

# 6.863J Natural Language Processing Lecture 2: Automata, Two-level phonology, & PC-Kimmo (the Hamlet lecture)

---

Instructor: Robert C. Berwick  
berwick@ai.mit.edu



## The Menu Bar

---

- Administrivia

web page: [www.ai.mit.edu/courses/6.863/](http://www.ai.mit.edu/courses/6.863/) now with  
Lecture 1, Lab components I (background), II (Lab  
1a) – due date is February 17 (holiday the 16<sup>th</sup>)

- *What* and *How*: word processing, or computational morphology
- What's in a word: morphology
- Modeling morphophonology by finite-state devices
- Finite-state automata vs. finite state transducers
- Some examples from English
- PC-Kimmo & Laboratory 1a :how-to




## Levels of language

---

- Phonetics/phonology/morphology: what words (or subwords) are we dealing with?
- Syntax: What phrases are we dealing with? Which words modify one another?
- Semantics: What's the literal meaning?
- Pragmatics: What should you conclude from the fact that I said something? How should you react?

6.863J/9.611J SP04 Lecture 2



## Start with words: they illustrate all the problems (and solutions) in NLP

---

- Parsing words  
Cats → CAT + N(oun) + PL(ural)
- Used in:
  - Traditional NLP applications
  - Finding word boundaries (e.g., Latin, Chinese)
  - Text to speech (*boathouse*)
  - Document retrieval (example next slide)
- In particular, the problems of *parsing*, *ambiguity*, and *computational efficiency* (as well as the problems of *how people do it*)

6.863J/9.611J SP04 Lecture 2

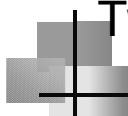


## Terminology

---

- Morpheme, etc..
- See webpage  
<http://pandora.cii.wvu.edu/vajda/ling201/test1materials/Morphologyoverhead.htm>

6.863J/9.611J SP04 Lecture 2



## Two parts to the “what”

---

1. Which units can glue to which others (roots and affixes) (or stems and affixes), eg,
2. What ‘spelling changes’ (orthographic changes) occur – like dropping the *e* in ‘chase + ed’

OK, let’s tackle these one at a time

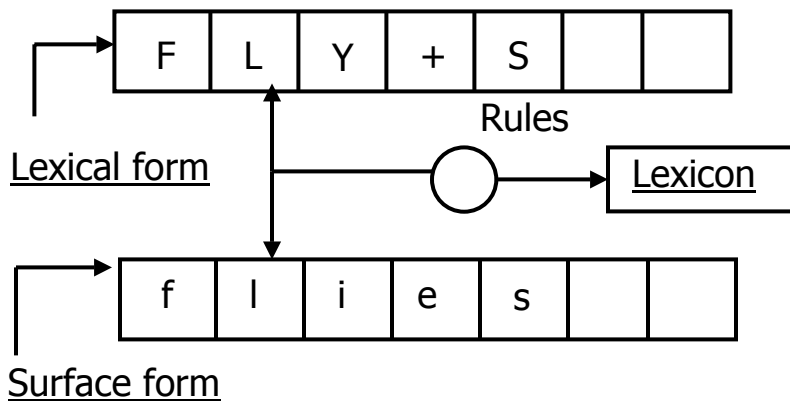
6.863J/9.611J SP04 Lecture 2

## The Knowledge

- ✓ Which units can glue to which others (roots and affixes) (or stems and affixes)
- ✓ What 'spelling changes' (orthographic changes) occur – like dropping the *e* in 'chase + ed'

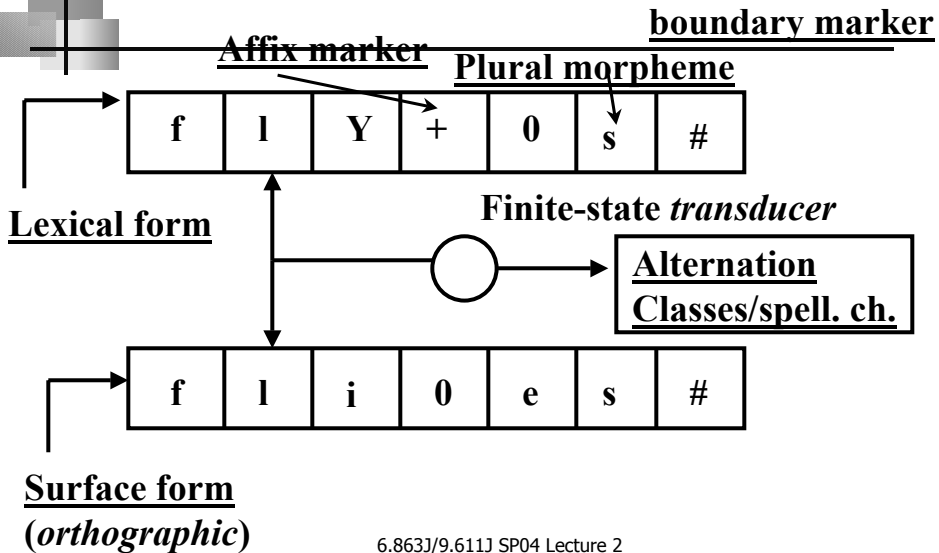
6.863J/9.611J SP04 Lecture 2

## Our goal: PC-Kimmo



6.863J/9.611J SP04 Lecture 2

## Terminology



## What is the Knowledge?

- Morpheme classes
- Like beads on a string
- Computational model:
  - Finite-state automata



## Two parallel finite-state machines

---

- Machine 1: order of morphemes
- Machine 2: spelling changes

6.863J/9.611J SP04 Lecture 2



## Why finite-state machines?

---

- Minimal model
- Fast
- What makes up a finite-state automaton?
- Linear concatenation of equivalence classes

6.863J/9.611J SP04 Lecture 2

## Definition of finite-state automaton (fsa)

---

- A finite-state automaton (FSA) is a quintuple  $(Q, \Sigma, \delta, q_0, F)$  where
  - $Q$  is a *finite* set of states
  - $\Sigma$  is a finite set of terminal symbols, the alphabet
  - $q_0 \in Q$  is the initial state
  - $F \subseteq Q$ , the set of final states
  - $\delta$  is a mapping from  $Q \times \Sigma \rightarrow 2^Q$ , the transition mapping

6.863J/9.611J SP04 Lecture 2

## Application to morphology

---

- Purely concatenative morphemes like a rooted tree (but what about prefixes?)

6.863J/9.611J SP04 Lecture 2

## English morphology: what states do we need for the fsa?

---

- As an example, consider adjectives

*Big, bigger, biggest*

*Cool, cooler, coolest, coolly*

*Red, redder, reddest*

*Clear, clearer, clearest, clearly, unclear, unclearly*

*Happy, happier, happiest, happily*

*Unhappy, unhappier, unhappiest, unhappily*

*Real, unreal, silly*

6.863J/9.611J SP04 Lecture 2

## FSA states

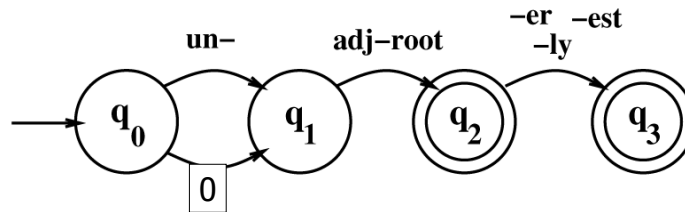
---

- Equivalence classes under the notion of 'substitution': elements that 'behave alike'
- Consider: \_\_\_er What goes in the space?
  - Cool, big, happy (Adjectives that can be comparative)
- If we need to make more refinements, we need a new class

6.863J/9.611J SP04 Lecture 2



Will this fsa work?



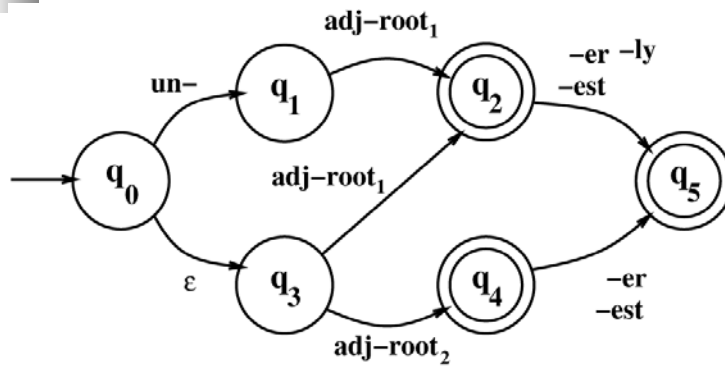
6.863J/9.611J SP04 Lecture 2

Ans: no!

- Accepts all adjectives above, but
- Also accepts *unbig*, *readly*, *realest*
- Common problem: overgeneration
- Solution?

6.863J/9.611J SP04 Lecture 2

## Revised picture



6.863J/9.611J SP04 Lecture 2

## Extension to politics: an alien language

- Bush could win the election
- Bush will win the election
- Bush did win the election
- Bush could have won the election
- Bush will have won the election

...

6.863J/9.611J SP04 Lecture 2

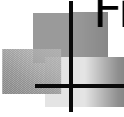


## Structural linguistics circa 1940s

---

“[Words] are assigned to classes on the basis of the environments in which they occur. Each environment determines one and only one class, namely the class of all [words] occurring in that environment... A word *A* belongs to the class determined by the environment \_\_\_\_X if X is either an utterance or occurs as part of some utterance.” Wells, 1947, pp. 81-82.

6.863J/9.611J SP04 Lecture 2



## Finite transition networks (FTNs)

---

- 1-1 picture of fsa-ftn
- Easy to see, easy to prove certain properties: closure under concatenation, intersection
- Conversion of nondeterministic to deterministic fsa – always possible

6.863J/9.611J SP04 Lecture 2

## Limits of FTNs

---

- Finite # of 'memory states'
- Can distinguish only a finite # of classes (bins) i.e., the states
- This sets limits on the patterns fsa's can recognize
- FSA's cannot even describe all possible human words...

6.863J/9.611J SP04 Lecture 2

## Then can be *indescribable* words (for an fst)

---

- Can we even do all natural languages?
- Example: Bambarra (African language in Mali)
- Words in form *Noun+o+Noun*, as in *wuluowulo* = 'whichever dog'
- Also have repeated endings (like anti-anti...)  
*wulu+nyini+la* = 'dog searcher'  
*wulunyinina+ nyini+la* = 'one who searches for dog searchers'
- Fatal bite: combine with word *o word* formation:  
*wulunyininyinila o wulunyininyinila* (arbitrarily long!)

6.863J/9.611J SP04 Lecture 2

## Why is this not describable by an fsa?

---

- Intuition:  $w^n$  o  $w^n$  language
- Need arbitrary # of bins to keep track of  $w^n$  string to match it up with  $w^n$  after the 'o'
- Must be able to count to arbitrary  $n$  to keep track of # of copies...
- But, only a finite # of bins...so....

6.863J/9.611J SP04 Lecture 2

---

## What does fsa machine for English look like?

6.863J/9.611J SP04 Lecture 2

## Next: what about the spelling changes? That's harder!

---

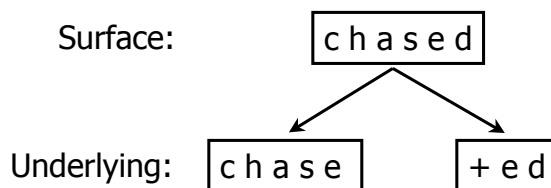
✓ Which units can glue to which others (roots and affixes) (or stems and affixes)

2. What 'spelling changes' (orthographic changes) occur – like dropping the *e* in 'chase + ed'

6.863J/9.611J SP04 Lecture 2

## Mapping between surface form & underlying form

---



But clearly this can go either way – given the underlying form, we can *generate* the surface form – so we really have a *relation* betw. surface & underlying form, viz.:

6.863J/9.611J SP04 Lecture 2

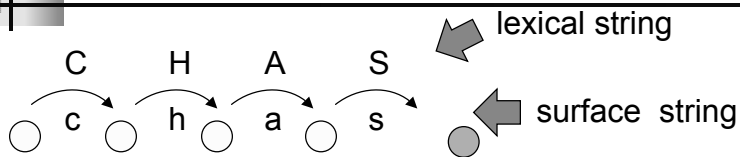
## Conventional notation

Lexical (underlying) form: c h a s e + e d  
Surface form: c h a s 0 0 e d

The 0's "line up" the lexical & surface strings  
This immediately suggests a finite-state automaton  
'solution' : an extension known as a  
*finite-state transducer*

6.863J/9.611J SP04 Lecture 2

## Finite-state transducers: a *pairing* between lexical/surface strings



- Or more carefully

6.863J/9.611J SP04 Lecture 2

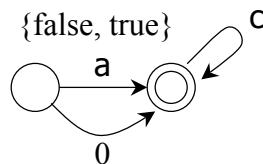
## Definition of finite-state transducer

- state set  $Q$
- initial state  $q_0$
- set of final states  $F$
- input alphabet  $S$  (also define  $\Sigma^*$ ,  $\Sigma^+$ )
- output alphabet  $D$
- transition function  $\delta : Q \times \Sigma \rightarrow 2^Q$
- output function  $\sigma : Q \times \Sigma \times Q \rightarrow D^*$

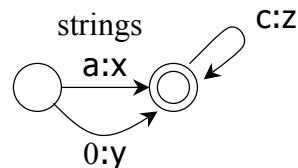
6.863J/9.611J SP04 Lecture 2

## The difference between (familiar) fsa's and fst's: functions from...

### Acceptors (FSAs)



### Transducers (FSTs)



6.863J/9.611J SP04 Lecture 2

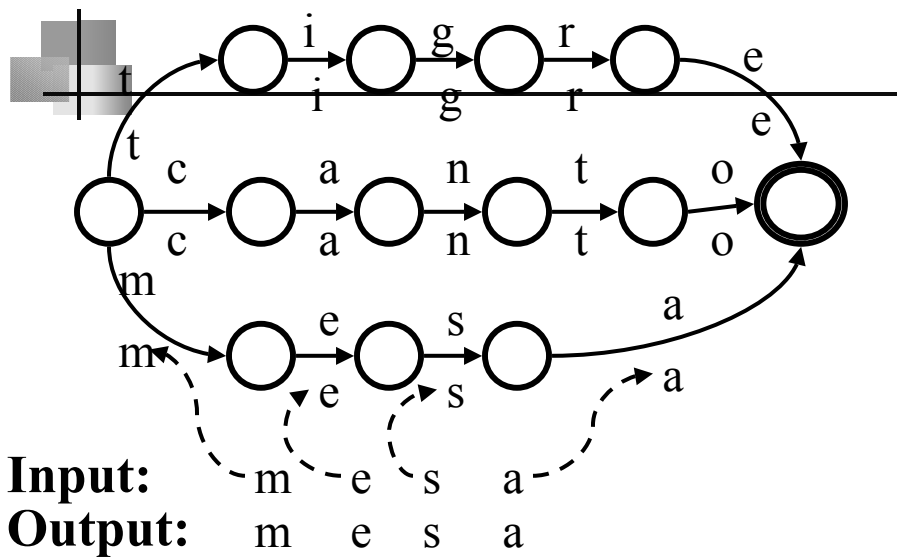


## Properties of fst's – compare to fsa's

- Closed under concatenation – get another fst if wired together
- NOT closed under intersection
- NOT always able to make deterministic

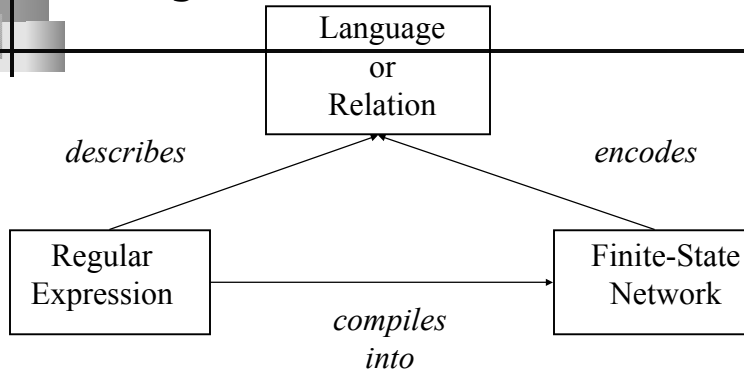
6.863J/9.611J SP04 Lecture 2

## A Two-Level Transducer



6.863J/9.611J SP04 Lecture 2

## The Big Picture



6.863J/9.611J SP04 Lecture 2

## Defining an fst for a spelling-change rule

- Suggests all we need to do is build an FST for a spelling-change rule that 'matches' lexical and surface strings
- Example: fox+s, foxes; buzz+s, buzzes
- Rule: e before non initial x,s,z
- Instantiation as an FST:

<b>F</b>	<b>O</b>	<b>X</b>	<b>+</b>	<b>0</b>	<b>S</b>	<b>#</b>	<b>lexical</b>
f	o	x	0	e	s	#	<b>surface</b>

6.863J/9.611J SP04 Lecture 2

## Implication

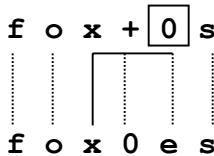
---

- 0:e can occur only in this context
- Must write this as a constraint
- Write an FTN that accepts only strings of this form, e.g., dafjakjdx0:es#

6.863J/9.611J SP04 Lecture 2

## FTN imposes a Constraint

---



0:e correspondence requires  
a preceding x on the  
lexical side, s:s following

In this context, all other  
possible realization of a  
0:s are prohibited.

6.863J/9.611J SP04 Lecture 2

## Turning this into an fst

- Write down the left, center, and right context
- In this case:

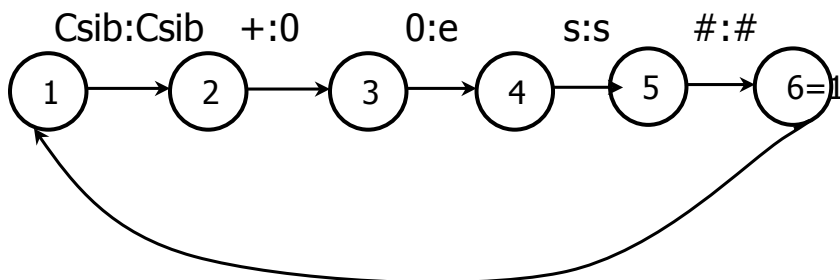
x:x    +:0    0:e    s:s    #:#

Csib:Csib

- Pad out with nulls (0's)
- Write an fsa (ftn) that accepts exactly this string

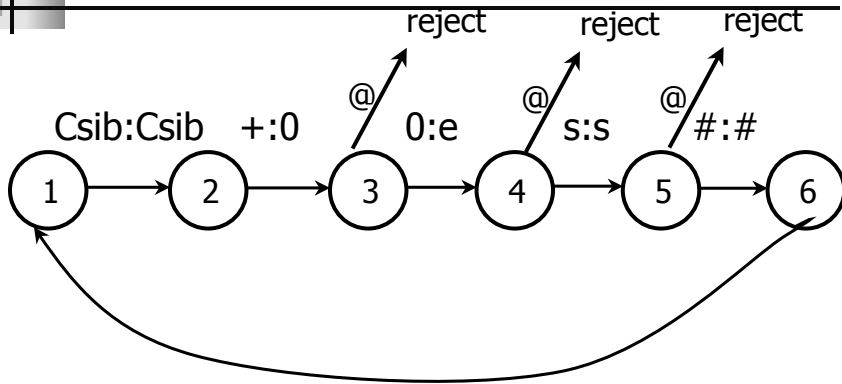
6.863J/9.611J SP04 Lecture 2

## Start with straightline fst



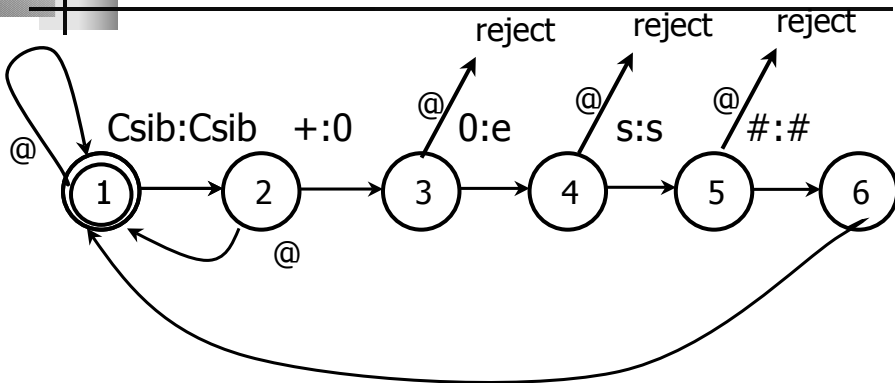
6.863J/9.611J SP04 Lecture 2

Now add rejection notices...

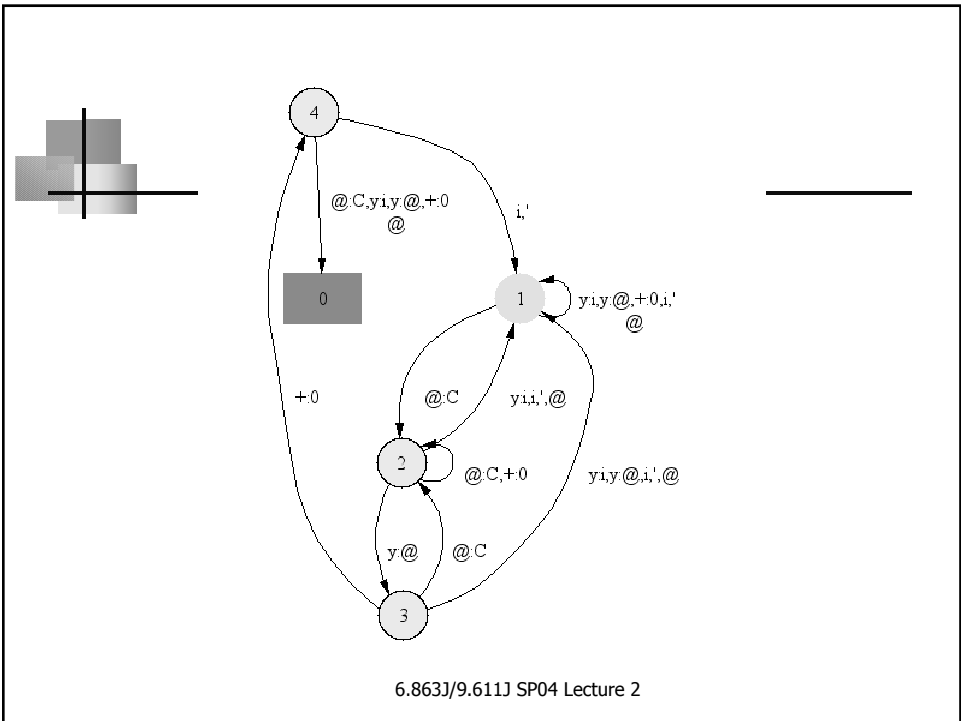
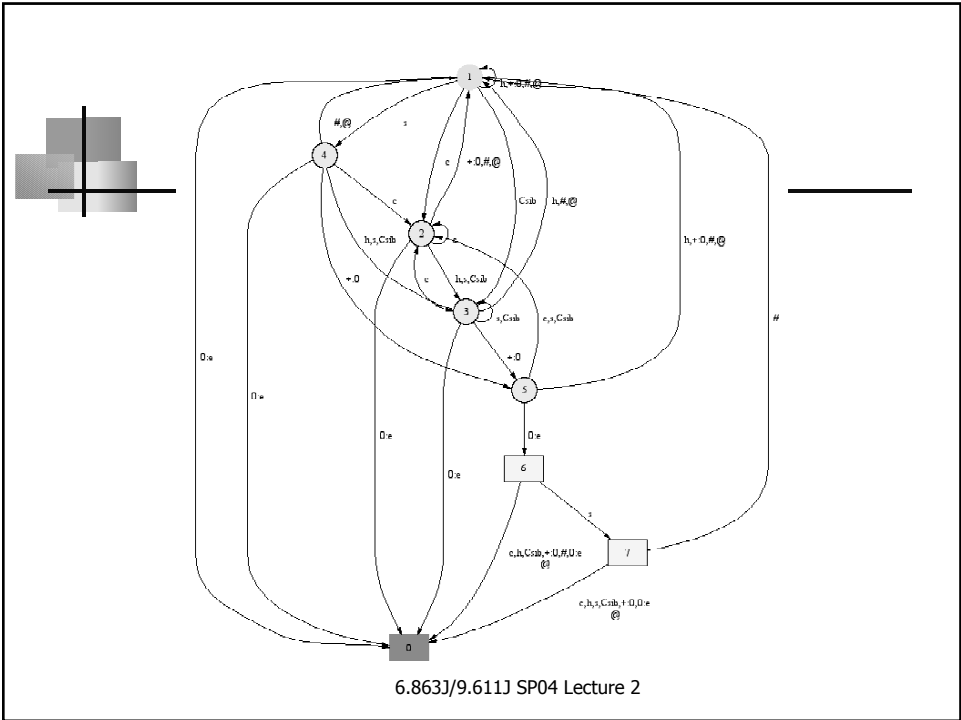


6.863J/9.611J SP04 Lecture 2

And acceptance (cook until done)



6.863J/9.611J SP04 Lecture 2



## Tabular format

---

```
RULE "3 Epenthesis, 0:e => [Csib|ch|sh|] +:0__s [#]"  
7 8
```

```
  c h s Csib + # 0 @  
  c h s Csib 0 # e @  
1: 2 1 4 3    1 1 0 1  
2: 2 3 3 3    1 1 0 1  
3: 2 1 3 3    5 1 0 1  
4: 2 3 3 3    5 1 0 1  
5: 2 1 2 2    1 1 6 1  
6: 0 0 7 0    0 0 0 0  
7: 0 0 0 0    0 1 0 0
```

6.863J/9.611J SP04 Lecture 2

## And that's (almost) all folks...

---

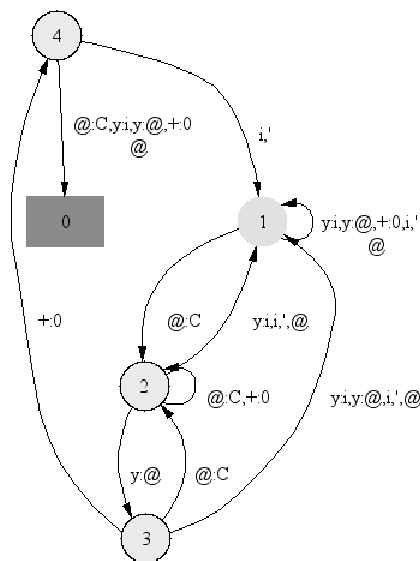
- Except...
- There's more than one rule...
  
- Spy+s → spies
- Quiz+s → quizzes
- Make+ing → making

6.863J/9.611J SP04 Lecture 2

# Spelling change rules

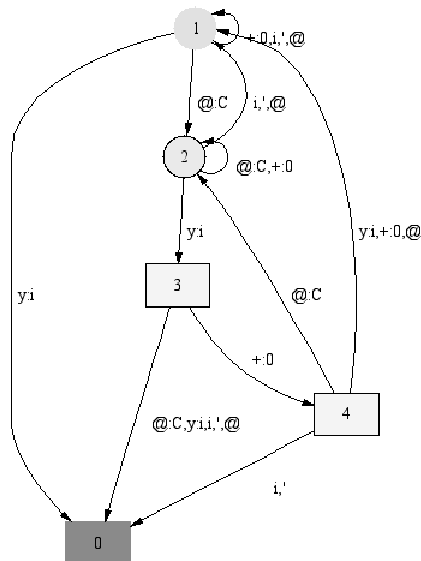
Name	Description	Example
Consonant Doubling (gemination, G)	1-letter consonant doubled before <i>-ing/ed</i>	beg/begging
E deletion (elision, EL),	Silent <i>e</i> dropped before <i>-ing, -ed</i>	make/making
E insertion (epenthesis, EP)	<i>e</i> added after <i>-s, -z, -ch, -sh</i> before <i>-s</i>	fox/foxes
Y replacement (Y)	<i>-y</i> changes to <i>-ie</i> before <i>-ed</i>	try/tries
I spelling (I)	<i>I</i> goes to <i>y</i> before vowel	lie/lying

# Set of 5 spelling change FTNs





## Another...name that automaton...



## So, we're done, right?

- So, not so fast...!!!!
- Sometimes, more than 1 spelling change rule applies. Example: spy+s, spies: y
- $y$  goes to  $i$  before an inserted  $e$  (compare, "spying")
- $e$  inserted at affix +s
- Here's the picture:

## Two-Level Constraints 2

s p y 0 + s  
 .....  
 .....  
 s p i e 0 s

s p y 0 + s  
 .....  
 .....  
 s p i e 0 s

y:i <=> \_ 0:e

0:e <=> Cons: y: \_ +:0 s:s

6.863J/9.611J SP04 Lecture 2

## Another Example from English ("gemination")

underlying quiz + s



Rule A: s -> es after z

intermediate quiz + es

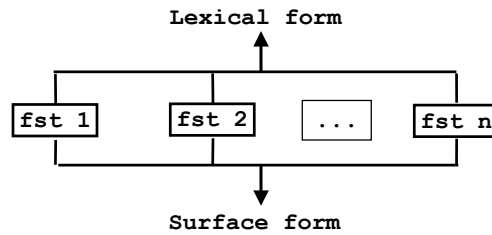


Rule B: z doubles before  
 Suffix beginning with  
 vowel

surface quizzes

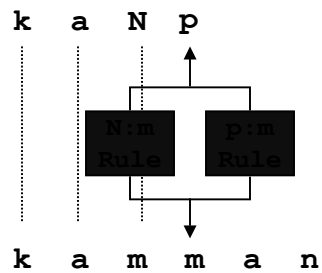
6.863J/9.611J SP04 Lecture 2

## Run transducers in parallel?



6.863J/9.611J SP04 Lecture 2

## Parallel application – how?



6.863J/9.611J SP04 Lecture 2

# Sequential Application

k a N p a n



k a m p a n



k a m m a n

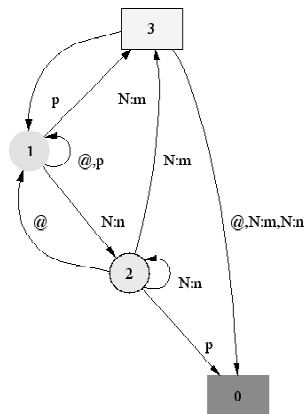
$N \rightarrow m / \_ p$

$p \rightarrow m / m \_$

6.863J/9.611J SP04 Lecture 2

# Machine Rule 1 ("N goes to m")

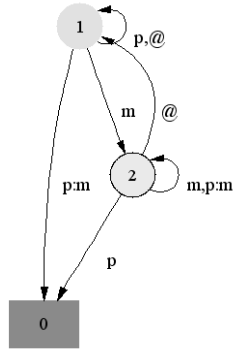
Rule 1:  $N \rightarrow m \mid \_ p$



6.863J/9.611J SP04 Lecture 2

# Machine Rule 2 ("p goes to m")

Rule 2:  $p \rightarrow m \mid m \_$

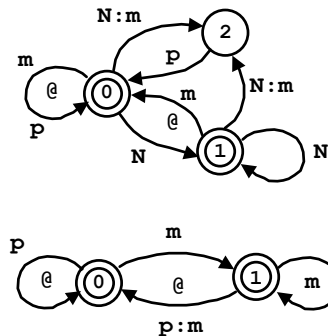


6.863J/9.611J SP04 Lecture 2

# Sequential Application in Detail

```

    k a N p a n
      |
    0 0 0 2 0 0 0
      |
    k a m p a n
      |
    0 0 0 1 0 0 0
      |
    k a m m a n
  
```



6.863J/9.611J SP04 Lecture 2

## Constraints on both sides

k a N **p** a n  
 .....  
 k a m m a n

N:m correspondence  
 requires a following p on  
 the lexical side.

In this context, all other  
 possible realization of a  
 lexical N are prohibited.

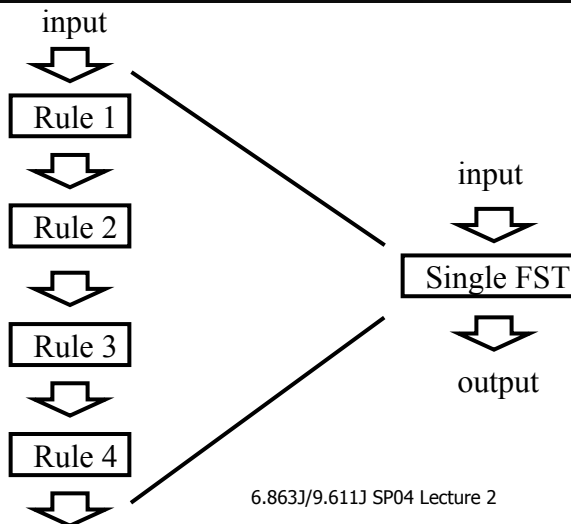
k a N p a n  
 .....  
 k a **m** m a n

p:m correspondence  
 requires a preceding m  
 on the surface side.

In this context, all other  
 possible realization of a  
 lexical p are prohibited.

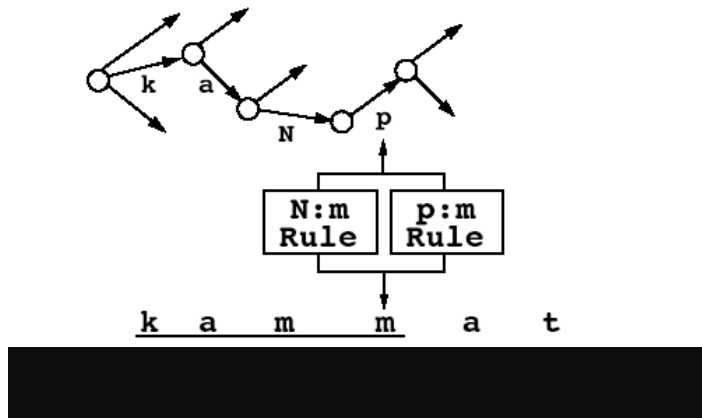
6.863J/9.611J SP04 Lecture 2

## When is this possible?



6.863J/9.611J SP04 Lecture 2

Plus lexicon – lexical forms always constrained by the path we're following through the lexicon tree



6.863J/9.611J SP04 Lecture 2

We trace through the finite-state devices in tandem

