

# Autómatas 2

Escuela Nacional de Antropología e Historia (ENAH)

Agosto – diciembre de 2015

# Análisis (Parsing)

- Obtener una estructura lingüística que describa una entrada.
- going → VERB-go + GERUND-ing
- Análisis morfológico

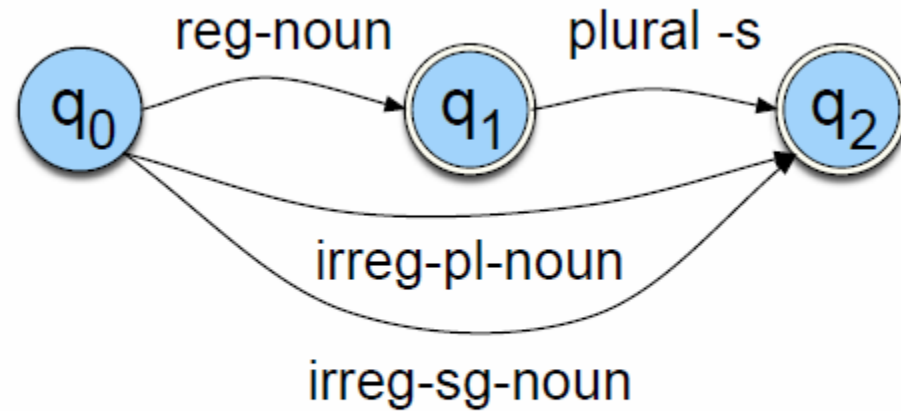
# Análisis (Parsing)

- En lugar de analizar, ¿por qué no mejor almacenar todas las formas posibles de todas las palabras de una lengua?

# Lexicón (Lexicon)

- Lista de palabras.
- Lista de bases y afijos.

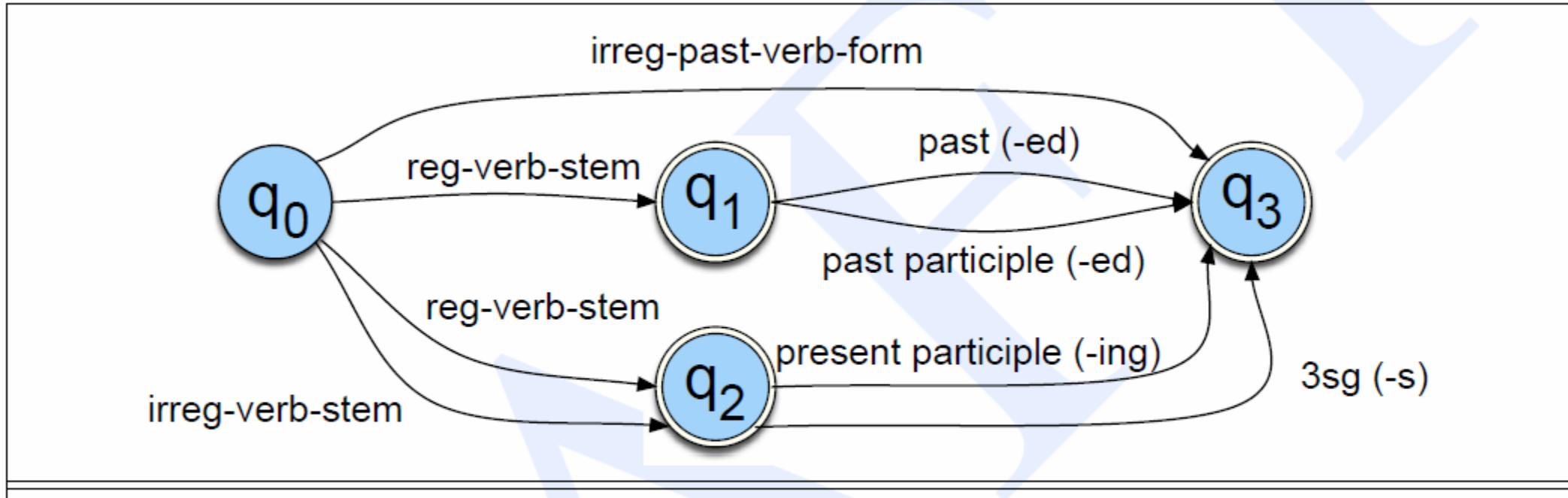
# Autómata morfológico



**Figure 3.3** A finite-state automaton for English nominal inflection.

| <b>reg-noun</b>        | <b>irreg-pl-noun</b>   | <b>irreg-sg-noun</b>    | <b>plural</b> |
|------------------------|------------------------|-------------------------|---------------|
| fox<br>cat<br>aardvark | geese<br>sheep<br>mice | goose<br>sheep<br>mouse | -s            |

# Autómata morfológico



| <b>reg-verb-stem</b>           | <b>irreg-verb-stem</b> | <b>irreg-past-verb</b>         | <b>past</b> | <b>past-part</b> | <b>pres-part</b> | <b>3sg</b> |
|--------------------------------|------------------------|--------------------------------|-------------|------------------|------------------|------------|
| walk<br>fry<br>talk<br>impeach | cut<br>speak<br>sing   | caught<br>ate<br>eaten<br>sang | -ed         | -ed              | -ing             | -s         |

# Autómata morfológico

big, bigger, biggest,

happy, happier, happiest, happily

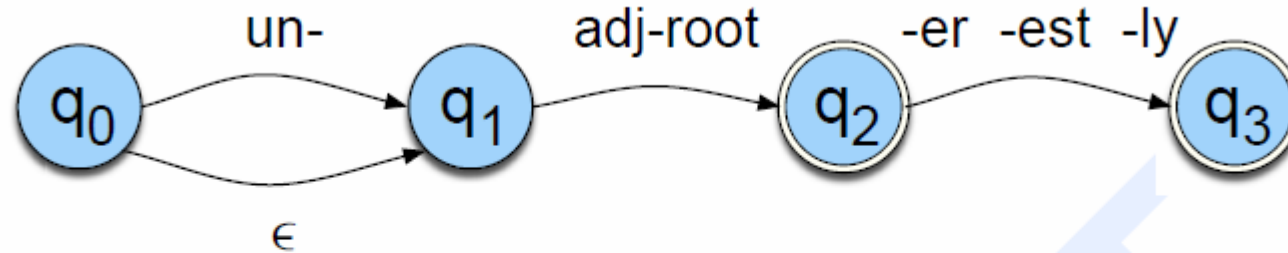
unhappy, unhappier, unhappiest, unhappily

clear, clearer, clearest, clearly, unclear, unclearly

cool, cooler, coolest, coolly

red, redder, reddest

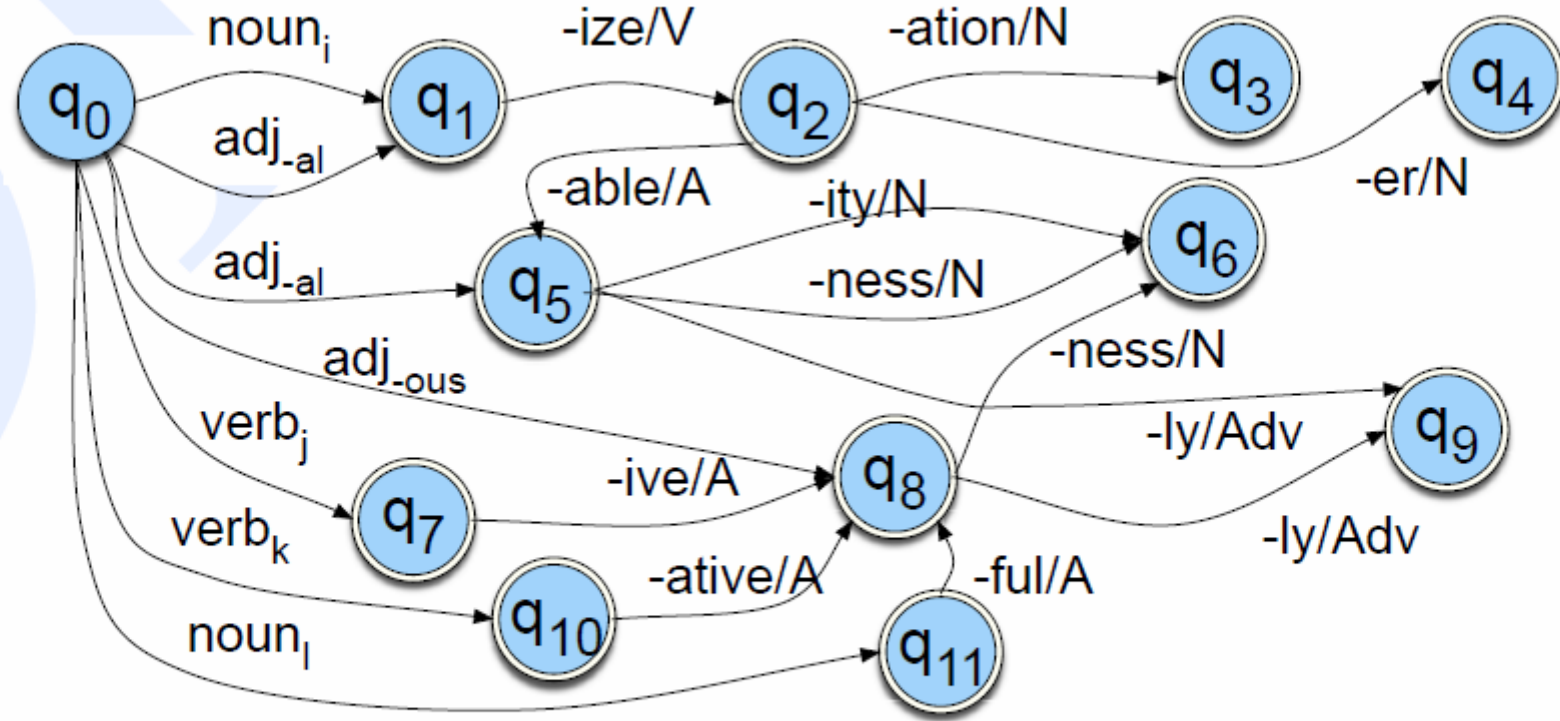
real, unreal, really



**Figure 3.5**  
posal #1.

An FSA for a fragment of English adjective morphology: Antworth's Pro-

# Autómata morfológico



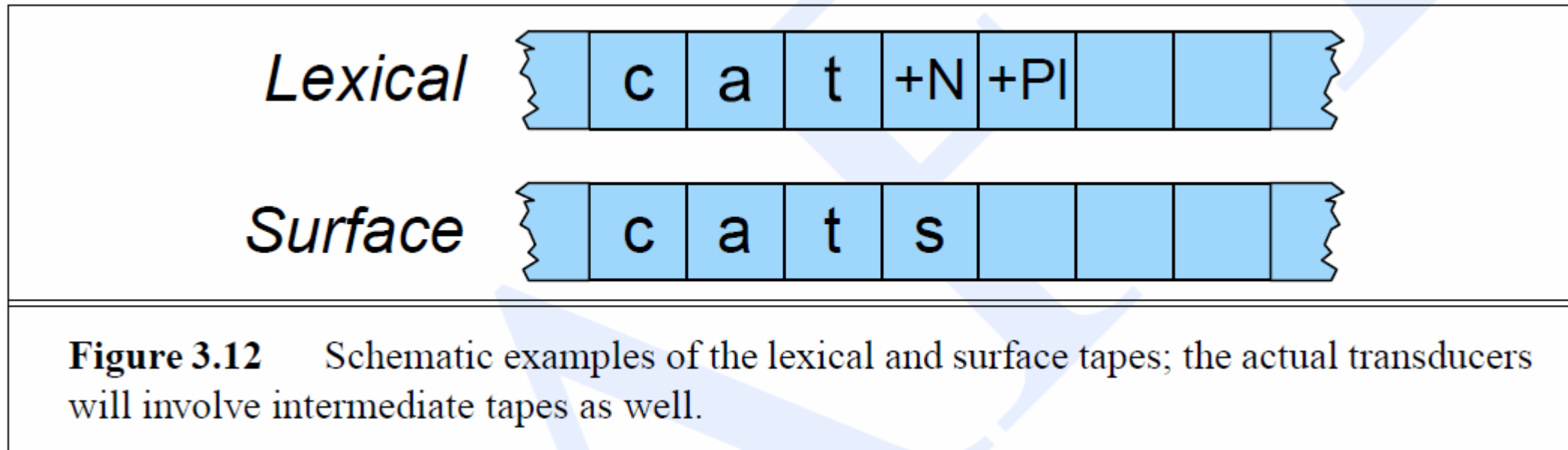
**Figure 3.6** An FSA for another fragment of English derivational morphology.



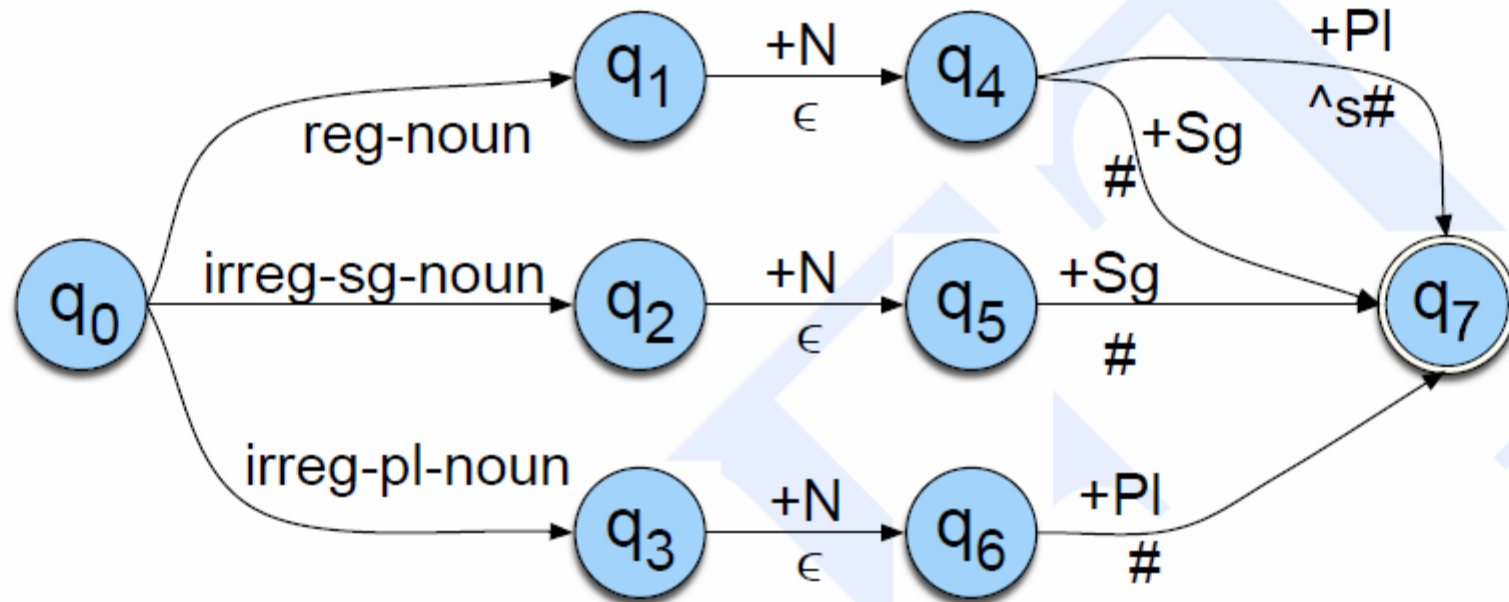
# Transductores de estados finitos (FST)

- Tipo de autómeta de estados finitos que hace correspondencia entre dos conjuntos de símbolos.
- Aceptador: recibe pares de cadenas y los acepta.
- Generador: genera pares de cadenas.
- Traductor: recibe una cadena y genera otra de salida.

# Transductores de estados finitos (FST)



# Transductores de estados finitos (FST)

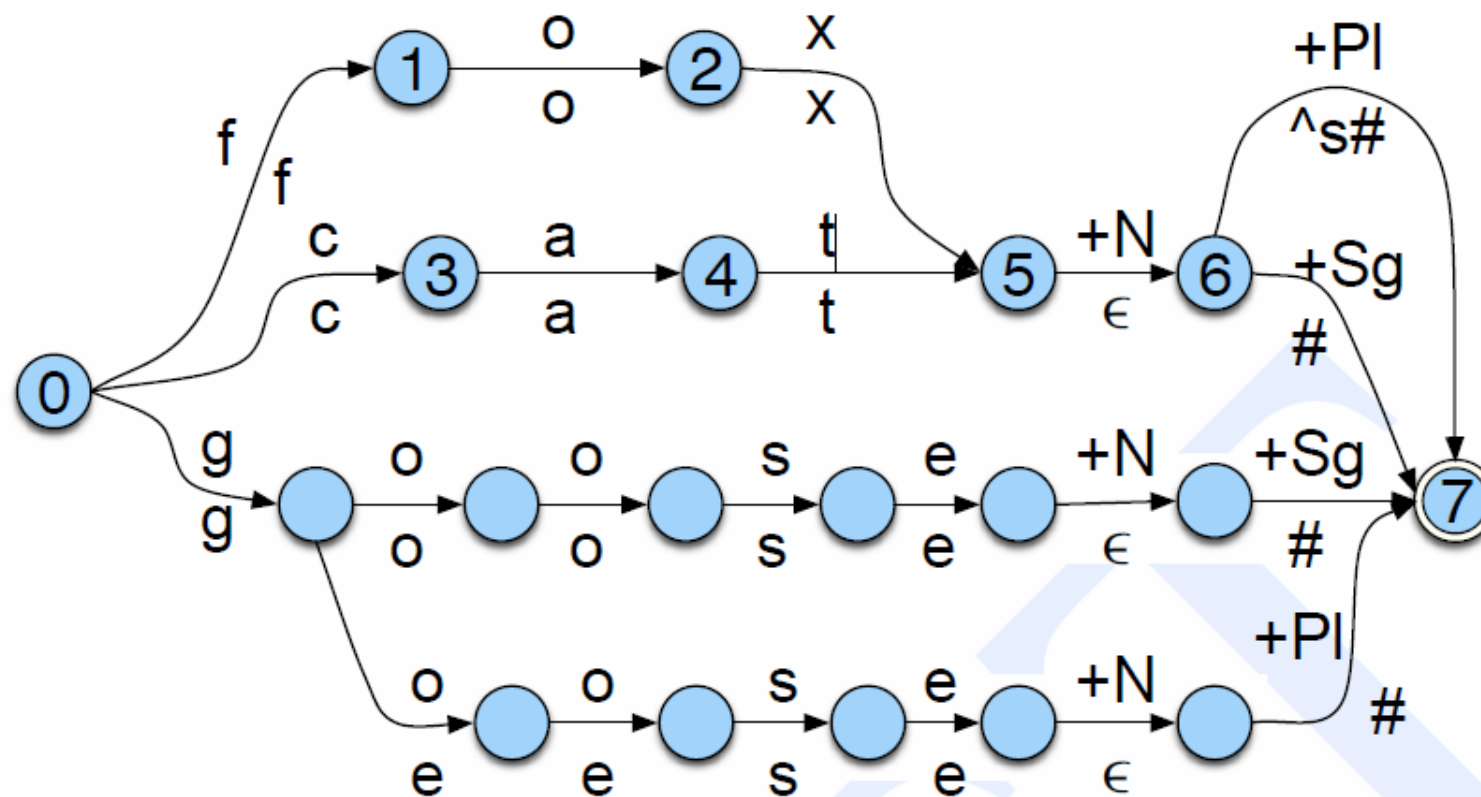


**Figure 3.13** A schematic transducer for English nominal number inflection  $T_{num}$ . The symbols above each arc represent elements of the morphological parse in the lexical tape; the symbols below each arc represent the surface tape (or the intermediate tape, to be described later), using the morpheme-boundary symbol  $\wedge$  and word-boundary marker  $\#$ . The labels on the arcs leaving  $q_0$  are schematic, and need to be expanded by individual words in the lexicon.

# Transductores de estados finitos (FST)

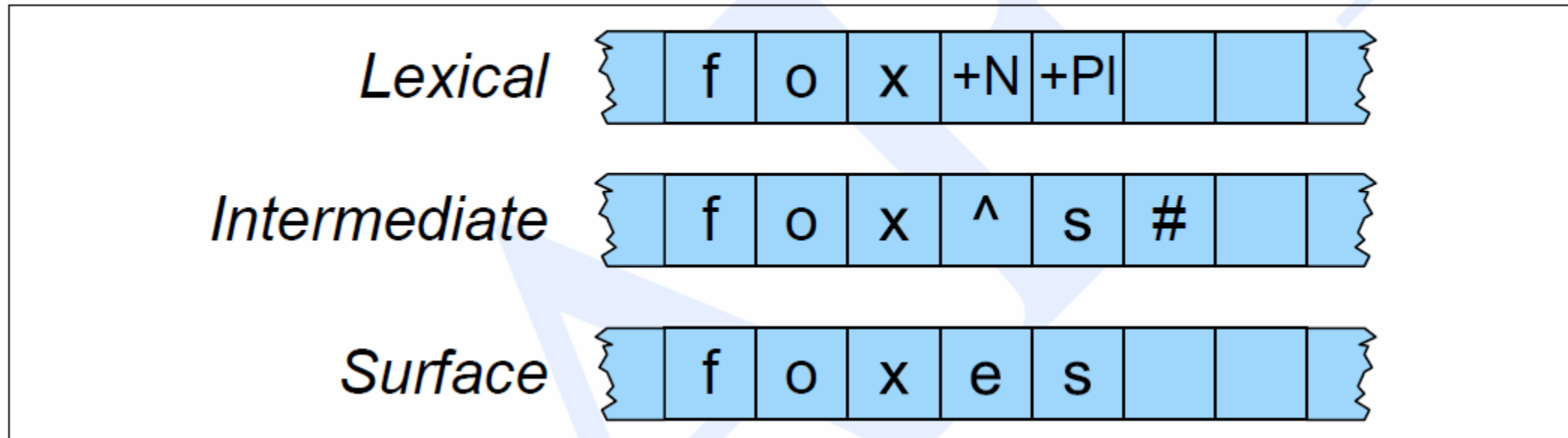
| <b>reg-noun</b> | <b>irreg-pl-noun</b> | <b>irreg-sg-noun</b> |
|-----------------|----------------------|----------------------|
| fox             | g o:e o:e s e        | goose                |
| cat             | sheep                | sheep                |
| aardvark        | m o:i u:ε s:c e      | mouse                |

# Transductores de estados finitos (FST)

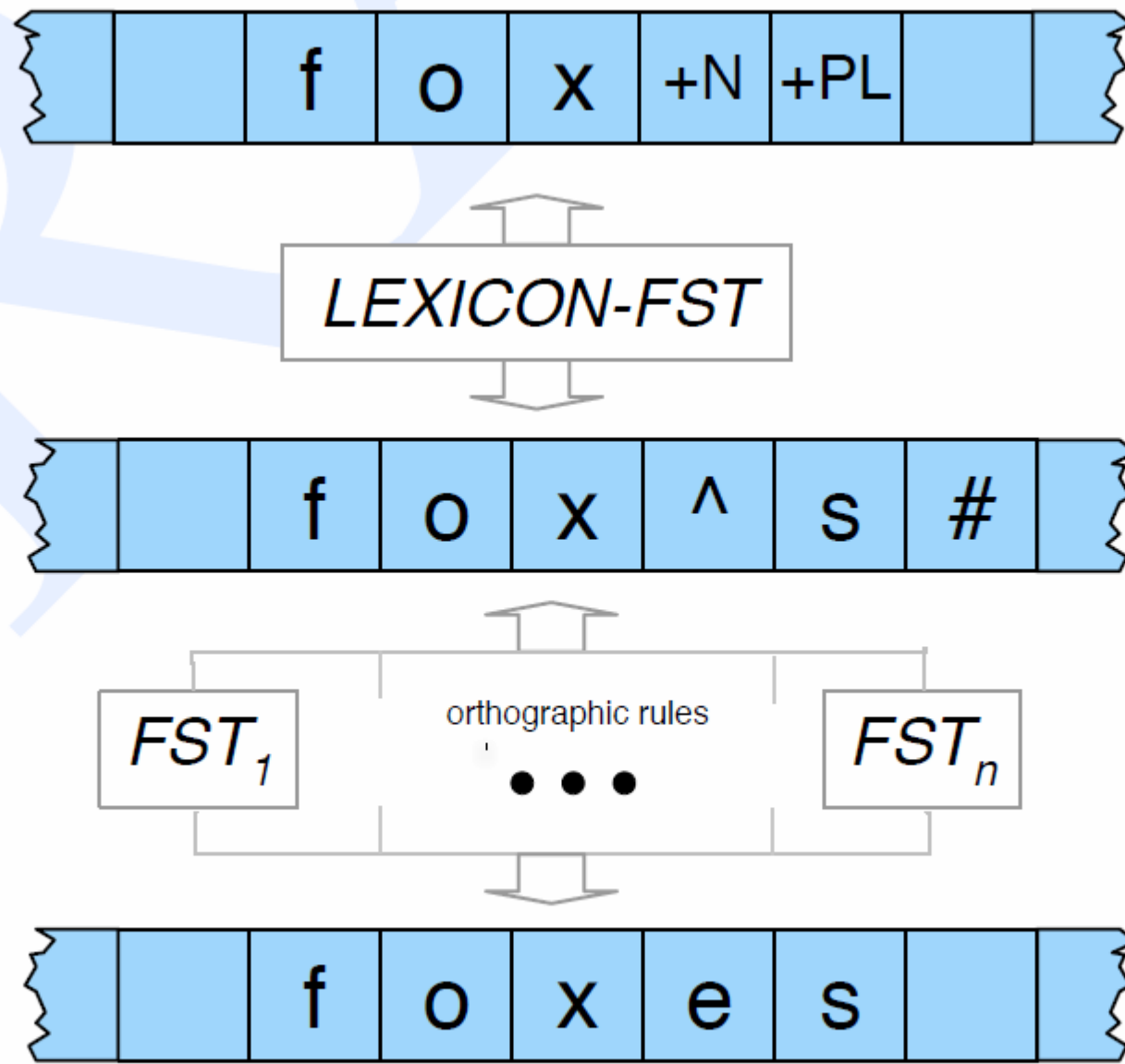


**Figure 3.14** A fleshed-out English nominal inflection FST  $T_{lex}$ , expanded from  $T_{num}$  by replacing the three arcs with individual word stems (only a few sample word stems are shown).

# Transductores de estados finitos (FST)



**Figure 3.16** An example of the lexical, intermediate, and surface tapes. Between each pair of tapes is a two-level transducer; the lexical transducer of Fig. 3.14 between the lexical and intermediate levels, and the E-insertion spelling rule between the intermediate and surface levels. The E-insertion spelling rule inserts an *e* on the surface tape when the intermediate tape has a morpheme boundary  $\wedge$  followed by the morpheme *-s*.



**Figure 3.19** Generating or parsing with FST lexicon and rules

# Descubrir morfología

- Sucesor frecuente (Harris, 1955)

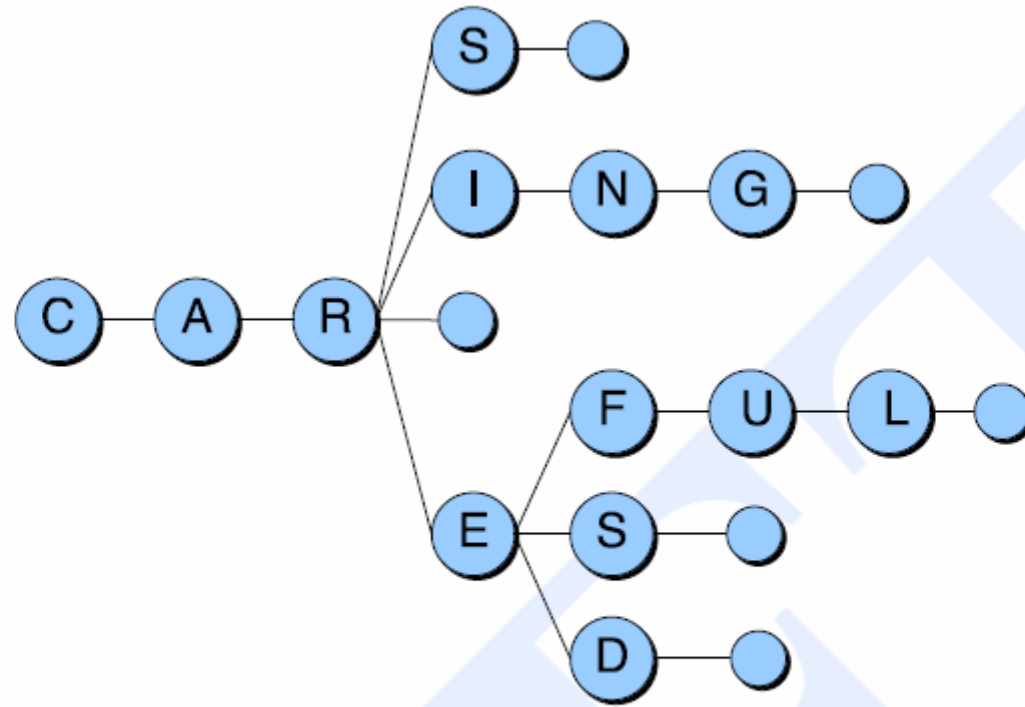
|   |  |
|---|--|
| gover~n,<br>gover~ned,<br>gover~ning,<br>gover~nment,<br>gover~nor,<br>gover~ns | govern,<br>govern~ed,<br>govern~ing,<br>govern~ment,<br>govern~or,<br>govern~s |
| Sucesor frecuente<br>(gover) = 1 (n).   | Sucesor frecuente<br>(govern) = 6 (e, i, m, o,<br>s, espacio/signo).           |



# Descubrir morfología

- Trie: estructura de árbol que guarda cadenas. Cada cadena es un camino desde la raíz hasta la última hoja.

# Descubrir morfología



**Figure 11.17** Example of a letter trie. A Harris style algorithm would insert morpheme boundaries after *car* and *care*. After Schone and Jurafsky (2000).

# Ejercicio

- xooko'
- xook
- xooki'
- xooke'
- xoot
- xooko'obo'
- xooknak
- xookna'ako'on
- xooka'
- xooko'obe'

# TAREA

- wiliko'
- wilike'ex
- wila'aj
- wilej
- wilike'
- wilik
- wilibe'
- wilaj
- wil
- wili'
- wile'
- wile'ex
- wilmaj
- wilajo'
- wilmaji'
- wilmi'
- wili'i
- wilo'ob
- wilike'exe'
- wilaje'
- wilme'exi'
- wilmaje'ex
- wilmaje'exi'
- wilike'exo'
- wila'e'

Fin