

# **NATURAL LANGUAGE PHONOLOGY & MORPHOLOGY**

CS662: Natural Language Processing

2015-01-29

*Elision:* Say \_\_\_\_\_. Now say \_\_\_\_\_ without saying \_\_\_\_\_.

*Matching:* The word \_\_\_\_\_ starts with the \_\_\_\_\_ sound.  
Which one of these words starts with the \_\_\_\_\_ sound like  
\_\_\_\_\_: \_\_\_\_\_ or \_\_\_\_\_?

*Isolation:* The word \_\_\_\_\_ has four sounds: \_\_\_\_\_, \_\_\_\_\_, \_\_\_\_\_,  
and \_\_\_\_\_. What is the second sound in \_\_\_\_\_?

*Blending:* \_\_\_\_\_ and \_\_\_\_\_ makes ...?

## Arabic *binyanim* (K-T-B):

<i>yu-CCiC-u</i>	<i>yukt<u>i</u>b<u>u</u></i>	‘he dictates’
<i>CaCaC-a</i>	<i>kat<u>a</u>ba</i>	‘he wrote’
<i>ya-CCuC-u</i>	<i>yakt<u>u</u>b<u>u</u></i>	‘he writes’
<i>CaaCiC</i>	<i>kaat<u>i</u>b</i>	‘writer’
<i>ma-CCaC</i>	<i>ma<u>k</u>ta<u>b</u></i>	‘letter’

## Turkish agglutination:

*çöp-lük-ler-imiz-de-ki-ler-den-mi-y-di*

‘Was it from those that were in our garbage cans?’

## English variable bracketing:

We could open the door—it’s [*un-lock*]able.

We can’t keep people out—the door’s *un*[*lock-able*].

# WHY MODEL IT?

- Most languages have richer morphology than English:
  - In Archi, *every single verb* has 1.5 million forms
  - In Turkish, there is *no known upper bound* on word length
  - And Zipf's Law applies within morphological paradigms (see slides from 2015-01-20)
- The identities and boundaries of morphological units are aggressively obscured by phonological processes

[Source: Kibrik 1998, Hankamer 1992]

# OUTLINE

- Phonology:
  - Phonemes, allophones, phones, and features
  - Phonological processes as FSTs
- Morphology:
  - Morphemes and allomorphs
  - Morphological typology
  - Morphological analysis

# PHONOLOGY

# CONSONANTS

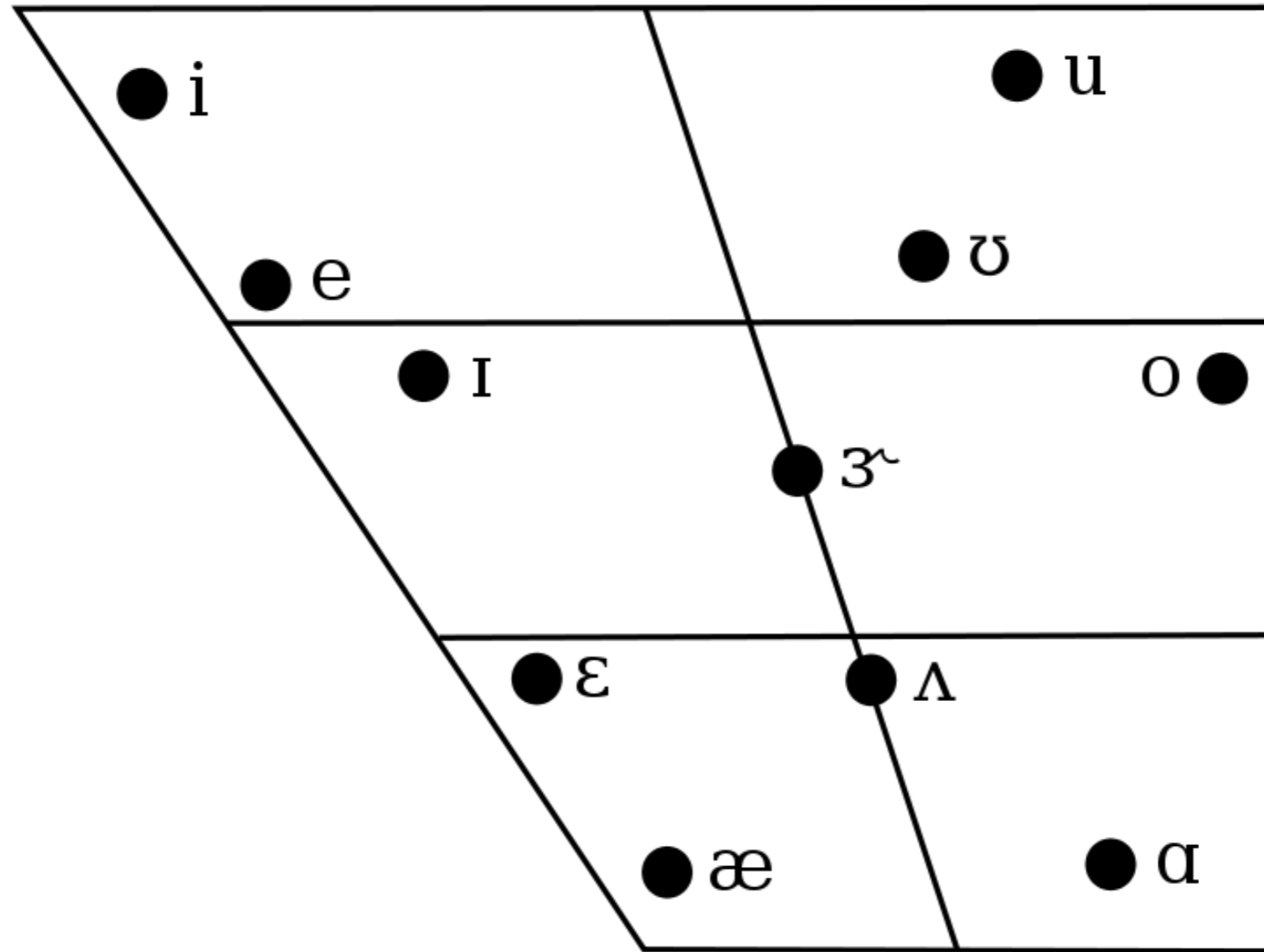
	Bilabial	Labiodental	Dental	Alveolar	Post-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Glottal
Plosive	p b			t d		ʈ ɖ	c ɟ	k ɡ	q ɢ		ʔ
Nasal	m	ɱ		n		ɳ	ɲ	ŋ	ɴ		
Trill	ʙ			ɾ					ʀ		
Tap or flap		ɸ		ɾ		ɽ					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	h ɦ
Lateral Fricative				ɬ ɮ							
Approximate	ʋ			ɹ		ɻ	j	ɰ			
Lateral Approximant				l		ɭ	ʎ	ʟ			

# IN ENGLISH...

	Bilabial	Labiodental	Dental	Alveolar	Post-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Glottal
Plosive	<u>p</u> <u>b</u>			<u>t</u> <u>d</u>		<u>ʈ</u> <u>ɖ</u>	<u>c</u> <u>ɟ</u>	<u>k</u> <u>g</u>	q ɢ		ʔ
Nasal	<u>m</u>	ɱ		<u>n</u>		ɳ	ɲ	ŋ	ɴ		
Trill	ʙ			ɾ					ʀ		
Tap or flap		ɸ		ɾ		ɽ					
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ʝ	x ɣ	χ ʁ	ħ ʕ	h ɦ
Lateral Fricative				ɬ ɮ							
Approximant	ʋ			ɹ		ɻ	j	ɰ			
Lateral Approximant				l		ɭ	ʎ	ʟ			



# VOWELS IN ENGLISH



# DEFINITIONS

*Phoneme*: a contrastive unit of sound

e.g., *pat* /pæt/ vs. *bat* /bæt/

*Allophone*: a surface contextual variant of a phoneme

e.g., regular past tense *-ed* and regular plural *-s*

*Phone*: a single sound unit/segment (contrastive or not)

e.g., *train* [tɹeɪn] has four or five phones (depending on whether you count [eɪ] as one or two)

***/-d/***

***/-z/***

[næp-t]

'napped'

[læp-s]

'laps'

[næb-d]

'nabbed'

[læb-z]

labs'

[saɪt-əd]

'sighted'

[li:s-əz]

'leases'

[saɪd-əd]

'sided'

[tʃi:z-əz]

'cheeses'

# DISTINCTIVE FEATURES & NATURAL CLASSES

- *Laryngeal* features specify the state of the glottis:

[±VOI] distinguishes [b, d, g] from [p, t, k]

- *Place* features specify place of articulation:

[p] is [LAB], [t] is [COR]

- *Manner* features specify manner of articulation:

[s] is [+CONT], [d] is [-CONT]

- Conjunctions of features define *natural classes*:

[-VOI, +CONT] ↔ [f, θ, s, ʃ, h]

# REWRITE RULE FOR ENGLISH PAST TENSE

A rewrite rule is a four-tuple  $\phi \rightarrow \psi / \lambda \_ \rho$  read as “ $\phi$  is replaced by  $\psi$  when preceded by  $\lambda$  and followed by  $\rho$ ”

ASSIMILATION:  $[+OBS] \rightarrow [-VOI] / [-VOI, +OBS] \_$

$\phi \rightarrow \psi$ :  $[p, b \rightarrow p], [t, d \rightarrow t], [k, g \rightarrow k], \dots$

$\lambda$ :  $[p, t, k, \dots]$

$\rho$ : (null)

$/næp-d/ \rightarrow [næpt]$

# EFFICIENT COMPILATION OF REWRITE RULES

- Rule  $\phi \rightarrow \psi / \lambda \_ \rho$  can be represented as the composition of 5 finite state transducers  $r \circ f \circ R \circ L_1 \circ L_2$  such that
  - $r$ : insert  $>$  before every  $\rho$
  - $f$ : insert  $<_1$  and  $<_2$  before  $\phi$  followed by  $>$  (marking just those  $\phi$  immediately before a  $\rho$ )
  - $R$ : replace  $\phi$  with  $\psi$  in the context  $<_1 \_ >$  and delete  $>$
  - $L_1$ : delete  $<_1$  preceded by a  $\lambda$
  - $L_2$ : delete  $<_2$  not preceded by a  $\lambda$

A → B / C — D

C A C A D C D D C A A C A D

insert > before every D:

C A C A> D C> D> D C A A C A> D

insert 21 before A>:

C A C21A> D C> D> D C A A C21A> D

replace 1A> with 1B:

C A C21B D C> D> D C A A C21B D

clean up:

C A C B D C D D C A A C B D

# RULE INTERACTIONS

In most dialects of English, the consonants /t, d/ undergo FLAPPING when preceded by a stressed vowel; as a result the word *coder* thus sounds the same as *coater*, though *code* and *coat* do not.

For Canadian speakers (a.o.), the stressed vowel in words like *price* undergoes RAISING before /t/ but not before /d/; as a result, *write* and *ride* have different vowels.

Despite the fact that FLAPPING eliminates the distinction between /t ~ d/, *writer* and *rider* still have different vowels for these speakers.

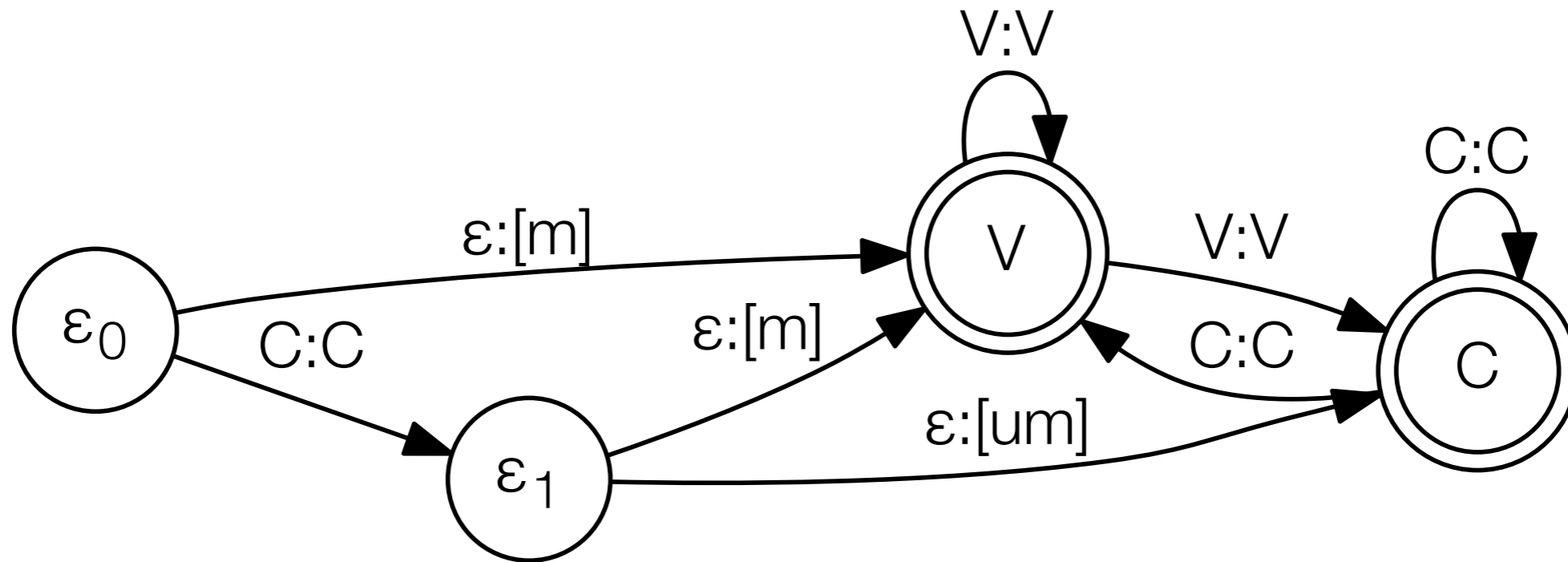
This shows that RAISING applies before FLAPPING; i.e., the output of FLAPPING is the input to RAISING.

[Source: Joos 1942, Chomsky 1957]



# INFIXATION IN THAO

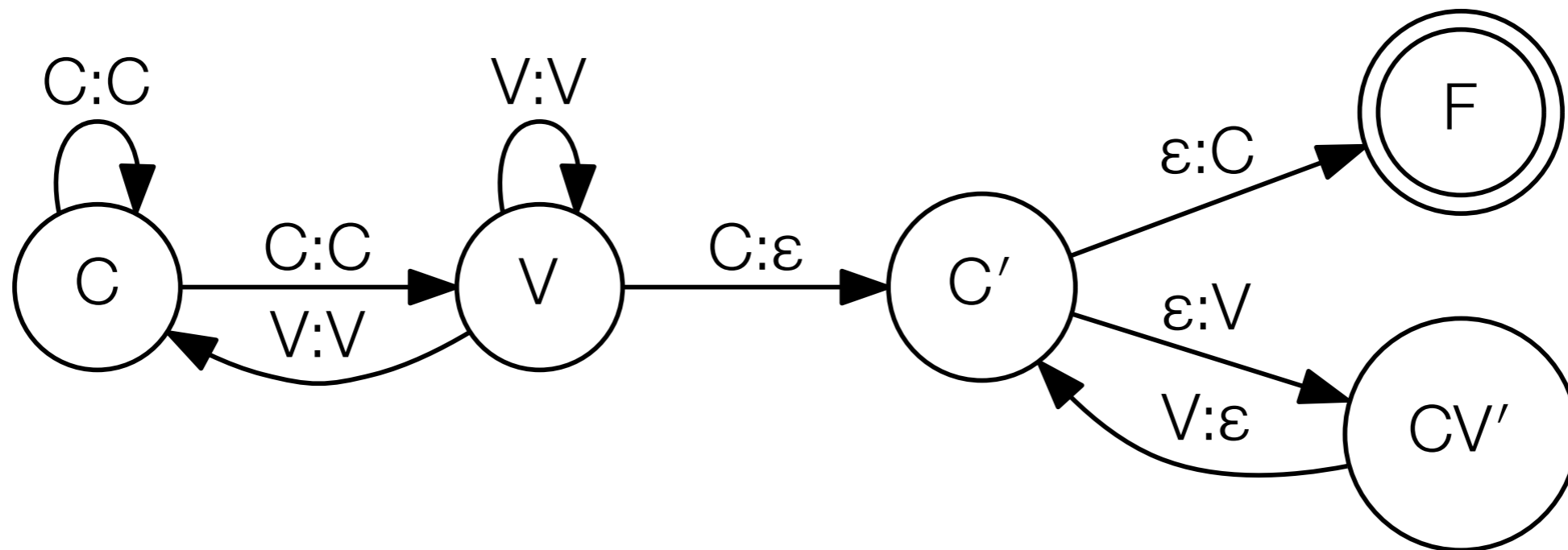
<i>hafuj</i>	<i>h<u>m</u>afuj</i>	‘chant’
<i>tqir</i>	<i>t<u>u</u>mqir</i>	‘protest’
<i>iup</i>	<i><u>m</u>iup</i>	‘blow with the mouth’



[Source: Lu 2011]

# METATHESIS IN ROTUMAN

<i>hosa</i>	<i>hoas</i>	‘flower’
<i>pure</i>	<i>puer</i>	‘to rule’
<i>parofita</i>	<i>parofiat</i>	‘prophet’



[Source: Churchward 1940, Chandlee et al. 2012]

# REDUPLICATION IN TAGALOG

<i>ól</i>	<i>olól</i>	‘mad’
<i>súlat</i>	<i>susúlat</i>	‘write’
<i>magka-útang</i>	<i>magkaka-útang</i>	‘owe’

<i>áraw</i>	<i>araw-aráw</i>	‘day’
<i>sangpówo</i>	<i>sangposangpówo</i>	‘ten’

Total reduplication is very difficult to model with FSTs.

Many NLP systems really work with written (orthographic) forms. That's an interesting subject, but beyond the scope of this class.

# MORPHOLOGY

## **Morphological differences distinguishing *H. neanderthalensis* from *H. sapiens*:**

Projecting mid-face

Low, flat, elongated skull

Lack of a protruding chin

No groove on canine teeth

Barrel-shaped rib cage

Large kneecaps

Long collar bones, wider shoulders

# DEFINITIONS

- *Morpheme*: a unit of meaning defined such that each word contains at least one

e.g., *dog* has one morpheme, *dogs* has two

- *Allomorph*: a surface contextual variant of a morpheme

e.g., the English regular past tense /-d/ has three allomorphs: [-d] as in *nabbed*, [-t] as in *napped*, and [-əd] as in *sighted*

# TYPES OF AFFIXATION

- Prefixation: *un-lock*
- Suffixation: *excite-ment*
- Infixation: *saxo-ma-phone, Minne-bloody-sota*
- Circumfixation: *in-vigor-ate, in-toxic-ate*
- Templatic affixation: Arabic *yuktibu vs. kaatib*



# FUSIONAL (INFLECTING) VS. AGGLUTINATIVE LANGUAGES

## **Fusional:**

(Latin)

*gladi-us*

sword-MASC.NOM.SG.

‘my sword’

*me-us*

1POSS.SG.-MASC.

## **Agglutinative:**

(Turkish)

*oda-lar-im-iz*

room-PL.-1POSS.-PL.

‘our rooms’

# ISOLATING (ANALYTIC) VS. SYNTHETIC LANGUAGES

**Isolating:**

(Vietnamese)

*Khi tôi đến nhà bạn tôi, chúng tôi bắt đầu làm bài.*

‘When I came to my friend’s house, we began to do the lesson.’

**Synthetic:**

(Turkish)

*Çöp-lük-ler-imiz-de-ki-ler-den-mi-y-di?*

‘Was it from those that were in our garbage cans?’

# INFLECTION VS. DERIVATION

*Inflection*: affixation that preserves the major part of speech, adding syntactically required information

e.g., *dog vs. dogs*

*Derivation*: affixation that (usually) changes the part of speech, adding additional information

e.g., *destroy vs. destruction*

# AUTOMATIC ANALYSIS OF MORPHOLOGY

- Stemming (e.g., Porter II):

*In linguistics, a morpheme is the smallest grammatical unit in a language.*

- Morphological segmentation (e.g., Morfessor):

*In linguis-tics, a morphem-e is the small-e-s-t gramma-tical uni-t in a language.*

- Morphological analysis (e.g., KIMMO):

*In linguistic-s[n.pl.], a morpheme is the small-est[comparative] gramma-tical[adj.] unit in a language.*