

國立臺灣大學電機資訊學院資訊工程學系

博士論文

Department of Computer Science and Information Engineering

College of Electrical Engineering & Computer Science

National Taiwan University

doctoral dissertation

台語文處理技術：以變調及詞性標記為例

Processing Techniques for Written Taiwanese

-- Tone Sandhi and POS Tagging

楊允言

Iûnn Ún-giân

指導教授：高成炎博士 陳克健博士

Advisors: Gao Cheng-yan Ph.D. Chen Keh-jiann Ph.D.

中華民國 98 年 1 月

January, 2009

國立臺灣大學博士學位論文
口試委員會審定書

台語文處理技術：以變調及詞性標記為例
Processing Techniques for Written Taiwanese
-- Tone Sandhi and POS Tagging

本論文係楊允言君（學號 D93922001）在國立臺灣大學資訊工程學系完成之博士學位論文，於民國 98 年 1 月 16 日承下列考試委員審查通過及口試及格，特此證明

口試委員：

高成炎

陳宜成

(指導教授)

陳言希

高勳彰

高照明

張俊盛

劉昭麟

李淑娟

余明興

呂育道

系主任

Preface

Taiwanese is my mother tongue. I lived in Taipei for 30 years, since I was born. I spoke my mother tongue before I entered elementary school. I could also understand Mandarin, but could not speak it well. However, after entering elementary school, I gradually forgot my mother tongue. I was able to speak Mandarin fluently, and I could understand Taiwanese but found it difficult to express my thoughts in my mother tongue. In school, they told me that Taiwanese was only a dialect, and that we should not speak it in school. In fact, if we spoke Taiwanese, we could be punished. I accepted this argument until I was 20 years old.

At that time, I enrolled at National Taiwan University (NTU), during the period when Taiwan was still under martial law. A political event occurred on campus that gave me a big shock. I could not concentrate on my studies for several months. I was gradually moving out of the shadows, establishing my own independent thinking and values. I realized that cultural enlightenment is more important than political enlightenment, and that our mother tongue is a solid foundation for Taiwanese culture.

I began to relearn my mother tongue. Then, I realized that it would not be enough to just be able to understand and speak the language. The ability to read and write the Taiwanese language could contribute to its promotion.

I began to read Taiwanese material, which is difficult to find in bookstores.

I also wrote letters to my friends in Taiwanese, mixing various kinds of scripts, including Han characters, phonetic symbols, KK phonetics, and Kana. Afterwards, I borrowed a book from the NTU library entitled, *The speech structure and transcription method of Taiwan Hokkianese* ‘台灣福建話的語音結構及標音法,’ which was written by Tēⁿ, Liông-úi and published in 1978 . It took me four days to learn POJ. I also founded the NTU Puppet Shows Association ‘台大掌中劇團’ with my classmate Tiuⁿ, Lân-sék ‘張蘭石’ during my last year of university life. I strongly advocated the use of written Taiwanese to write the drama scripts, but finally failed due to a lack of time

Based on my computer science background, my study interest shifted to natural language processing during my graduate years, from 1991 to 1993. I processed Mandarin. At the same time, I was still engaged in the written Taiwanese movement. My friends and I founded the “Students Taiwanese Promotion Association,” (STAPA) ‘學生台灣語文促進會,’ published the “Taiwanese Student” magazine ‘台語學生’ (4 trial issues and 22 issues in total), wrote articles in Taiwanese, *etc.*

I entered the Dahan Institute of Technology CSIE Department as a lecturer and began to do Taiwanese processing related research work in 1999. Due to the lack of Taiwanese language related resources, I tried to establish some online written Taiwanese tools, including the Online Taiwanese-Mandarin Dictionary (OTMD), Online Taiwanese Syllable Dictionary (OTSD), and Online Taiwanese Concordancer System (OTCS), with more than 9-million-syllable written Taiwanese materials, Taiwanese word frequency reports with both

Han-Romanization mixed script and Taiwanese Romanization script, *etc.*

In my opinion, the computational linguistics related academic society is a minority in Taiwan. It is difficult for a researcher to master both computer science and the language itself. Obviously, Taiwanese processing related research is less popular in this minority, although I cannot understand why most researchers abandon their mother tongues. I selected Taiwanese language processing as my major field for study, tried to create this “new” field of research, and have achieved some preliminary results. I had the computer science background, and have been engaged in written Taiwanese work for more than 20 years.

I love my mother tongue and want to promote its social status. This is a difficult but, I believe, an important task.

Acknowledgments

I started reading and writing Taiwanese in 1987 when I was an undergraduate student, began written Taiwanese processing related research in 1999, and re-entered NTU CSIE as an adult student in 2004. I have always believed that working in this field is a worthwhile endeavor.

I am very grateful to my advisor, Professor Cheng-yan Kao, who encouraged me to select this topic and offered me some related research opportunities. My appreciation also goes to research fellow Keh-jiann Chen, my co-advisor, who generously provided me with many Chinese language processing resources. Many thanks to my committee members, Professor Jason S. Chang, Hsi-Jian Lee, Hsin-hsi Chen, Chao-Lin Liu, Ming-shing Yu, Khîn-huānn Lí, and Zhao-ming Gao, for their many constructive suggestions.

Thanks to my academic partner Hák-khiam Tiunn. Thanks also to Professor Yuh-dauh Lyu and Doctor Hók-chû Tiunn for correcting my written English in this dissertation. Thanks to Professor Jū-hông Tiuⁿ for helping me to understand some linguistic terms that confused me.

OK, I will leave space for my mother tongue. The following content will be written in Taiwanese.

Iōng ka-kī ê bú-gú siá, siá liáu khah kín, thang khah kín pit-giáp ☺

Ùi 1993 nî gián-kiù-só' pit-giáp liáu-āu khai-sí kóng. Hit-chūn tui lâu tī hák-sút-kài sit-chāi bô chhù-bī, kám-kah chòe Tâi-gú-bûn ūn-tōng khah sit-chāi. Góa seng kè-siòk ku tī Tiong-gián-īⁿ Sū-khò' sió-cho' chiáh thâu-lō', ná tih chòe Tâi-bûn kang-chok-chiá hóng-tâm. 1994 nî góa khi ín Tâi-tâi tiān-sòan tiong-sim ê thâu-lō', kám-siā Tân Bûn-chìn lâu-su, thiaⁿ-kóng góa ùi 60 kúi ê èng-cheng-chiá hông keng--tiòh, lâu-su ū tau-saⁿ-kāng. Tī hiá chòe ê khah tiōng-iàu ê tãi-chì, tō-sī siá tãi-hák liân-khó hun-hoat thêng-sek, che sī Khó' Sùn-khim lâu-su hō' góa ê ki-hōe.

Góa mā khai-sí peh-soaⁿ, chham-ka Siā Tiōng-têng teng-san hōe, koh ū chòe

kàn-pō, sèk-sāi chin chōe hó pêng-iú. Kám-siā Kì Chhun-heng ê kài-siāu, in-ūi peh-soaⁿ, góa sèk-sāi Giòk-lêng, liáu-āu kiat-hun. Bit-góat ê tē it chām, sī tè Siā Tiông-têng teng-san-hōe ê pêng-iú chòe-tīn peh Giòk-san, chòe chit kiāⁿ tiông-iàu ê tãi-chì, chit hāng tãi-chì, it-tit kàu 12 tang āu chiah ũ mui-thé pò-tō.

¹ Góa koh ták lé-pái kò-têng chit àm tī Tâi-ôan ê Tiàm chòe gī-kang, thâu-ke A-sam-ko kah Biāu-lêng-chí tui Tâi-ôan ê ài, sī góa chiok kèng-pōe ê.

In-ūi peh-soaⁿ, jú lâi jú kám-kah tōa tī Tâi-pak sit-chāi sī chiok kan-khó' ê tãi-chì, lāng chōe, chhia chōe, khong-kan óeh, khong-khì báí. Góa khai-sí chhē só-chāi, ũ soaⁿ koh lí to'-chhi khah hng ê, ũ Hoa-liân, Tâi-tang kah Lâm-tâu. 1996 nî, in-ūi chhē-tiòh Hoa-liân ê thâu-lō', tī lāu-pē lāu-bú kah Giòk-lêng lóng bô chi-chhi ê chêng-hêng hā, kā Tâi-tâi ê thâu-lō' sí-tiāu, cháu lâi Hoa-liân.

1996 nî, Hoa-liân iah-koh ũ pèh-sek khióng-pò' ê khi-hun, góa bô kā lí phiàn!

Hó ka-chài chit tang āu Giòk-lêng tiàu lâi Hoa-liân, só-í seng-òah thang khah ún-têng. Gún tī Hoa-liân bóe chhù, Tan-hông kah Ka-chhái mā sio liân-sòa chhut-sì.

Seng-òah ún-têng bô piáu-sī thâu-lō' ún-têng, góa tī Hoa-liân, saⁿ tang ōaⁿ gō' ê thâu-lō'. Góa iah-sī hoat-kak, nā-sī beh chòe Tâi-gú-bún, tī hák-sút-kài iah-sī khah hó-sè. 1999 nî, góa jip-lâi Tâi-hàn ki-sút hák-íⁿ.

Kám-siā Hák-khiam hiaⁿ, i chin khòaⁿ-tiông góa, ũ-sī-á hō' góa ki-hōe khi i ê pan ián-káng/siông-khò, mā ũ chiap kè-ék, pun keng-hùi hō' góa tau chòe bāng-chām.

Góa khai-sí kiàn-lip Tâi-gú-bún gián-kiù ê ki-chhó', Tâi-hôa sū-tián (kám-siā Tēⁿ Liông-úi kàu-siū kah chin chōe pêng-iú), Tâi-gú jī-tián, gú-liáu-khò (kám-siā Chùn-iók-hiaⁿ, Pek-nî hiaⁿ, Tek-hôa, Lē-soat, ...), Siók-gú (kám-siā Siau Pêng-tī lāu-su), Koa-iâu (kám-siā Gô' Jîn-sek lāu-su), Chhiò-khe-tâi (kám-siā Iūⁿ Chín-jū hāu-tiúⁿ, Siau lāu-su, Pek-nî hiaⁿ) téng-téng.

Góa siūⁿ-beh kái-tō', 1999 nî pat thak kè Tong-hôa Chòk-kún só' chāi-chit choan-pan chit hák-kí, 2002 nî khi khó' Sêng-tâi Tâi-bún-hē tē it kài phok-sū-pan,

¹ <http://www.libertytimes.com.tw/2007/new/feb/13/today-fo5.htm>

kám-siā Līm Sìn-kian bók-su thê-kiong sok-sià hō' góa tóa, gún sī chho'-chhù kìn-bīn. Sêng-tāi bô khó tiâu, tū-liáu ka-kī bô chài-tiâu, thiaⁿ-kong koh-khah chú-iàu ê gōan-in sī lāi-tóe chit kho' chiok chheh chòe Tâi-gú ê. Kám-siā Lī Heng-chhiong kàu-siū, i khó-lêng kám-kah tit-chiap kā góa kóng khó phok-sū-pan bô bāng tui góa ê táⁿ-kek siuⁿ tōa, só-í kiàn-gī góa ùi góa gōan-lāi ê choan-giap chit-pêng lāi seng-téng khòaⁿ-māi. Kám-siā Tēⁿ Liông-úi kàu-siū, i hiòng góa iau kó, siá chit phiⁿ lūn-būn tâu khì IJSL, khó-sioh āu--lāi chit pún Special Issue on Taiwanese pêng bô chhut-pán, kám-siā Hông-giâu (Henry), i tap-èng chòe-tīn siá chit phiⁿ lūn-būn, i ê Eng-būn iáu-siū chán.

In-ūi ū chit phiⁿ lūn-būn, góa tō sin-chhéng seng-téng, kám-siā Kán Lip-hong lāu-su, i siông-sòe khòaⁿ góa seng-téng ê chu-liāu, hō' góa chin chōe kiàn-gī, i kā góa kóng, in-ūi chin-chōe lāng bô liáu-kái góa tih chòe ê mih-kiāⁿ, só-í seng-téng ê chu-liāu ài lēng-gōa siá chit hūn Proposal, kau-tài Tâi-gú-būn chū-jiân gú-giân chhú-lí gián-kiù ê tiōng-iàu-sèng. Seng-téng sàng-sím ê sī-kan, tán chin kú chiah ū siau-sit, thiaⁿ-kóng ū úi-ôan bôai sím, thòe tng-khì Hák-sím-hōe. Góa chin hó-ūn, sūn-lī seng-téng chòe chō'-lí kàu-siū, m̄-koh bú seng-téng bú kah chin thiám-thâu, khai-sí ū pèh-thâu-mo' a.

Kán Lip-hong lāu-su koh ū kā góa kóng chit ê chin tiōng-iàu ê koan-liām, i kā góa kóng, sī-kan kàu--a, tō ài chòe kai chòe ê tã-chì, chit kù òe góa ū kā kì tī sim-koaⁿ tóe. Só-í, tī seng-téng kiát-kó chhut-lāi chìn-chêng, góa khai-sí sin-chhéng kok-kho-hōe kè-ék. Lēng-gōa, Tēⁿ Liông-úi kàu-siū hit-chūn tī Kau-tāi, i hi-bōng góa ē-tàng khì Chheng-hôa kè-siok thak phok-sū-pan.

Góa khai-sí liân-lók, pau-koah Tiuⁿ Chùn-sēng lāu-su, āu-lāi mā koh liân-lók Lí Sek-kian lāu-su kah Tân Sìn-hi lāu-su. In piáu-sī hoan-gêng, m̄-koh kám-kah chòe Tâi-gú gián-kiù beh pít-giap ū khùn-lân. Chha-put-to tī chit ê sī-chūn, Ko Sêng-iām lāu-su pān chit tiūⁿ Tâi-gú-būn siong-koan ê gián-thó-hōe, hōe-āu chòe-tīn chiáh-png ê sī, i kiò góa tng-lāi Tâi-tāi Chu-sìn-hē thak, tō chòe Tâi-gú-būn hong-bīn. Góa tng-khì siūⁿ chit lé-pài, che kī-kan koh ùi Hoa-liân cháu-khì Koe-lāng chhiú khan chhiú, koh kah Ko lāu-su liân-lók, khak-jīm, liáu-āu chōaⁿ-á koat--lòh-lāi, khai-sí lāi chún-pí.

2004 nî, chha-put-to kāng chit ê sî-chūn, seng-téng thong-kè, phok-sú-pan lók-chhú, Kok-kho-hōe kè-ék sin-chhéng--tiòh, Tan-hông mā beh jip Hù-sió. Ūi-tiòh Giòk-lêng chhōa gín-á hong-piān, gún ùi Kiat-an poaⁿ lâi Bí-lūn-á, tī Hù-sió āu-piah-m̄ng bóe chit keng a-phà-toh. Góa siūⁿ-kóng choan-sim thak khah khòai pit-giap, só-í Tâi-hàn chit pêng sin-chheng liū-chit thêng-sin.

Ko lāu-su chit ê sî-chūn chhéng-tiòh Tâi-būn-kóan ê kè-ék, beh chòe Tâi-gú ê gú-im hák-seng, pún-lâi chhiáⁿ góa ê tông-òh Sēng-an tàu-tīn chòe, khó-sioh tòng bô chit kò gèh, tō koh kâ Sēng-an tiàu khi chòe Bioinfo, Ko lāu-su ê pún-tō. Pún-lâi koh chhiáⁿ Hank chòe-tīn lâi chòe Tâi-gú gú-im piān-sek, bô gī-gō i āu-lâi bô thong-kè chu-keh-khó chit koan. Kè-ék ū chit ê choan-jīm ê giáh, chhiáⁿ Kiát-gák hiaⁿ lâi tàu-saⁿ-kāng. Góa ka-kī ê kè-ék leh ? Sui-jīan sin-chhéng--tiòh, chóng--sī pó-chō kim-giáh siuⁿ chió, só-í chú-chhî-jīn hùi poah chhut-lâi chhiáⁿ Lē-soat lâi Hoa-liân chòe kè-ék ê khang-khè.

Kám-siā lāu-pē lāu-bú, ūi-tiòh góa beh tng-khì Tâi-pak thak-chu, in tèk-piát kâ chhù koh hoan chit piàn. Góa chit lé-pài tòa Tâi-pak gō-kang, m̄-koh lāu-su kau-tài--lòh-lâi ê khang-khè li-li khok-khok, mā bô-hoat-tō choan-sim thak-chu, koh-khah hāi ê, Giòk-lêng ùi 11 gèh khai-sí, khó-lêng sī seng-òah siuⁿ kín-tiuⁿ, liân-sòa n̄ng kò gèh ūi-tng bô sòng-khòai. Só-í góa liū-chit thêng-sin ūi-chhî chit hák-kí, liáu-āu tō koh tng-lâi kà-chu, thang kâ lāu-su kóng góa bōe-tàng tiāⁿ-tiāⁿ tī Tâi-pak.

Chu-keh-khó sī chit kiāⁿ khióng-pòⁱ ê tã-chì, pêng-kun chit kai ū 1/3 ê tông-òh in-ūi bô thong-kè chu-keh-khó hông thòe-hák, chhiūⁿ góa chit-khóan nî-hè khah tōa ê lāu hák-seng, ap-lèk iū-kí tōa. 2005 nî 3 gèh, tē it kho Complexity án-nóa kè ê ? In-ūi tè tiòh thak-chu-hōe ê kî-tiong kúi pái, kám-siā lī-hāi ê Ông Hông-lūn tông-òh. 2005 nî 10 gèh koh khó kè Algorithm kah AI, góa kau bōe chió hák-hùi hioh-jóah tak lé-pài cháu Tâi-pak chham-ka thak-chu-hōe. Ē-kì-e khó-chhì chêng ê lé-pái, Lêng-ông hong-thai lâi, Kiat-an ê chhù, chhù-téng sià chúi lòh-lâi, Bí-lūn-á ê chhù, po-lê-thang-á phò--khì, thêng-chúⁱ n̄ng kang, iah bōe-hù chhú-lí, thang-á seng iōng chóa kō--khí-lâi, pau-hók-á khóan--leh tō kín piāⁿ khi Tâi-pak chún-pī khó-chhì, sim-koaⁿ-thâu phit-phók-chhái. 2006 nî 3

gêh chòe-āu chit kho Architecture bô kè, ap-lèk koh piàn tōa, nî kè 40, hiông-hiông kám-kah chiân-tô' bông-bông, 2006 nî 10 gêh chòe-āu chit pái ê ki-hōe, chóng-sng ôan-sêng tē sî kho chu-keh-khó, ná-chhiūⁿ í-keng keng-lèk chit-sì-lâng hiah-nî kú.

Thák phok-sū-pan ê kî-kan, tú-tiòh 3 pái gin-á in-ūi hì-iām tōa īⁿ, Giòk-lêng siông-pan cháu bōe khui kha, só-í sî góa khi pīⁿ-īⁿ kò', ē-kì-ê ū chit-piàn, chai-iaⁿ beh tōa īⁿ, kín cháu-khì tô'-su-kóan chioh chit pún “Chhèng-phàu, pīⁿ-khún kah kng-thih”, lí-iông tōa tī pīⁿ-īⁿ ê sî-kan thák ôan, che sî Khîn-hōaN hiaⁿ kài-siāu góa thák ê. Chit tōaⁿ sî-kan, góa koh ū thák “Tiūⁿ-niū sè-kài”, “chit thág ò-á”, “Hun-lân keng-iām”, “Tâi-ôan chhú-hun” ... téng-téng ê chu, lâi pâi-kái sim-thâu ê ut-chut, tâi-pō-hūn sî tī hé-chhia téng thák ôan ê. Koh ū chit tōaⁿ sî-kan, chit thâu-chá khí-lâi seng sng chit ê Sudoku, chiām-sî kâ chòe bōe ôan ê khang-khè khng chit pīⁿ.

Pâi-kái ut-chut, góa koh ū chòe pát-hāng, pau-koah khi chham-ka Google sió kang-khū pí-sài. Thiaⁿ Kán Lip-hong lâu-su kóng, góa ê Tâi-hôa sū-tián sió kang-khū ū jip-úi, chóng--sî, āu--lâi su hō' sng-miā ê sió kang-khū. Tâi-ôan-lâng ài sng-miā, khah iāⁿ kè òh bú-gú, sit-chāi chin bú-nāi. Kám-siā Lí Chì-kiông lâu-su, i chio gún khiâ khóng-bêng-chhia, 2008 nî 3 gêh góa ū tī 20 tiám-cheng í-lâi ôan-sêng 300 khí-lò, iân-lō' Seven bōe ê sio ê chiáh-mih lóng chhiúⁿ bô, hit àm thiaⁿ chhân-kap-á ê siaⁿ thiaⁿ kah pá, khiâ kah kha kiông-beh tng--khi, kha-chhng óah-beh thiàⁿ--sî, tng--lâi koh sūi sán 3 kong-kun, m̄-koh tit-tiòh chit tiuⁿ cheng-su.² Kám-siā Tō' Chèng-sèng pō-tiúⁿ, i kâ Kok-gú-bún keng-sài cheng-ka Tâi-gú pheng-im pí-sài, hō' góa ū ki-hōe tit-tiòh siā-hōe-cho' tē it miā,³ i koh sak thui-tián bú-gú kiát-chhut kòng-hiàn-chióng, góa ū tiòh-chióng, mā in-ūi tiòh-chióng, hō' góa ê siau-sit ū ki-hōe khan tī Chū-iū sî-pò ê tōe-hng-pán.⁴

Kè-ék ê pō-hūn, thák phok-sū-pan ê chit tōaⁿ kî-kan, chin hó-ūn, chip-hêng sî ê Kok-kho-hōe kè-ék, kám-siā Kiát-gák hiaⁿ, Tek-hôa kah Tōa-thâu-liân kâ góa tàu chòe chin chōe khang-khè. Lêng-gōa, kám-siā Lí Heng-chhiong lâu-su, tī i ê

² <http://ungian.pixnet.net/blog/post/16038270>

³ <http://iug.csie.dahan.edu.tw/iug/Ungian/engu/2007/phengimpisai/phengimpisai.asp>

⁴ <http://iug.csie.dahan.edu.tw/iug/Ungian/engu/2008/bugujit/bugujit.asp>

chhui-chiàn hā, góa chip-hêng chit ê Tâi-bûn-kóan ê kè-ék, mā kám-siā Tiuⁿ Iām-hiàn lāu-su, i ê chhui-chiàn, hō' góa sin-chhéng tiòh Kok-sú-kóan ê kè-ék, chòe Tâi-gú-bûn ūn-tōng hóng-tâm, kám-siā Bí-chhin, i thòe góa kā siōng khùn-lân ê hóng-tâm pō'-hūn chhú-lí hó, Kok-sú-kóan chit pún chu, 2008 nî 4 gèh chhut-pan 800 pún, 8 gèh tō bōe ôan a. Chòe che bô pí siá phok-sū lūn-bûn khah khin-sang, ūi-tiòh chit ê hóng-tâm, liân-sòa kúi-lō kò-gèh lóng òaⁿ khùn, kè-nî ê sí koh mo' h kúi tháh ê kó-kiāⁿ tih khòaⁿ.⁵

Mā kám-siā Hák-khiam hiaⁿ, Khîn-hōaⁿ hiaⁿ, Ūi-bûn, Gō' Chiau-sin i-su, Si-chong hiaⁿ kah Sò'-eng, in iau-chhiáⁿ góa chòe in ê ke-ék ê kiōng-tōng chú-chhi-jîn, só'-í góa ê keng-lék ē-tàng siá kah chiok súi-khùi ☺

Kám-siā Ū Bêng-heng lāu-su, i kā góa iau-kó, hō' góa hoat-piáu góa ê tē it phiⁿ kok-chè kî-khan lūn-bûn, kám-siā Lâu Chiau-lîn kah Tân Pek-lím lāu-su, in ê iau-kó sī tē jī phiⁿ. Kám-siā Tân Sìn-hi lāu-su, in-ūi i ê iau-chhiáⁿ, hō' góa ū ki-hōe chòe Chu-sìn jī hák-mŋg ê kui-ōe úi-ôan, thê-chhut góa tui pún-thó' gú-giân chu-sìn chhú-lí ê kiàn-gī.

Kám-siā Tâi-hàn ki-sút hák-īⁿ, góa tī chiá chiáh thâu-lō', chit-chūn sī tē 10 tang, tī chiá, góa khai-sí ē-tàng choan-sim lâi chòe Tâi-gú-bûn ê hák-sút gián-kiù, mā tī chiá seng-téng chòe chō'-lí kàu-siū, sui-jiân chit keng hák-hāu í-keng put-chí-á hui-hiám a.

Kám-siā Tâi-tang tãi-hák Tâi-gú-só' ê hák-seng, sui-jiân hioh-jòah tàk lé-pái cháu Tâi-tang put-chí-á sin-khó', m̄-koh in jîn-chin kah chun-tiōng lāu-su ê thāi-tō', hō' góa kám-siū tiòh chòe lāu-su ê chun-giâm. Bók-chiân ū 3 ê góa chí-tō' ê hák-seng pit-giáp, hō' góa chin ū sêng-chiū-kám.

Kám-siā Chhit-chhiⁿ-thâm, Chhit-chhiⁿ-thâm m̄ sī thâm, sī tōa-hái, kè-khì chit tang, góa tiāⁿ-tiāⁿ khiá khóng-bêng-chhia siong-hā-pan, khòaⁿ tōa-soaⁿ tōa-hái ê hó keng-tì, koh kiam liân sin-thé.

Khó-lêng koh ū chit kóa lāng góa bô kóng-tiòh, m̄-sī lín bô tiōng-iàu, sī góa kì-tì báí.

Tē it pái jip Tâi-tâi, sī 1984 nî khó tiòh Chu-sìn-hē, tō'-kè jîn-seng

⁵ http://writtentaiwanese-movement.blogspot.com/2008/04/blog-post_16.html

cheng-chhái ê saⁿ tang gōa, kiàn-lip kè-tát-koan kah Tâi-ôan ì-sek. Tē jī pái jip Tâi-tâi sī 1994 nī tī Tâi-tâi tiān-sòan tiong-sim chiáh thâu-lō̍, hit-chūn thé-hōe Tâi-ôan tōa-soaⁿ ê súi, mā ôan-sêng chiong-sin tâi-sū. Tē saⁿ pái jip Tâi-tâi, tō-sī chit pái thak phok-sū-pan, hi-bōng ē-tàng hō̍ góa chhē tiòh khah ún-tēng ê thâu-lō̍.

Ka-tēng, khang-khè, gián-kiù, hák-giap, ... Nī-hè khah tōa, tâi-chì chōe, thak-chu khí-lâi chin sng sin-thé, chit sī tang pòⁿ í-lâi, bák-chiu ū ló-hoe, koh tiòh péh-lâi-chiong, chhùi-khí chhut bün-tōe, koaⁿ-kong-lêng chí-sò̍ chhèng kōan, siⁿ phê-chōa, kha-tóe mā ē thiàⁿ, kui sin-khu pīⁿ liáu-liáu. Góa kī-sit m̄ chai-íⁿ thak phok-sū-pan sī m̄-sī cheng-khak ê koat-tēng, chóng--sī, Tâi-gú-bün chiáⁿ-chòe góa it-seng ê chì-giap, che sī phah sí bô thòe ê.

Chòe-āu, beh kám-siā Tâi-gú-bün-kài chōe-chōe ê pēng-iú, in-ūi ū tak-ke ê phah-piàⁿ, góa chiah ū ki-hōe chòe chit ê tōe-bók. Tâi-gú-bün ün-tōng ê lō̍ iah-koh chin kú-tng, lán chòe-tīn kè-siok kiáⁿ lóh-khì, òng-bāng Tâi-gú-bün chhut-thâu-thiⁿ !

ps: Lün-bün kàu-sī ē chiūⁿ-bāng, chit phiⁿ siā-sū ē khng tī góa ê Blog.⁶

ps2: Acknowledgments ê chho̍-kó ū khng khi Tâi-gú-bāng, kám-siā Siau Pēng-tī lāu-su, Tiuⁿ Hók-chū i-su, Ngō̍ Chiau-sin i-su, Tân Phek-siū bs, Lâm-hēng-hiaⁿ, Kim-ēng lāu-su, Jīn-sek lāu-su, Khái--ko, Siok-hūi, Tsùn-tsiu ê liáh-pau kah chàn-siaⁿ. Chit pō̍-hūn siá liáu siōng sóng-khòai, bián chhin-chhiūⁿ siá lün-bün an-ne, tak jī tak kù ài chim-chiok.

ps3: 2009/1/16 sui-jiân thong-kè kháu-chhì, m̄-koh ài chun-chiàu kháu-chhì úi-ôan ê ì-kiàn siu-kái, in-ūi kin-nī kè-nī chá, hák-hāu iau-kiū siōng òaⁿ 2/2 ài kau, siu-kái ê sí-kan bōe-hù, phah-sng iân chit hák-kī chiah thang pit-giap, bô gī-gō̍ 1/20 siu-tiòh hák-hāu ê thong-ti, kóng kau lün-bün ê kī-hān kái-chòe 2/15. Che ke--chhut-lâi ê nng lé-pài, ná-chhiūⁿ sī chiam-tùi góa ê tiâu-khóaⁿ. Hō̍ góa thang lī-iōng kè-nī hiòh-khùn ê sí-kan kóaⁿ siu-kái ê khang-khè, kī-hān lâi ôan-sêng siu-kái, kiám láp chit hák-kī kùi-som-som ê chù-chheh-hùi. Kám-siā Tâi-tâi, mā kám-siā Giok-lêng ê thé-thiap.

⁶ <http://ungian.pixnet.net/blog/post/25236012>

摘要

台語是世界上重要的語言，可惜沒有受到應有的重視。在某些方面，台語文的特性與華文或英文相當不同。本論文主要討論台語文處理技術。

白話字（台語羅馬字）是台語文的重要書寫系統。我們先介紹白話字的字元編碼，提及白話字數字調號做為不同白話字字元編碼的內部表示法。針對白話字文本搜尋，我們提出兩階段搜尋策略，並提出白話字音節近似搜尋的方法。我們還描述白話字顯示方法、白話字文字處理相關應用程式以及漢羅台語文斷詞方法。

我們提出以規則方法處理變調問題的演算法。先將每個台語詞翻成華語詞，找出其詞類標記訊息，以詞類標記和變調規則來決定變調後的聲調。我們實作出台語變調系統。此系統在訓練資料及測試資料分別達到 97.4%和 89.0%的變調正確率。

此外，我們提出詞類標記方法。我們先開發語詞對齊檢查程式將逐段對齊的兩種台語文本做語詞對齊，之後利用 HMM 機率模型挑選最適當的華語對應詞，再利用 MEMM 分類器挑選出其詞性標記。我們的方法達到 91.5%的正確率。

過去幾年，我們建立了一些有用的線上台語文工具。希望這些工具以及我們所做的初步研究成果，能讓台語文處理相關研究更加蓬勃發展。

關鍵詞： 台語文，變調，詞類標記，白話字，自然語言處理。

Tiah-iàu

Tâi-gú sī sè-kài siōng chin tiōng-iàu ê gú-giân, khó-sioh soah bô hông khòaⁿ-tiōng. Tī bó chit-kóa hong-bîn, Tâi-gú-bûn ê tèk-sèng kah Hôa-bûn iah-sī Eng-bûn chin bô kâng-khóan. Pún lûn-bûn chú-iàu beh thó-lûn Tâi-gú-bûn chhú-lí ki-sút.

Pêh-ōe-jī (Tâi-gú Lô-má-jī, POJ) sī Tâi-gú-bûn ê tiōng-iàu su-siá hē-thóng. Gún seng kài-siâu POJ ê jī-gôan pian-bé, thê-chhut iōng POJ sò-jī tiâu-hō chiâⁿ-chòe bô kâng POJ jī-gôan pian-bé ê lōe-pōⁿ piáu-sī-hoat. Chiam-tùi PÔJ bun-pún chhiau-chhê, gún thê-chhut n̄g kai-tōaⁿ chhiau-chhê chhek-liók, koh thê-chhut POJ im-chiat lûi-sū chhiau-chhê ê hong-hoat. Gún koh ū kài-siâu POJ hián-sī ê hong-hoat, POJ ê bûn-jī chhú-lí siong-koan èng-iōng thêng-sek, kah Hàn-lô t̄ng-sū ê hong-hoat.

Gún thê-chhut iōng kui-chek hong-hoat chhú-lí piàn-tiâu bûn-tôe ê ián-sòan-hoat. Seng kā múi chit ê Tâi-gú sū hoan chòe Hôa-gú sū, chhê chhut i ê sū-lūi piau-kì, iōng sū-lūi piau-kì kah piàn-tiâu kui-chek lâi koat-têng piàn-tiâu liáu-āu ê siaⁿ-tiâu. Gún khai-hoat chhut Tâi-gú piàn-tiâu hē-thóng. Chit ê hē-thóng tī hùn-liân chu-liâu kah chhì-giām chu-liâu hun-piát tát-kàu 97.4% kah 89.0% ê piàn-tiâu chêng-khak-lút.

Lēng-gōa, gún thê-chhut sū-lūi piau-kì hong-hoat. Gún seng khai-hoat gú-sū tùi-chòe kiám-cha thêng-sek, kā chit tōaⁿ chit tōaⁿ tùi-chòe ê n̄g-khóan Tâi-gú bûn-pún chòe gú-sū ê tùi-chòe, liáu-āu lī-iōng HMM ki-lút bô-hêng keng chhut siōng tàu-tah ê Hôa-gú tùi-èng-sū, koh lī-iōng MEMM hun-lūi-khì keng chhut sū-lūi piau-kì. Gún ê hong-hoat tát-kàu 91.5% ê chêng-khak-lút.

Kè-khì kúi-lō tang, gún kiàn-lip chit-kóa chin ū lō-ēng ê sòⁿ-téng Tâi-gú-bûn kang-khū. Hi-bōng chiâ ê kang-khū kah gún sóⁿ chòe ê chhoⁿ-pōⁿ gián-kiù sêng-kó, thang hōⁿ Tâi-gú-bûn chhú-lí siong-koan gián-kiù koh-khah heng.

Koan-kiân-sū : Tâi-gú-bûn, Piàn-tiâu, Sū-lūi piau-kì, Pêh-ōe-jī, Chū-jiân gú-giân chhú-lí.

Abstract

Taiwan Southern Min (Taiwanese) is an important language that has received very little attention in the world. The characteristics of written Taiwanese are quite different from Mandarin or English in some respects. This dissertation focuses on Taiwanese processing techniques.

POJ is an important Taiwanese script. We introduce the POJ character code and mention numbered POJ as the interchange code for various POJ encodings. Then, we propose a two-stage search strategy for a POJ text search, and propose a POJ syllable query expansion. We also describe the display method for POJ, POJ word processing utilities, and a word segmentation method for HR mixed script.

We propose a rule-based tone sandhi algorithm. We translated every word into Mandarin, and obtained the POS information. Using the POS data and tone sandhi rules, we then tagged each syllable with its post-sandhi tone marker. Finally, we implemented a Taiwanese tone sandhi processing system. Our system achieved accuracy rates of 97.4% and 89.0% with training and test data, respectively.

Additionally, we propose a POS tagging method. We developed a word alignment checker to assist with the word alignment of the two Taiwanese scripts, selected the most adequate Mandarin word using a Hidden Markov probabilistic model, and finally tagged the word using a Maximal Entropy Markov Model classifier. We achieved an accuracy rate of 91.5% in the Taiwanese POS tagging work.

Over the past several years, we have established some useful online written Taiwanese tools. Based on these tools and our preliminary research results, we hope that written Taiwanese processing related research can be promoted.

Keywords: Written Taiwanese, Tone Sandhi, POS Tagging, Peh-Oe-Ji, Natural Language Processing

Table of Contents

Preface	i
Acknowledgments	iv
摘要	xi
Abstract	xiii
Abbreviations	xxiii
Chapter 1 Introduction	1
1.1 Background	1
1.1.1 Language Population in Taiwan	1
1.1.2 Southern Min Language Population.....	2
1.1.3 Another Investigation: the Taiwan Southern Min Viewers	3
1.1.4 The Confusing Name of This Language.....	5
1.2 Different Types of Written Taiwanese Scripts	6
1.2.1 The Han Characters Script	7
1.2.2 The Romanized Scripts	9
1.2.3 The Han-Romanization Mixed Script.....	10
1.2.4 Other Scripts	10
1.2.5 Target Scripts in This Dissertation	11
1.3 Issues Related to Written Taiwanese Processing	11
1.4 Organization of This Dissertation	12
Chapter 2 Resources and Survey of Written Taiwanese Processing.....	19
2.1 Digital Resources for Written Taiwanese	19
2.1.1 Fonts	19
2.1.2 Dictionary	20
2.1.3 Text Corpora.....	23
2.1.4 Electronic Books	27
2.2 Survey of Written Taiwanese Processing Techniques	28
2.2.1 Input Method	28

2.2.2	Word Segmentatation	29
2.2.3	POS Tagging	30
2.2.4	Scripts Conversion	30
2.2.5	Text-to-Speech	30
2.2.6	Translation.....	33
2.2.7	Parsing.....	33
2.3	Summary.....	33
Chapter 3 Coding, I/O for POJ, and Text Processing		35
3.1	Character Code of POJ.....	35
3.2	Two Kinds of POJ Representation.....	39
3.3	Search Problem with POJ Text	41
3.3.1	Issues with POJ Text Search	41
3.3.2	Two-Stage Search Method: String Matching Then Filtering	42
3.3.3	Query Expansions: Toneless, Glottal Stop, Checked Syllable, and Vowel	44
3.3.4	Examples of Search Results.....	47
3.4	POJ Text Display.....	49
3.4.1	Issues with POJ Text Display	49
3.4.2	POJ and Numbered POJ Conversion Method	50
3.4.3	POJ Graph Display.....	52
3.4.4	Examples of Display Results	53
3.5	Some Text Processing Utilities for POJ.....	55
3.5.1	POJ Phoneme Segmentation and Spelling Checker	55
3.5.2	POJ Syllable/Word/Sentence Count.....	57
3.6	Word Segmentation for HR Mixed Script.....	58
Chapter 4 Tone Sandhi Problem and Algorithm.....		63
4.1	Tone Sandhi Problem of the Taiwanese Language.....	63
4.1.1	Types of the Taiwanese Language Tone Sandhi.....	64

4.1.2	Boundary of Tone Sandhi Group	68
4.2	Implementation of the Taiwanese Pronunciation System	68
4.2.1	System Diagram.....	68
4.2.2	Observation Data and Test Data	70
4.2.3	POS Tagging Set	71
4.2.4	Tone Sandhi Marks	73
4.3	Rule-based Tone Sandhi Algorithm.....	73
4.4	Results, Accuracy Rate and Discussion	78
4.4.1	Experiment Results	78
4.4.2	Accuracy Rate and Related Analysis.....	80
4.4.3	Discussion	83
4.5	Summary and Possible Direction	85
Chapter 5	POS Tagging Method	87
5.1	Problems of POS Tagging.....	87
5.2	POS Tagging Methods.....	88
5.2.1	Origin of the Corpus	89
5.2.2	Word for Word Alignment	89
5.2.3	Searching for the Corresponding Mandarin Candidate Words.	90
5.2.4	Selecting the Best Mandarin Translation	91
5.2.5	Selecting the Most Appropriate POS According to the Corresponding Mandarin Word.....	92
5.3	Results.....	94
5.4	Error Analysis	99
5.4.1	Incorrect Corresponding Mandarin Word Selection.....	99
5.4.2	Absence of Appropriate Mandarin Words in the OTMD	100
5.4.3	Unknown Words from the Viewpoint of Mandarin.....	101
5.4.4	Propagation Error	101

5.4.5	Other Cases.....	101
5.4.6	Summary of Error Conditions.....	102
5.5	Discussion.....	103
5.5.1	Is Improvement Possible ?.....	103
5.5.2	Hyphen Problems, Distinction between Taiwanese and Mandarin.....	104
5.5.3	The Distinction between Different Eras or Different Genres....	105
5.6	Summary.....	106
Chapter 6	Conclusion and Future Work.....	109
6.1	Our Contributions to Written Taiwanese Resources and Processing	109
6.2	Future Work and Prospects for Written Taiwanese Processing Research.....	112
Reference.....		117
Appendix.....		127
A.1	Brief Introduction to The Phoneme of Taiwanese.....	127
A.1.1	Initials.....	127
A.1.2	Vowels.....	128
A.1.3	Tones.....	129
A.1.4	Compared with Mandarin.....	130
A.2	Examples of Written Taiwanese.....	132
A.3	Terminologies.....	136
A.4	Webpages Made by Author.....	138
A.5	Differences between POJ and TL.....	139

List of Tables

Table 1 - 1 Living Languages in Taiwan.....	1
Table 1 - 2 Southern Min Language Populations around the World.....	2
Table 1 - 3 Visits to the Taiwan Southern Min Contents Website	4
Table 1 - 4 The Names of “Taiwan Southern Min”	6
Table 3 - 1 POJ Unicode Encoding.....	36
Table 3 - 2 POJ and Numbered POJ Synopsis	40
Table 3 - 3 Example of Taiwanese Article in 3 Scripts	40
Table 3 - 4 Toneless Search for “hoe-chhia”.....	48
Table 3 - 5 Checked Syllable Search for “cha”	48
Table 3 - 6 Vowel Search for “uiN”	49
Table 4 - 1 The Taiwanese Tone Sandhi Phenomena	67
Table 4 - 2 Observation Data Sources	70
Table 4 - 3 Test Data Sources	71
Table 4 - 4 POS Classes	72
Table 4 - 5 Tone Sandhi Marks.....	73
Table 4 - 6 Tone Sandhi Marking Algorithm.....	73
Table 4 - 7 Number of Errors and Accuracy Rate for Each Paragraph	80
Table 4 - 8 Affected and Accurately Affected Syllables of Each Rule.....	81
Table 4 - 9 Number of Dominant Rule, Accurate Dominate Rule and Accuracy Rate.....	82
Table 4 - 10 Additional Rules to Obtain Higher Accuracy Rate	83
Table 4 - 11 Error Conditions and Possible Solutions.....	84
Table 5 - 1 Test Data List.....	95
Table 5 - 2 Tagging Accuracy Rate for the Test Data	96
Table 5 - 3 Example of POS Tagging Result	97
Table 5 - 4 The Incorrect Mandarin Words Selected and Their Respective POS	100

Table 5 - 5 Errors Caused by the Absence of Appropriate Mandarin Words in the OTMD.....	100
Table 5 - 6 Unknown Words from the Viewpoint of Mandarin	101
Table 5 - 7 The Reasons for the POS Tagging Errors.....	102
Table 5 - 8 Tagging Accuracy Rates for Different Genres.....	106
Table 5 - 9 Tagging Accuracy Rates for Different Eras.....	106
Table A - 1 Consonants.....	127
Table A - 2 Single Vowels.....	128
Table A - 3 Compound Vowels	128
Table A - 4 Nasal and Glottal	128
Table A - 5 Nasal Finals and Final Stops	129
Table A - 6 Syllabic Consonants	129
Table A - 7 Tones	130
Table A - 8 Terminologies	136
Table A - 9 Differences between POJ and TL	139

List of Figures

Fig 3 - 1 POJ Text Match Algorithm (Target: Word).....	43
Fig 3 - 2 POJ Toneless Search Algorithm (Target: Syllable).....	45
Fig 3 - 3 POJ Checked Syllable Search Algorithm (Target: Syllable)	46
Fig 3 - 4 POJ Vowel Search Algorithm (Data: Syllable)	47
Fig 3 - 5 POJ to Numbered POJ Algorithm	50
Fig 3 - 6 Numbered POJ to POJ Algorithm	51
Fig 3 - 7 Numbered POJ to POJ Graph Algorithm	53
Fig 3 - 8 Unicode Display of POJ.....	54
Fig 3 - 9 Graph Display of POJ.....	54
Fig 3 - 10 Check If a Legal POJ Syllable Algorithm.....	56
Fig 3 - 11 Result from Online Syllable/Word/Sentence Count System for POJ	58
Fig 3 - 12 Backward Maximal Matching Algorithm for HR Mixed Script.....	59
Fig 4 - 1 Taiwanese Tone Sandhi System Diagram.....	69
Fig 5 - 1 Taiwanese Language POS Tagging System Architecture Diagram	89
Fig A - 1 Han script (Taiwanese Folk Song)	132
Fig A - 2 Han script (Taiwanese Textbook in Japanese-Ruled Period).....	133
Fig A - 3 POJ script (Taiwan Prefectural City Church News).....	134
Fig A - 4 HR Mixed Script (Taiwanese Writing Forum).....	135

Abbreviations

Abbrev.	Taiwanese (POJ)	English	Mandarin
CCA	Bûn-kiàn-hōe	Council for Cultural Affairs	文建會
CCTPLD	Tâi-ôan POJ bûn-hák chu-liâu so̍-chip	The Collection and Cataloging of Taiwanese POJ Literature Data	台灣白話字文 學資料蒐集
CED	Tiong-bûn tiān-chú sû-tián	Chinese electronic dictionary	中文電子辭典
CKIP	Tiong-bûn sû tì-sek-khò̍ sió-cho̍	Chinese Knowledge and Information Processing	中文詞知識庫 小組
CWSTS	Tiong-bûn tng-sû kap Sû- lûi Phiau-kì hē-thóng	Chinese Word Segmentation and Tagging System	中文斷詞及詞 類標記系統
DADWT	Tâi-gú-bûn sò-ūi tián-chông chu-liâu-khò̍	Digital Archive Database for Written Taiwanese	台語文數位典 藏資料庫
HMM	Ún Markov mô-hêng	Hidden Markov Model	隱馬可夫模型
LDC	Gú-giân Chu-liâu Tjong- sim	Language Data Consortium	語言資料中心
MEMM	Siōng tōa lōan-tō̍ Markov mô-hêng	Maximal Entropy Markov Model	最大熵馬可夫 模型
MOE	Kàu-iòk-pō̍	Ministry of Education	教育部
NMTL	Tâi-ôan bûn-hák-kóan	National Museum of Taiwan Literature	台灣文學館
NSC	Kok-kho-hōe	National Science Council	國科會
OTCS	Tâi-gú-bûn gú-sû kiám-sek hē-thóng	Online Taiwanese Concordancer System	台語文語詞檢 索系統
OTJDDT	Tâi-Jit tōa sû-tián Tâi-gú ék-pún	Online Taiwanese- Japanese Dictionary with Taiwanese Translation	台日大辭典台 語譯本

Abbrev.	Taiwanese (POJ)	English	Mandarin
OTMD	Tâi-hôa sòa ⁿ -téng sū-tián	Online Taiwanese- Mandarin Dictionary	台華線上辭典
OTSD	Tâi-gú sòa ⁿ -téng jī-tián	Online Taiwanese Syllable Dictionary	台語線上字典
POJ	Peh-ōe-jī	Taiwanese vernacular writing	白話字
POS	Sū-sèng	Part of speech	詞性
SMHLA	Bân-Kheh-gú tián-chông	Southern Min and Hakka Language Archive”	閩客語典藏
TAICORP	Tâi-ôan jī-tông gú-liâu-khò	Taiwan Child Language Corpus	台灣兒童語料庫
THMLW	Tâi-gú kap Kheh-gú hián-tâi bûn-hák choan-tôe bāng-chām	Taiwanese and Hakka Modern Literature Website	台語及客語現代文學專題網站
TL	Tâi-ôan Bân-lâm-gú Lô-má-jī pheng-im hē-thóng hong-àn	Taiwan Southern Min Lô-má-jī phonetic scheme	臺灣閩南語羅馬字拼音方案
TVLA	Tâi-ôan Peh-ōe-jī bûn-hiàn chu-liâu-kóan	Taiwanese Vernacular Literature Archive	台灣白話字文獻資料館