:0]$
Epenthesis before -s suffix.

7. $y \rightarrow i / C_ + :0$
Spell *y* as *i* after consonant and before suffix.
8. $y \rightarrow y / C_ + :0 [i:i | ':]$
Exception to Rule 7; cf. "trying", "fly's".
9. $s \rightarrow 0 / [+0 | +e]s + :0 _$
Delete possessive 's after plural suffix.
10. $e \rightarrow 0 / VCC^+ _ + :0 V$
Elision.⁵
11. $e \rightarrow e / VC^+ Cpa1 _ + :0 Vbk$
Exception to Rule 10.
12. $i \rightarrow y / _ e:0 + :0 i$
Spell *i* as *y* before elided *e* before *i*-initial suffix.
13. $+ \rightarrow l / ' : 0 C^+ V l : l _ [V | y] :$
 $l \in \{b, d, g, l, m, n, p, r, t\}$
Gemination.

• Arcs in $A_{s,i}$:

14. $\# \rightarrow 0 / _$
Delete right word boundary.

• Illustrations

I. The derivation that relates #kiss+s# to Okisses0 proceeds as follows.

1. Begin in state *i* looking at #:0.
2. Follow arc 2 to *s*, recognizing k:k. (This is the only applicable arc.)
3. Follow arc 3 to *s*, recognizing i:i. (This is the only applicable arc.)
4. Follow arc 3 to *s*, recognizing s:s. (This is the only applicable arc.)
5. Follow arc 3 to *s*, recognizing s:s. (This is the only applicable arc.)
6. Follow arc 6 to *s*, recognizing +:e. (Arc 2 is also applicable here; but see the next illustration.)
7. Follow arc 3 to *s*, recognizing s:s. (This is the only applicable arc.)
8. Follow arc 14 to *f*, recognizing #:0. (This is the only applicable arc.)

II. No derivation relates #kiss+s# to Okiss0s0. Any such derivation would have to proceed like the above derivation through Step 5. At the next step, the conditions for two arcs are met: arc 2 (replacing + with 0) and arc 6 (replacing + with e). Since the context of the latter

⁵Here, C^+ can be any string of no more than four consonants.

arc is more specific, it must apply; there is no derivation from this point using arc 2.

III. The derivation that relates #try+ing# to Otry0ing0 proceeds as follows.

1. Begin in state *i* looking at #:0.
2. Follow arc 2 to *s*, recognizing t:t. (This is the only applicable arc.)
3. Follow arc 3 to *s*, recognizing r:r. (This is the only applicable arc.)
4. Follow arc 8 to *s*, recognizing y:y. (There are three applicable arcs at this point: arc 3, arc 7, and arc 8. However, arcs 3 and 7 are illegal here, since their contexts are both less specific than arc 8's.)
5. Follow arc 2 to *s*, recognizing +:0. (This is the only applicable arc.)
6. Follow arc 3 to *s*, recognizing i:i. (This is the only applicable arc.)
7. Follow arc 3 to *s*, recognizing n:n. (This is the only applicable arc.)
8. Follow arc 3 to *s*, recognizing g:g. (This is the only applicable arc.)
9. Follow arc 14 to *f*, recognizing #:0. (This is the only applicable arc.)

IV. No derivation relates #try+ing# to Otri0ing0. Any such derivation would have to proceed like the above derivation through Step 3. At the next step, arc 7 cannot be traversed, since arc 8 is also applicable and its context is more specific. Therefore, no arc is minimally satisfied and the derivation halts at this point.

COMPUTATIONAL COMPLEXITY

We now consider the complexity of using DFSM's to create one side of a US-string, given the other side as input. There are basically two tasks to be analyzed:

- **DFSM GENERATION:** Given a DFSM, D , over an alphabet, \mathcal{L} , and an underlying form, u , does D generate a surface form, s , from u ?
- **DFSM RECOGNITION:** Given a DFSM, D , over an alphabet, \mathcal{L} , and a surface form, s , does D generate an underlying form, u , from s ?

These two tasks are related to the tasks of KIMMO GENERATION and KIMMO RECOGNITION, the various versions of which Barton et al. (1987) proved to be NP-complete or worse.

Relationship to Kimmo

The DFSM is not a generalization of KIMMO; it is an alternative architecture for two-level rules.

KIMMO takes a programming approach; it provides a declarative rule formalism, which can be related to a very large FS automaton or to a system of parallel FS automata. The automata are in general too unwieldy to be pictured or managed directly; they are manipulated using the rules. By integrating rules into the automata, the DFSM approach provides a procedural formalism that is compact enough to be diagrammed and manipulated directly.

DFSM rules are procedural; their meaning depends on the role that they play in an algorithm. In a DFSM with many states, the effect achieved by a rule (where a rule is a context-dependent replacement type) will in general depend on how the rule is attached to states. In practice, however, the proceduralism of the DFSM approach can be limited by allowing only a few states, which have a natural morphonemic interpretation. The English spelling example that we presented in the previous section illustrates the idea. There are only four states. Of these, two of them delimit word processing; one of them begins processing by traversing a left word boundary, the other terminates processing after traversing a final word boundary. Of the remaining two states, one processes the word; all of the rules concerning possible replacements are attached to arcs that loop from this state to itself. The other is a nonterminal state with no arcs leading from it. In the example, the only purpose of this state is to render certain insertions or deletions obligatory, by "trapping" all US-strings in which the operation is not performed in the required context.

In cases of this kind, where the ways in which rules can be attached to arcs are very restricted, the proceduralism of the DFSM formalism is limited. The uses of rules in such cases correspond roughly to two traditional types of phonological constructs: rules that allow certain replacements to occur, and constraints that make certain replacements obligatory.

Although DFSM's are less declarative than KIMMO, we believe that it may be possible to interpret at least some DFSM's (those in which the roles that can be played by states are limited) using a nonmonotonic formalism that provides for prioritization of defaults, such as prioritized default logic; see (Brewka, 1993). In this way, DFSM's could be equated to declarative, axiomatic theories with a nonmonotonic consequence relation. But we have not carried out the details of this idea.

Though it is desirable to constrain the number of states in a DFSM, there may be applications in which we may want more states than in the English example. For instance, one natu-

ral way to process vowel harmony would multiply states by creating a word-processing state for each vowel quality. Multiple modes of word-processing could also be used to handle cases (as in many Athabaskan languages) where different morphophonemic processes occur in different parts of the word.

If they are desired, local translations of the four varieties of KIMMO rules⁶ into DFSM's are available, by using only one state plus a sink state. The following correspondences provide translations, in polynomial time, to one or more DFSM arcs:

- *Exclusion*, $u : s / \Leftarrow LC_RC$: an arc $u \rightarrow s / LC_RC$ from the state to a sink state.
- *Context Restriction*, $u : s \Rightarrow LC_RC$: a loop $u \rightarrow s / LC_RC$, and an arc $u \rightarrow s / _$ to a sink state.
- *Surface Coercion*, $u : s \Leftarrow LC_RC$: a loop $u \rightarrow s / LC_RC$, and for each surface character $s' \in L$, an arc $u \rightarrow s' / LC_RC$ to a sink state.
- *Composite*, $u : s \Leftrightarrow LC_RC$: all of the arcs mentioned in Context Restriction or Surface Coercion.

Extended DFSM's

The differences between KIMMO and DFSM's prohibit the complexity analysis for the corresponding two KIMMO problems from naturally extending to an analysis of DFSM generation and recognition. In fact, we can define an *extended DFSM* (EDFSM), which drops the finite encodability requirement that KIMMO lacks, for which we have the following result:

Theorem 1. EDFSM GENERATION is PSPACE-hard

Proof by reduction of REGULAR EXPRESSION NON-UNIVERSALITY (see Figure 1). Given an alphabet Σ , and a regular expression, $a \neq \phi$, over Σ , we define an EDFSM over the alphabet, $\Sigma \cup \{\$\}$, where $\$ \notin \Sigma$. We choose one non-empty string $\alpha \in L(a)$ of length n . The EDFSM first recognizes each character in α , completing the task at state s_0 :

$$\begin{aligned} \alpha_1 &\rightarrow \alpha_1 / (L:L)^* _ (L:L)^* \text{ }^7 \\ &\vdots \\ \alpha_n &\rightarrow \alpha_n / (L:L)^* _ (L:L)^* \end{aligned}$$

From s_0 , there are two arcs, which map to different states:

⁶Sproat (1992), p. 145.

⁷Unlike with normal DFSM's, we will use regular expressions for the contexts themselves in EDFSM's, not their encodings, since they may be infinite anyway.

$$\begin{aligned} \$ &\rightarrow \$ / \Sigma^* _ \Sigma^* \\ \$ &\rightarrow \$ / (a + \$) _ (a + \$) \end{aligned}$$

where the latter rule traverses to some state s_1 , with a being the expression which replaces each atom, b , in a by its constant replacement, $b:b$, and likewise for Σ .

From s_1 , the EDFSM then recognizes α again, terminating at the only final state. We provide this EDFSM, along with the input $\alpha\$ \alpha$ to EDFSM GENERATION. This EDFSM can accept $\alpha\$ \alpha$ if and only if, at state s_0 , the context $\langle \Sigma^*, \Sigma^* \rangle$ is not more specific than the context $\langle (a + \$), (a + \$) \rangle$. So, we have:

$$\begin{aligned} \langle \Sigma^*, \Sigma^* \rangle &\not\prec \langle (a + \$), (a + \$) \rangle \\ \Leftrightarrow \langle \Sigma^*, \Sigma^* \rangle &\not\leq \langle (a + \$), (a + \$) \rangle \\ \text{or } \langle \Sigma^*, \Sigma^* \rangle &\equiv \langle (a + \$), (a + \$) \rangle \\ \Leftrightarrow \Sigma^* &\not\subseteq L(a + \$) \\ \text{or } \Sigma^* &= L(a + \$) \\ \Leftrightarrow \Sigma^* &\not\subseteq L(a + \$), \text{ since } \$ \notin \Sigma, \\ \Leftrightarrow \Sigma^* &\not\subseteq L(a) \cup \{ \$ \} \\ \Leftrightarrow \Sigma^* &\not\subseteq L(a), \text{ since } \$ \notin \Sigma, \\ \Leftrightarrow \Sigma^* &\not\subseteq L(a) \\ \Leftrightarrow L(a) &\neq \Sigma^* \text{ (we know } L(a) \subseteq \Sigma^*) \end{aligned}$$

The translation function is linear in the size of the input. \square

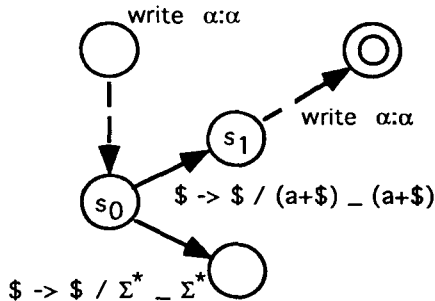


Figure 1. EDFSM constructed in Theorem 1.

The Complexity of DFSM GENERATION

Finite encodability foils the above proof technique, since one can no longer express arbitrary regular expressions over pairs in the contexts of rules. In fact, as we demonstrated above, there is a polynomial-time algorithm for comparing the specificities of finitely-encodable contexts. Finite encodability does not, however, restrict the complexity of DFSM's enough to make DFSM GENERATION polynomial time:

Theorem 2. DFSM GENERATION is NP-complete.

Proof DFSM GENERATION is obviously in NP. The proof of NP-hardness is a reduction of

3-SAT. Given an input formula, w , we construct a DFSM consisting of one state over an alphabet consisting of 0, 1, #, one symbol, u_i , for each variable in w , and one symbol, c_j , for each conjunct in w . Let m be the number of variables in w , and n , the number of conjuncts. For each variable, u_i , we add four loops:

$$\begin{aligned} u_i &\rightarrow 1 / \# : \# u_1 : \mathcal{L} \dots u_{i-1} : \mathcal{L} _ \\ u_i &\rightarrow 0 / \# : \# u_1 : \mathcal{L} \dots u_{i-1} : \mathcal{L} _ \\ u_i &\rightarrow 1 / u_i : 1 u_{i+1} : \mathcal{L} \dots u_m : \mathcal{L} \mathcal{L} : \mathcal{L} \\ &\quad u_1 : \mathcal{L} \dots u_{i-1} : \mathcal{L} _ \\ u_i &\rightarrow 0 / u_i : 0 u_{i+1} : \mathcal{L} \dots u_m : \mathcal{L} \mathcal{L} : \mathcal{L} \\ &\quad u_1 : \mathcal{L} \dots u_{i-1} : \mathcal{L} _ \end{aligned}$$

The first two choose an assignment for a variable, and the second two enforce that assignment's consistency. For each conjunct, $l_{j_1} \vee l_{j_2} \vee l_{j_3}$, where the l 's are literals, we also add three loops, one for each literal. The loops enforce a value of 1 on the symbol u_{j_i} if l_{j_i} is a positive literal, or 0, if it is negative. For example, for the conjunct $u_1 \vee \neg u_3 \vee u_4$, we add the following three rules:

$$\begin{aligned} c_j &\rightarrow c_j / u_1 : 1 u_2 : \mathcal{L} \dots u_m : \mathcal{L} _ \\ c_j &\rightarrow c_j / u_3 : 0 u_4 : \mathcal{L} \dots u_m : \mathcal{L} _ \\ c_j &\rightarrow c_j / u_4 : 1 u_5 : \mathcal{L} \dots u_m : \mathcal{L} _ \end{aligned}$$

Thus, the input to DFSM GENERATION is the above DFSM plus an input string created by iterating the substring $u_1 \dots u_m c_j$ for each conjunct. The input string corresponding to the formula, $(\neg u_1 \vee u_2 \vee u_4) \wedge (\neg u_2 \vee u_3 \vee \neg u_4) \wedge (u_1 \vee u_2 \vee u_3)$, would be $\# u_1 u_2 u_3 u_4 c_1 u_1 u_2 u_3 u_4 c_2 u_1 u_2 u_3 u_4 c_3$. The DFSM accepts this input string if and only if the input formula is satisfiable; and this translation is linear in $m + n$. \square

Compilation

Of course, we should consider whether the complexity of DFSM GENERATION can be compiled out, leaving a polynomial-time machine which accepts input strings. This can be formalized as the separate problem:

- **FIXED-DFSM-GENERATION:** For some DFSM, D , over alphabet, \mathcal{L} , given an underlying form, u , does D generate a surface form, s , from u ?

Whether or not FIXED DFSM GENERATION belongs to P remains an open problem. It is, of course, no more difficult than the general DFSM GENERATION problem, and thus no more difficult than NP-complete. The method used in the proof given above, however, does not naturally extend to the case of FIXED DFSM GENERATION, since we cannot, with a fixed DFSM, know in advance

how many variables to expect in a given input formula, without which we cannot use the same trick with the left context to preserve the consistency of variable assignment.

Even more interestingly, the technique used in the proof of PSPACE-hardness of EDFSM GENERATION does not naturally extend to fixed EDFSM's either; thus, whether or not FIXED DFMSM GENERATION belongs to P is an open question as well⁸. Dropping finite encodability, of course, affects the compilation time of the problem immensely.

Nulls

The two proofs we have given remain valid if we switch all of the underlying forms with their surface counterparts. Thus, without nulls, EDFSM RECOGNITION is PSPACE-hard, DFMSM RECOGNITION is NP-complete, and, if FIXED DFMSM GENERATION is in P, then we can presumably use the same compilation trick with the roles of underlying and surface strings reversed to show that FIXED DFMSM RECOGNITION is in P as well.

If nulls are permitted in surface realizations, however, DFMSM RECOGNITION becomes much more difficult, even with finite encodability enforced:

Theorem 3. DFMSM RECOGNITION with nulls is PSPACE-hard.

Proof by reduction of CONTEXT-SENSITIVE LANGUAGE MEMBERSHIP (see Figure 2). Given a context-sensitive grammar and an input string of length m , we let the input surface form to the DFMSM RECOGNITION problem be the same as the input string. We then design a DFMSM with an alphabet equal to $\Sigma \cup \{\$, !\}$, where Σ is the set of non-terminals plus the set of terminals. The DFMSM first copies each surface input symbol to the corresponding position in the underlying form, and then adds the pair $\$:0$, completing the task in a state s_0 .

Having copied the string onto the underlying side of the pair, the remainder of the recognized underlying form will consist of rewritings of the string for each rule application, and will be paired with surface nulls at the end of the input string. Each rewriting will be separated by a $\$$ symbol, and, as the string length changes, it will be padded by $!$ symbols. For each rule $\alpha \rightarrow \beta$, we add a cycle to the DFMSM, emanating from state s_0 , which first

writes j copies of the $!$ symbol to the underlying form, where $j = b - a$, $b = |\beta|$, and $a = |\alpha|$:

$$! \rightarrow 0 / \Sigma : \mathcal{L} (\mathcal{L} : \mathcal{L} \dots_m \mathcal{L} : \mathcal{L}) _$$

$j \geq 0$ since the rules are context-sensitive.

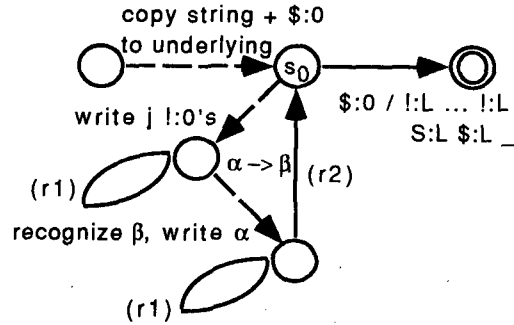


Figure 2. DFMSM constructed in Theorem 3.

The cycle then copies part of the most recent $\$$ -bounded string of symbols with a family of loops of the form:

$$\sigma \rightarrow 0 / \sigma : \mathcal{L} (\mathcal{L} : \mathcal{L} \dots_{m+j} \mathcal{L} : \mathcal{L}) _ \quad (r1)$$

for each $\sigma \in \Sigma$. It then recognizes β , and $:$ writes α , with:

$$\alpha_1 \rightarrow 0 / (\beta_1 : \mathcal{L} \dots_b \beta_b : \mathcal{L}) \\ (\mathcal{L} : \mathcal{L} \dots_{m+j+1-b} \mathcal{L} : \mathcal{L}) _$$

followed by:

$$\alpha_2 \rightarrow 0 / _ \\ \vdots \\ \alpha_a \rightarrow 0 / _$$

It then copies the rest of the most recent $\$$ -bounded string, using copy of the family of loops in (r1), and then adds a new $\$$ with a rule that also ensures that this second loop has iterated the appropriate number of times by checking that the length has been preserved:

$$\$ \rightarrow 0 / \$: \mathcal{L} (\mathcal{L} : \mathcal{L} \dots_m \mathcal{L} : \mathcal{L}) _ \quad (r2)$$

The DFMSM also has a loop emanating from s_0 which adds more $!$ symbols:

$$! \rightarrow 0 / ! : \mathcal{L} (\mathcal{L} : \mathcal{L} \dots_m \mathcal{L} : \mathcal{L}) _$$

All of the rule-cycles will use this to copy previously-added $!$ symbols, as the string shrinks in size. The proper application of this loop is also ensured by the length-checking of (r2).

Finally, we add one arc to the DFMSM from s_0 to the only final state which checks that the final copy of the string contains only the distinguished symbol, S :

$$\$ \rightarrow 0 / (! : \mathcal{L} \dots_{m-1} ! : \mathcal{L}) S : \mathcal{L} \$: \mathcal{L} _$$

⁸It is quite unlikely, however, since the reduction can probably be made with a different PSPACE-complete problem, from which the NP-completeness of FIXED EDFSM GENERATION would follow as a corollary.

Thus, the DFSM recognises the surface form if and only if there is a series of rewritings from the input string to S using the rules of the grammar, and the translation is linear in the size of the input string times the number of rules. \square

Since there exist fixed context-sensitive grammars for which the acceptance problem is NP-hard⁹, the NP-hardness of FIXED DFSM RECOGNITION with nulls follows as a corollary.

CONCLUSION

We claimed that DFSM's provide an approach to rules that is likely to seem more natural and intuitive to phonologists. Bridging the gap between linguistically adequate formalisms and computationally useful formalisms is a long-term, community effort, and we feel that it would be premature to make claims about the linguistic adequacy of the approach; this depends on whether two-level approaches can be developed and deployed in a way that will satisfy the theoretical and explanatory needs of linguists. A specific claim on which our formalism depends is that all natural two-level phonologies can be reproduced using DFSM's with finitely encodable rules. We feel that this claim is plausible, but it needs to be tested in practice.

Computationally, our complexity work so far on DFSM's does not preclude the possibility that compilers for generation and recognition (without nulls) exist which will allow for polynomial-time behavior at run-time. Although this question must eventually be resolved, we feel that any implementation is likely to be simpler than that required for KIMMO, and that even a direct implementation of DFSM's can prove adequate in many circumstances. We have not constructed an implementation as yet.

Like other two-level approaches, we have a problem with surface nulls. It is possible in most realistic recognition applications to bound the number of nulls by some function on the length of the overt input; and it remains to be seen whether a reasonable bound could sufficiently improve complexity in these cases.

We have dealt with the problem of underlying nulls by simply ruling them out. This simplifies the formal situation considerably, but we do not believe that it is acceptable as a general solution; for instance, we can't expect all cases of epenthesis to occur at morpheme boundaries. If underlying nulls are allowed, though, we will somehow need to limit the places where underlying nulls can occur; this is another good reason to pay attention to a phonotactic level of analysis.

⁹Garey and Johnson, (1979), p. 271.

ACKNOWLEDGEMENTS

This material is based upon work supported under a National Science Foundation Graduate Research Fellowship. This work was funded by National Science Foundation grant IRI-9003165. We thank the anonymous referees for helpful comments.

REFERENCES

- Evan Antworth. 1990. *PC-KIMMO: a two-level processor for morphological analysis*. Dallas, Texas: Summer Institute of Linguistics.
- Edward Barton, Robert Berwick, and Eric Ristad. 1987. *Computational Complexity and Natural Language*. Cambridge, Massachusetts: MIT Press.
- Gerhard Brewka. 1993. *Adding priorities and specificity to default logic*. DMG Technical Report. Sankt Augustin, Germany: Gesellschaft für Mathematik und Datenverarbeitung.
- Mary Dalrymple, Ronald Kaplan, Lauri Karttunen, Kimmo Koskenniemi, Sami Shaio, and Michael Wescoat. 1987. *Tools for Morphological Analysis*. Stanford, California, 1987: CSLI Technical Report CSLI-87-108.
- Michael Garey and David Johnson. 1979. *Computers and Intractability: A Guide to the Theory of NP-completeness*. New York, New York: Freeman and Co.
- Lauri Karttunen. 1991. "Finite-state constraints." *International conference on current issues in computational linguistics*. Penang, Malaysia.
- Lauri Karttunen and Kent Wittenburg. 1983. "A two-level morphological analysis of English." *Texas Linguistic Forum* 22 pp. 217-228.
- Graeme Ritchie, Graham Russell, Alan Black and Stephen Pulman. 1992. *Computational morphology*. Cambridge, Massachusetts: MIT Press.
- Richard Sproat. 1992. *Morphology and computation*. Cambridge, Massachusetts: MIT Press.