

HỌ TÊN BỐ

a,b

, LÊ NG CHI MAI

a

Tóm tắt

: Bài viết này nhằm giới thiệu những khái niệm cơ bản và tình hình nghiên cứu

1. Mở đầu

Gần đây khi có dịp nói chuyện về xử lý ngôn ngữ (XLNN) và xử lý tiếng Việt (XLTV) trong công nghệ

2. Những khái niệm cơ bản

Tiếng nói và chữ viết là hai yếu tố cơ bản nhất của bất kỳ ngôn ngữ nào. Trong sự phát triển của công

(a) Trắc nghiệm là các bộ gõ chữ Việt và thành công của việc chuyển mã chữ Việt vào bộ
<http://nomfoundation.org>

(b) Tiếp theo có thể kể đến các công nghệ nhận diện (OCR: optical character recognition), như hệ VnDoc

(c) Các phần mềm hỗ trợ việc sử dụng các công cụ ngoài, tiêu biểu là các trình soạn thảo văn bản

(d) Các nghiên cứu trong việc phân tích các dịch Anh-Việt, Việt-Anh như các hệ dịch EVTRAN và VETRAN.

(e) Một loạt các nghiên cứu là Việt hoá các phần mềm mà gần đây tiêu biểu là bộ công cụ Việt hoá Windows

Vấn đề xử lý tiếng Việt trong công nghệ thông tin

Written by Administrator

Wednesday, 08 August 2012 00:05 - Last Updated Wednesday, 08 August 2012 00:10

Tuy liên quan đến ngôn ngữ Việt, không phải tất cả các vấn đề trên đều thuộc về lĩnh vực xử lý ngôn ngữ tự nhiên.

Để làm sáng tỏ điều này, chúng ta sẽ xem xét các ví dụ trong bài báo của ông Nguyễn Văn Khoa, một chuyên gia về xử lý ngôn ngữ tự nhiên.

Các dữ liệu liên quan đến ngôn ngữ viết (văn bản) và nói (tiếng nói) đang dần trở nên khó khăn hơn vì chúng ta đang phải xử lý các dữ liệu chính xác hơn.

Giống như chúng ta có các câu sau trong các tiếng nước ngoài:

- "We meet here today to talk about Vietnamese language and speech processing."

- "Aujourd'hui nous nous réunissons ici pour discuter le traitement de langue et de parole vietnamienne."

- "Мы встречаемся здесь сегодня, чтобы говорить о вьетнамском языке и обработке речи."

- "私たちはここに集まりベトナムについて話します。"

- "Chúng ta gặp nhau ở đây để bàn về xử lý ngôn ngữ và tiếng nói Việt."

Và giống như chúng ta không ai biết cả năm thì tiếng trên, nhưng tò mò muốn biết các câu đó nói gì. Nếu chúng ta nhìn vào các câu này để tìm hiểu về những câu tiếng Anh, Pháp, Nga, Nhật, Hàn và Việt như ta nhìn thấy ở trên.

- "Hôm nay chúng ta gặp nhau ở đây để bàn về xử lý ngôn ngữ và tiếng nói Việt."

Nếu các câu này được đưa ra trước những câu tiếng Anh, Pháp, Nga, Nhật, Hàn và Việt như ta nhìn thấy ở trên.

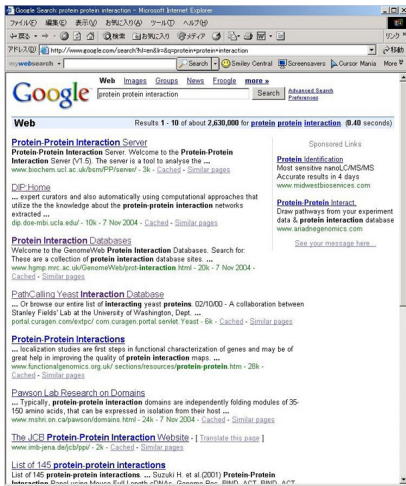


Tuy nhiên, một văn bản thô (một bài báo khoa học chẳng hạn) có thể có đến hàng nghìn câu, và ta

Có thể hình dung phần mềm gõ chữ Việt cho phép ta trực tiếp tạo ra một tệp văn bản trên máy tính (nh

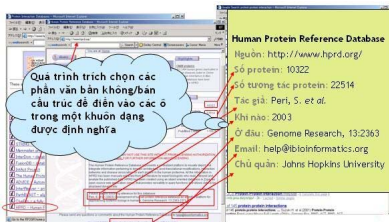
1. Nhận dạng tiếng nói (speech recognition), nhận biết và chuyển chúng thành dữ liệu số (digital data).
2. Tổng hợp tiếng nói (speech synthesis) bản, phân tích và chuyển thành tiếng nói.
3. Nhận dạng chữ viết (optical character recognition, OCR) nhận biết từng chữ cái và từ.
4. Dịch tự động (machine translation) dữ liệu văn bản trong một ngôn ngữ (tiếng Anh) sang một ngôn ngữ khác (tiếng Việt).
5. Tóm tắt văn bản (text summarization) bản dài (mỗi trang chẳng hạn) máy tóm tắt.
6. Tìm kiếm thông tin (information retrieval) tìm kiếm Google bản hay tiếng nói.

[1]



Hình 1

7. Trích ch n thông tin (information extraction) hi u đ p v b n b r m h a y t i ng n o i t i s t
[2]



Hình 2

8. Phát hi n tri th c và khai th á h i ng u n g i n b r m (knowledge discovery and text d

Còn nhi u bài toán và công nh x lý ngôn ngữ khác, nh giao di n ng i i máy b ng ngôn ngữ t nhi

ng d ng c a công nh <http://pagefault.havsta.com>) g p h n. Ở t h i n t y N a i t h i n / g i v i e t e t a t i n t r i n g t h
.co.jp/world/english
http://www.google.com/language_tools?hl=en
, hay

Có thể phân loại các bài toán:

– 1-3 thuộc lĩnh vực xử lý tín hiệu nói và xử lý ảnh (speech and image processing),

– 4-5 thuộc lĩnh vực xử lý văn bản (text processing),

– 6-8 thuộc lĩnh vực khai phá văn bản và Web (text and Web mining).

Phân loại này là tổng quát và không có nghĩa là processing and audio processing không liên quan đến xử lý tín hiệu.

Các bài toán 1-3, 5-8 liên quan đến xử lý hai ngôn ngữ khác nhau. Khi đó

3. Vấn đề phát triển của xử lý ngôn ngữ và tin học nói trong CNTT

Có thể nói xử lý ngôn ngữ tin học trên máy tính là một trong những vấn đề khó nhất của CNTT. Cái khó

Mấu chốt ở đây là bản chất của ngôn ngữ (có của coc sense) trong ngôn ngữ, định nghĩa và định hình ngôn ngữ.

Công nghệ ngôn ngữ, nhất là xử lý văn bản, với đời sống bao gồm các bước (tầng, layer) cơ bản sau đây

1. Tầng ngữ âm (phonetic and phonological) và ngữ âm (linguistic sounds), như mô hình

2. Tầng hình thái (morphological) và ngữ pháp (grammar): các thành phần có nghĩa của từ (word),

Vấn đề lý thuyết Việt trong công nghệ thông tin

Written by Administrator

Wednesday, 08 August 2012 00:05 - Last Updated Wednesday, 08 August 2012 00:10

Chức năng kinh nghiệm và các mô hình đã có thể nhưng thái độ của thị trường 80 đến 90% là chấp nhận

X lý văn bản và tiếng nói là phần giao diện mà vẫn qua mô hình thông kê và tiếp cận dựa vào

Như đã trình bày sơ bộ trên, x lý ngôn ngữ là một việc khó, phức tạp, chỉ có thể làm lâu dài theo nhu

Trên thế giới, nhiều tổ chức (ví dụ như <http://www.cic.nyu.edu>) về x lý ngôn ngữ tự nhiên đã được thành lập với các

Nhiều chính phủ đã đầu tư vào <http://www.till.edu> trong CNTT (Mỹ, Nhật Bản, Trung Quốc, Singapore,

Là người đi sau trong lĩnh vực này, các nhà nghiên cứu của chúng ta cũng đang cố gắng nghiên cứu và đưa ra các giải pháp

4. Tình hình và những vấn đề chính trong x lý ngôn ngữ Việt

Hãy thử nhìn lại tình hình của chúng ta. Ngoài những việc đã làm và bước đầu làm được có phần

Bên ngoài Việt Nam, cũng có những nhà x lý ngôn ngữ Việt, như nhóm dịch Anh-Việt của tiến sĩ PH

Ngoài những kết quả ban đầu, sau đây có thể là một vài điểm chính về hoạt động x lý ngôn ngữ

- Thứ nhất tập trung vào làm các sản phẩm công nghệ như các phần mềm dịch máy, một số

- Ít các nghiên cứu nền tảng về lý thuyết "hệ thống" và tài nguyên: thiếu

Vấn đề lý thuyết và thực tiễn trong công nghệ thông tin

Written by Administrator

Wednesday, 08 August 2012 00:05 - Last Updated Wednesday, 08 August 2012 00:10

- Phần đông là các nghiên cứu ứng dụng và định hướng ứng dụng thực tiễn và sử dụng kiến thức. Đa số mô hình nghiên cứu là mô hình nghiên cứu định lượng.
- Rất có thể nhiều nhóm nghiên cứu đã có ý định hành công việc nhưng còn thiếu kiến thức chuyên môn.
- Đáng băn khoăn hơn cả là các nhà nghiên cứu chưa chú trọng, thậm chí bỏ qua các kỹ thuật nghiên cứu định lượng.

5. Kết luận

Chúng ta ai cũng biết rằng công nghệ thông tin có thể mang lại nhiều lợi ích cho xã hội và con người.

1. Xây dựng và phát triển một số lý thuyết và kỹ thuật thông tin trên Internet và các ứng dụng khác.
2. Xây dựng các công cụ và kỹ thuật nghiên cứu định lượng, thực hiện mục tiêu 1 của đề tài.

Mặc dù việc phân tích các số liệu cho người dùng cuối là cấp bách và thiết yếu cùng ta cần nghiên cứu, nhưng

Lời cảm ơn

Xin chân thành cảm ơn các ý kiến đóng góp của bạn đọc tài liệu này và các đồng nghiệp: Ngô Văn Tuấn

Tài liệu tham khảo chính

Allen, J. (1994). Natural Language Understanding. The Benjamin/Cummings Publishing Co.

Baeza-Yates, R., Ribeiro-Neto, B. (1999). Modern Information Retrieval, Addison Wesley.

Bao, H.T., Thang, N.T., Chien, N.P., Mai, L.C. (2001). Towards a Practical Framework for Vietnamese I

Bao, H.T., Funakoshi, K. (1998). Information Retrieval Using Rough Sets, Journal of Japanese Society

Bao, H.T., Binh, N.N. (2002). Nonhierarchical Document Clustering by a Tolerance Rough Set Model, I

Bao, H.T., Tuan, N.A., Son, N.C. (2003). Issues in Construction of a Vietnamese Machine Tractable Dic

Berry, M.W. (2004). Survey of Text Mining: Clustering, Classification, and Retrieval, Springer.

Chakrabarti, S. (2003). Mining the Web, Morgan Kaufmann Publishers.

Cohen, W., McCallum, A. (2003). Information Extraction from the World Wide Web, Tutorial in ACM Con

Cole, R., Mariani, J., Uszkoreit, H., Varile, G., Zaenen, A., Zampolli, A., Zue, V. (1997). Language Tech

Dale, R., Moisl, H., Somers, H. (2000). Handbook of Natural Language Processing, Marcel Dekker.

Dien, D., Kiem, H., Toan, N.V. (2001). Vietnamese Word Segmentation, Proceedings of the Sixth Natur

Dien, D. (2002). Building a Training Corpus for Word Sense Disambiguation in English-to-Vietnamese M

Dien, D., Kiem, H. (2003). POS-Tagger for English-Vietnamese Bilingual Corpus, Proceedings of HLT-N

Door, B.J., Jordan, P.W., Benoit, J.W. (2000). A Survey of Current Paradigms in Machine Translation.

EDR Electronic Dictionary Technical Guide (1993). Japan Electronic Dictionary Research Institute.

Jelinek, F. (1998). Statistical Methods for Speech Recognition. The MIT Press.

Jurafsky, D., Martin, J. H. (2000). Speech and Language Processing. An Introduction to Natural Language

Hieu, P.X., Horiguchi, S., Bao, H.T. (2005). Conditional Models for Automatic Data Integration from the

Huong, L.T. (2004). Investigation into an Approach to Automatic Text Summarisation, Ph.D. dissertation

Huyen, N.T.M., Laurent Romary, Luong, V.X. (2003). A Case Study in POS Tagging of Vietnamese Text

Khang, B.H. et al. (2004). Báo cáo T ng k t Khoa h c và K thu t Đ tài Nghiên c u Phát tri n Công n

Mani, I., Maybury, M.T. (1999). Advanced in Automatic Text Summarization, The MIT Press.

Manning C. D., Schutze H. (1999). Foundations of Statistical Natural Language Processing. MIT Press.

Minh, N.L., Shimazu, A., Horiguchi, S., Bao, H.T. (2004). Example-Based Sentence Reduction Using H

Minh, N.L., Shimazu, A., Horiguchi, S., Bao, H.T., Fukushi, M. (2004). Probabilistic Sentence Reduction

Minh, N.L. (2004). Statistical Machine Learning Approaches to Cross-language Text Summarization, PhD

Nagao, M. (1989). Machine translation: how far can it go? Oxford University.

Oracle Text – An Oracle [White Paper \(2001\)](http://www.oracle.com/technology/products/text/pdf/text_bwp.pdf)

Sirmakessis S. (2004). Text Mining and Its Applications, Springer.

a,b

Viện Công nghệ Thông tin,

Viện Khoa học và Công nghệ T

[1] Chú thích của Vietlex.

[2] Chú thích của Vietlex.