# CMPT 413
# Computational Linguistics

Anoop Sarkar

`http://www.cs.sfu.ca/~anoop`

# Finite-state transducers

- Many applications in computational linguistics

- Popular applications of FSTs are in:
  - Orthography
  - Morphology
  - Phonology

- Other applications include:
  - Grapheme to phoneme
  - Text normalization
  - Transliteration
  - Edit distance
  - Word segmentation
  - Tokenization
  - Parsing

# Orthography and Phonology

- Orthography: written form of the language (affected by morpheme combinations)

  move + ed $\rightarrow$ moved

  swim + ing $\rightarrow$ swimming <u>S W IH1 M IH0 NG</u>

- Phonology: change in pronunciation due to morpheme combinations (changes may not be confined to morpheme boundary)

  intent <u>IH2 N T EH1 N T</u> + ion

  $\rightarrow$ intention <u>IH2 N T EH1 N CH AH0 N</u>

# Orthography and Phonology

- Phonological alternations are not reflected in the spelling (orthography):
  - Newton  Newtonian
  - maniac  maniacal
  - electric  electricity

- Orthography can introduce changes that do not have any counterpart in phonology:
  - picnic    picnicking
  - happy    happiest
  - gooey    gooiest

# Segmentation and Orthography

- To find entries in the lexicon we need to segment any input into morphemes
- Looks like an easy task in some cases:

  *looking* → look + ing

  *rethink* → re + think

- However, just matching an affix does not work:

  \**thing* → th + ing

  \**read* → re + ad

- We need to store valid stems in our lexicon

  what is the stem in *assassination* (*assassin* and not *nation*)

# Porter Stemmer

- A simpler task compared to segmentation is simply stripping out all affixes (a process called **stemming**, or finding the stem)
- Stemming is usually done without reference to a lexicon of valid stems
- The Porter stemming algorithm is a simple composition of FSTs, each of which strips out some affix from the input string
  - input=..*ational*, produces output=..*ate* (*relational* → *relate*)
  - input=..V..*ing*, produces output=ε (*motoring* → *motor*)

# Porter Stemmer

- False positives (stemmer gives incorrect stem): *doing* → *doe*, *policy* → *police*

- False negatives (should provide stem but does not): *European* → *Europe*, *matrices* → *matrix*

*I'm a rageaholic. I can't live without rageahol.*

    Homer Simpson, from *The Simpsons*

- Despite being linguistically unmotivated, the Porter stemmer is used widely due to its simplicity (easy to implement) and speed

# Segmentation and orthography

- More complex cases involve alterations in spelling
  - *foxes* → fox + s     [ ***e*-insertion** ]
  - *loved* → love + ed [ ***e*-deletion** ]
  - *flies* → fly + s      [ ***i* to *y*, *e*-deletion** ]
  - *panicked* → panic + ed [ **k-insertion** ]
  - *chugging* → chug + ing [ **consonant doubling** ]
  - **singging* → sing + ing
  - *impossible* → in + possible [ ***n* to *m*** ]
- Called *morphographemic* changes.
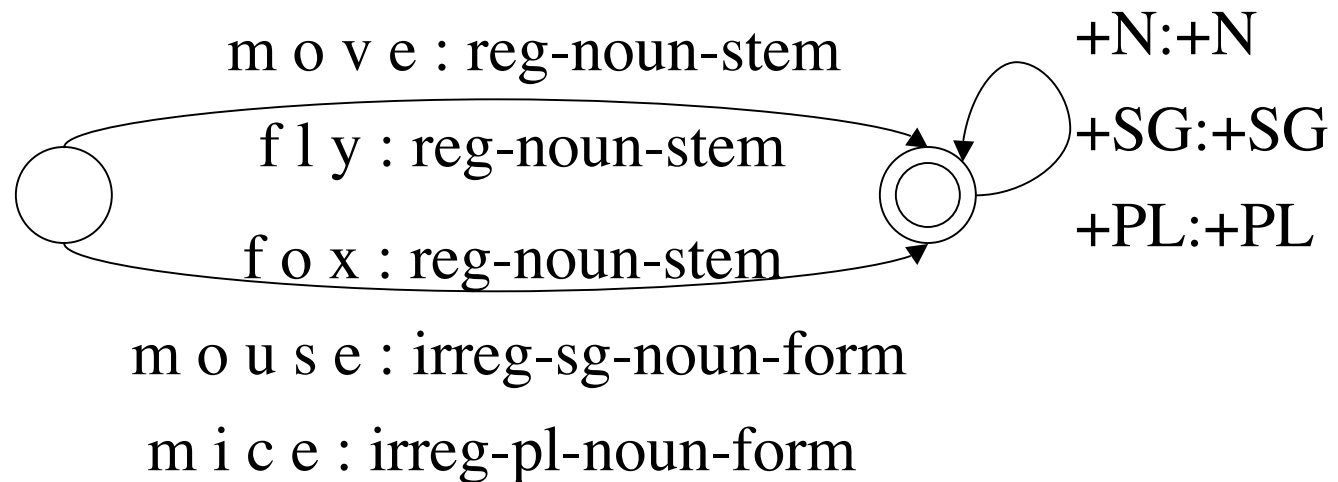- Similar to but not identical to changes in pronunciation due to morpheme combinations

# Morphological Parsing with FSTs

- Think of the process of decomposing a word into its component morphemes in the reverse direction: as *generation* of the word from the component morphemes
- Start with an abstract notion of each morpheme being simply combined with the stem using concatenation
  - Each stem is written with its part of speech, e.g. cat+N
  - Concatenate each stem with some suffix information, e.g. cat+N+PL
  - e.g. cat+N+PL goes through an FST to become *cats* (also works in reverse!)
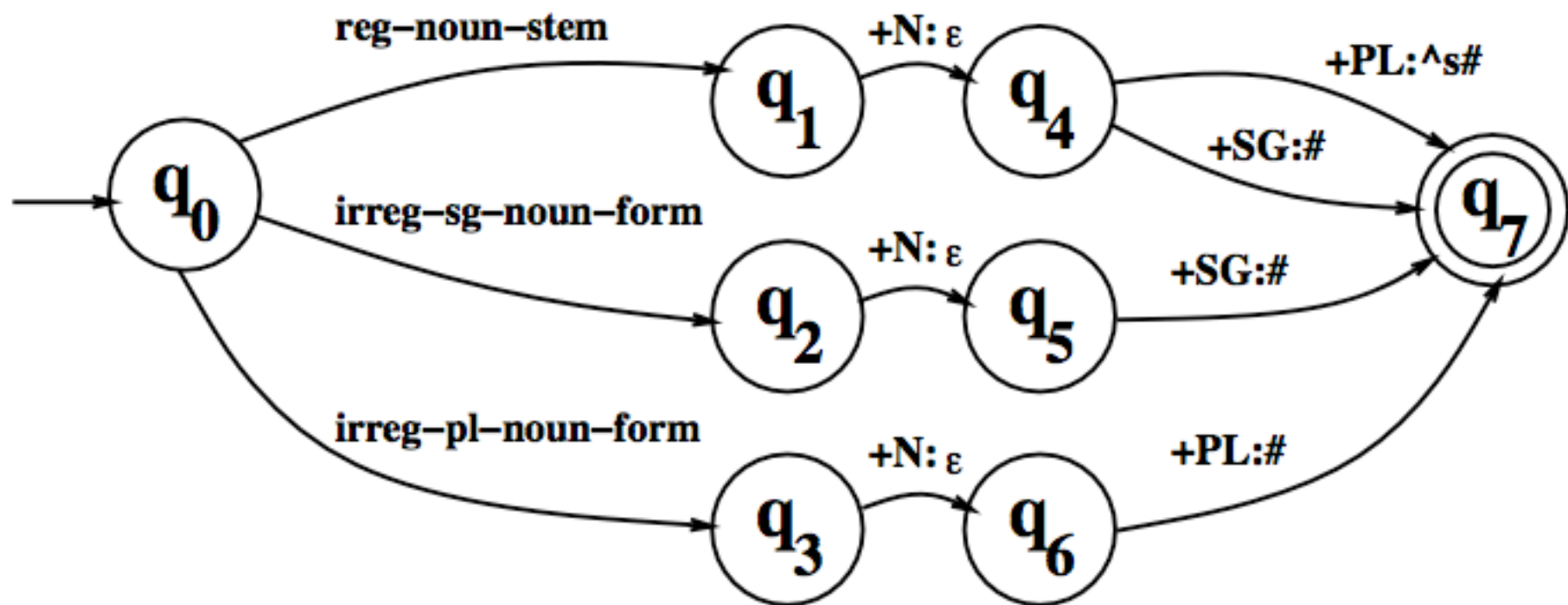
# Morphological Parsing with FSTs

- Retain simple morpheme combinations with the stem by using an intermediate representation:
  - e.g. cat+N+PL becomes *cat^s#*
- Separate rules for the various spelling changes. Each spelling rule is a different FST
- Write down a separate FST for each spelling rule

  *foxes* → fox^s#     [ ***e*-insertion FST** ]

  *loved* → love^ed# [ ***e*-deletion FST** ]

  *flies* → fly^s#       [ ***i* to *y*, *e*-deletion FST** ]

  *panicked* → panic^ed# [ **k-insertion FST** ]

  *etc.*

# Lexicon FST (stores stems)



m o v e : reg-noun-stem

f l y : reg-noun-stem

f o x : reg-noun-stem

+N:+N

+SG:+SG

+PL:+PL

m o u s e : irreg-sg-noun-form

m i c e : irreg-pl-noun-form

Compose the above lexicon FST with some inflection FST

This machine relates intermediate forms like `fox^s#` to underlying lexical forms like `fox+N+PL`

| Lexical | f | o | x | +N | +PL | | | |

| Intermediate | f | o | x | ^ | s | # | | |

*e*-insertion FST

- The label *other* means pairs not use anywhere in the transducer.
- Since # is used in a transition, $q_0$ has a transition on # to itself
- States $q_0$ and $q_1$ accept default pairs like (*cat^s#, cats#*)
- State $q_5$ rejects incorrect pairs like (*fox^s#, foxs#*)

# *e*-insertion FST

- Run the e-insertion FST on the following pairs:

  *(fir#, fir#)*                 *(fizz^s#, fizzs#)*

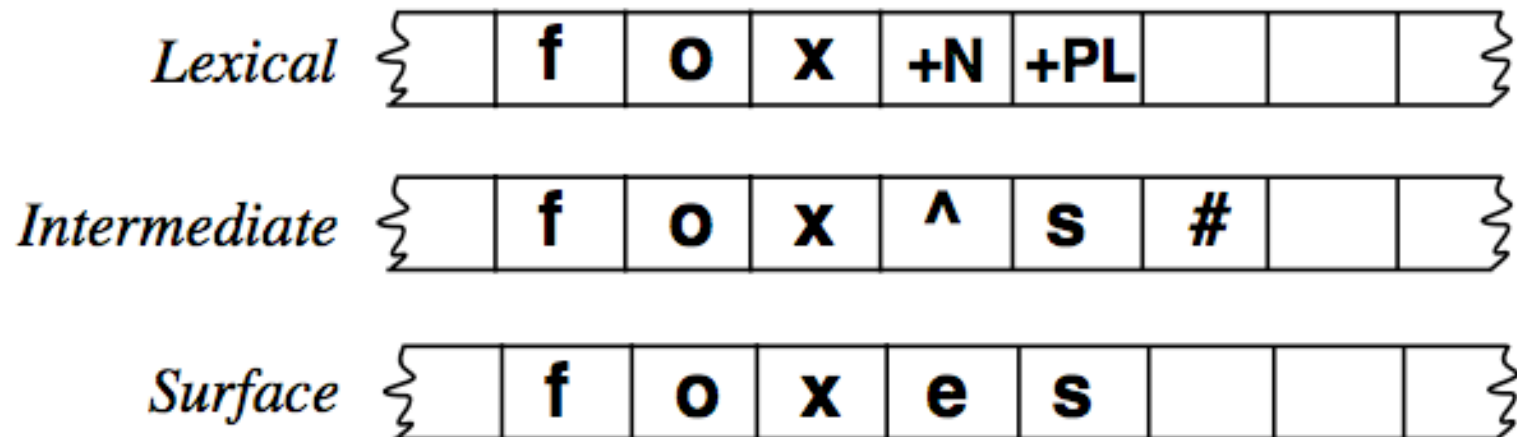  *(fir^s#, firs#)*              *(fizz^s#, fizzes#)*

  *(fir^s#, fires#)*            *(fizz^ing#, fizzing#)*

- Find the state the FST reaches after attempting to accept each of the above pairs

- Is the state a final state, i.e. does the FST accept the pair or reject it

- We first use an FST to convert the lexicon containing the stems and affixes into an intermediate representation
- We then apply a spelling rule that converts the intermediate form into the surface form
- **Parsing**: takes the surface form and produces the lexical representation
- **Generation**: takes the lexical form and produces the surface form
- But how do we handle multiple spelling rules?

| Lexical | f | o | x | +N | +PL | | | |
|---|---|---|---|---|---|---|---|---|

| Intermediate | f | o | x | ^ | s | # | | |
|---|---|---|---|---|---|---|---|---|

| Surface | f | o | x | e | s | | | |
|---|---|---|---|---|---|---|---|---|

# Method 1: Composition

**FST composition**: creates one FST for all rules

.. y+s

Lexicon

FST$_1$

FST$_2$

.
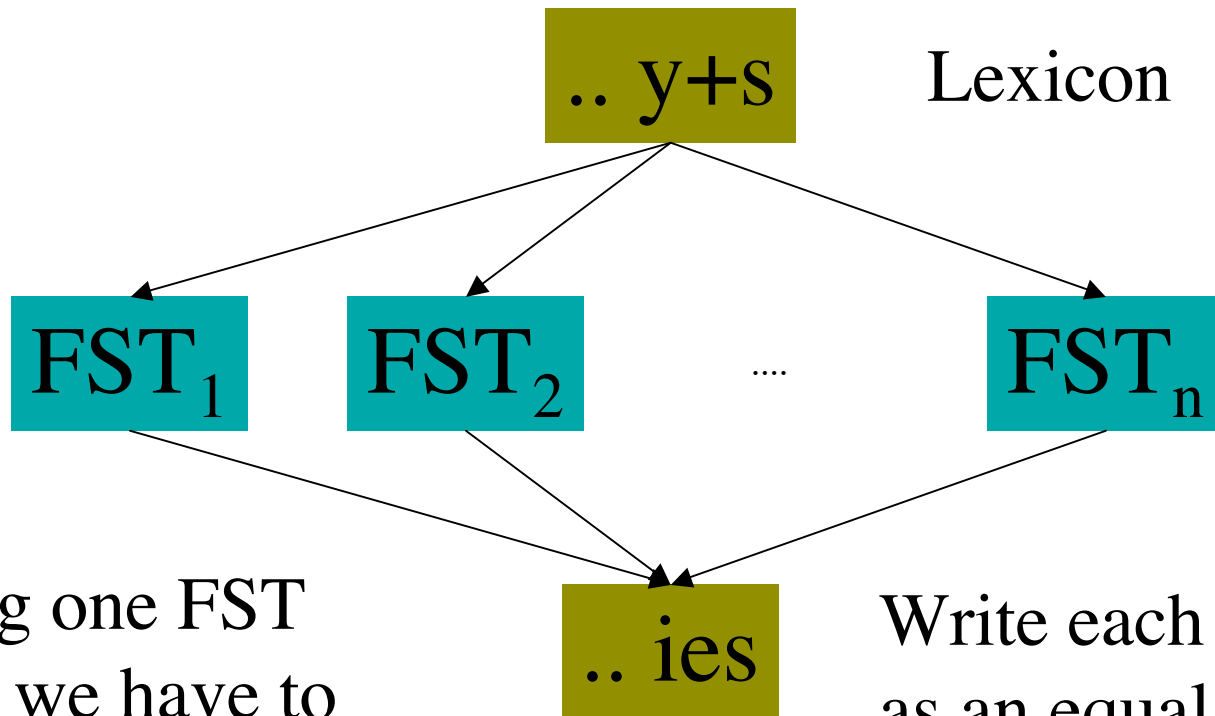.

FST$_n$

.. ies

write one FST for each spelling rule: each FST has to provide input to next stage

# Method 2: Intersection

.. y+s  Lexicon

FST$_1$    FST$_2$    ....    FST$_n$

.. ies

Creating one FST
implies we have to
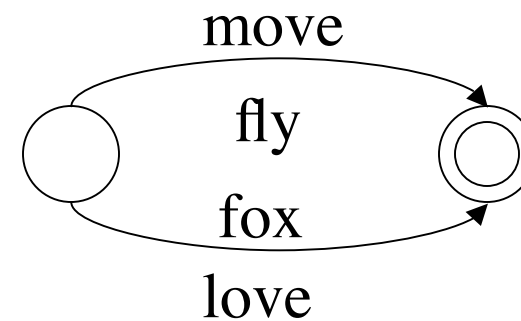do **FST intersection**
(but there's a catch:
*what is it?*)

Write each FST
as an equal length
mapping (ε is taken
to be a real symbol)

# Intersecting/Composing FSTs

- Implement each spelling rule as a separate FST
- We need slightly different FSTs when using Method 1 (composition) vs. using Method 2 (intersection)
  - In Method 1, each FST implements a spelling rule if it matches, and transfers the remaining affixes to the output (composition can then be used)
  - In Method 2, each FST computes an equal length mapping from input to output (intersection can then be used). Finally compose with lexicon FST and input.
- In practice, composition can create large FSTs
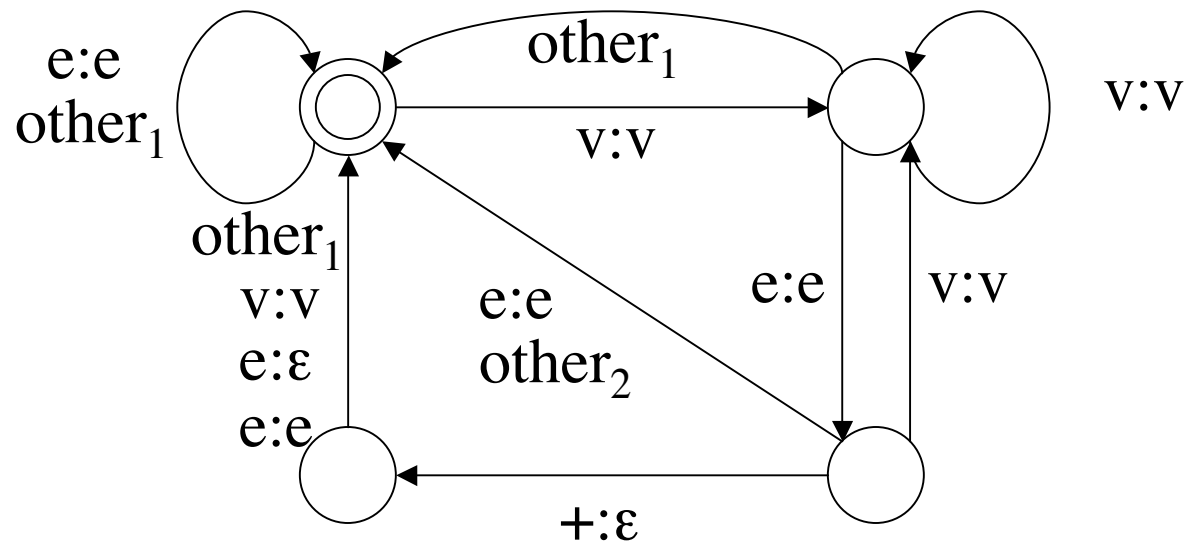
# Length Preserving "two-level" FST for *e-deletion*

Stems/Lexicon



move + ed

move ε εd

$other_1 = \Sigma - \{e, v\}$

$other_2 = \Sigma - \{e, v, +\}$

# Rewrite Rules

left        right
context    context

- Context dependent rewrite rules: $\alpha \rightarrow \beta / \lambda$ __ $\rho$
  - $(\lambda\ \alpha\ \rho \rightarrow \lambda\ \beta\ \rho$; that is $\alpha$ becomes $\beta$ in context $\lambda$ __ $\rho$)
  - $\alpha, \beta, \lambda, \rho$ are regular expressions, $\alpha$ = input, $\beta$ = output
- How to apply rewrite rules:
  - Consider rewrite rule: $a \rightarrow b / ab$ __ $ba$
  - Apply rule on string *ababababa*
  - Three different outcomes are possible:
    - *abbbabbbaba*  (left to right, iterative)
    - *ababbbabbba*  (right to left, iterative)
    - *abbbbbbbba*  (simultaneous)

# Rewrite Rules

u → i / i C* __

(u → i / Σ* i C* __ Σ*)

Input: kikukuku

# Rewrite Rules

u → i / i C* __

kikukuku
kik<span style="color:red">u</span>kuku
kik<span style="color:red">i</span>kuku
kik<span style="color:red">i</span>k<span style="color:red">u</span>ku
kik<span style="color:red">i</span>k<span style="color:red">i</span>ku
kik<span style="color:red">i</span>k<span style="color:red">i</span>k<span style="color:red">u</span>
kik<span style="color:red">i</span>k<span style="color:red">i</span>k<span style="color:red">i</span>

output of one application *feeds* next application

*left to right application*

# Rewrite Rules

u → i / i C* __

kikukuku

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

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

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

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

kik<span style="color:red">ik</span><span style="color:red">ik</span><span style="color:red">u</span>

kik<span style="color:red">ik</span><span style="color:red">ik</span><span style="color:red">i</span>

← *right to left application*

# Rewrite Rules

u → i / i C* __   kikukuku

kikukuku

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

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

*simultaneous application*
(context rules apply to input
string only)

# Rewrite Rules

- Example of the e-insertion rule as a rewrite rule:

  $\varepsilon \rightarrow e \: / \: (x \mid s \mid z)\char`\^ \underline{\phantom{xx}} \: s\#$

- Rewrite rules can be optional or obligatory
- Rewrite rules can be ordered wrt each other
- This ensures exactly one output for a set of rules

# Rewrite Rules

- Rule 1: iN → im / __ (p | b | m)
- Rule 2: iN → in / __
- Consider input *iNpractical* (N is an abstract nasal phoneme)
- Each rule has to be obligatory or we get two outputs: *impractical* and *inpractical*
- The rules have to be ordered wrt to each other so that we get *impractical* rather than *inpractical* as output
- The order also ensures that *intractable* gets produced correctly

# Rewrite Rules

- Under some conditions, these rewrite rules are equivalent to FSTs
- We cannot apply output of a rule as input to the rule itself iteratively:

  $\varepsilon \rightarrow ab \;/\; a \;\underline{\quad}\; b$

  If we allow this, the above rewrite rule will produce $a^n\, b^n$ for $n >= 1$ which is not regular

  Why? Because we rewrite the $\varepsilon$ in $a\varepsilon b$ which was introduced in the previous rule application

  Matching the a__b as left/right context in $a\varepsilon b$ is ok

# Rewrite Rules

- In a rewrite rule: $\alpha \rightarrow \beta \,/\, \lambda \,\underline{\quad}\, \rho$
- Rewrite rules are interpreted so that the **input** $\alpha$ does not match something introduced in the previous rule appliction
- However, we are free to match the **context** either $\lambda$ or $\rho$ or both with something introduced in the previous rule application (see previous examples)
- In this case, we can convert them into FSTs

# Rewrite rules to FSTs

u → i / Σ* i C* __ Σ*     (example from R. Sproat's slides)

- Input: kikukupapu (use left-right iterative matching)
- Mark all possible right contexts

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

- Mark all possible left contexts

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

- Change u to i when delimited by <>

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

- But the next u is not delimited by <> and so cannot be changed even though the rule matches

# Rewrite rules to FSTs

u → i / Σ* i C* __ Σ*

- Input: kikukupapu
- Mark all possible right contexts

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

- Mark all *u* followed by > with $<_1$ and $<_2$

  k > i > k $<_1$ > u > k $<_1$ > u > p > a > p $<_1$ > u >

  $<_2$  u    $<_2$  u        $<_2$  u

- Change all *u* to *i* when delimited by $<_1$ >

  k > i > k $<_1$ > i > k $<_1$ > i > p > a > p $<_1$ > i >

  $<_2$  u    $<_2$  u        $<_2$  u

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

# Rewrite rules to FSTs

$$k > i > k <_1 > i > k <_1 > i > p > a > p <_1 > i >$$
$$\quad\quad <_2 \quad u \quad\quad <_2 \; u \quad\quad\quad\quad <_2 \quad u$$

- Delete >

$$k\; i\; k <_1 i\; k <_1 i\; p\; a\; p <_1 i$$
$$\quad\quad <_2 u \quad <_2 u \quad\quad <_2 u$$

- Only allow $i$ where $<_1$ is preceded by $iC^*$, delete $<_1$

$$k\; i\; k \quad i\; k \quad i\; p\; a\; p$$
$$\quad\quad <_2 u \quad <_2 u \quad\quad\quad <_2 u$$

- Allow only strings where $<_2$ is **not** preceded by $iC^*$, delete $<_2$

$$k\; i\; k\; i\; k\; i\; p\; a\; p\; u$$

# Rewrite rules to FST

- For every rewrite rule: $\alpha \rightarrow \beta \ / \ \lambda \ \underline{\quad} \ \rho$:
  - FST $r$ that inserts > before every $\rho$
  - FST $f$ that inserts $<_1$ & $<_2$ before every $\alpha$ followed by >
  - FST $replace$ that replaces $\alpha$ with $\beta$ between $<_1$ and > and deletes >
  - FST $\lambda_1$ that only allows all $<_1 \beta$ preceded by $\lambda$ and deletes $<_1$
  - FST $\lambda_2$ that only allows all $<_2 \beta$ **not** preceded by $\lambda$ and deletes $<_2$
- Final FST = $r \ _o \ f \ _o \ replace \ _o \ \lambda_1 \ _o \ \lambda_2$
- This is only for left-right iterative obligatory rewrite rules: similar construction for other types

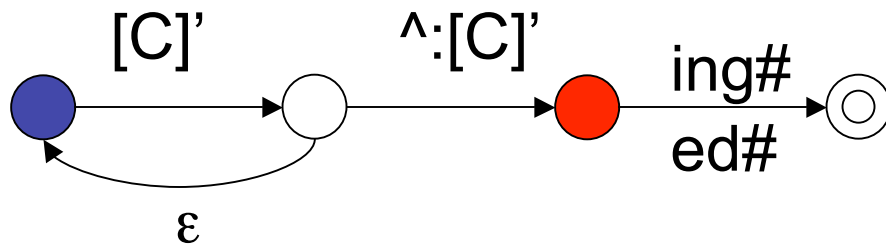# Rewrite Rules to FST

## FST for *replace*

$\Sigma: \Sigma$

$<_2:<_2$

$>: \varepsilon$          $<_1:<_1$

$[\alpha\times\beta]$

$<_1:\varepsilon, <_2:\varepsilon, >:\varepsilon$

$>:\varepsilon$

Create a new FST by taking the cross product of the languages $\alpha$ and $\beta$ and each state of this new FST: $[\alpha\times\beta]$ has loops for the transitions $<_1:\varepsilon, <_2:\varepsilon, >:\varepsilon$
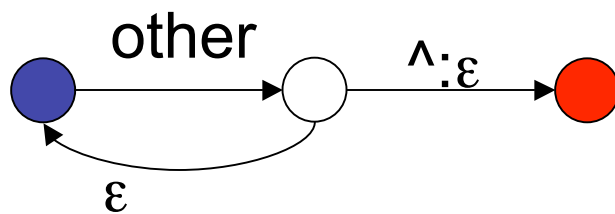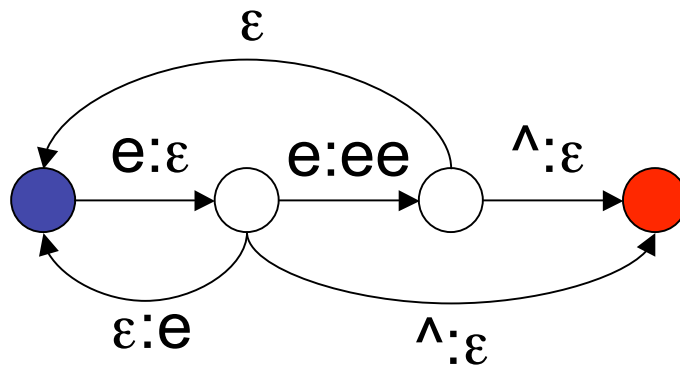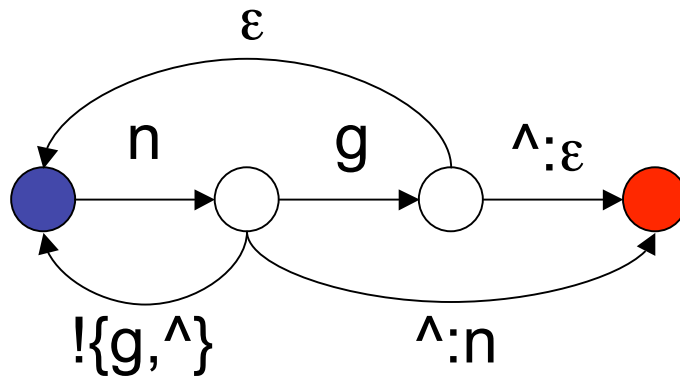
# Ambiguity (in parsing)

- Global ambiguity: (de+light+ed *vs.* delight+ed)

  *foxes* → fox+N+PL (*I saw two foxes*)

  *foxes* → foxes+V+3SG (*Clouseau foxes them again*)

- Local ambiguity:

  *assess* has a prefix string *asses* that has a valid analysis:
  *asses* → ass+N+PL

- Global ambiguity results in two valid answers, but local ambiguity returns only one.

- However, local ambiguity can also slow things down since two analyses are considered partway through the string.

# Summary

- FSTs can be applied to creating lexicons that are aware of morphology
- FSTs can be used for simple stemming
- FSTs can also be used for morphographemic changes in words (spelling rules), e.g. fox+N+PL becomes foxes
- Multiple FSTs can be composed to give a single FST (that can cover all spelling rules)
- Multiple FSTs that are length preserving can also be run in parallel with the intersection of the FSTs
- Rewrite rules are a convenient notation that can be converted into FSTs automatically
- Ambiguity can exists in the lexicon: both global & local

[C]'   ^:[C]'   ing#   ed#   [C]' = [C]-{n}

ε
n   g   ^:ε
!{g,^}   ^:n

ε
e:ε   e:ee   ^:ε
ε:e   ^:ε

other   ^:ε
ε
other = Σ-[C]'-{n,e}