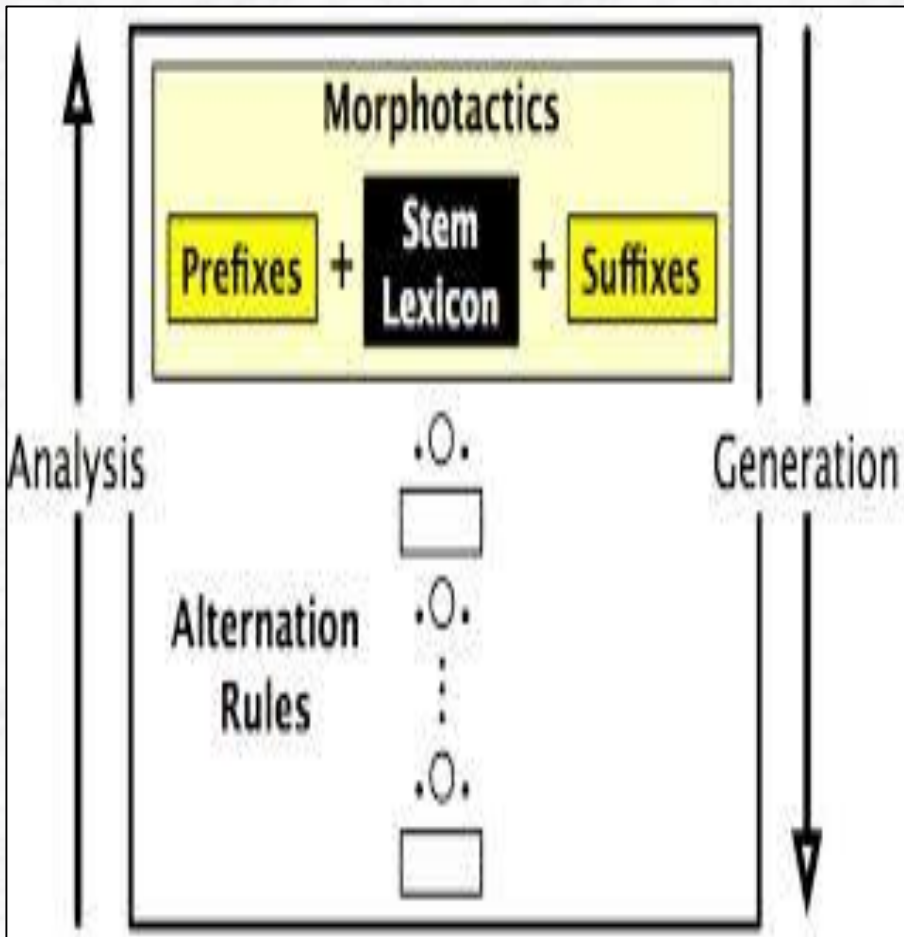


Morphological analysis



Words and Transducers

Instructor : Dr. Hanaa Bayomi Ali
Mail : h.mobarz @ fci-cu.edu.eg

Stages of language processing

1- Phonetics and phonology

Speech sound

2- Lexical Analysis

Dividing the whole chunk of txt into paragraphs, sentences, and words

3- Morphology & Lexicon

Words & their forms

4- Syntactic Analysis

Structure of sentences

5- Semantic Analysis

Meaning of words & sentences

6- Pragmatics

Meaning in context & for a purpose

7- Discourse

Connected sentence processing in a larger body of text

Words

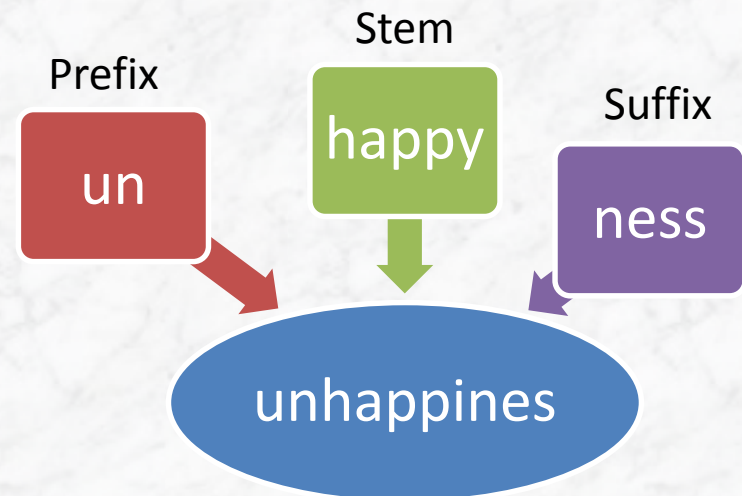
- In formal languages, **words** are arbitrary strings
- In natural languages, words are made up of meaningful subunits called **morphemes**
- morpheme = "minimal meaning-bearing unit in a language"
- Morphemes may be

Stems: the main morpheme of the word.

Affixes: convey the word's role, number, gender, etc.

"Unhappiness" Contains 3 morphemes

- 1- "un" means "not"
- 2- "ness" means "being in a state or condition"
- 3- "Happy" is a **free morpheme** (as a "word" in its own right)



Morphemes definitions

- ***Affix***

- A bound **morpheme** that is joined before, after, or within a root or stem

- **Suffixes:** follow the root/stem “eat, eats”

- **Prefixes:** precede the root/stem “happy, unhappy”

- **Infixes:** inserted into the root/stem “s(plural), mothers-in-law”

- Affixes maybe derivational affixes, inflectional affixes or both

- ***Root***

- The portion of the word that:

- is common to a set of derived or inflected forms, if any, when all affixes are removed

- is not further analyzable into meaningful elements

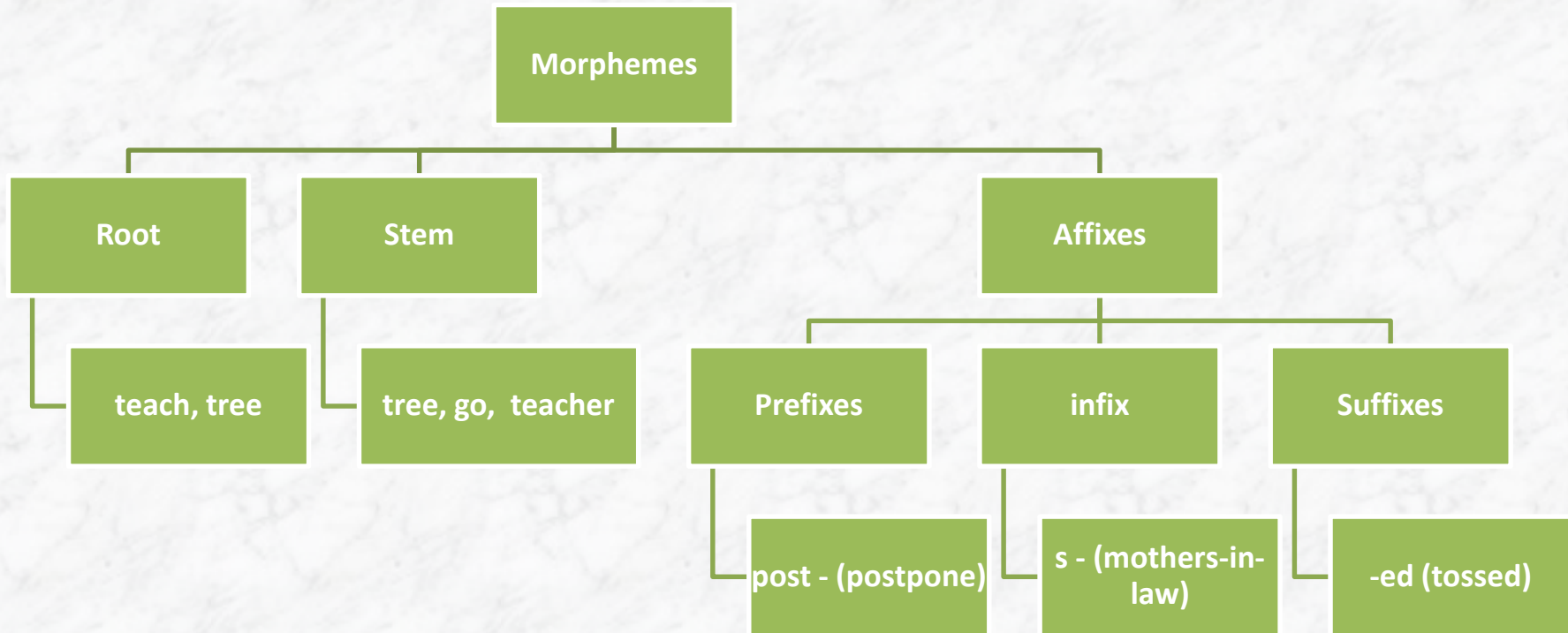
- carries the principle portion of meaning of the words

- ***Stem***

- The root of a word, together with any derivational affixes

- the morpheme that forms the central meaning unit in a word

Morphemes definitions



Free and Bound Morpheme

Analysis at a morphological level is concerned with structural elements of meaning called **morphemes**. Morphemes are classified into two types:

1- Free Morphemes: girl, boy, mother, etc.

- These are words with a complete meaning, so they can stand alone as an independent word in a sentence.

2- Bound Morphemes: These are lexical items incorporated into a word as a dependent part. They cannot stand alone, but must be connected to another morpheme.

- Bound morphemes operates in the connection processes by means of **derivation**, **inflection**, and **compounding**.

EX. Bound morpheme: book**s**

- Books = +morpheme +morpheme
- Books= +book +s

Morphological Analysis

Why do we need to do Morphological Analysis?

- determine meanings of individual word.
- To break word down into component morphemes and build a structured representation.
- Important for many tasks:
 - machine translation
 - information retrieval
 - Parsing
 - text simplification
 - Spelling correction

Morphological Analysis

- How words are constructed out of basic meaning unit
 - **Inflection** (number, person, tense , Gender ...): same core meaning (don't change syntactic category) (**Stemmer**)
 - love + past tense → loved
 - **Derivation** (adjective-adverb, noun-verb): change meaning
 - un + friend + ly → unfriendly
 - **Compounding** (separate words or single word)
 - black + bird → blackbird
 - Foot + ball → football
- Arabic words is highly inflected

Morphological Analysis

- ***Inflection*** is the form variation of a word under certain **grammatical** conditions, these conditions consist notably of the number, gender, conjugation, or tense

stem + gram. morpheme → same class

Ex: help + ed → helped

- ***Derivation*** combines affixes to an existing root or stem to form a **new word**. ***Derivation is more irregular and complex than inflection.*** It often results in a change in the **part of speech** for the derived word.

stem + gram. morpheme → different class

Ex: civilization

Can you divide these words?

- teaches
- Happiness
- Unbelievable
- Teacher
- Monster
- Rattlesnake

Can you divide these words?

• Teaches	teach-s	Inflection
• Happiness	happi-ness	<i>Derivation</i>
• Unbelievable	un-believ-able	<i>Derivation</i>
• Teacher	teach-er	<i>Derivation</i>
• Monster	monster	Inflection
• Rattlesnake	rattle-snake	Compound

Distinguishing inflection from derivation

main criteria:

1- Category change: Inflection does not change grammatical category; derivation sometimes does (thereby creating new words).

Form 1	Form 2	Category change	Der/Inf
car	cars	No	(s) Inflection
Play	plays	no	(s) Inflection
Formal(adj)	Formalize (V)	yes	(ize) Derivation
Read (V)	Readable(adj)	yes	(able) Derivation

Distinguishing inflection from derivation

Three main criteria:

2- Order: Derivational affixes must combine with the base before an inflectional affix does (root - aff_{der} - aff_{inf} " teach_{root} - er_{der} - s_{inf}).

- Inflectional affixes always have a regular meaning.

- the plural 's *in word-forms like* bicycles, dogs, shoes, tins, trees, *and so on*

'more than one'

- Derivational affixes may have irregular meaning.

- write :: writer (one who writes)
- paint :: painter (one who paints)
- cut :: cutter? (an instrument used to cut)
- Mix :: Mixer (an instrument used to mix)

More on Inflection

Noun inflectional suffixes	<ul style="list-style-type: none">• Plural marker -s• Possessive marker 's
Verb inflectional suffixes	<ul style="list-style-type: none">• Third person present singular marker -s• Past tense marker -ed• Progressive marker -ing• Past participle markers -en or -ed
Adjective inflectional suffixes	<ul style="list-style-type: none">• Comparative marker -er• Superlative marker -est

Morphological Parsing

- **Parsing** means taking an input and producing some sort of structure for it.
- **Morphological parsing** means breaking down a word form into its constituent morphemes.
 - e.g. Colourful \Rightarrow colour+ ful
- Mapping of a word form to its baseform is called **stemming**.
 - e.g. Teacher \Rightarrow Teach

Finite-State Morphological Parsing

- In order to build a parser we need the following:
 - a **lexicon** containing the stems and affixes,
 - **morphotactics**, i.e. the model of morpheme ordering, e.g. *un + friend + ly* instead of *ly + friend + un*
 - a set of **rules** (orthographic, etc.), i.e. the model of changes that occur in a word, usually when two morphemes combine, e.g. *city + s* \Rightarrow *cities*.

Finite-State Automaton for Inflection of English Verbs

Verbal Inflection

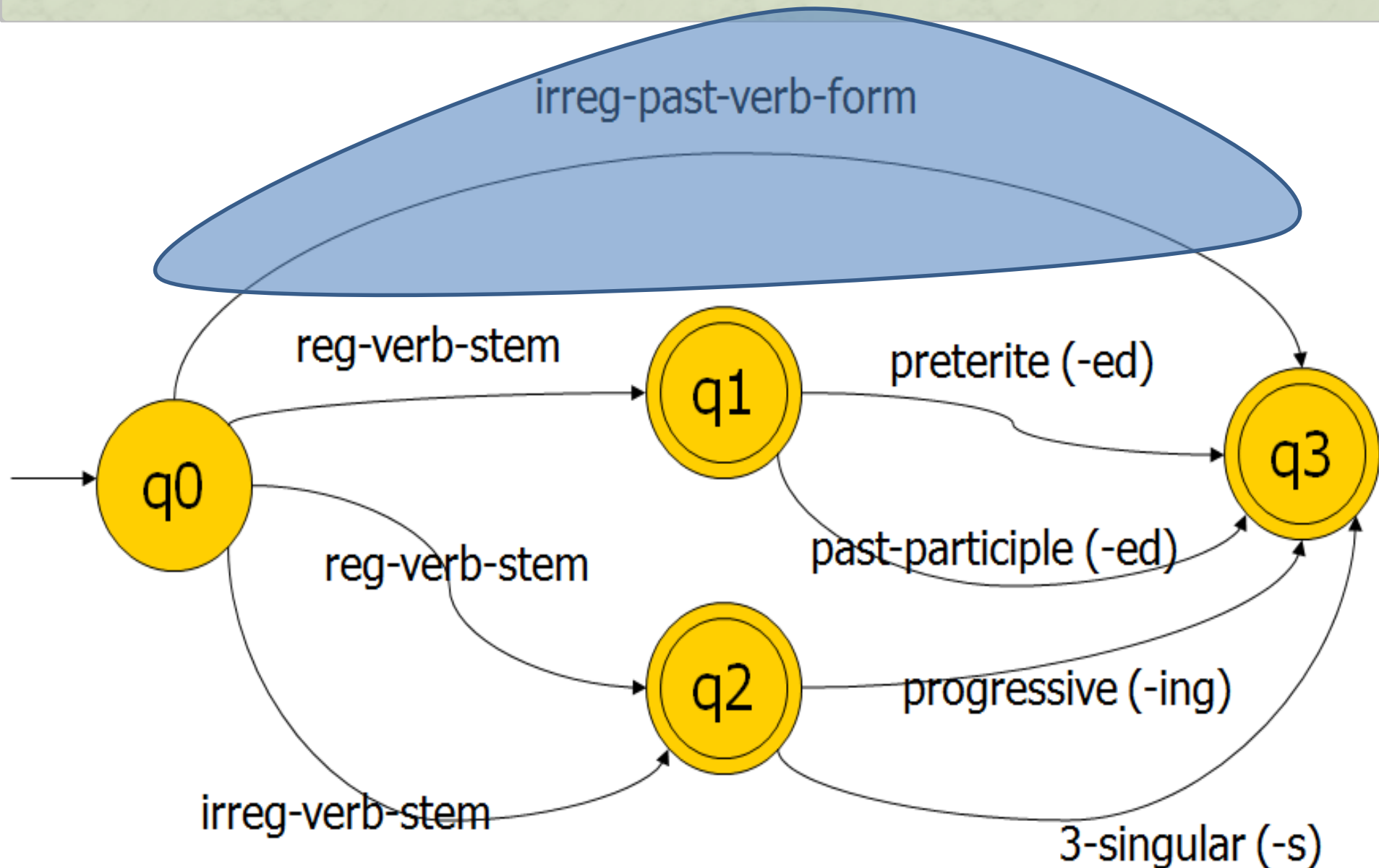
Morphological Form Classes	Regularly Inflected Verbs				Irregularly Inflected Verbs		
Stem	Jump	Parse	Fry	Sob	Eat	Bring	Cut
-s form	Jumps	Parses	Fries	Sobs	Eats	Brings	Cuts
-ing participle	Jumping	Parsing	Frying	Sobbing	Eating	Bringing	Cutting
Past form	Jumped	Parsed	Fried	Sobbed	Ate	Brought	Cut
–ed participle	Jumped	Parsed	Fried	Sobbed	Eaten	Brought	Cut

Finite-State Automaton for Inflection of English Verbs

Verbal Inflection

Morphological Form Classes	Regularly Inflected Verbs				Irregularly Inflected Verbs		
Stem	Jump	Parse	Fry	Sob	Eat	Bring	Cut
-s form	Jumps	Parses	Fries	Sobs	Eats	Brings	Cuts
-ing participle	Jumping	Parsing	Frying	Sobbing	Eating	Bringing	Cutting
Past form	Jumped	Parsed	Fried	Sobbed	Ate	Brought	Cut
–ed participle	Jumped	Parsed	Fried	Sobbed	Eaten	Brought	Cut

Finite-State Automaton for Inflection of English Verbs

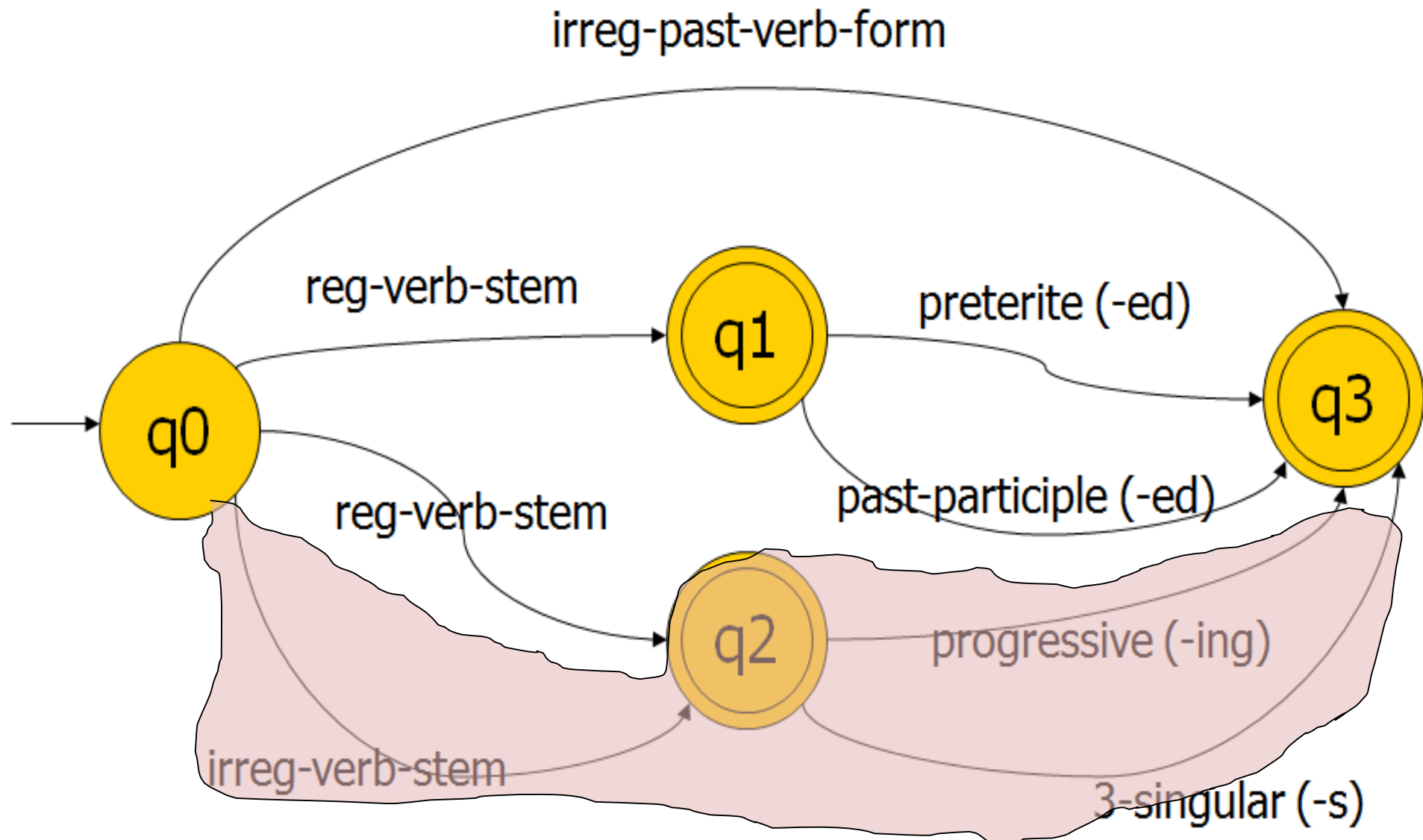


Finite-State Automaton for Inflection of English Verbs

Verbal Inflection

Morphological Form Classes	Regularly Inflected Verbs				Irregularly Inflected Verbs		
Stem	Jump	Parse	Fry	Sob	Eat	Bring	Cut
-s form	Jumps	Parses	Fries	Sobs	Eats	Brings	Cuts
-ing participle	Jumping	Parsing	Frying	Sobbing	Eating	Bringing	Cutting
Past form	Jumped	Parsed	Fried	Sobbed	Ate	Brought	Cut
–ed participle	Jumped	Parsed	Fried	Sobbed	Eaten	Brought	Cut

Finite-State Automaton for Inflection of English Verbs

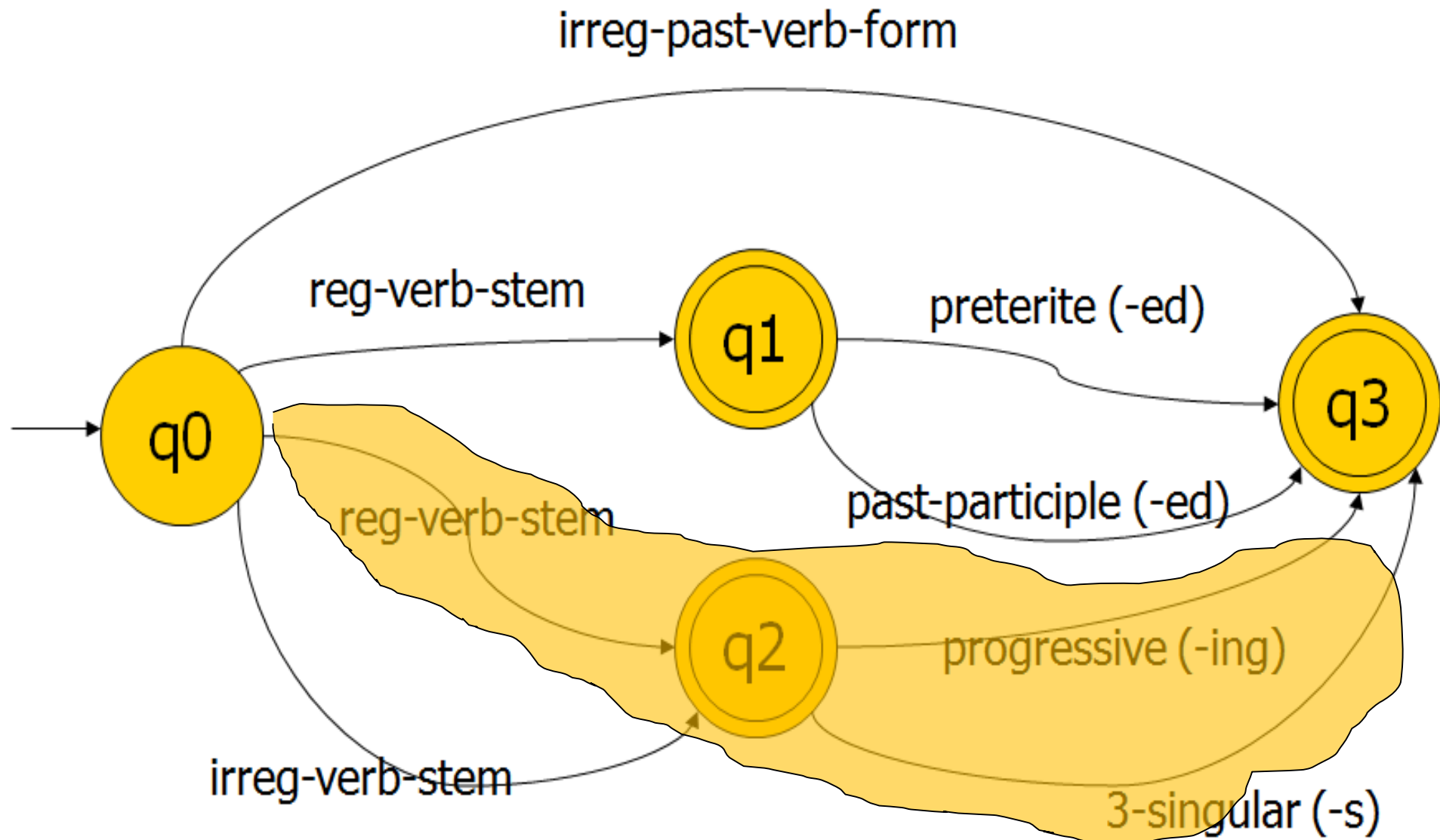


Finite-State Automaton for Inflection of English Verbs

Verbal Inflection

Morphological Form Classes	Regularly Inflected Verbs				Irregularly Inflected Verbs		
Stem	Jump	Parse	Fry	Sob	Eat	Bring	Cut
-s form	Jumps	Parses	Fries	Sobs	Eats	Brings	Cuts
-ing participle	Jumping	Parsing	Frying	Sobbing	Eating	Bringing	Cutting
Past form	Jumped	Parsed	Fried	Sobbed	Ate	Brought	Cut
–ed participle	Jumped	Parsed	Fried	Sobbed	Eaten	Brought	Cut

Finite-State Automaton for Inflection of English Verbs

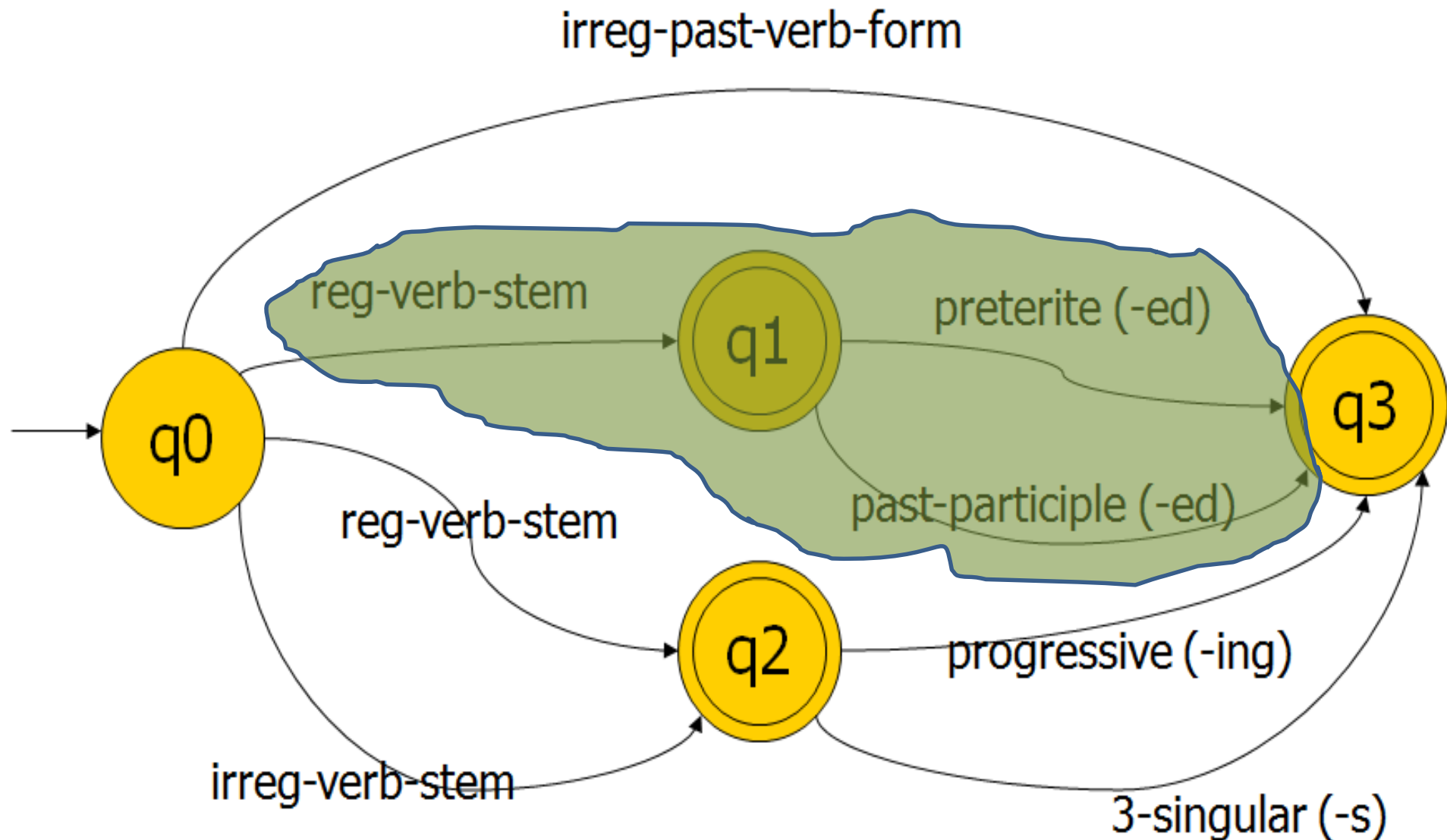


Finite-State Automaton for Inflection of English Verbs

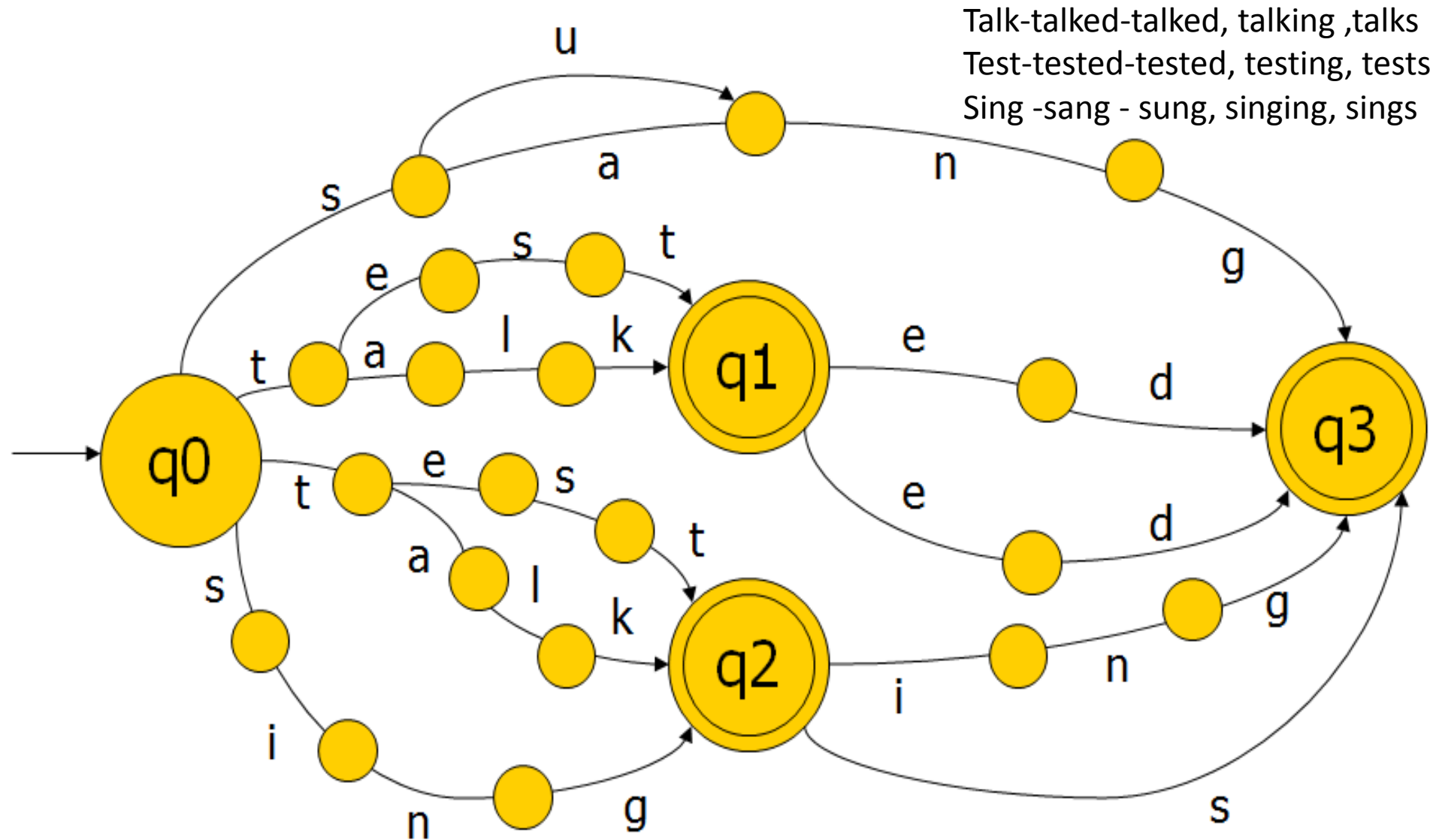
Verbal Inflection

Morphological Form Classes	Regularly Inflected Verbs				Irregularly Inflected Verbs		
Stem	Jump	Parse	Fry	Sob	Eat	Bring	Cut
-s form	Jumps	Parses	Fries	Sobs	Eats	Brings	Cuts
-ing participle	Jumping	Parsing	Frying	Sobbing	Eating	Bringing	Cutting
Past form	Jumped	Parsed	Fried	Sobbed	Ate	Brought	Cut
–ed participle	Jumped	Parsed	Fried	Sobbed	Eaten	Brought	Cut

Finite-State Automaton for Inflection of English Verbs

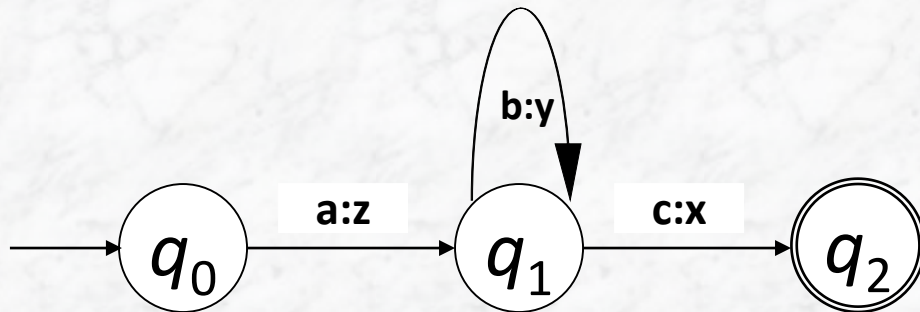


Finite-State Automaton for Inflection of the Verbs 'talk', 'test' and 'sing'



Finite State Transducer (FST)

- Transduce: to convert from a form to another
- The input symbol is transduced into the output symbol as a transition occurs on the arc
- FST accepts and generates string (to define **internal structure** of the word <Analysis>)



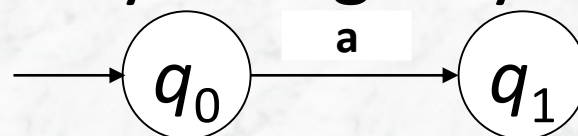
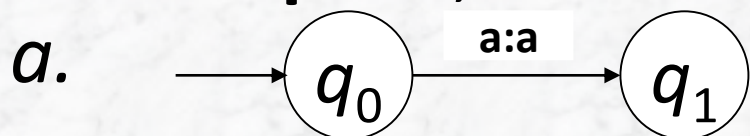
Accepts
abbbc
and generates
zyyyx

Finite State Transducer (Cont.)

- Finite-state transducers have a formal definition, which is similar to that of finite state automata. A FST consists of five components $(Q, \Sigma, q_0, F, \delta)$ where:
 - Q is a finite set of states
 - Σ is a finite set of symbol or character pairs $i : o$, where i is a symbol of the input alphabet and o of the output alphabet
 - q_0 is the start state, $q_0 \in Q$
 - F is the set of final states, $F \subseteq Q$
 - δ is a relation from states and symbols to states.
$$\delta: Q \times \Sigma \times Q$$

Finite-State Transducer

- A transducer maps between one set of symbols and another; a finite state transducer does this via a finite automaton.
- Where an FSA accepts a language stated over a finite alphabet of single symbols, e.g. $\Sigma = \{a, b, c, \dots\}$, an FST accepts a language stated over **pairs of symbols**, e.g. $\Sigma = \{a:a, b:b, a:c, a:\varepsilon, \varepsilon:\varepsilon, \dots\}$
- In two-level morphology, we call pairs like $a:a$ **default pairs**, and refer to them by a single symbol a .



- An FST can be seen as a **recognizer**, **generator**, **translator** or a **set relator**.

Two-Level Morphology

- Two-level morphology represents a word as a correspondence between :-
 - 1- a **lexical level**, which represents a simple **concatenation of morphemes** making up a word
 - 2- the **surface level**, which represents the **actual spelling** of the final word.

Lexical

	s	i	n	g	+V	+PROG				
--	---	---	---	---	----	-------	--	--	--	--

Surface

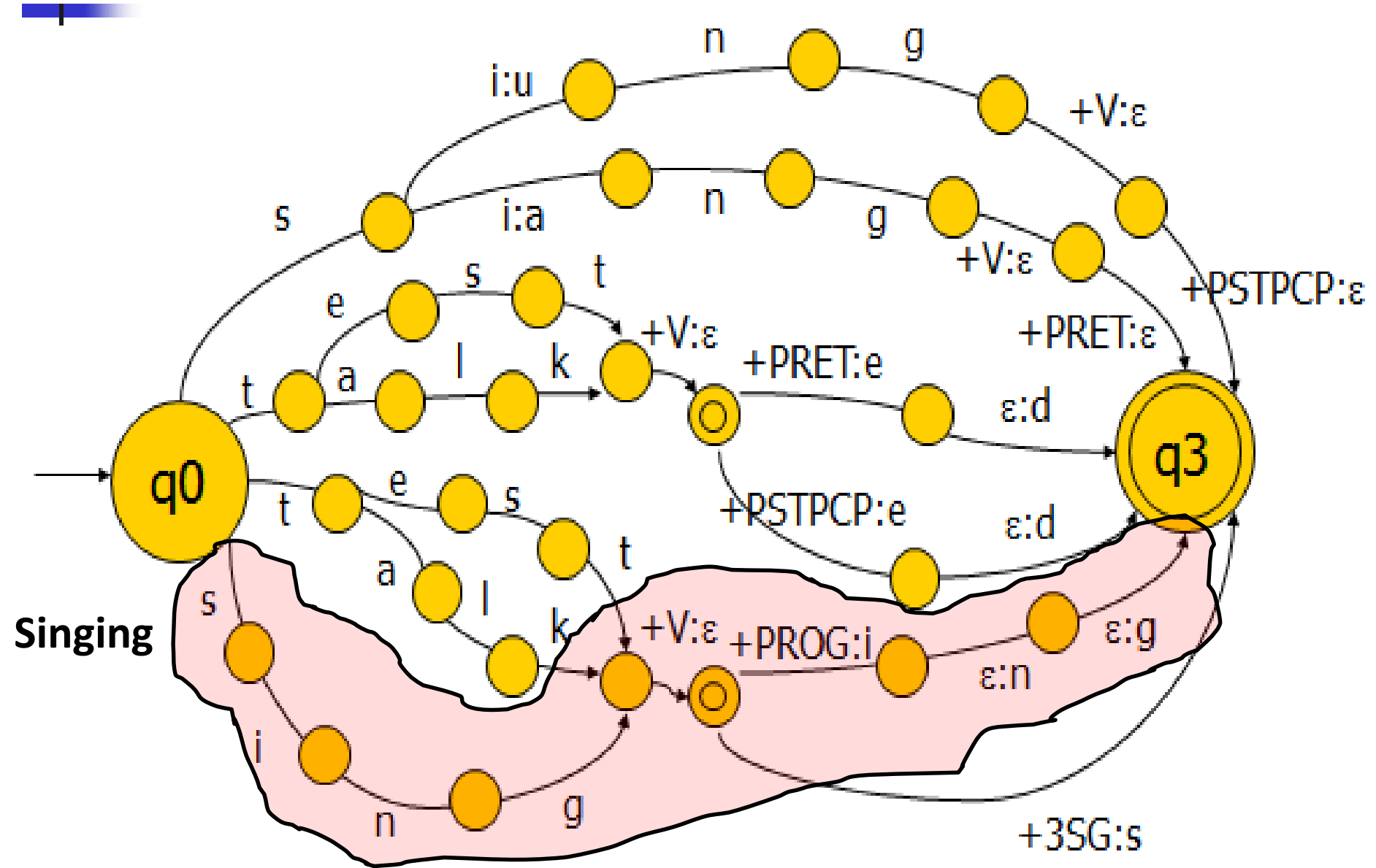
	s	i	n	g	i	n	g				
--	---	---	---	---	---	---	---	--	--	--	--

Examples

<i>Lexical form</i>	<i>Surface form</i>
sing +V +3SG	sings
test +V +PROG	testing
talk +V +PRET	talked
sing +V +PRET	sang
talk +V +PSTPCP	talked
sing +V +PSTPCP	sung

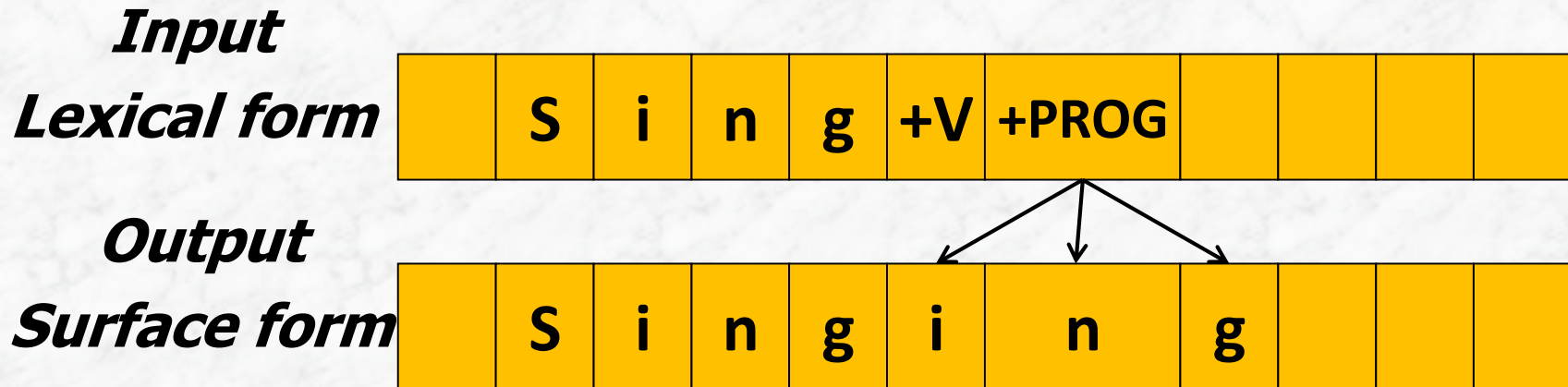
• ϵ epsilon  Do Nothing

Finite-State Transducer for Inflection of the Verbs 'talk', 'test' and 'sing'



Finite-State Transducer for Inflection of the Verbs 'talk', 'test' and 'sing'

Trace

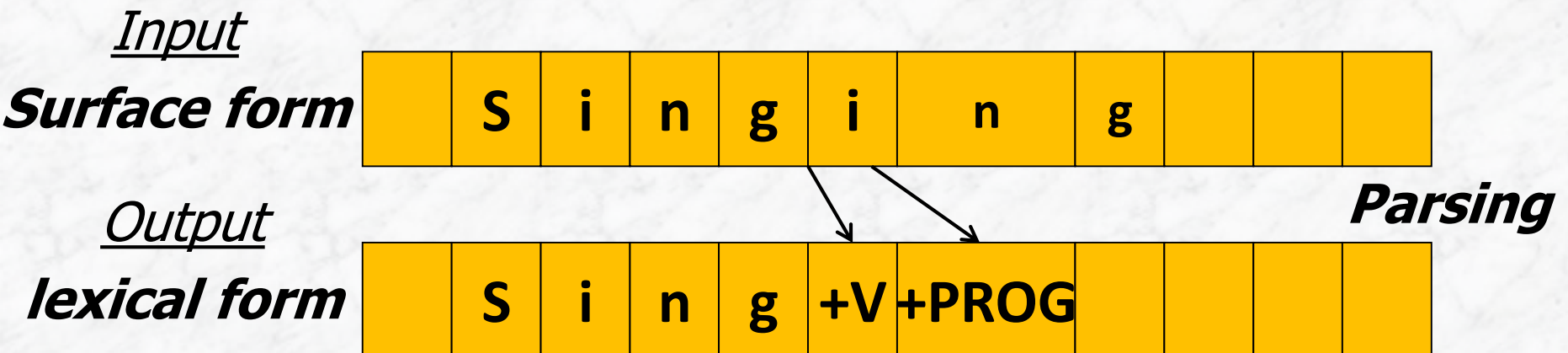
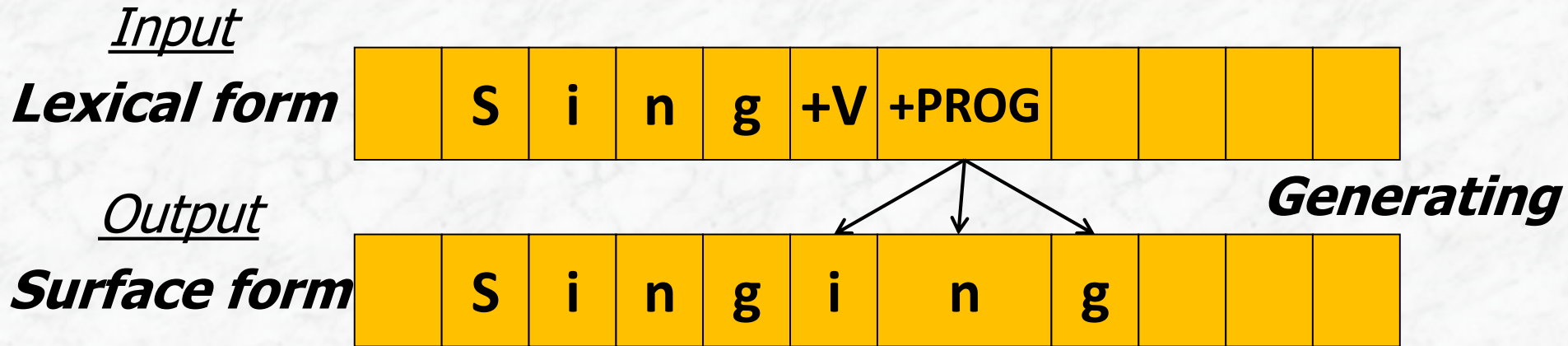


Useful FST Operations

- **Inversion:** Switch input and output labels.
 - e.g. $\Sigma(T) = \{a:b, c:d\} \Rightarrow \Sigma(\text{inv}(T)) = \{b:a, d:c\}$
 - This can be used to switch from generating words to parsing words.
- **Intersection:** Only sequences of pairs accepted by **both** transducer $T1$ **and** transducer $T2$ are accepted by transducer $T1 \wedge T2$.
- **Composition:** The output of transducer $T1$ serves as input to $T2$. This is marked as $T1 \circ T2$ or $T2(T1)$.

Finite-State Transducer for Inflection of the Verbs 'talk', 'test' and 'sing'

Trace Inverse



Finite-state transducers

A FST $T = L_{in} \times L_{out}$ defines a relation between two regular languages L_{in} and L_{out} :

$L_{in} = \{ \text{cat}, \text{cats}, \text{fox}, \text{foxes}, \dots \}$



$L_{out} = \{ \text{cat+N+sg}, \text{cat+N+pl}, \text{fox+N+sg}, \text{fox+N+PL} \dots \}$

$T = \{ \langle \text{cat}, \text{cat+N+sg} \rangle, \langle \text{cats}, \text{cat+N+pl} \rangle, \langle \text{fox}, \text{fox+N+sg} \rangle, \langle \text{foxes}, \text{fox+N+pl} \rangle \}$

English spelling rules

English spelling (orthography) is funny:
The underlying morphemes (*plural-s, etc.*) *can have different* orthographic surface realizations (-s, -es)

Spelling changes at morpheme boundaries:

- E-insertion: fox +s = fox**e**s
- E-deletion: make +**ing** = **making**

Intermediate representations

English plural -s: **cat** \Rightarrow **cats** **dog** \Rightarrow **dogs**

but: **fox** \Rightarrow **foxes**, **bus** \Rightarrow **buses** , **buzz** \Rightarrow **buzzes**

We define an intermediate representation which **captures**
- morpheme boundaries (^) and word boundaries (#).

<i>Lexicon:</i>	cat+N+PL	fox+N+PL
<i>Intermediate representation:</i>	cat^s#	fox^s#
<i>Surface string:</i>	cats	foxes

- *Intermediate-to-Surface Spelling Rule:*
If plural 's' follows a morpheme ending in 'x','z' or 's',
insert 'e'.

Spelling Rules and FSTs

Name	Description of Rule	Example
Consonant doubling	1-letter consonant doubled before <i>-ing/-ed</i>	beg/begging
E deletion	Silent e dropped before <i>-ing</i> and <i>-ed</i>	make/making
E insertion	e added after <i>-s, -z, -x, -ch, -sh</i> before <i>-s</i>	watch/watches
Y replacement	<i>-y</i> changes to <i>-ie</i> before <i>-s</i> , and to <i>-i</i> before <i>-ed</i>	try/tries
K insertion	verbs ending with <i>vowel + -c</i> add <i>-k</i>	panic/panicked

Three levels

- Add an intermediate level between the lexical and surface levels

Lexical

		k	i	s	s	+V	+3\$G			
--	--	---	---	---	---	----	-------	--	--	--

Intermediate

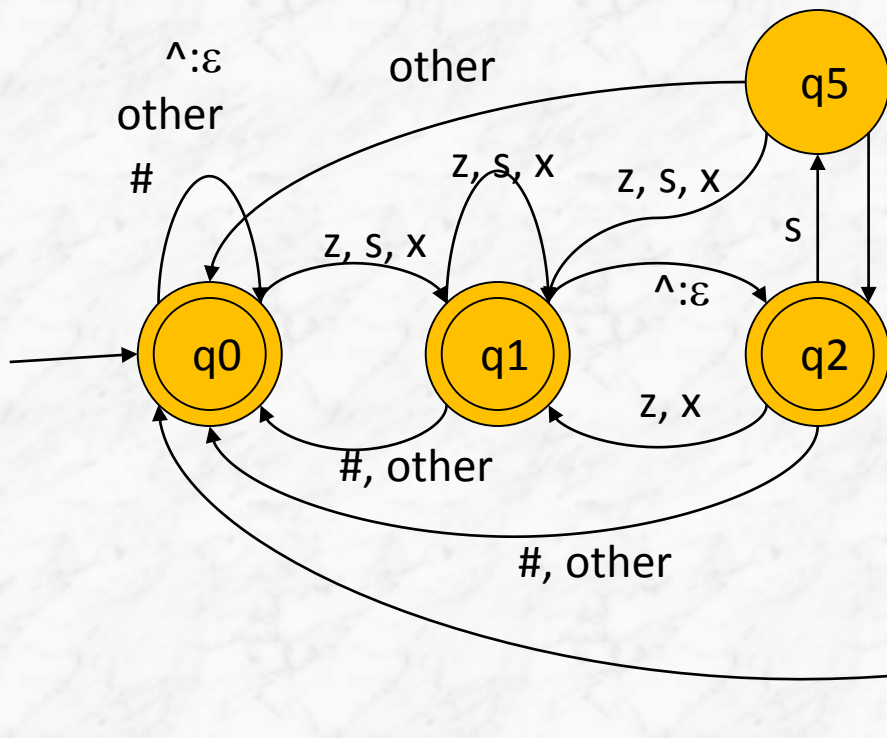
		k	i	s	s	^	s	#		
--	--	---	---	---	---	---	---	---	--	--

Surface

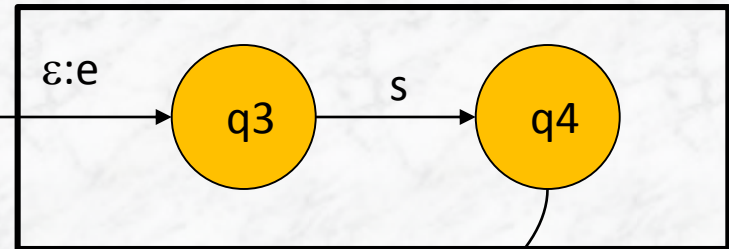
		k	i	s	s	e	s			
--	--	---	---	---	---	---	---	--	--	--

FST for the E-insertion Rule

fox^{^s}#

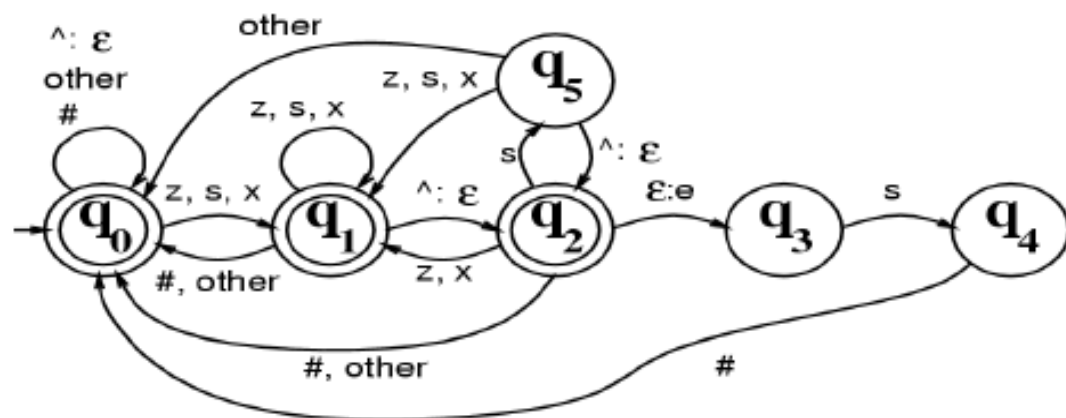


$$\varepsilon \rightarrow e / \left\{ \begin{array}{c} x \\ s \\ z \end{array} \right\}^{\wedge} _ s \#$$



The state-transition table for E-insertion rule

State \ Input	s : s	x : x	z : z	\wedge : ϵ	ϵ : e	#	other
q_0 :	1	1	1	0	-	0	0
q_1 :	1	1	1	2	-	0	0
q_2 :	5	1	1	0	3	0	0
q_3 :	4	-	-	-	-	-	-
q_4 :	-	-	-	-	-	0	-
q_5 :	1	1	1	2	-	-	0



The state-transition table for E-insertion rule

f	o	x	^	s	#
f	o	x	e	s	#

Intermediate Tape

Surface Tape

- Trace:
 - generating **foxes#** from **fox^s#**:
 $q_0\text{-f-}\rightarrow q_0\text{-o-}\rightarrow q_0\text{-x-}\rightarrow q_1\text{-^:}\epsilon\rightarrow q_2\text{-}\epsilon\text{:e-}\rightarrow q_3\text{-s-}\rightarrow q_4\text{-\#-}\rightarrow q_0$
 - generating **foxs#** from **fox^s#**:
 $q_0\text{-f-}\rightarrow q_0\text{-o-}\rightarrow q_0\text{-x-}\rightarrow q_1\text{-^:}\epsilon\rightarrow q_2\text{-s-}\rightarrow q_5\text{-\#-}\rightarrow \text{FAIL}$
 - generating **salt#** from **salt#**:
 $q_0\text{-s-}\rightarrow q_1\text{-a-}\rightarrow q_0\text{-l-}\rightarrow q_0\text{-t-}\rightarrow q_0\text{-\#-}\rightarrow q_0$
 - parsing **assess#**:
 $q_0\text{-a-}\rightarrow q_0\text{-s-}\rightarrow q_1\text{-s-}\rightarrow q_1\text{-^:}\epsilon\rightarrow q_2\text{-}\epsilon\text{:e-}\rightarrow q_3\text{-s-}\rightarrow q_4\text{-s-}\rightarrow \text{FAIL}$
 $q_0\text{-a-}\rightarrow q_0\text{-s-}\rightarrow q_1\text{-s-}\rightarrow q_1\text{-e-}\rightarrow q_0\text{-s-}\rightarrow q_1\text{-s-}\rightarrow q_1\text{-\#-}\rightarrow q_0$

State\Input	s:s	x:x	z:z	^:ε	ε:e	#	other
q_0 :	1	1	1	0	-	0	0
q_1 :	1	1	1	2	-	0	0
q_2 :	5	1	1	0	3	0	0
q_3 :	4	-	-	-	-	-	-
q_4 :	-	-	-	-	-	0	-
q_5 :	1	1	1	2	-	-	0