# LING 306: Introduction to Computational Linguistics

Richard Sproat

rws@uiuc.edu

http://www.staff.uiuc.edu/ rws/

URL for this course: http://www.staff.uiuc.edu/ rws/Courses/L306

Lecture 7: Rule Compilation

September 8, 2003

# Homework

- Problem 3 (= J&M, 2.8)

- Problem 5 (= J&M, 2.10–2.11)

- Problem 6

# Context-dependent rewrite rules

Given a rule of the form

$$\phi \rightarrow \psi / \lambda \underline{\quad} \rho$$

where $\phi$, $\psi$, $\lambda$, $\rho$ are regular expressions, how do you compile this rule into a transducer?

# An Illustrative Example

u → i / i C* __

(u → i / Σ* i C* __ Σ*)

Input: kikukuku

# An Illustrative Example

u → i / i C* __

# kikukuku

# An Illustrative Example

u → i / i C* __

kik**u**kuku

# An Illustrative Example

u → i / i C* __

kik**i**kuku

# An Illustrative Example

u → i / i C* __

# kikikuku

# An Illustrative Example

u → i / i C* __

kikikiku

# An Illustrative Example

u → i / i C* __

# kikikiku

# An Illustrative Example

u → i / i C* __

kikikiki

# A Note on Directionality

- The output of one application *feeds* the next application across the string.

- But I took it for granted that we were applying left to right

# Right-to-Left Application

u → i / i C* __

# kikukuku

# Right-to-Left Application

u → i / i C* __

# kikukuk<span style="color:red">u</span>

# Right-to-Left Application

u → i / i C* __

kikuk<span style="color:red">u</span>k<span style="color:red">u</span>

# Right-to-Left Application

u → i / i C* __

kikukuku

# Right-to-Left Application

u → i / i C* __

kikikuku

# Simultaneous Application

u → i / i C* __

# kikukuku

# Simultaneous Application

u → i / i C* __

# kikukuku

# Simultaneous Application

u → i / i C* __

# kikikuku

# Bookkeeping

- Let's say we are being extra careful and we want to make sure that we're applying the rule correctly.

- To help us, we might use some auxiliary annotations, such as angle brackets, to mark contexts, as well as the actual string to be changed.

# Bookkeeping: First Attempt

u → i / Σ* i C* __ Σ*

1. Mark all possible right contexts:

# kikukupapu

# Bookkeeping: First Attempt

u → i / Σ* i C* __ Σ*

1. Mark all possible right contexts:

$$_>k_>i_>k_>u_>k_>u_>p_>a_>p_>u_>$$

# Bookkeeping: First Attempt

u → i / Σ* i C* __ Σ*

2. Mark all possible left contexts:

$$>k>i>k>u>k>u>p>a>p>u>$$

# Bookkeeping: First Attempt

u → i / Σ* i C* __ Σ*

2. Mark all possible left contexts:

$$_>k_>i_{<>}k_{<>}u_>k_>u_>p_>a_>p_>u_>$$

# Bookkeeping: First Attempt

u → i / Σ* i C* __ Σ*

3. Change any *u* into a *i* when it's delimited by $_<$ and $_>$.

$$_>k_>i_{<>}k_{<>}u_>k_>u_>p_>a_>p_>u_>$$

# Bookkeeping: First Attempt

u → i / Σ* i C* __ Σ*

3. Change any *u* into a *i* when it's delimited by $_<$ and $_>$.

$$_>k_>i_{<>}k_{<>}\textcolor{red}{i}_>k_>u_>p_>a_>p_>u_>$$

# Oops

This is wrong: we are not allowing for the left-to-right feeding. So how about the following:

- Mark the beginnings of right contexts—$\rho$. That much was correct.

- Mark the left edge of all $\phi$ that have $>$ on the righthand side, with *two* different markers, $<_1$, $<_2$:

  ★ $<_1$ is to be used to trigger application of the change of $\phi$ to $\psi$
  ★ $<_2$ is to be used for *non*application of the change

- Change $\phi$ into $\psi$ wherever $\phi$ is delimited by $<_1$ and $>$.

  We can now delete $>$: we don't need it any more.

- This will give you a set of possible strings, some with the change, some without.

  We are now going to filter those outcomes.

- **Now we check $\lambda$. In particular:**

  - $\star$ Allow only strings where $<_1$ is preceded by $\lambda$; $<_1$ can then be deleted.
  - $\star$ Allow only strings where $<_2$ is *not* preceded by $\lambda$; $<_2$ can then be deleted.

# Bookkeeping: Second Attempt

u → i / Σ* i C* __ Σ*

1. Mark all possible right contexts:

    k    i    k    u    k    u    p    a    p    u

# Bookkeeping: Second Attempt

u → i / Σ* i C* __ Σ*

## 1. Mark all possible right contexts:

> k > i > k > u > k > u > p > a > p > u >

# Bookkeeping: Second Attempt

$u \rightarrow i / \Sigma^* i C^* \underline{\quad} \Sigma^*$

2. Mark all $u$ followed by $>$ with $<_1$ and $<_2$

$>$ k $>$ i $>$ k $>$ u $>$ k $>$ u $>$ p $>$ a $>$ p $>$ u $>$

# Bookkeeping: Second Attempt

u → i / Σ* i C* __ Σ*

2. Mark all $u$ followed by $>$ with $<_1$ and $<_2$

| > | k | > | i | > | k | $<_{1,>}$ | u | > | k | $<_{1,>}$ | u | > | p | > | a | > | p | $<_{1,>}$ |
|   |   |   |   |   |   | $<_2$ |   |   |   | $<_2$ |   |   |   |   |   |   |   | $<_2$ |

# Bookkeeping: Second Attempt

u → i / Σ* i C* __ Σ*

3. Change all *u* between $<_1$ and $>$, and delete $>$

$>$   k   $>$   i   $>$   k   $<1,>$   u   $>$   k   $<1,>$   u   $>$   p   $>$   a   $>$   p   $<1,>$

                                   $<2$                           $<2$                                       $<2$

# Bookkeeping: Second Attempt

u → i / Σ* i C* __ Σ*

3. Change all $u$ between $<_1$ and $>$, and delete $>$

| k | i | k | $<_1$ | i | k | $<_1$ | i | p | a | p | $<_1$ | i |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | $<_2$ | u | | $<_2$ | u | | | | $<_2$ | u |

# Bookkeeping: Second Attempt

u → i / Σ* i C* __ Σ*

4. Allow only strings where $<_1$ is preceded by *i C\**, and delete $<_1$:

| k | i | k | $<_1$ | i | k | $<_1$ | i | p | a | p | $<_1$ | i |
|---|---|---|-------|---|---|-------|---|---|---|---|-------|---|
|   |   |   | $<_2$ | u |   | $<_2$ | u |   |   |   | $<_2$ | u |

# Bookkeeping: Second Attempt

u → i / $\Sigma$* i C* __ $\Sigma$*

4. Allow only strings where $<_1$ is preceded by *i C\**, and delete $<_1$:

| k | i | k | i | k | i | p | a | p | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | $<_2$ u | | $<_2$ u | | | | $<_2$ u | |

# Bookkeeping: Second Attempt

u → i / $\Sigma$* i C* __ $\Sigma$*

5. Allow only strings where $<_2$ is *not* preceded by *i C\**, and delete $<_2$:

| k | i | k | i | k | i | p | a | p | |
|---|---|---|---|---|---|---|---|---|---|
| | | $<_2$ u | | $<_2$ u | | | | $<_2$ u | |

# Bookkeeping: Second Attempt

u → i / $\Sigma$* i C* __ $\Sigma$*

5. Allow only strings where $<_2$ is *not* preceded by *i C\**, and delete $<_2$:

<table>
<tr><td>k</td><td>i</td><td>k</td><td style="color:red">i</td><td>k</td><td style="color:red">i</td><td>p</td><td>a</td><td>p</td></tr>
<tr><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td style="color:red">u</td></tr>
</table>

# Now in fact . . .

Each of these five steps can be represented as a transducer:

- A transducer $r$ that inserts $>$ before every $\rho$.

- A transducer $f$ that inserts $<_1$ and $<_2$ before every $\phi$ followed by $>$.

- A transducer $replace$ that replaces $\phi$ with $\psi$ between $<_1$ and $>$, and also deletes $>$.

- A transducer $l_1$ that filters out all $<_1$ not preceded by $\lambda$, and deletes any $<_1$

- A transducer $l_2$ that filters out all $<_2$ preceded by $\lambda$, and deletes any $<_2$
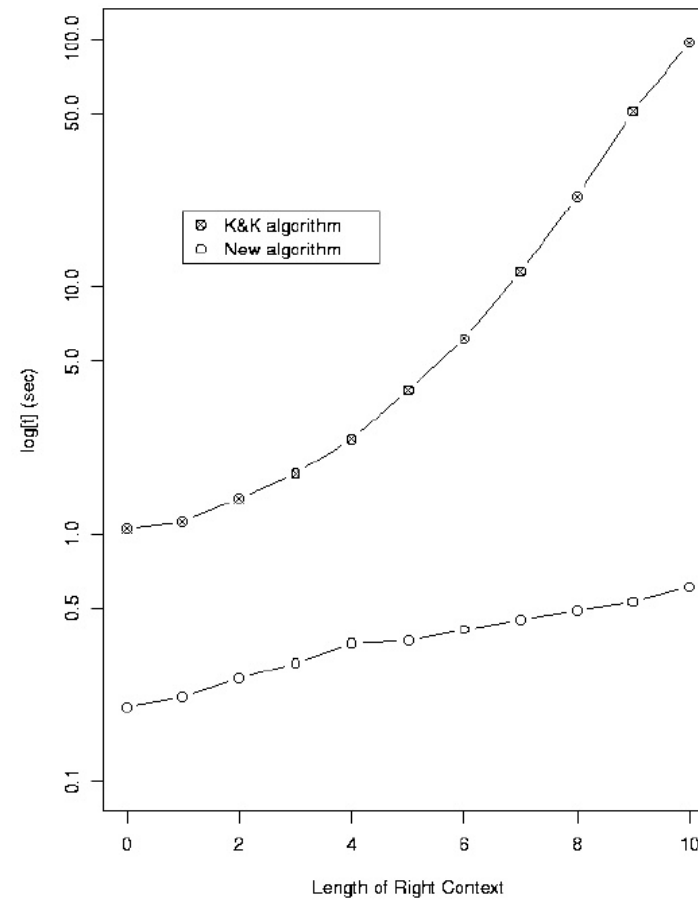
# So the Rule Transducer Is

$$r \circ f \circ replace \circ l_1 \circ l_2$$

# Literature

- The formulation I just sketched for left-to-right rules is from Mohri and Sproat (1996) "An Efficient Compiler for Weighted Rewrite Rules", *34th Annual Meeting of the Association for Computational Linguistics*, Santa Cruz, CA, available at http://acl.ldc.upenn.edu/P/P96/P96-1031.pdf

- It differs from the earlier Kaplan and Kay formulation in part because K&K introduce brackets everywhere, and subsequently filter, whereas M&S try to only introduce brackets where they are really needed.

- They are algebraically equivalent, but the M&S formulation is algorithmically better.

# Compilation Times

Rules of the form $a \rightarrow b/\underline{\quad} c^k$, $(k \in [0, 10])$

# Comparison in a Real Example

| Rules | KK | | | New | | |
|-------|------|--------|--------|------|--------|---------|
| | time (s) | space | | time (s) | space | |
| | | states | arcs | | states | arcs |
| <ö> | 62 | 412 | 50,475 | 47 | 394 | 47,491 |
| <a> | 284 | 1,939 | 215,721 | 240 | 1,927 | 213,408 |

- German text-to-speech system (B. Möbius)

  ⋆ Pronunciation rules for <ö> and <a>
  ⋆ ≈ 20 sets of rules (≈ 200 rules)

- Time for overall construction of each set of rules
  (SGI Indy 4000, 100 MHZ IP20, 128 Mbytes RAM)

# Reading

Nada: Have a Nice Weekend.