

Word segmentation

Word definition

Maximum matching

Word segmentation
Word definition
Lexicography

- Linguistics
 - Object: natural languages
- Lexicography
 - Sub domain of linguistics
 - Object: vocabulary

Word segmentation
Word definition
Vietnamese

- History of Vietnamese language
- Chinese
- Han-Nom
- Quoc-ngu

Word segmentation
Word definition
Vietnamese

- Categorization
 - Content words: noun, verb, adjective, adverb
 - Function words: pronouns, preposition, conjunction

Word segmentation
Word definition
Vietnamese

- Noun: cơm, gạo...
- Verb: chạy
- Adjective: (vấn đề) khó
- Adverb: (chạy) nhanh

Word segmentation
Word definition
Mono-syllable words

- Containing only one syllable
- Most of monosyllable words are thuan-Viet

Word segmentation
Word definition
Multi-syllable words

- Contains multiple syllables
- Từ ghép đẳng lập: to lớn, khỏe mạnh
- Từ ghép chính phụ: bằng khen, phòng ăn

Word segmentation
Word definition
Han-Viet words

- Chinese origin
- Composing of meaningful Chinese characters
- E.g.: học viên, nhân viên, công trường, đại học, phổ thông

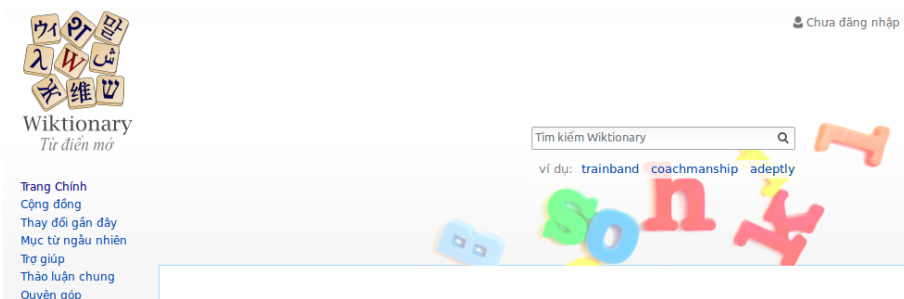
Word segmentation
Word definition
Foreign words

- Borrowed from other languages like English, French, Russia, Japanese, Korean.
- Domain terminologies
- E.g: cây (borrowed from 'case' in English), web, data.

Word segmentation
Word definition
Dictionaries

- <http://vlsp.org.vn/>
- Containing 40,000 common Vietnamese words

Word segmentation
Word definition
Dictionaries



Word segmentation
Word definition
Dictionaries



Word segmentation
Word definition
Proper names

- Persons: Hồ Chí Minh, Ngô Bảo Châu
- Locations: Hà Nội, Nhật

Word segmentation
Word definition
Dictionaries

Độ dài	# từ	%
1	6,303	15.69
2	28,416	70.72
3	2,259	5.62
4	2,784	6.93
5	419	1.04
Tổng	40,181	100

Word segmentation
Maximum matching

“Tách từ là bài toán nhận diện từ trong văn bản tiếng Việt.”

- WS is a problem in many Asian languages including Chinese, Japanese, Vietnamese, Thai, and Bummese.

Word segmentation
Maximum matching

- Input:
 - Dictionary
 - Unsegmented texts
- Algorithm:
 - Greedy algorithm
 - Move from left to right; Get the longest word;
 - Repeat until finish

Word segmentation
Maximum matching

-
- **START** initialize
 - (1) Input sequence $[w_0 w_1 \dots w_{n-1}]$
 - (2) $words \leftarrow []$
 - (3) $s \leftarrow 0$
-
- (4) $e \leftarrow n$ iteration
 - (5) When $[w_s \dots w_e]$ **has not been a word yet**: $e \leftarrow e - 1$
 - (6) $words \leftarrow words + [w_s \dots w_e]$
 - (7) $s \leftarrow e + 1$
 - (8) If $e < n$: Return to (4)
-
- (9) Get the segmented words finish
 - **END**
-

Word segmentation
Maximum matching
Examples

- Example 1:
“thời khóa biểu đang được cập nhật”
- Example 2:
“môn học xử lý ngôn ngữ tự nhiên”
- Example 3:
“con ngựa đá con ngựa đá”

18

Word segmentation
Maximum matching
Examples

- Example 1:
“thời khóa biểu đang được cập nhật”
→ *“thời_khóa_biểu đang được cập_nhật”*
- Example 2:
“môn học xử lý ngôn ngữ tự nhiên”
→ *“môn_học xử_ly ngôn_ngữ tự_nhiên”*
- Example 3:
“con ngựa đá con ngựa đá”
→ *“con_ngựa đá con_ngựa đá”*

Word segmentation
Maximum matching
Examples

- “học sinh học sinh học”*
- Left to right
“học_sinh học_sinh học”
 - Right to left
“học sinh_học sinh_học”

19

20

- Advantages:
 - Simple implementation
 - $O(V.n)$
- Shortcomings:
 - Require dictionaries
 - Segmentation ambiguities