

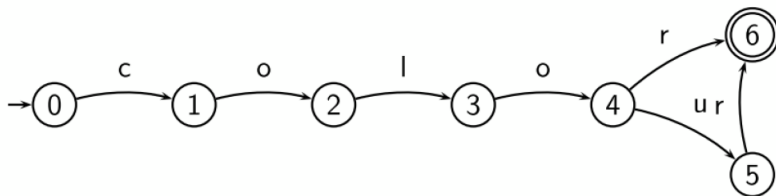
Finite-state methods for morphology

Pawan Goyal

CSE, IITKGP

Week 3: Lecture 3

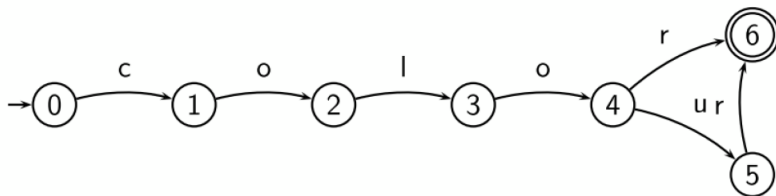
Finite State Automaton (FSA)



What is FSA?

- A kind of directed graph

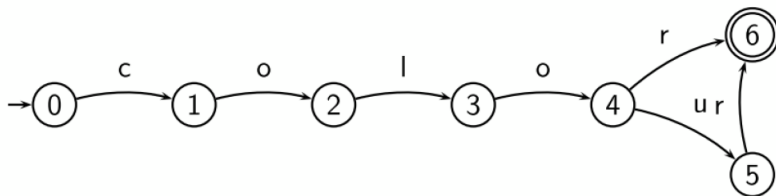
Finite State Automaton (FSA)



What is FSA?

- A kind of directed graph
- Nodes are called states, edges are labeled with symbols (possibly empty ϵ)

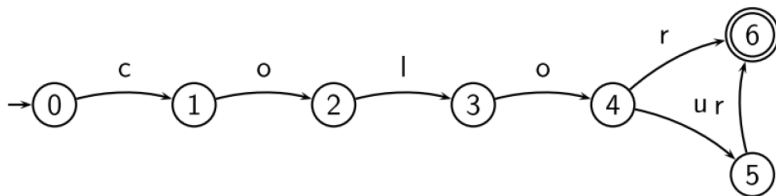
Finite State Automaton (FSA)



What is FSA?

- A kind of directed graph
- Nodes are called states, edges are labeled with symbols (possibly empty ϵ)
- Start state and accepting states

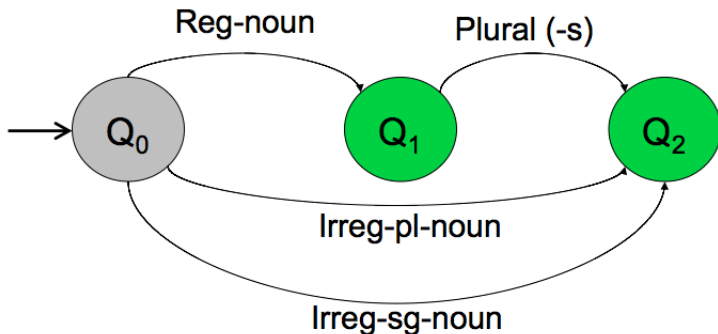
Finite State Automaton (FSA)



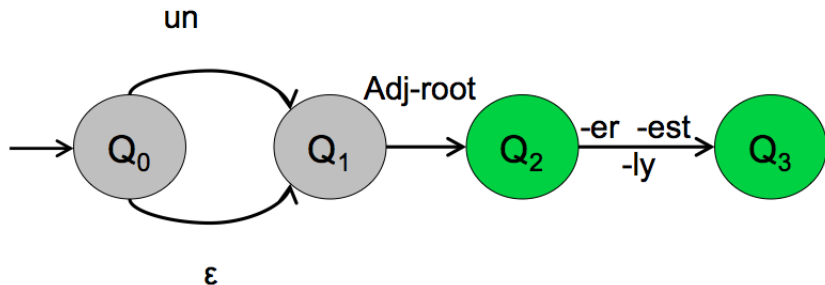
What is FSA?

- A kind of directed graph
- Nodes are called states, edges are labeled with symbols (possibly empty ϵ)
- Start state and accepting states
- Recognizes regular languages, i.e., languages specified by regular expressions

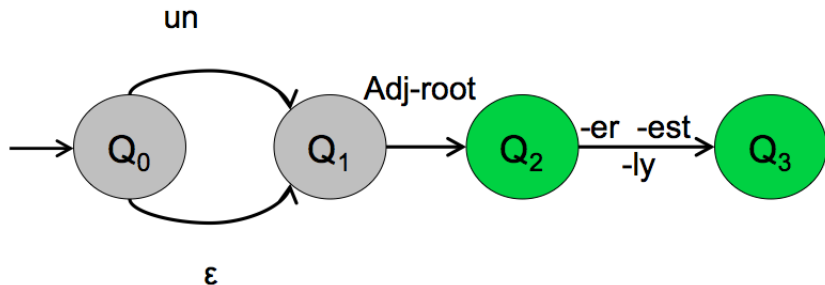
FSA for nominal inflection in English



FSA for English Adjectives



FSA for English Adjectives



Word modeled

happy, happier, happiest, real, unreal, cool, coolly, clear, clearly, unclear, unclearly, ...

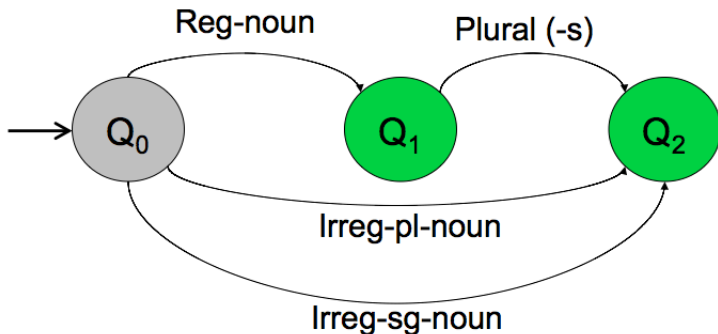
- The last two examples model some parts of the English morphotactics
- But what about the information about regular and irregular roots?

- The last two examples model some parts of the English morphotactics
- But what about the information about regular and irregular roots?

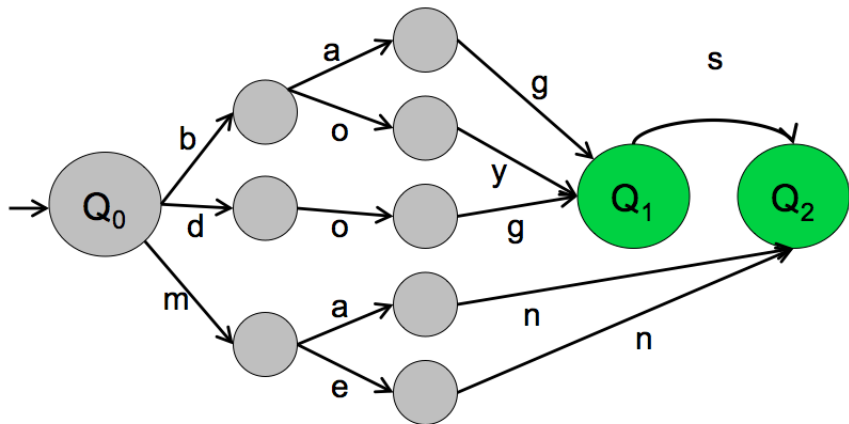
Lexicon

Can we include the lexicon in the FSA?

FSA for nominal inflection in English



After adding a mini-lexicon



Some properties of FSAs: Elegance

- Recognizing problem can be solved in linear time (independent of the size of the automaton)
- There is an algorithm to transform each automaton into a unique equivalent automaton with the least number of states
- An FSA is deterministic iff it has no empty (ϵ) transition and for each state and each symbol, there is at most one applicable transition
- Every non-deterministic automaton can be transformed into a deterministic one

But ...

FSAs are language recognizers/generators.

But ...

FSAs are language recognizers/generators.

We need transducers to build Morphological Analyzers

But ...

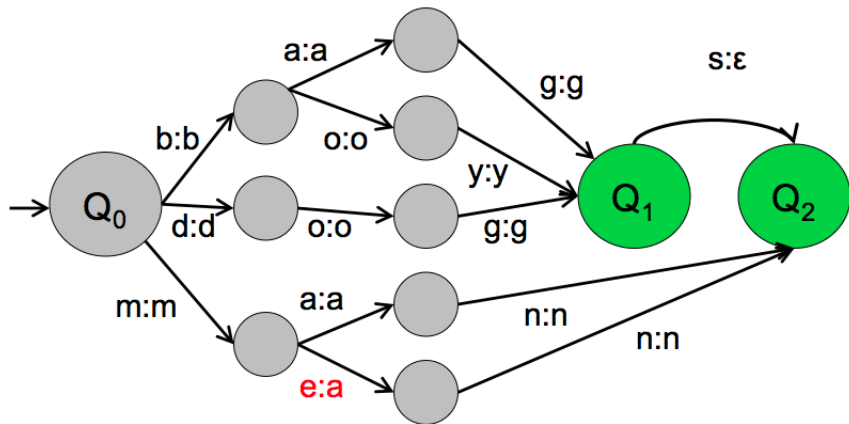
FSAs are language recognizers/generators.

We need transducers to build Morphological Analyzers

Finite State Transducers

- Translate strings from one language to strings in another language
- Like FSA, but each edge is associated with two strings

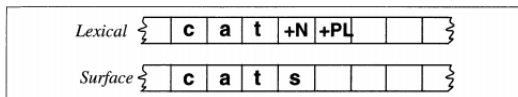
An example FST



Two-level morphology

Given the input *cats*, we would like to output *cat+N+PL*, telling us that cat is a plural noun.

We do this via a version of **two-level morphology**, a correspondence between a lexical level (morphemes and features) to a surface level (actual spelling).



Intermediate tape for Spelling change rules

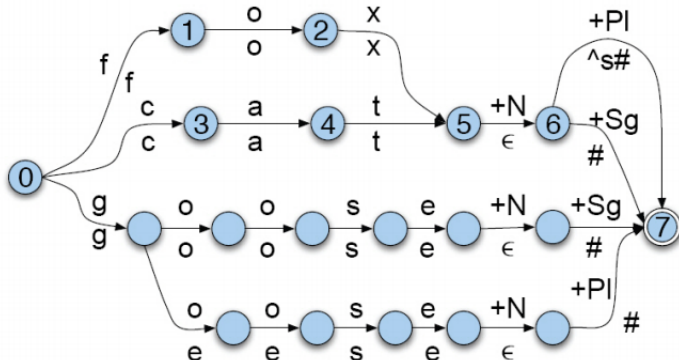
Lexical

f	o	x	+N	+PI			
---	---	---	----	-----	--	--	--

Intermediate

f	o	x	^	s	#		
---	---	---	---	---	---	--	--

English Nominal Inflection FST



Spelling Handling

A spelling change rule would insert an e only in the appropriate environment.

Lexical

	f	o	x	+N	+Pl			
--	---	---	---	----	-----	--	--	--

Intermediate

	f	o	x	^	s	#		
--	---	---	---	---	---	---	--	--

Surface

	f	o	x	e	s			
--	---	---	---	---	---	--	--	--

Rule Notation

$a \rightarrow b/c_d$: “rewrite a as b when it occurs between c and d .”

Morphological Analysis: Approaches

Two different ways to address phonological/graphemic variations

- Linguistic approach: A phonological component accompanying the simple concatenative process of attaching an ending
- Engineering approach: Phonological changes and irregularities are factored into endings and a higher number of paradigms

Different Approaches: Example from Czech

	woman	owl	draft	iceberg	vapor	fly
S1	žen-a	sov-a	skic-a	kr-a	pár-a	mouch-a
S2	žen-y	sov-y	skic-i	kr-y	pár-y	mouch-y
S3	žen-ě	sov-ě	skic-e	kř-e	pář-e	mouš-e
:						
P2	žen-0	sov-0	skic-0	ker-0	par-0	much-0

A linguistic approach

$$\begin{array}{cccccc}
 \text{žen} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix} & \text{sov} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix} & \text{skic} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix} & \text{kr} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix} & \text{pár} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix} & \text{mouch} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix}
 \end{array}$$

An engineering approach

$$\begin{array}{cccccc}
 \text{žen} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix} & \text{sov} + \begin{Bmatrix} a \\ y \\ \text{ě} \\ 0 \end{Bmatrix} & \text{skic} + \begin{Bmatrix} a \\ \text{i} \\ \text{e} \\ 0 \end{Bmatrix} & \text{k} + \begin{Bmatrix} \text{ra} \\ \text{ry} \\ \text{ře} \\ \text{er} \end{Bmatrix} & \text{p} + \begin{Bmatrix} \text{ára} \\ \text{áry} \\ \text{áře} \\ \text{ar} \end{Bmatrix} & \text{m} + \begin{Bmatrix} \text{oucha} \\ \text{ouchy} \\ \text{ouše} \\ \text{uch} \end{Bmatrix}
 \end{array}$$

- AT&T FSM Library and Lextools

<http://www2.research.att.com/~fsmtools/fsm/>

- OpenFST (Google and NYU)

<http://www.openfst.org/>