

# A Note on Sequential Rule-Based POS Tagging

Sylvain Schmitz

LSV, ENS Cachan & CNRS, Cachan, France  
schmitz@lsv.ens-cachan.fr

## Abstract

Brill’s part-of-speech tagger is defined through a cascade of leftmost rewrite rules. We revisit the compilation of such rules into a single sequential transducer given by Roche and Schabes (*Comput. Ling.* 1995) and provide a direct construction of the minimal sequential transducer for each individual rule.

**Keywords.** Brill Tagger; Sequential Transducer; POS Tagging

## 1 Introduction

Part-of-speech (POS) tagging consists in assigning the appropriate POS tag to a word in the context of its sentence. The program that performs this task, the *POS tagger*, can be learned from an annotated corpus in case of supervised learning, typically using hidden Markov model-based or rule-based techniques. The most famous rule-based POS tagging technique is due to Brill (1992). He introduced a three-parts technique comprising:

1. a *lexical tagger*, which associates a unique POS tag to each word from an annotated training corpus. This lexical tagger simply associates to each known word its most probable tag according to the training corpus annotation, i.e. a unigram maximum likelihood estimation;
2. an *unknown word tagger*, which attempts to tag unknown words based on suffix or capitalization features. It works like the contextual tagger, using the presence of a capital letter and bounded sized suffixes in its rules: for instance in English, a *-able* suffix usually denotes an adjective;
3. a *contextual tagger*, on which we focus in this paper. It consists of a cascade of string rewrite rules, called *contextual rules*, which correct tag assignments based on some surrounding contexts.

In this note, we revisit the proof that contextual rules can be translated into sequential transducers<sup>1</sup> proposed by Roche and Schabes (1995): whereas Roche and Schabes give a separate proof of sequentiality and exercise it to show that their constructed non-sequential transducer can be determinized (at the

---

<sup>1</sup>Historically, what we call here “sequential” used to be called “subsequential” (Schützenberger, 1977), but we follow the more recent practice initiated by Sakarovitch (2009).

expense of a worst-case exponential blow-up), we give a direct translation of a contextual rule into the minimal normalized sequential transducer, by adapting Simon (1994)’s string matching automaton to the transducer case. Our resulting sequential transducers are of linear size (before their composition). A similar construction can be found in (Mihov and Schultz, 2007), but no claim of minimality is made there.

## 2 Contextual Rules

### 2.1 Example

We start with an example by Roche and Schabes (1995): Let us suppose the following sentences were tagged by the lexical tagger (using the Penn Treebank tagset):

\*Chapman/NNP killed/VBN John/NNP Lennon/NNP  
 \*John/NNP Lennon/NNP was/VBD shot/VBD by/IN Chapman/NNP  
 He/PRP witnessed/VBD Lennon/NNP killed/VBN by/IN Chapman/NNP

with mistakes in the first two sentences: *killed* should be tagged as a past tense form “VBD”, and *shot* as a past participle form “VBN”.

The contextual tagger learns contextual rules over some tagset  $\Sigma$  of form  $uav \rightarrow ubv$  (or  $a \rightarrow b/u\_v$  using phonological rule notations (Kaplan and Kay, 1994)), meaning that the tag  $a$  rewrites to  $b$  in the context of  $u\_v$ , where the context is of length  $|uv|$  bounded by some fixed  $k$ ; in practice,  $k = 2$  or  $k = 3$  (Brill (1992) and Roche and Schabes (1995) use slightly different *templates* than the one parametrized by  $k$  we present here). For instance, a first contextual rule could be

$$\text{npv vbn} \rightarrow \text{npv vbd}$$

resulting in a new tagging

Chapman/NNP killed/VBD John/NNP Lennon/NNP  
 \*John/NNP Lennon/NNP was/VBD shot/VBD by/IN Chapman/NNP  
 \*He/PRP witnessed/VBD Lennon/NNP killed/VBD by/IN Chapman/NNP

A second contextual rule could be

$$\text{vbd in} \rightarrow \text{vbn in}$$

resulting in the correct tagging

Chapman/NNP killed/VBD John/NNP Lennon/NNP  
 John/NNP Lennon/NNP was/VBD shot/VBN by/IN Chapman/NNP  
 He/PRP witnessed/VBD Lennon/NNP killed/VBN by/IN Chapman/NNP

As stated before, our goal is to compile the entire sequence of contextual rules learned from a corpus into a single sequential function.

### 2.2 Cascade of Contextual Rules

Let us first formalize the semantics of Brill’s contextual rules. Let  $\mathcal{C} = r_1 r_2 \cdots r_n$  be a finite sequence of string rewrite rules in  $\Sigma^* \times \Sigma^*$  with  $\Sigma$  a POS tagset of *fixed size*. In practice the rules constructed in Brill’s contextual tagger are *length-preserving* and *1-change-bounded*, i.e. they modify a single letter, but this is not

a useful consideration for our transducer construction. Each rule  $r_i = u_i \rightarrow v_i$  defines a *leftmost rewrite relation*  $\xrightarrow[\text{lm}]{r_i}$  defined by

$$w \xrightarrow[\text{lm}]{r_i} w' \text{ iff } \exists x, y \in \Sigma^*, w = xu_iy \wedge w' = xv_iy \wedge (\forall z, z' \in \Sigma^*, w \neq zu_i z' \vee x \leq_{\text{pref}} z)$$

where  $x \leq_{\text{pref}} z$  denotes that  $x$  is a prefix of  $z$ . Note that the domain of  $\xrightarrow[\text{lm}]{r_i}$  is  $\Sigma^* \cdot u_i \cdot \Sigma^*$ . The *behavior* of a single rule is then the relation  $\llbracket r_i \rrbracket$  included in  $\Sigma^* \times \Sigma^*$  defined by

$$\llbracket r_i \rrbracket \stackrel{\text{def}}{=} \xrightarrow[\text{lm}]{r_i} \cup \text{Id}_{\Sigma^* \setminus (\Sigma^* \cdot u_i \cdot \Sigma^*)},$$

i.e. it applies  $\xrightarrow[\text{lm}]{r_i}$  on  $\Sigma^* \cdot u_i \cdot \Sigma^*$  and the identity on its complement  $\Sigma^* \setminus (\Sigma^* \cdot u_i \cdot \Sigma^*)$ .

The behavior of  $\mathcal{C}$  is then the composition

$$\llbracket \mathcal{C} \rrbracket \stackrel{\text{def}}{=} \llbracket r_1 \rrbracket \mathbin{\text{\$}} \llbracket r_2 \rrbracket \mathbin{\text{\$}} \cdots \mathbin{\text{\$}} \llbracket r_n \rrbracket.$$

Note that this behavior does *not* employ the transitive closure of the rewriting rules.

A naive implementation of  $\mathcal{C}$  would try to match each  $u_i$  at every position of the input string  $w$  in  $\Sigma^*$ , resulting in an overall complexity of  $O(|w| \cdot \sum_i |u_i|)$ . One often faces the problem of tagging a *set* of sentences  $\{w_1, \dots, w_m\}$ , which yields  $O((\sum_i |u_i|) \cdot (\sum_j |w_j|))$ . As shown in Roche and Schabes' experiments, compiling  $\mathcal{C}$  into a single sequential transducer  $\mathcal{T}$  results in practice in huge savings, with overall complexities in  $O(|w| + |\mathcal{T}|)$  and  $O(|\mathcal{T}| + \sum_j |w_j|)$  respectively.

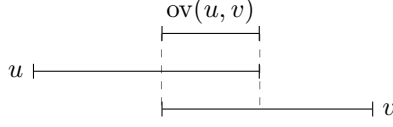
Each  $\llbracket r_i \rrbracket$  is a rational function, being the union of two rational functions over disjoint domains. Let  $|r_i|$  be the length  $|u_i v_i| \leq k$ . Roche and Schabes (1995, Section 8.2) provide a construction of an exponential-sized transducer  $\mathcal{T}_{r_i}$  for each  $\llbracket r_i \rrbracket$ , and compute their composition  $\mathcal{T}_{\mathcal{C}}$  of size  $|\mathcal{T}_{\mathcal{C}}| = O(\prod_{i=1}^n 2^{|r_i|})$ . As they show that each  $\llbracket r_i \rrbracket$  is actually a sequential function, their composition  $\llbracket \mathcal{C} \rrbracket$  is also sequential, and  $\mathcal{T}_{\mathcal{C}}$  can be determinized to yield a sequential transducer  $\mathcal{T}$  of size doubly exponential in  $\sum_{i=1}^n |r_i| \leq nk$  (see Roche and Schabes, 1995, Section 9.3). By contrast, our construction directly yields linear-sized minimal sequential transducers for each  $\llbracket r_i \rrbracket$ , resulting in a final sequential transducer of size  $O(\prod_{i=1}^n |r_i|) = O(2^{n \log k})$ .

### 3 Sequential Transducer of a Rule

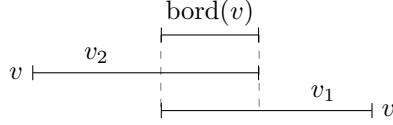
Intuitively, the sequential transducer for  $\llbracket r_i \rrbracket$  is related to the *string matching automaton* (Simon, 1994; Crochemore and Hancart, 1997) for  $u_i$ , i.e. the automaton for the language  $\Sigma^* u_i$ . This insight yields a *direct* construction of the minimal sequential transducer of a contextual rule, with at most  $|u_i| + 1$  states. Let us recall a few definitions:

#### 3.1 Preliminaries

**Overlaps, Borders** (see e.g. Crochemore and Hancart, 1997, Section 6.2). The *overlap*  $\text{ov}(u, v)$  of two words  $u$  and  $v$  is the longest suffix of  $u$  which is simultaneously a prefix of  $v$ .



A word  $u$  is a *border* of a word  $v$  if it is both a prefix and a suffix of  $v$ , i.e. if there exist  $v_1, v_2$  in  $\Sigma^*$  such that  $v = uv_1 = v_2u$ . For  $v \neq \varepsilon$ , the longest border of  $v$  different from  $v$  itself is denoted  $\text{bord}(v)$ .



**Fact 3.1.** For all  $u, v$  in  $\Sigma^*$  and  $a$  in  $\Sigma$ ,

$$\text{ov}(ua, v) = \begin{cases} \text{ov}(u, v) \cdot a & \text{if } \text{ov}(u, v) \cdot a \leq_{\text{pref}} v \\ \text{bord}(\text{ov}(u, v) \cdot a) & \text{otherwise.} \end{cases}$$

**Sequential Transducers** (see e.g. Sakarovitch, 2009, Section V.1.2). Formally, a sequential transducer from  $\Sigma$  to  $\Delta$  is a tuple  $\mathcal{T} = \langle Q, \Sigma, \Delta, q_0, \delta, \eta, \iota, \rho \rangle$  where  $\delta : Q \times \Sigma \rightarrow Q$  is a partial transition function,  $\eta : Q \times \Sigma \rightarrow \Delta^*$  a partial transition output function with the same domain as  $\delta$ , i.e.  $\text{dom}(\delta) = \text{dom}(\eta)$ ,  $\iota \in \Delta^*$  is an initial output, and  $\rho : Q \rightarrow \Delta^*$  is a partial final output function.  $\mathcal{T}$  defines a partial *sequential function*  $\llbracket \mathcal{T} \rrbracket : \Sigma^* \rightarrow \Delta^*$  with

$$\llbracket \mathcal{T} \rrbracket(w) \stackrel{\text{def}}{=} \iota \cdot \eta(q_0, w) \cdot \rho(\delta(q_0, w))$$

for all  $w$  in  $\Sigma^*$  for which  $\delta(q_0, w)$  and  $\rho(\delta(q_0, w))$  are defined, where  $\eta(q, \varepsilon) = \varepsilon$  and  $\eta(q, wa) = \eta(q, w) \cdot \eta(\delta(q, w), a)$  for all  $w$  in  $\Sigma^*$  and  $a$  in  $\Sigma$ .

Let us note  $\mathcal{T}_{(q)}$  for the sequential transducer with  $q$  for initial state. We write  $u \wedge v$  for the longest common prefix of strings  $u$  and  $v$ ; the longest common prefix of all the outputs from state  $q$  can be written formally as  $\bigwedge_{v \in \Sigma^*} \llbracket \mathcal{T}_{(q)} \rrbracket(v)$ . A sequential transducer is *normalized* if this value is  $\varepsilon$  for all  $q \in Q$  such that  $\text{dom}(\llbracket \mathcal{T}_{(q)} \rrbracket) \neq \emptyset$ ; any sequential transducer can be normalized.

The *translation* of a sequential function  $f$  by a word  $w$  in  $\Sigma^*$  is the sequential function  $w^{-1}f$  with

$$\text{dom}(w^{-1}f) \stackrel{\text{def}}{=} w^{-1}\text{dom}(f) \quad w^{-1}f(u) \stackrel{\text{def}}{=} \left( \bigwedge_{v \in \Sigma^*} f(wv) \right)^{-1} \cdot f(wu)$$

for all  $u$  in  $\text{dom}(w^{-1}f)$ . As in the finite automata case where minimal automata are isomorphic with residual automata, the minimal sequential transducer for a sequential function  $f$  is defined as the *translation transducer*  $\langle Q, \Sigma, \Delta, q_0, \delta, \eta, \iota, \rho \rangle$ , where

- $Q \stackrel{\text{def}}{=} \{w^{-1}f \mid w \in \Sigma^*\}$  (which is finite),
- $q_0 \stackrel{\text{def}}{=} \varepsilon^{-1}f$ ,

- $\iota \stackrel{\text{def}}{=} \bigwedge_{v \in \Sigma^*} f(v)$  if  $\text{dom}(f) \neq \emptyset$  and  $\iota = \varepsilon$  otherwise,
- $\delta(w^{-1}f, a) \stackrel{\text{def}}{=} (wa)^{-1}f$ ,
- $\eta(w^{-1}f, a) \stackrel{\text{def}}{=} \bigwedge_{v \in \Sigma^*} (w^{-1}f)(av)$  if  $\text{dom}((wa)^{-1}f) \neq \emptyset$  and  $\eta(w^{-1}f, a) = \varepsilon$  otherwise, and
- $\rho(w^{-1}f) \stackrel{\text{def}}{=} (w^{-1}f)(\varepsilon)$  if  $\varepsilon \in \text{dom}(w^{-1}f)$ , and is otherwise undefined.

### 3.2 Main Construction

Here is the definition of our transducer for a contextual rule:

**Definition 3.2** (Transducer of a Contextual Rule). The sequential transducer  $\mathcal{T}_r$  associated with a contextual rule  $r = u \rightarrow v$  with  $u \neq \varepsilon$  is defined as  $\mathcal{T}_r \stackrel{\text{def}}{=} \langle \text{Pref}(u), \Sigma, \Sigma, \varepsilon, \delta, \eta, \varepsilon, \rho \rangle$  with the set of prefixes of  $u$  as state set,  $\varepsilon$  as initial state and initial output, and for all  $a$  in  $\Sigma$  and  $w$  in  $\text{Pref}(u)$ ,

$$\delta(w, a) \stackrel{\text{def}}{=} \begin{cases} wa & \text{if } wa \leq_{\text{pref}} u \\ w & \text{if } w = u \\ \text{bord}(wa) & \text{otherwise} \end{cases}$$

$$\rho(w) \stackrel{\text{def}}{=} \begin{cases} \varepsilon & \text{if } w \leq_{\text{pref}} (u \wedge v) \\ (u \wedge v)^{-1} \cdot w & \text{if } (u \wedge v) <_{\text{pref}} w <_{\text{pref}} u \\ \varepsilon & \text{otherwise, i.e. if } w = u \end{cases}$$

$$\eta(w, a) \stackrel{\text{def}}{=} \begin{cases} a & \text{if } wa \leq_{\text{pref}} (u \wedge v) \\ \varepsilon & \text{if } (u \wedge v) <_{\text{pref}} wa <_{\text{pref}} u \\ (u \wedge v)^{-1} \cdot v & \text{if } wa = u \\ a & \text{if } w = u \\ \rho(w)a \cdot \rho(\text{bord}(wa))^{-1} & \text{otherwise.} \end{cases}$$

For instance, the sequential transducer for the rule  $ababb \rightarrow abbbb$  is shown in Figure 1 (one can check that  $ababb \wedge abbbb = ab$ ,  $\text{bord}(b) = \varepsilon$ ,  $\text{bord}(aa) = a$ ,  $\text{bord}(abb) = \varepsilon$ ,  $\text{bord}(abaa) = a$ , and  $\text{bord}(ababa) = aba$ ). The intuition behind the definition of  $\eta(w, a)$  is to decompose the rewriting according to  $u \rightarrow v$  into four phases:

1. while in the common prefix  $u \wedge v$  of  $u$  and  $v$ , implement the identity function (states  $\varepsilon$ ,  $a$ , and  $ab$  in Figure 1),
2. as soon as we start reading a symbol of  $u$  that does not match that of  $v$  (upon reading  $a$  in state  $ab$  in Figure 1), we stop outputting symbols and wait for the whole of  $u$  to be read,
3. if  $u$  has been read, we output the remaining rewritten string  $(u \wedge v)^{-1} \cdot v$  we had been saving (upon reading  $b$  in state  $abab$  in Figure 1),
4. after having read the first occurrence of  $u$  in full, we merely implement the identity again (state  $ababb$  in Figure 1).

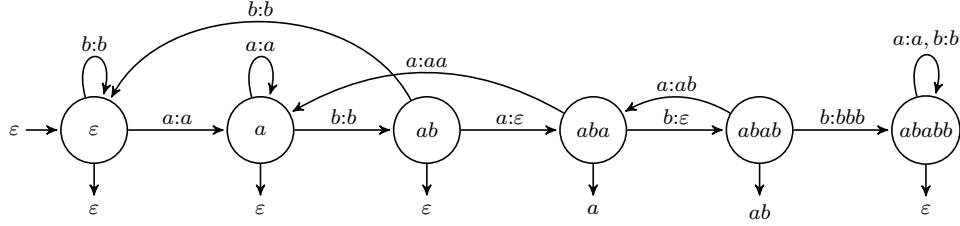


Figure 1: The sequential transducer constructed for  $ababb \rightarrow abbbb$ .

5. If on the other hand we realize that  $u$  cannot be read after all in some state  $w$  upon reading some  $a$  (e.g. transition on  $a$  in state  $aba$  in Figure 1), we need to flush the missing output  $(u \wedge v)^{-1} \cdot w = \rho(w)$ , minus the saved output if the state we reach is itself in phase 2.

It remains to show that this sequential transducer is indeed the minimal normalized sequential transducer for  $\llbracket r \rrbracket$ .

**Proposition 3.3** (Correctness). *Let  $r = u \rightarrow v$  with  $u \neq \varepsilon$ . Then  $\llbracket \mathcal{T}_r \rrbracket = \llbracket r \rrbracket$ .*

*Proof.* Let us first consider the case of input words in  $\Sigma^* \setminus (\Sigma^* \cdot u \cdot \Sigma^*)$ :

*Claim 3.3.1.* For all  $w$  in  $\Sigma^* \setminus (\Sigma^* \cdot u \cdot \Sigma^*)$ ,

$$\delta(\varepsilon, w) = \text{ov}(w, u) \quad \eta(\varepsilon, w) = w \cdot \rho(\text{ov}(w, u))^{-1}.$$

*Proof of the claim.* By induction on  $w$ : since  $u \neq \varepsilon$ , the base case is  $w = \varepsilon$  with

$$\delta(\varepsilon, \varepsilon) = \varepsilon = \text{ov}(\varepsilon, u) \quad \eta(\varepsilon, \varepsilon) = \varepsilon = \varepsilon \cdot \varepsilon^{-1} = \varepsilon \cdot \rho(\varepsilon)^{-1}.$$

For the induction step, we consider  $wa$  in  $\Sigma^* \setminus (\Sigma^* \cdot u \cdot \Sigma^*)$  for some  $w$  in  $\Sigma^*$  and  $a$  in  $\Sigma$ :

$$\begin{aligned} \delta(\varepsilon, wa) &= \delta(\delta(\varepsilon, w), a) && \text{(by def.)} \\ &= \delta(\text{ov}(w, u), a) && \text{(by ind. hyp.)} \end{aligned}$$

where by definition of  $\delta$ , we have  $\delta(\text{ov}(w, u), a) = \text{ov}(w, u) \cdot a$  if  $\text{ov}(w, u) \cdot a \leq_{\text{pref}} u$  and  $\delta(\text{ov}(w, u), a) = \text{bord}(\text{ov}(w, u) \cdot a)$  otherwise (the case  $\text{ov}(w, u) = u$  is impossible since  $w$  in  $\Sigma^* \setminus (\Sigma^* \cdot u \cdot \Sigma^*)$ ). In all cases:

$$\begin{aligned} \delta(\varepsilon, wa) &= \text{ov}(wa, u) && \text{(by Fact 3.1)} \\ \eta(\varepsilon, wa) &= \eta(\varepsilon, w) \cdot \eta(\delta(\varepsilon, w), a) && \text{(by def.)} \\ &= w \cdot \rho(\delta(\varepsilon, w))^{-1} \cdot \eta(\delta(\varepsilon, w), a) && \text{(by ind. hyp.)} \\ &= w \cdot \rho(w')^{-1} \cdot \eta(w', a); && \text{(by setting } w' = \delta(\varepsilon, w)) \end{aligned}$$

we need to do a case analysis for this last equation:

**Case  $w'a \not\leq_{\text{pref}} u$**  Then  $\eta(w', a) = \rho(w') \cdot a \cdot \rho(\text{border}(w'a))^{-1}$ , which yields

$$\begin{aligned} \eta(\varepsilon, wa) &= w \cdot \rho(w')^{-1} \cdot \rho(w') \cdot a \cdot \rho(\delta(\varepsilon, wa))^{-1} && \text{(by Fact 3.1)} \\ &= wa \cdot \rho(\delta(\varepsilon, wa))^{-1}. \end{aligned}$$

**Case  $w'a <_{\text{pref}} u$**  Then  $\delta(\varepsilon, wa) = w'a$ , and we need to further distinguish between several cases:

$w'a \leq_{\text{pref}} (\mathbf{u} \wedge \mathbf{v})$  then  $\rho(w') = \varepsilon$ ,  $\eta(w', a) = a$ , and  $\rho(w'a) = \varepsilon$ , thus

$$\eta(\varepsilon, wa) = wa = wa \cdot \varepsilon^{-1} = wa \cdot \rho(w'a)^{-1},$$

$w' = (\mathbf{u} \wedge \mathbf{v})$  then  $\rho(w') = \varepsilon$ ,  $\eta(w', a) = \varepsilon$ , and  $\rho(w'a) = (u \wedge v)^{-1} \cdot w'a = a$ ,

$$\eta(\varepsilon, wa) = w = wa \cdot a^{-1} = wa \cdot \rho(w'a)^{-1},$$

$(\mathbf{u} \wedge \mathbf{v}) <_{\text{pref}} w'$  then  $\rho(w') = (u \wedge v)^{-1} \cdot w'$ ,  $\eta(w', a) = \varepsilon$ , and  $\rho(w'a) = (u \wedge v)^{-1} \cdot w'a$ , thus

$$\begin{aligned} \eta(\varepsilon, wa) &= w \cdot ((u \wedge v)^{-1} \cdot w')^{-1} \\ &= wa \cdot a^{-1} \cdot ((u \wedge v)^{-1} \cdot w')^{-1} \\ &= wa \cdot \rho(w'a)^{-1}. \end{aligned} \quad \square$$

The claim yields that  $\llbracket \mathcal{T}_r \rrbracket$  coincides with  $\llbracket r \rrbracket$  on words in  $\Sigma^* \setminus (\Sigma^* \cdot u \cdot \Sigma^*)$ , i.e. is the identity over  $\Sigma^* \setminus (\Sigma^* \cdot u \cdot \Sigma^*)$ . Then, since  $u \neq \varepsilon$ , a word in  $\Sigma^* \cdot u \cdot \Sigma^*$  can be written as  $waw'$  with  $w$  in  $\Sigma^* \setminus (\Sigma^* \cdot u \cdot \Sigma^*)$ ,  $a$  in  $\Sigma$  with  $wa$  in  $\Sigma^* \cdot u$ , and  $w'$  in  $\Sigma^*$ . Let  $u = u'a$ ; the claim implies that

$$\delta(\varepsilon, w) = u' \qquad \eta(\varepsilon, w) = w \cdot \rho(u')^{-1}.$$

Thus, by definition of  $\mathcal{T}_r$ ,  $\delta(\varepsilon, wa) = u'a = u$  and thus

$$\eta(\varepsilon, wa) = \eta(\varepsilon, w) \cdot \eta(u', a) = w \cdot \rho(u')^{-1} \cdot (u \wedge v)^{-1} \cdot v;$$

**if**  $(\mathbf{u} \wedge \mathbf{v}) <_{\text{pref}} u'$

$$\eta(\varepsilon, wa) = w \cdot ((u \wedge v)^{-1} \cdot u')^{-1} \cdot (u \wedge v)^{-1} \cdot v = w \cdot u'^{-1} \cdot v = wa \cdot u^{-1} \cdot v;$$

**otherwise** i.e. if  $u' = (u \wedge v)$ :

$$\eta(\varepsilon, wa) = w \cdot u'^{-1} \cdot v = wa \cdot u^{-1} \cdot v.$$

Thus in all cases  $\llbracket \mathcal{T}_r \rrbracket(wa) = \llbracket r \rrbracket(wa)$ , and since  $\mathcal{T}_r$  starting in state  $u$  (i.e.  $\mathcal{T}_{r(u)}$ ) implements the identity over  $\Sigma^*$ , we have more generally  $\llbracket \mathcal{T}_r \rrbracket = \llbracket r \rrbracket$ .  $\square$

**Lemma 3.4** (Normality). *Let  $r = u \rightarrow v$ . Then  $\mathcal{T}_r$  is normalized.*

*Proof.* Let  $w \in \text{Prefix}(u)$  be a state of  $\mathcal{T}_r$ ; let us show that  $\bigwedge \llbracket \mathcal{T}_{r(w)} \rrbracket(\Sigma^*) = \varepsilon$ .

**If**  $(\mathbf{u} \wedge \mathbf{v}) <_{\text{pref}} w <_{\text{pref}} \mathbf{u}$  let  $u' = w^{-1}u \in \Sigma^+$ , and consider the two outputs

$$\begin{aligned} \llbracket \mathcal{T}_{r(w)} \rrbracket(u') &= \eta(w, u')\rho(u) = (u \wedge v)^{-1}v \\ \llbracket \mathcal{T}_{r(w)} \rrbracket(\varepsilon) &= \rho(w) = (u \wedge v)^{-1}w. \end{aligned}$$

Since  $(u \wedge v) <_{\text{pref}} u$  we can write  $u$  as  $(u \wedge v)au''u'$ , and either  $v = (u \wedge v)bv'$  or  $v = u \wedge v$ , for some  $a \neq b$  in  $\Sigma$  and  $u'', v'$  in  $\Sigma^*$ ; this yields  $w = (u \wedge v)au''$  and thus  $\llbracket \mathcal{T}_{r(w)} \rrbracket(u') \wedge \llbracket \mathcal{T}_{r(w)} \rrbracket(\varepsilon) = \varepsilon$ .

**otherwise**  $\rho(w) = \varepsilon$ , which yields the lemma.  $\square$

**Proposition 3.5** (Minimality). *Let  $r = u \rightarrow v$  with  $u \neq \varepsilon$  and  $u \neq v$ . Then  $\mathcal{T}_r$  is the minimal sequential transducer for  $\llbracket r \rrbracket$ .*

*Proof.* Let  $w <_{\text{pref}} w'$  be two different states in  $\text{Prefix}(u)$ ; we proceed to prove that  $\llbracket w^{-1}\mathcal{T}_r \rrbracket \neq \llbracket w'^{-1}\mathcal{T}_r \rrbracket$ , hence that no two states of  $\mathcal{T}_r$  can be merged. By Lemma 3.4 it suffices to prove that  $\llbracket \mathcal{T}_r(w) \rrbracket \neq \llbracket \mathcal{T}_r(w') \rrbracket$ , thus to exhibit some  $x \in \Sigma^*$  such that  $\llbracket \mathcal{T}_r(w) \rrbracket(x) \neq \llbracket \mathcal{T}_r(w') \rrbracket(x)$ . We perform a case analysis:

if  $w' \leq_{\text{pref}} (u \wedge v)$  then  $w <_{\text{pref}} (u \wedge v)$  thus  $\llbracket \mathcal{T}_r(w) \rrbracket(x) = x$  for all  $x \notin w^{-1} \cdot \Sigma^* \cdot u \cdot \Sigma^*$ ; consider

$$\llbracket \mathcal{T}_r(w) \rrbracket(w'^{-1}u) = w'^{-1}u \neq w'^{-1}v = \llbracket \mathcal{T}_r(w') \rrbracket(w'^{-1}u) ;$$

if  $w \leq_{\text{pref}} (u \wedge v)$  and  $w' = u$  then  $\llbracket \mathcal{T}_r(w') \rrbracket(x) = x$  for all  $x$  and we consider

$$\llbracket \mathcal{T}_r(w) \rrbracket(w^{-1}u) = w^{-1}v \neq w^{-1}v = \llbracket \mathcal{T}_r(w') \rrbracket(w^{-1}u) ;$$

otherwise that is if  $w \leq_{\text{pref}} (u \wedge v)$  and  $(u \wedge v) <_{\text{pref}} w' <_{\text{pref}} u$ , or  $(u \wedge v) <_{\text{pref}} w <_{\text{pref}} w' \leq_{\text{pref}} u$ , we have  $\rho(w) \neq \rho(w')$  thus

$$\llbracket \mathcal{T}_r(w) \rrbracket(\varepsilon) \neq \llbracket \mathcal{T}_r(w') \rrbracket(\varepsilon) . \quad \square$$

## 4 Conclusion

The results of the previous section yield (the cases  $u = \varepsilon$  and  $u = v$  are trivial):

**Theorem 4.1.** *Given a contextual rule  $r = u \rightarrow v$ , one can construct directly the minimal normalized sequential transducer  $\mathcal{T}_r$  of size  $O(|r|)$  for  $\llbracket r \rrbracket$ .*

The remaining open question is whether we can obtain good upper bounds on the size of the sequential transducer  $\mathcal{T}_{\mathcal{C}}$  for a cascade  $\mathcal{C} = r_1 \cdots r_n$ : string matching automata enjoy nice combinatorial properties, but it seems unlikely for better bounds than  $O(\prod_{i=1}^n |r_i|)$  to exist: indeed, the membership problem, i.e. given  $u, v$  in  $\Sigma^*$  and a cascade  $\mathcal{C}$  whether  $\llbracket \mathcal{C} \rrbracket(u) = v$ , is easily seen to be PTIME-complete for length-preserving rules, but a LOGSPACE construction of a sequential transducer for  $\llbracket \mathcal{C} \rrbracket$  would result in a LOGSPACE algorithm for the problem.

## References

- Brill, E., 1992. A simple rule-based part of speech tagger. In *ANLP '92, third Conference on Applied Natural Language Processing*, pages 152–155. ACL Press. doi:10.3115/974499.974526.
- Crochemore, M. and Hancart, C., 1997. Automata for matching patterns. In Rozenberg, G. and Salomaa, A., editors, *Handbook of Formal Languages*, volume 2. Linear Modeling: Background and Application, chapter 9, pages 399–462. Springer. ISBN 3-540-60648-3.



- Kaplan, R.M. and Kay, M., 1994. Regular models of phonological rule systems. *Computational Linguistics*, 20(3):331–378. <http://www.aclweb.org/anthology/J94-3001.pdf>.
- Mihov, S. and Schultz, K.U., 2007. Efficient dictionary-based text rewriting using subsequential transducers. 13(4):353–381. doi:10.1017/S1351324905004092.
- Roche, E. and Schabes, Y., 1995. Deterministic part-of-speech tagging with finite-state transducers. *Computational Linguistics*, 21(2):227–253. <http://www.aclweb.org/anthology/J95-2004.pdf>.
- Sakarovitch, J., 2009. *Elements of Automata Theory*. Cambridge University Press. ISBN 978-0-521-84425-3.
- Schützenberger, M.P., 1977. Sur une variante des fonctions séquentielles. *Theoretical Computer Science*, 4(1):47–57. doi:0.1016/0304-3975(77)90055-X.
- Simon, I., 1994. String matching algorithms and automata. In Karhumäki, J., Maurer, H., and Rozenberg, G., editors, *Results and Trends in Theoretical Computer Science: Colloquium in Honor of Arto Salomaa*, volume 812 of *Lecture Notes in Computer Science*, pages 386–395. Springer. ISBN 978-3-540-58131-4. doi:10.1007/3-540-58131-6.61.