
UNIVERSITI SAINS MALAYSIA

First Semester Examination
Academic Session 2002/2003

September 2002

CCS503 – Intelligent Document Processing

CSC514 – Natural Language Processing

Duration : 3 hours

INSTRUCTION TO CANDIDATE:

- Please ensure that this examination paper contains **FOUR** questions in **SIX** printed pages before you start the examination.
 - Answer **ALL** questions.
 - This is an Open Book Examination.
 - You can choose to answer either in Bahasa Malaysia or English.
-

ENGLISH VERSION OF THE QUESTION PAPER

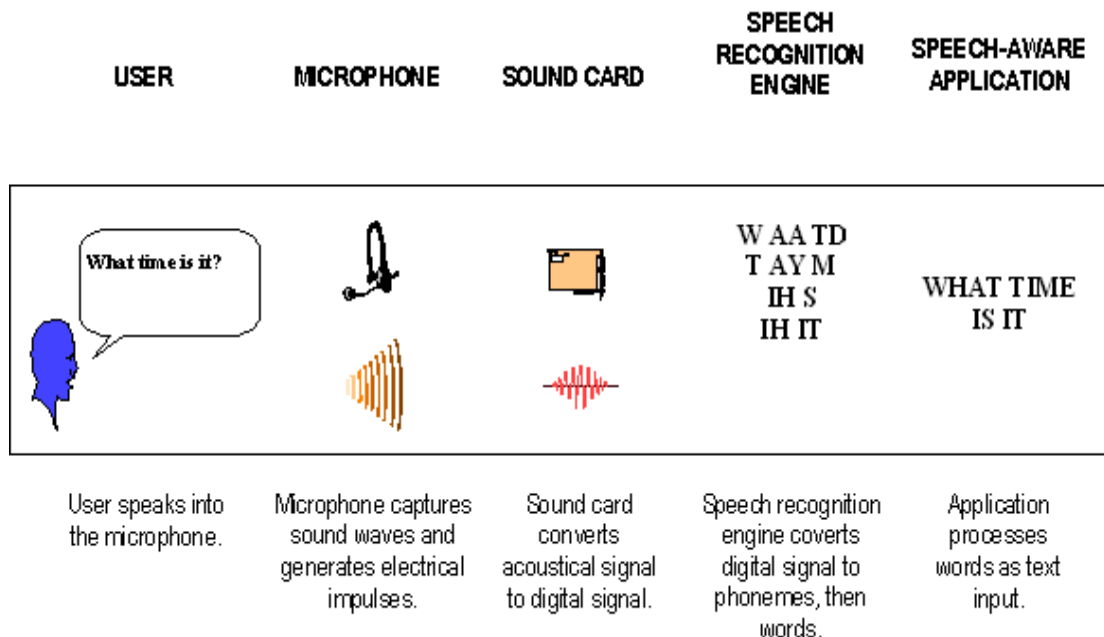
1. Given the following sentence:

"What time is it?"

- (a) How many phonemes are there in the given sentence? Write down the pronunciation of the sentence given above in IPA symbols.

[25/100]

- (b) Speech recognition, or speech-to-text, involves capturing and digitizing the sound waves, converting them to basic language units or phonemes, constructing words from phonemes, and contextually analyzing the words to ensure correct spelling for words that sound alike (such as write and right). The figure below illustrates this high-level description of the process.



Discuss in detail how the various NLP tools introduced in this course can be applied to improve the performance of the speech recognition engine as mentioned above.

[50/100]

- (c) In English, the regular way of forming the plural of a noun is to add an 's' to the end of the word. Such 's' is pronounced differently (i.e. /z/ or /s/) in different contexts. For example, the plural of *bid* is [bɪdz] and the plural of *bit* is [bɪts]. Construct the appropriate phonological rule(s) for this phenomenon.

[25/100]

2. This question has two parts. Both parts MUST be answered.

- (a) A noun may be expressed in the singular or plural form, e.g. *cat ~ cats*, *house ~ houses*. In English, plurality is expressed by various inflections.

| | |
|---------------------------|--------------------------|
| <i>baby ~ babies</i> | <i>maze ~ mazes</i> |
| <i>box ~ boxes</i> | <i>rash ~ rashes</i> |
| <i>bush ~ bushes</i> | <i>shelf ~ shelves</i> |
| <i>buzz ~ buzzes</i> | <i>spy ~ spies</i> |
| <i>cat ~ cats</i> | <i>switch ~ switches</i> |
| <i>chicken ~ chickens</i> | <i>tax ~ taxes</i> |
| <i>echo ~ echoes</i> | <i>toad ~ toads</i> |
| <i>frog ~ frogs</i> | <i>tomato ~ tomatoes</i> |
| <i>glass ~ glasses</i> | <i>watch ~ watches</i> |
| <i>house ~ houses</i> | <i>wife ~ wives</i> |
| <i>knife ~ knives</i> | <i>wolf ~ wolves</i> |

- (i) Using ONLY the data provided above, determine the various possible inflections in English.

[5/100]

- (ii) The inflections are used only under certain conditions. Now write rules to generate the appropriate plural form in English. Your rules should be in the following format. Remember to give the condition(s) under which the rule may apply.

singular_ending → plural_ending | condition(s)

[40/100]

- (iii) Indicate the order in which the rules are to occur.

[5/100]

- (b) In Spanish, nouns too may be expressed in the singular or plural form, e.g. *libro* ~ *libros*, *iglú* ~ *iglúes*. As with English, plurality is expressed by various inflections. [In Spanish, accent on a syllable is indicated by ' , as in *el bambú* 'bamboo', *el ladrón* 'robber' and *el lápiz* 'pencil'.]

| | | | |
|-----------------------|--------------|--------------------------|----------------|
| <i>el bambú</i> | 'bamboo' | <i>los bambúes</i> | 'bamboos' |
| <i>el corredor</i> | 'runner' | <i>los corredores</i> | 'runners' |
| <i>el hombre</i> | 'man' | <i>los hombres</i> | 'men' |
| <i>el iglú</i> | 'igloo' | <i>los iglúes</i> | 'igloos' |
| <i>el joven</i> | 'young man' | <i>los jóvenes</i> | 'young men' |
| <i>el ladrón</i> | 'robber' | <i>los ladrones</i> | 'robbers' |
| <i>el lápiz</i> | 'pencil' | <i>los lápices</i> | 'pencils' |
| <i>el libro</i> | 'book' | <i>los libros</i> | 'books' |
| <i>el lunes</i> | 'Monday' | <i>los lunes</i> | 'Mondays' |
| <i>el mujer</i> | 'woman' | <i>las mujeres</i> | 'women' |
| <i>el nivel</i> | 'level' | <i>los niveles</i> | 'levels' |
| <i>el padre</i> | 'parent' | <i>los padres</i> | 'parents' |
| <i>el papel</i> | 'paper' | <i>los papeles</i> | 'papers' |
| <i>el paraguas</i> | 'umbrella' | <i>los paraguas</i> | 'umbrellas' |
| <i>el pollo</i> | 'chicken' | <i>los pollos</i> | 'chickens' |
| <i>el régimen</i> | 'regime' | <i>los regímenes</i> | 'regimes' |
| <i>el rubí</i> | 'ruby' | <i>los rubíes</i> | 'rubies' |
| <i>la bici</i> | 'bike' | <i>las bicis</i> | 'bikes' |
| <i>la casa</i> | 'house' | <i>las casas</i> | 'houses' |
| <i>la ciudad</i> | 'city' | <i>las ciudades</i> | 'cities' |
| <i>la flor</i> | 'flower' | <i>las flores</i> | 'flowers' |
| <i>la nación</i> | 'nation' | <i>las naciones</i> | 'nations' |
| <i>la pluma</i> | 'pen' | <i>las plumas</i> | 'pens' |
| <i>la universidad</i> | 'university' | <i>las universidades</i> | 'universities' |
| <i>una vez</i> | 'one time' | <i>algunas veces</i> | 'a few times' |

- (i) Using ONLY the data provided above, determine the various possible inflections in Spanish.

[5/100]

- (ii) The inflections are used only under certain conditions. Now write rules to generate the appropriate plural form in Spanish. Your rules should be in the following format. Remember to give the condition(s) under which the rule may apply.

singular_ending → plural_ending | condition(s)

[40/100]

- (iii) Indicate the order in which the rules are to occur.

[5/100]

3. Given the following grammar:

$S \rightarrow \text{noun}(\text{SEM0}) \text{ VP}(\text{SEM0})$
 $\text{VP}(\text{SEM0}) \rightarrow \text{verb}(\text{Comps}, \text{SEM0}, \text{SEM1}) \text{ VerbComps}(\text{Comps}, \text{SEM1})$
 $\text{VerbComps}(\text{none}, \text{SEM1}) \rightarrow \emptyset$
 $\text{VerbComps}(\text{noun}, \text{SEM1}) \rightarrow \text{noun}(\text{SEM1})$
 $\text{noun}(\text{SEM0}) \rightarrow \{\text{noun}, \text{SEM0}\}$
 $\text{noun}(\text{SEM1}) \rightarrow \{\text{noun}, \text{SEM1}\}$
 $\text{verb}(\text{Comps}, \text{SEM0}, \text{SEM1}) \rightarrow \{\text{verb}, \text{Comps}, \text{SEM0}, \text{SEM1}\}$

lexicon:

cows : noun, +ANIMATE
 grass : noun, -ANIMATE
 eat : verb, noun, +ANIMATE, -ANIMATE
 sleeps : verb, none, +ANIMATE, -

- (a) List all sentences generated by the above grammar which are semantically valid but grammatically incorrect (in English). [20/100]
- (b) Extend the grammar and lexicon given above so that only grammatically correct (in English) sentences will be generated. [30/100]
- (c) List all sentences generated by the grammar given in 3(b) which are grammatically correct and semantically valid (in English). [15/100]
- (d) Construct a chart produced by the bottom-up chart parsing technique for one of the longest sentences listed in 3(c) above. Values of all variables involved are required to be shown in the chart. [35/100]

4. (a) For each of the following NLP tools, describe its functionality and give an example input/output pair.

- (i) Spelling checker
- (ii) Inflectional morphological analyzer
- (iii) Part of speech tagger
- (iv) Parser
- (v) Word sense tagger

[40/100]

(b) Discuss in detail how the NLP tools in 4(a) can be applied to each of the following NLP applications.

- (i) Search engine
- (ii) Grammar checker
- (iii) Intelligent helpdesk system (or call center)

[60/100]