

附加詞類訊息的台語語詞搭配佇教學上的應用

Hù-ka sù-luī sìn-sit ê Tâi-gí gí-sû tah-phuê tī kàu-hák siōng ê ìng-iōng

楊允言，大漢技術學院、資訊工程系，助理教授。

Iûnn Ún-giân, Tuā-hàn Kì-sút Hák-īnn, Tsu-sìn Kang-tīng Hē, tsōo-lí kàu-siū.

劉杰岳，拋荒台文工作室，負責人。

Lâu Kiát-gák, Pha-hng Tâi Bûn Kang-tsok-sik, hū-tsik-jîn.

陳鄭弘堯，台語文工作者。

Tân-Tinn Hông-giâu, tâi gí-bûn kang-tsok-tsiá.

陳柏中，國立清華大學、物理系，助理教授。

Tân Peh-tiong, Kok-lip Tshing-huâ Tâi-hák, Bûn-lí Hē, tsōo-lí kàu-siū.

摘要

Tiah-iàu

語詞搭配(collocation)是學習語詞按怎使用真好的工具，

Gí-sû tah-phuê(collocation) sī hák-sip gí-sû án-nuá sú-iōng tsin hó ê kang-khū,

這會當透過 對語料的互訊息(mutual information) 恰 相關度(correlation) 的統計來得著，毋過嘛會產生袂少無必要的資料(noise)。

tse ē-tàng thau-kuè tuì gí-liâu ê Hōo-sìn-sit (mutual information) kah Siong-kuan-tōo (correlation) ê thóng-kè lâi tit-tiōh, m̄-koh mā ē sán-sing bē-tsió bô pit-iàu ê tsu-liâu (noise).

本文提出利用詞類的訊息來進行篩選 ê 方法，自動共無必要的資料提掉，提升自動台語語詞搭配的品質，幫贊台語語詞的學習。

Pún-bûn thê-tshut lī-iōng sù-luī sìn-sit lâi tsin-hîng thai-suán ê hong-huat, tsū-tōng kā bô pit-iàu ê tsu-liâu thêh-tiâu, thê-sing tsū-tōng Tâi-gí gí-sû tah-phuê ê phín-tsit, pang-tsân Tâi-gí gí-sû ê hák-sip.

實驗結果顯示，利用詞類ê篩選，動詞—名詞(V—N)詞組比形容詞—名詞(A—N)詞組會當得著較好的結果；

Sit-giām kiat-kó hián-sī, lī-iōng sū-luī ê thai-suán, tōng-sū — bīng-sū (V—N) sū-tsoo pí hīng-iōng-sū — bīng-sū (A—N) sū-tsoo ē-tàng tit-tiōh khah-hó ê kiat-kó;

查詢一个動詞正片的名詞 抑是一个名詞倒片的動詞的語詞搭配，所揣出來的資料嘛真有學習參考的價值。

tshâ-sūn tsit-ê tōng-sū tsiānn-pīng ê bīng-sū iah-sī tsit-ê bīng-sū tò-pīng ê tōng-sū ê gí-sū tah-phuè, sóo tshuē--tshut-lâi ê tsu-liāu mā tsin ū hák-síp tsham-khó ê kè-tát.

毋過若欲進一步提升自動語詞搭配的品質，咱猶閣有誠濟空課愛做。

M̄-koh nā beh tsin-tsit-pōo thê-sing tsū-tōng gí-sū tah-phuè ê phín-tsit, lán iah-koh ū tsiānn-tsē khang-khuè ài tsò.

關鍵詞：台語文、語詞搭配、互訊息、相關度、語文教學

Kuan-kiān-sū: Tâi-gí-bûn, Gí-sū Tah-phuè, Hōo-sìn-sit, Siong-kuan-tōo, Gí-bûn Kàu-hàk

1. 踏話頭：語料庫、語詞搭配

It. Tàh uē-thâu: Gí-liāu-khòo, Gí-sū Tah-phuè

「搭配」這個詞，台語閣會使講「參佱」、「照佱」、「sann 佱」、「配搭」、...等等，毋過為著佱華語詞對應，所以採用「搭配」。

"Tah-phuè" tsit-ê sū, Tâi-gí koh ē-sái kóng "tsham-kah", "tsiàu-kah", "sann-kah", "phuè-tah", ... tít-tít, m̄-koh uī-tiōh kah Huā-gí-sū tui-ìng, sóo-í tshái-iōng "tah-phuè".

語料庫就是因語言材料的倉庫，伊會使講是人類所講、所寫的語言文字的一个樣本，

Gí-liāu-khòo tō-sī khng gí-giân tsai-liāu ê tshng-khòo, i ē-sái kóng sī jîn-luī sóo kóng, sóo siá ê gí-giân bûn-jī ê tsit-ê iūnn-pún,

透過這個樣本，咱通透過語料庫了解這個語言、文字的使用情形。

Thàu-kuè tsit ê iunn-pún, lán thang thàu-kuè gí-liâu-khò liáu-kái tsit-ê gí-giân, bûn-jī ê sú-iōng tsing-hîng.

二十世紀 60 年代，英文的語料庫開始建立。

Jī-tsáp sè-kí liók-khòng nî-tāi, Ing-bûn ê gí-liâu-khò khai-sí kiàn-lip.

90 年代，華文的語料庫嘛開始佇台灣、中國等所在建立。

kiú-khòng nî-tāi, Huâ-bûn ê gí-liâu-khò mā khai-sí tī Tâi-uân, Tiong-kok tóng sóo-tsāi kiàn-lip.

台文的起步較晏，一直到這世紀(21 世紀)才開始。

Tâi-bûn ê khí-pō khah uann, it-tit kàu tsit sè-kí (jī-tsáp-it sè-kí) tsiah khai-sí.

對強勢語言來講，語料庫的建立代表這個語言的重要性，因為建立了後會帶動這的語言的相關研究，嘛宣示伊佇學術界的地位；

Tuì kiông-sè gí-giân lâi kóng, gí-liâu-khò ê kiàn-lip tī tsit-ê gí-giân ê tiōng-iàu-sing, in-uī kiàn-lip liáu-āu ē tài-tōng tsit-ê gí-giân ê siōng-kuan gián-kiù, mā suan-sī i tī hák-sút-kài ê tē-uī;

對弱勢語言來講，語料庫的建立是這個語言欲復興的一個重要基礎，因為利用著主導世界行向的電腦科技，而且有機會予這個語言恰其它的強勢語言做連結。

tuì jiók-sè gí-giân lâi kóng, gí-liâu-khò ê kiàn-lip sī tsit-ê gí-giân beh hók-hing ê tsit-ê tiōng-iàu ki-tshóo, in-uī lī-iōng-tiōh tsú-tō sè-kài kiann-hiòng ê tiān-nóo kho-ki, jī-tshiánn ū ki-huê hōo tsit-ê gí-giân kah kī-tha ê kiông-sè gí-giân tò liân-kiat.

語料庫的應用研究有誠濟，佇遮，阮針對語文教學的語詞搭配做較詳細的探討。

Gí-liâu-khò ê ìng-iōng gián-kiù ū tsiann tsē, tī tsia, gún tsiam-tuì gí-bûn kàu-hák ê gí-sū tah-phuê tò khah siōng-suê ê thàm-thó.

語詞搭配(collocation) 就是「語詞用會當預測的方式組合做伙的方式」，
嘛有人講是「全一个文本內底，佇真近的所在共同出現的兩個抑是兩個以上的語詞」。

Gí-sû tah-phuè(collocation) tō-sī "Gí-sû iōng ē-tàng ū-tshik ê hong-sik tsoo-háp tsò-hué ê hong-sik", mā ū-lâng kóng sī "Kāng tsit-ê bûn-pún lâi-tué, ī tsin kūn ê sóo-tsāi kiōng-tông tshut-hiān ê n̄g ê iah-sī n̄g ê í-siōng ê jí-sû".

語詞搭配會使學習著一个語詞 ê 慣用方法，

Gí-sû tah-phuè ē-sái hák-síp-tiōh tsit-ê jí-sû ê kuàn-iōng hong-huat,

比一个例，

pí tsit-ê lē,

咱訂火車單， 閣去車頭付錢提著車單， 咱會使講「車單拍好矣」，

lán tīng hé-tshia-tuann, koh khi tshia-thâu hù-tsīnn thêh-tiōh tshia-tuann, lán ē-sái kóng "tshia-tuann phah hó--ah",

「拍」恰「損」的意思相全， 毋過咱袂使講「車單損好矣」，

"phah" kah "kòng" ê í-sù sio-kâng, m̄-koh lán bē-sái kóng "tshia-tuann kòng hó--ah",

所以，「拍」kah「車單」就有搭配 ê 關係， 其中一个詞若是用同義詞共換掉， 這個搭配關係無一定會繼續存在。

sóo-í, "phah" kah "tshia-tuann" tō ū tah-phuè ê kuan-hē, kî-tiong tsit-ê sū nā-sī iōng tông-gī-sū kā uānn-tiāu, tsit-ê tah-phuè kuan-hē bô-it-tīng ē kè-siók tsūn-tsāi.

兩個語詞欲按怎鬥陣使用， 一般 ê 辭典並無特別 kā 註明， 這對語言 ê 學習加減有阻礙。

N̄g-ê jí-sû beh án-nuá tau-tīn sú-iōng, it-puann ê sū-tián pīng-bô tīk-pià̍t kā tsù-bīng, tse tuì jí-giân ê hák-síp ke-kiám ū tsóo-gāi.

台灣人學英文， 硬記英文語詞，

Tâi-uân-lâng òh Ing-bûn, ngē-kì Ing-bûn gí-sû,

欲用英文佮人溝通 ê 時陣， 利用遮的語詞鬥出一句咱掠做會通 ê 英文句，
毋過講英語 ê 人可能消無貓仔門，

beh iōng Ing-bûn kah lâng koo-thong ê sī-tsūn, lī-iōng tsia ê gí-sû tau tshut tsit-kù lán liàh-tsò ē thong ê Ing-bûn-kù, m̄-koh kóng Ing-gí ê lâng khó-lîng sa bô niau-á-m̄g,

因為咱 ê 語詞搭配錯誤， 錯誤可能包括語法 kah 慣用語。

in-uī lán ê gí-sû tah-phuè tshò-gōo, tshò-gōo khó-lîng pau-kuah gí-huat kah kuàn-iōng-gí.

所以語詞搭配成做語言學習真重要 ê 一部分。

Sóo-í gí-sû tah-phuè tsiānn-tsò gí-giân hāk-sîp tsin tiōng-iàu ê tsit-pōo-hūn.

英文 ê 部分， 出版袂少語詞搭配辭典， 嘛是為著濟濟英語 ê 學習者。

Ing-bûn ê pōo-hūn, tshut-pán bē-tsió gí-sû tah-phuè sū-tián, mā-sī uī-tiōh tsuē-tsuē Ing-gí ê hāk-sîp-tsiá.

台語 ê 學習敢有需要注意語詞搭配？

Tâi-gí ê hāk-sîp kám ū su-iàu tsù-ì gí-sû tah-phuè?

當然有需要， 包括對想 beh 學台語 ê 外國人， 嘛包括咱這代 kah 下一代。

Tong-jiân ū su-iàu, pau-kuah tuī siūnn-beh òh Tâi-gí ê guā-kok-lâng, mā pau-kuah lán tsit-tâi kah ē-tsit-tâi.

受著強勢 ê 華語影響， 誠濟自認為會曉講台語 ê 人， 台語講起來 mā 是 li-li-lak-lak,

Siū-tiōh kiōng-sè ê Huâ-gí íng-hióng, tsiānn tsē tsū-jīm-uī ē-hiáu kóng Tâi-gí ê lâng, Tâi-gí kóng--khí-lâi li-li-lak-lak,

嘛因為用華語思考， 有 ê 人「一尾魚」煞變做「一條魚」。

mā in-uī iōng Huâ-gí su-khó, ū-ê-lâng "tsit-bué hī" suah pinn-tsò "tsit-tiâu hī".

建立台語 ê 語詞搭配辭典是台語教學真重要 ê 一部份，

Kiàn-líp Tâi-gí ê gí-sû tah-phuê sū-tián sī Tâi-gí kàu-hák tsin tiōng-iàu ê tsit-pōo-hūn,

可惜目前猶無，這是語言學界 ê 責任。

khó-sioh bók-tsiân iah bô, tse-sī gí-giân-hák-kài ê tsik-jīm.

本文是倚佇資訊科學 ê 角度，提出現時另外一個通好運用 ê 方式：利用台語文語料庫這個資源來建立台語語詞搭配 ê 資料。

Pún-bûn sī khiā tī tsu-sin kho-hák ê kak-tōo, thê-tshut hiân-sī lîng-guā tsit-ê thang-hó ūn-iōng ê hong-sik: Lī-iōng Tâi-gí-bûn gí-liâu-khò tsit-ê tsu-guân lâi kiàn-líp Tâi-gí gí-sû tah-phuê ê tsu-liâu.

2. 自動語詞搭配做法

Jī. Tsū-tōng Gí-sû Tah-phuê Tsò-huat

根據語料庫，欲共共同出現而且互相關係密切 ê 語詞揣出來，是利用統計 ê 方法。

Kin-kù gí-liâu-khò, beh kā kiōng-tōng tshut-hiân jī-tshiánn hōo-siong kuan-hē bít-tshiat ê gí-sû tshuē--tshut-lâi, sī lī-iōng thóng-kè ê hong-huat.

有兩個公式會使做，

Ū n̄ng-ê kong-sik ē-sái tsò,

第一個是互訊息(mutual information，以下簡稱 MI)，設使 A、B 是語詞，這兩個語詞 ê MI ê 公式是寫做按呢。

tē-it-ê sī hōo-sin-sit (mutual information, í-hā kán-tshing MI), siat-sú A, B sī gí-sû, tsit n̄ng-ê gí-sû ê MI ê kong-sik sī siá-tsò án-ne.

其中，「P of A」是語詞 A ê 機率，

Kī-tiong, "P of A" sī gí-sû A ê ki-lút,

設使 A 攏總出現「Frequency of A」擺，所有 ê 語料攏總有 N 个詞 (word tokens)，「P of A」等於 N 分之「Frequency of A」

siat-sú A lóng-tsóng tshut-hiān "Frequency of A" pái, sóo-ū ê jí-liāu lóng-tsóng ū N-ê sū (word tokens), "P of A" tít-î N hun tsi "Frequency of A".

全款，「P of B」是語詞 B ê 機率，「P of AB」是詞組 AB ê 機率。

Kāng-khuán, "P of B" sī jí-sū B ê ki-lút, "P of AB" sī sū-tsoo AB ê ki-lút.

咱分三種情形來討論：

Lán hun sann-tsióng tsîng-hîng lâi thó-lūn:

(a) 設使 A 後壁干焦會出現 B，B 頭前干焦會出現 A，這兩個詞一定成對出現，

(a) Siat-sú A āu-piah kan-na ē tshut-hiān B, B thâu-tsíng kan-na ē tshut-hiān A, tsit n̄ng-ê sū it-tīng sīng-tui tshut-hiān,

這款情形下，「P of A」近倚「P of B」近倚「P of AB」，「MI of AB」近倚「負 log P of AB」，是一個較大 ê 正數；

tsit-khuán tsîng-hîng hā, "P of A" kīn-uá "P of B" kīn-uá "P of AB", "MI of AB" kīn-uá "hū log P of AB", sī tsit-ê khah tuā ê tsiann-sò;

若是 AB 定定成對出現，毋過 A 後壁會接 B 以外其它 ê 詞，B 頭前嘛可能會接 A 以外其它 ê 詞，

nā-sī AB tiānn-tiānn sīng-tui tshut-hiān, m̄-koh A āu-piah ē tsia p B í-guā ê sū, B thâu-tsíng mā khó-līng ē tsia p A í-guā ê sū,

這個時陣，「MI of AB」可能嘛是正數，毋過會較細；

tsit-ê sī-tsūn, "MI of AB" khó-līng mā-sī tsiann-sò, m̄-koh ē khah suè;

(b) 設使 A 恰 B 無關係（獨立事件），

(b) Siat-sú A kah B bô kuan-hē (tók-líp sū-kiānn),

「P of AB」近倚"P of A"乘"P of B"，所以「MI of AB」近倚 0；

"P of AB" kīn-uá "P of A" sîng "P of B" , sóo-í "MI of AB" kīn-uá 0;

(c) 設使語詞 A 出現致使 B 較袂出現，

(c) Siat-sú gí-sû A tshut-hiān tì-sú B khah buē tshut-hiān,

「P of AB」遠小於"P of A" 乘 "P of B" , 所以「MI of AB」是負數。

"P of AB" uán sió í "P of A" sîng "P of B" , sóo-í "MI of AB" sī hū-sòo .

基本上，語料愈大，統計 ê 結果愈有代表性。

Ki-pún siōng, gí-liāu jú tuā, thóng-kè ê kiát-kó jú ū tāi-piáu-sìng.

第二个是相關度(correlation，以下簡稱 CR)，這嘛是一个統計 ê 公式：

Tē-jī-ê sī siong-kuan-tōo (correlation, í-hā kán-tshing CR), tse mā-sī tsít-ê thóng-kè ê kong-sik:

簡單講，欲算「CR of AB」，就共語料內底 ê 詞組分做四个部分，透過頂面 ê 公式來計算。

Kán-tan kóng, beh sng "CR of AB", tō kā gí-liāu lâi-té ê sū-tsoo hun-tsò sì ê pōo-hūn, thau-kè tít-bīn ê kong-sik lai kè-sng.

算出來 ê 數字攏是正數，

Sng--tshut-lâi ê sòo-jī lóng-sī tsiānn-sòo,

AB 兩個語詞若定定做陣出現，「CR of AB」可能超過 10,000，甚至超過 100,000。

AB nn̄g-ê gí-sû nā tiānn-tiānn tsò-tīn tshut-hiān, "CR of AB" khó-ling tshiau-kuè tsít-bān, sīm-tsì tshiau-kuè tsáp-bān.

愛注意 ê 是，有時仔若拍字錯誤，因為錯字 ê 詞頻真低，會致使 MI 恰 CR 的分數變足懸。

Ài tsù-ì ê sī, ū-sī-á nā phah-jī tshò-gōo, in-uī tshò-jī ê sū-pîn tsin kē, ē tì-sú MI kah CR ê hun-sòo piàn tsiok kuân.

為著避免這個情形，會使設詞頻 ê 限制，

Uī-tio̍h phiah-bián tsit-ê tsīng-hīng, koh ài siat sū-pîn ê hān-tsè,

就是講，欲計算 MI(AB)，會使限制詞組 AB 的頻率超過一個數目（可比出現 10 擺以上），才來計算，若無就共提掉。

tō-sī kóng, beh kè-sng MI(AB), ē-sái hān-tsè sū-tsoo AB ê pîn-lùt tshiau-kè tsit-ê sòo-bók (khó-pí tshut-hiān tsáp-pái í-siōng), tsiah lâi kè-sng, nā-bô tō kā thêh-tiāu.

3. 利用詞類訊息提升品質

Sann. Lī-iōng sū-luī sìn-sit thê-sing phín-tsit

用統計 ê 方式來揣語詞搭配，恰語言專家利用伊專業智識揣出來 ê，當然加減會無全款。

Iōng thóng-kè ê hong-sik lâi tshuē gí-sū tah-phuè, kah gí-giân tsuan-ka lī-iōng i tsuan-giap tì-sik tshuē--tshut-lâi--ê, tong-jiân ke-kiám ē bô kāng-khuán.

用統計方式，咱有簡化，限定佇相倚 ê 兩個語詞，

Iōng thóng-kè hong-sik, lán ū kán-huà, hān-tiānn tī sio-uá ê n̄g-ê gí-sū,

實際上有 ê 互相搭配 ê 語詞並無相倚，「拍車單」有可能講做「拍兩張車單」，

sit-tsè siōng ū-ê hōo-siong tah-phuè ê gí-sū pīng-bô sio-uá, "phah tshia-tuann" ū khó-līng kóng tsò "phah n̄g tiunn tshia-tuann",

嘛有可能是三个語詞互相搭配，可比「除了……以外，猶閣……」。

mā ū khó-lîng sī sann-ê gí-sû hōo-siong tah-phuè, khó-pí "tû-liáu ... í-guā, iah-koh ...".

當然咱會使修正統計方式，共無相倚毋過離無遠 ê 語詞嘛考慮入來，抑是修改公式來計算三个語詞 ê 互相關係。

Tong-jiân lán ē-sái siu-tsing thóng-kè hong-sik, kā bô sio-uá m̄-koh lī bô hñg ê gí-sû mā khó-lū--jíp-lâi, iah-sī siu-kái kong-sik lâi kè-sng sann-ê gí-sû ê hōo-siong kuan-hē.

總是按呢做加足厚工，增加誠濟計算量，改進 ê 成果嘛真有限。

Tsóng--sī an-ne tsò ke tsiok kâu-kang, tsing-ka tsiân tsē kè-sng-liōng, kái-tsìn ê sîng-kó mā tsin iú-hān.

用統計方式做，並毋是欲取代專家，

Iōng thóng-kè ê hong-sik tsò, pīng m̄-sī beh tshú-tāi tsuan-ka,

毋過，咱利用電腦會當處理大量資料 ê 特性，提供電腦計算出來 ê 結果，對學習者恰專家攏有參考作用，

m̄-koh, lán lī-iōng tiān-nóo ē-tàng tshú-lí tuā-liōng tsu-liáu ê tik-sing, thê-kiong tiān-nóo kè-sng--tshut-lâi ê kiát-kó, tuì hák-síp-tsiá kah tsuan-ka lóng ū tsham-khó tsok-iōng,

專家所列出來 ê 語詞搭配，有 ê 因為實際上較少使用，致使統計結果並無表現出來；

tsuan-ka sóo liát--tshut-lâi ê gí-sû tah-phuè, ū-ê in-uī sít-tsè siōng khah tsió sú-iōng, tì-sú thóng-kè kiát-kó pīng-bô piáu-hiān--tshut-lâi;

統計結果有，毋過專家無列出來 ê，有可能並無什麼意義，

thóng-kè kiát-kó ū, m̄-koh tsuan-ka bô liát--tshut-lâi ê, ū khó-lîng pīng bô siānn-mé ì-gī,

總是伊佇實際使用上有較懸的牽連，嘛提供專家參考恰研究 ê 方向。

tsóng--sī i tī sít-tsè sú-iōng siōng ū khah kuān ê khan-liân, mā thê-kiong tsuan-ka tsham-khó kah gián-kiù ê hong-hiōng.

另外，根據英文 ê 語料庫做出來 ê 語詞搭配，介詞(in, at,...)、指示代名詞(that)、冠詞(a, the, ...) 因為頻率真懸，遮 ê 語詞 ê 語詞搭配("of the", "that the", "to be", "in a", ...) 嘛攏排佇誠頭前。

Līng-guā, kun-kù Ing-bûn ê gí-liāu-khò tsò--tshut-lâi ê gí-sû tah-phuè, kài-sû (in, at,...), tsí-sī tãi-bīng-sû (that), kuàn-sû (a, the,...) in-uī pîn-lút tsin kuân, tsia ê gí-sû ê gí-sû tah-phuè ("of the", "that the", "to be", "in a", ...) mā lóng pài tī tsiānn thâu-tsīng.

對遮咱發現，詞類會使提供咱一个真好 ê 線索，

Uì tsia lán huat-hiān, sū-luī ē-sái thê-kiong lán tsit-ê tsiānn hó ê suānn-soh,

假使咱會當利用詞類 ê 訊息，對統計出來 ê 語詞搭配結果做分類整理，對學習者來講應該閣較有幫贊。

ká-sú lán ē-tàng lī-iōng sū-pîn ê sìn-sit, tuì thóng-kè--tshut-lâi ê gí-sû tah-phuè kiāt-kó tsò hun-luī tsíng-lí, tuì hāk-sip-tsiá lâi kóng ìng-kai koh-khah ū pang-tsān.

有詞類訊息 ê 台語辭典，包括日本天理大學 ê 《現代閩南語辭典》(村上嘉英編，1981 年)、台北中華語文研習所 ê 《台英辭典》(Embree 編，1984)、...等，可惜並無電子檔案通好運用。

Ū sū-luī sìn-sit ê Tâi-gí sū-tián, pau-kuah Jit-pún Thian-lí tãi-hák ê "Hiān-tâi Bân-lâm-gí sū-tián"(Murakami Yoshihide pian, 1981 nî), Tâi-pak Tiong-huâ gí-bûn gián-sip-sóo ê "Tâi-ing sū-tián" (Embree pian, 1984), ... tít, khó-sioh pīng-bô tiān-tsú tóng-àn thang-hó ūn-iōng.

所以阮只好撻一輶，以線頂台文華文辭典做基礎，將台文 ê 詞條透過華文對譯，對到中研院詞庫小組整理 ê 八萬目華文詞條，對華文詞條得著詞類訊息。

Sóo-í gún tsí-hó seh tsit-lín, í suānn-tít Tâi-bûn Huâ-bûn sū-tián tsò ki-tshóo, tsiong Tâi-bûn ê sū-tiáu thau-kuè Huâ-bûn tuì-lk, tuì kàu Tiong-gián-īnn Sū-khò sió-tsoo tsíng-lí ê peh-bān-bók Huâ-bûn sū-tiáu, uì Huâ-bûn sū-tiáu tit-tiòh sū-luī sìn-sit.

詞庫小組 ê 詞類分較幼，有 46 个詞類，包括動詞、名詞等攏繼續有分細類落去，

Sū-khò sió-tsoo ê sū-luī pun khah iù, ū 46-ê sū-luī, pau-kuah tōng-sū, bīng-sū tít lóng kè-siok ū pun sè-luī lòh--khì,

毋過計算語詞搭配，咱無需要分遐爾幼，用大類來分就好，包括名詞、動詞、形容詞、副詞、連接詞、介詞、語氣詞、感嘆詞等等。

m̄-koh kè-sng̤ gí-sû tah-phuè, lán bô su-iàu pun hiah-nī iù, iōng tuā-luī lâi pun tō hó, pau-kuah bīng-sû, tōng-sû, hīng-iōng-sû, hù-sû, liân-tsiap-sû, kài-sû, gí-khì-sû, kám-thàn-sû tít-tít.

4. 實驗步數

Sit-giām pōo-sòo

下面是本實驗 ê 做法：

Ē-bīn sī pún sit-giām ê tsò-huat:

(a) 資料包括台語文語料、台文華文線頂辭典（以下簡稱台華辭典）佾中研院詞庫小組八萬目詞(華文)。

(a) Tsu-liâu pau-kuah Tâi-gí-bûn gí-liâu, Tâi-bûn Huâ-bûn suann-tíng sū-tián (í-hā kán-tshing Tâi-huâ sū-tián) kah Tìong-gián-īnn Sū-khò-sió-tsoo pueh-bān-bák sū (Huâ-bûn).

其中，台語文語料庫 ê 來源是台語文界 ê 朋友，楊允言負責整理，

Kî-tiong, Tâi-gí-bûn gí-liâu-khò ê lâi-guan sī Tâi-gí-bûn-kài ê pīng-iú, Iūnn Ún-giân hū-tsik tsíng-lí,

包括全羅 3,462,367 个音節 佾漢羅 5,568,057 个音節；

pau-kuah tsuân-lô sann-pah sì-tsáp-lák-bān nīng-tshing sann-pah lák-tsáp-tshit ê im-tsiat kah Hàn-lô gōo-pah gōo-tsáp-lák-bān peh-tshing khòng gōo-tsáp-tshit ê im-tsiat;

台華辭典主要貢獻者是鄭良偉，閣有台語文界朋友鬥補充詞條，目前有 6 萬 2 千外个詞條，每一个詞條包括漢羅、全羅、華文對譯三个欄位。

Tâi-huâ sū-tián tsú-iàu kòng-hiàn-tsiá sī Tēnn Liông-uí, koh ū Tâi-gí-bûn-kài pīng-iú tau póo-tshiong sū-tiâu, bók-tsiân ū lák-bān nīng-tshing-guā-ê sū-tiâu, muí tsit-ê sū-tiâu pau-kuah Hàn-lô, tsuân-lô, Huâ-bûn tui-ik sann-ê nuâ-uī.

(b) 全羅 ê 台語文語料無需要經過斷詞，毋過漢羅 ê 需要，

(b) Tsuân-lô ê Tâi-gí-bûn gí-liâu bô su-iàu keng-kè tñg-sû, m̄-koh Hàn-lô--ê su-iàu,

斷詞 ê 做法是根據台華辭典，採用「倒頭上大比對」(backward maximal matching, BMM)演算法。

tñg-sû ê tsò-huat sī kun-kù Tâi-huâ sū-tián, tshái-iōng "tò-thâu siōng-tuā pí-tuì" (backward maximal matching, BMM) ián-suàn-huat.

語詞 ê 數量，全羅有 2,436,599 个，漢羅有 4,051,195 个。

Gí-sû ê sòo-liōng, tsuân-lô u nñg-pah sì-tsáp-sann-bān lāk-tshing gōo-pah káu-tsáp-káu ê, Hàn-lô ū sì-pah khòng gōo-bān tsít-tshing tsít-pah káu-tsáp-gōo ê.

(c) 每一个語詞，透過台華辭典 ê 華文對譯，對到中研院詞庫小組八萬目詞 ê 詞條，揣出詞類，詞類可能有幾 lō 个，暫時無做處理。

(c) Muí tsít-ê gí-sû, thàu-kè Tâi-huâ sū-tián ê Huâ-bûn tuì-ik, tuì kàu Tìong-gián-īnn sū-khò sió-tsoo peh-bān-bāk sū ê sū-tiâu, tshuē-tshut sū-luī, sū-luī khó-līng ū kuí-lō ê, tsiām-sī bô tsò tshú-lí.

因為無處理含糊性問題，所以會影響著實驗結果 ê 品質。

In-ūi bô tshú-lí hām-hōo-sing ê bûn-tuê, ē íng-hióng-tiōh sít-giām kiāt-kó ê phín-tsit.

另外，本實驗採用簡化 ê 詞類，干焦分大類：動詞、名詞、形容詞、...等等，

Līng-guā, pún sít-giām tshái-iōng kán-huà ê sū-luī, kan-na pun tuā-luī: tōng-sū, bīng-sū, hīng-iōng-sū, ... tít-tít,

其中，動詞 VH 小類是「狀態不及物述詞」，親像「浪漫」、「特別」、「辛苦」、「豐富」、「心酸」、「感動」，阮共改做形容詞。

kī-tiong, tōng-sū VH sió-luī sī "Tsōng-thài put-kiap-bút sūt-sū", tshin-tshiunn "lōng-bān", "tīk-piāt", "sin-khóo", "hong-hù", "sim-sng", "kám-tōng", gún kā kái-tsò hīng-iōng-sū.

按呢修改，有 ê 無問題，部分可能有問題。

An-ne siu-kái, ũ-ê bô bûn-tuê, pōo-hûn khó-lîng ũ bûn-tuê.

(d) 計算兩個相倚 ê 語詞的 MI 佻 CR，閣經過詞類 ê 篩選。

(d) Kè-sng nîg-ê sio-uá ê gí-sû ê MI kah CR, koh king-kè sū-luī ê thai-suán.

因為資料量誠大，這擺實驗干焦揣 A-N（形容詞——名詞）、V-N（動詞——名詞）詞組。

In-ūi tsu-liâu-liōng tsiānn tuā, tsit-pái sít-giām kan-na tshuē A-N (hîng-iōng-sû —— bîng-sû), V-N (tōng-sû —— bîng-sû) sū-tsoo.

5. 實驗結果

Gōo. Sít-giām kiát-kó

本實驗 ê 程式用 Java 寫，共結果輸出到資料庫，資料量誠大。

Pún sít-giām ê thîng-sik iōng Java siá, kā kiát-kó su-tshut kàu tsu-liâu-khò, tsu-liâu-liōng tsiānn tuā.

本節利用表格佻文字說明來展示實驗成果。

Pún-tsiat lī-iōng pió-keh kah bûn-jī suat-bîng lâi tián-sī sít-giām sîng-kó.

5.1 MI 佻 CR 共同詞組數量

Gōo tiám it MI kah CR kiōng-tōng sū-tsoo sòo-liōng

表一 共 MI 佻 CR 頭前 30 / 100 / 500 / 1,000 / 3,000 / 5,000 个共同詞組 ê 數量列出來，漢羅佻全羅分開算，待討論：(這個詞組，佇 MI 佻 CR 兩個統計表攏有出現)

Pió-it kā MI kah CR thâu-tsîng 30 / 100 / 500 / 1,000 / 3,000 / 5,000 ê kiōng-tōng sū-tsoo ê sòo-liōng liát--tshut-lâi, Hàn-lô kah tsuân-lô hun-khui sng, Thāi thó-lūn: (tsit-ê sū-tsoo , tī MI kah CR nîg-ê thóng-kè-pió lóng-ū tshut-hiān)

其中，漢羅 ê 部分，詞組至少出現 10 擺，全羅 ê 部分，詞組至少出現 5 擺：

kî-tiong, Hàn-lô ê pōo-hūn, sū-tsoo tsì-tsió tshut-hiān 10-pái, tsuân-lô ê pōo-hūn, sū-tsoo tsì-tsió tshut-hiān 5-pái:

表一 互訊息 kah 相關度 共同詞組數量比例

Piáu tsit hōo-sìn-sit kah siang-kuan-tōo kiōng-tōng sū-tsoo sòo-liōng pí-lē

對表一看起來，漢羅 ê 部分，MI 佻 CR 算出來 ê 結果較一致，

Uì pió-it khuànn--khí-lâi, Hàn-lô ê pōo-hūn, MI kah CR sng--tshut-lâi ê kiát-kó khah it-tì,

這應該是因為，漢羅 ê 語料量較濟，統計出來 ê 結果較有可信度。

tse ìng-kai sī in-uī, Hàn-lô ê gí-liâu-liōng khah tsē, thóng-kè--tshut-lâi ê kiát-kó khah ū khó-sìn-tōo.

5.2 A-N V-N 詞組

Gōo-tiám jī A-N V-N sū-tsoo

本實驗用詞類是 A-N 佻 V-N 做篩選，

Pún sít-giām iōng sū-luī sī A-N kah V-N tsò thai-suán,

只要這個詞組有可能是 A-N 抑是 V-N 就共掠出來，

tsí-iàu tsit-ê sū-tsoo ū khó-líng sī A-N iah-sī V-N tō kā liáh--tshut-lâi,

毋過因為雜訊缺少，表 2 佻表 3 是對漢羅佻全羅，MI ê 統計表對懸到低分別用人工搵出是 A-N 抑是 V-N 的詞組：

m̄-koh in-uī tsàp-sìn bē tsíó, pió-jī kah pió-sann sī uì hàn-lô kah tsuân-lô, MI ê thóng-kè-pió uì kuân kàu kē hun-piát iōng lāng-kang tshuē-tshut sī A-N iah-sī V-N ê sū-tsoo:

A-N 詞組較僥揣，揀出來 ê 有 ê 看起來誠勉強。

A-N sū-tsoo khah oh tshuē, kóng--tshut-lâi--ê ū-ê khuànn--khí-lâi tsiānn bián-kióng.

親像漢羅部分 ê 「旺梨」、「ām 瓜」，自本應該是一個詞，

Tshin-tshiūnn Hàn-lô pōo-hūn ê "ōng-lâi", "ām kue", tsū-pún ìng-kai sī tsit-ê sū,

毋過台華辭典內底有「王梨」無「旺梨」，有「醃瓜」無「ām 瓜」，
致使這幾個詞予斷詞系統切做兩個詞，

m̄-koh Tâi-huâ sū-tián lâi-tué ū "ōng-lâi" bô "ōng-lâi", ū "am-kue" bô "ām-kue", tì-sú tsit kuí-ê sū hōo tñg-sū
hē-thóng tshiat-tsò nñg-ê sū,

煞閣因為切開 ê 「這兩個詞」ê 結合性懸，佇 A-N 詞組表現出來。

suah koh in-uī tshiat-khui ê "tsit nñg-ê sū" ê kiat-háp-sìng kuân, tī A-N sū-tsoo piáu-hiân--tshut-lâi .

對這個角度看，自動語詞搭配統計表，對辭典詞條 ê 收錄，會當提供一
寡建議。

Uì tsit-ê kak-tōo khuànn, tsū-tōng gí-sū tah-phuè thóng-kè-pió, tuì sū-tián sū-tiâu ê siu-lók, ē-tàng thê-kiong
tsit-kuá kiàn-gī.

V-N 詞組 ê 部分，全羅 ê 部分真好，對 MI 統計表頭前 64 個詞組內底就通
揀出真正符合 V-N ê 詞組，毋過漢羅 ê 部分就誠 bái。

V-N sū-tsoo ê pōo-hūn, tsuân-lô ê pōo-hūn tsin hó, uì MI thóng-kè-pió thâu-tsing 64-ê sū-tsoo lâi-té tō
thang kóng-tshut tsin-tsiann hù-háp V-N ê sū-tsoo, m̄-koh Hàn-lô ê pōo-hūn tō tsiann bái.

主要 ê 原因是斷詞，因為語料內底 ê 寫法佮辭典無一致。

Tsú-iàu ê guân-in sī tñg-sū, in-uī gí-liâu lâi-té ê siá-huat kah sū-tián bô it-tì.

5.3 名詞 ê 動詞搭配佮動詞 ê 名詞搭配

Gōo-tiám sann Bêng-sū ê tōng-sū tah-phuè kah tōng-sū ê bing-sū tah-phuè

表四 查詢「V-舌」

Pió-4 Tshâ-sūn "V-tsih"

「MI 小於零」ê 部分無列出來

"MI sió-î ling5" ê pōo-hūn bô liat--tshut-lâi

表五 查詢「損-N」

Piô-5 Tshâ-sûn "kòng-N"

假使咱欲知影一个名詞頭前會使用 siánn-mé 動詞（倒ㄟ搭配），抑是一个動詞後壁會用使用 siánn-mé 名詞（正ㄟ搭配），用 MI 抑是 CR 來查詢，查出來 ê 結果真有參考價值，

Ká-sú lán beh tsai-iánn tsít-ê bîng-sû thâu-tsîng ē-sái iōng siánn-mé tōng-sû (tò-pîng tah-phuè), iah-sī tsít-ê tōng-sû āu-piah ē-sái iōng siánn-mé bîng-sû (tsiánn-pîng tah-phuè), iōng MI iah-sī CR lâi tshâ-sûn, tshâ--tshut-lâi ê kiát-kó tsin ū tsham-khó kè-tát,

舉一个例，表四用「V-舌」查詢，表五用「kòng-N」查詢：

kú tsít ê lē, pió-si iōng "V-tsih" tshâ-sûn, pió-gōo iōng "kòng-N" tshâ-sûn:

這個查詢方式，對語言學習應該有真大 ê 幫贊。

Tsit-ê tshâ-sûn hong-sik, tuì jí-giân hák-sip ìng-kai ū tsin tuā ê pang-tsân.

5.4 主要問題

5.4 Tsú-iàu bûn-tuê

實驗 ê 結果並無原來按算 ê 遐爾好，下面分析問題出佇佗位：

Sit-giām ê kiát-kó pîng bô guân-lâi àn-sng ê hiah-nī hó, ē-bîn hun-sik bûn-tê tshut tī tó-uī:

(1) 漢羅用字無一致，致使斷詞效果無好：

(1) Hàn-lô iōng-jī bô it-tì, tì-sú tñg-sû hâu-kó bô hó:

V-N ê 部分，漢羅 ê 結果真 bái，

V-N ê pōo-hūn, Hàn-lô ê kiap-kó tsin bái,

人工 tih 看遮的詞組 ê 時，發現有誠濟詞組是斷詞錯誤 ê，

lâng-kang tih khuànn tsia-ê sū-tsoo ê sī, huat-hiān ū tsiānn-tsē sū-tsoo sī tñg-sū tshò-gōo--ê,

有 ê 是台華辭典內無這個詞，

ū-ê sī Tâi-huâ sū-tián lâi bô tsit-ê sū,

有 ê 是語料 ê 漢羅寫法 kah 台華辭典內 ê 無仝，實際上，語料內底 ê 寫法就無一致，

ū-ê sī gí-liāu ê Hàn-lô siá-huat kah Tâi-huâ sū-tián lâi ê bô-kâng, sít-tsè siōng, gí-liāu lâi-té ê siá-huat tō bô it-tì,

辭典無 ê 詞，若 beh 改進，包括加詞、專有名詞辨識、定量詞處理等等，寫法無一致 ê 問題 tō 真僉解決，kā 所有 ê 寫法 lóng kā lok tī 辭典內底並 m̄ 是好 ê 解決方式；

sū-tián bô ê sū, nā beh kái-tsin, pau-kuah ka sū, tsuan-iú-bîng-sū piān-sik, tīng-niū-sū tshú-lí tít-tít, siá-huat bô it-tì ê bûn-tuê tō tsin oh kái-kuat, kā sóo-ū ê siá-huat lóng kā lok tī sū-tián lâi-tuê pīng m̄-sī hó ê kái-kuat hong-sik;

(2) 詞類問題：

(2) Sū-luī bûn-tuê:

語言翻譯是多對多，咱透過華語對譯詞掠詞類，這個過程就有可能會 tiòng 一寡無必要 ê 詞類出來；

Gí-giân huan-ik sī to tuì to, lán thau-kè Huâ-gí tuì-ik-sū liáh sū-luī, tsit-ê kè-thîng tō ū khó-lîng ē tiòng tsit-kuá bô pit-iàu ê sū-luī tshut-lâi;

另外，一个詞有至少兩個詞類 mā 是真普遍 ê 情形，因為無處理，造成雜訊誠濟；

līng-guā, tsit-ê sū ū tsi-tsió n̄ng-ê sū-luī mā-sī tsin phóo-phiàn ê tsing-hing, in-ūi bô tshú-lí, tsō-sing tsáp-sin tsiann tsē;

建立台語詞 ê 詞類，恰確認正確詞類，這兩件代誌攏是大工程，毋過早慢愛做，咱才有法度閣進一步。

kiàn-lip Tâi-gí-sū ê sū-luī, kah khak-jīm tsing-khak sū-luī, tsit n̄ng-kiann tsi-tsió lóng-sī tuā-kang-thing, m̄-koh tsá-bān ài tsò, lán tsiah ū-huat-tōo koh tsin-tsit-pōo .

6. 未來方向

Liók Bī-lâi hong-hiòng

本實驗雖然有一寡成果，總是，值得拍拚 ê 所在猶閣誠濟，包括：

Pún sít-giām sui-jiân ū tsit-kuá sing-kó, tsóng--sī, tát-tit phah-piann ê sóo-tsāi iah-koh tsiann tsē, pau-kuah:

(a) 開發線頂台語語詞搭配查詢系統，幫贊台語相關研究 ê 進行；

(a) Khai-huat suann-ting Tâi-gí gí-sū tah-phuè tshâ-sūn hē-thóng, pang-tsân Tâi-gí siong-kuan gián-kiù ê tsin-hing;

(b) 除了 A-N V-N 詞組以外，閣整理其它詞組，親像 P-N(介詞—名詞)、N-V(名詞—動詞)、D-A(副詞—形容詞)、D-V(副詞—動詞)、V-R(動詞—代名詞)、……等；

(b) Tû-liáu A-N V-N sū-tsoo í-guā, koh tsing-lí kī-tha sū-tsoo, tsin-tshiūnn P-N (kài-sū — bing-sū), N-V (bing-sū — tōng-sū), D-A (hù-sū — hing-iōng-sū), D-V (hù-sū — tōng-sū), V-R (tōng-sū — tai-bing-sū), ... ting;

(c) 維護台華辭典，共 MI 抑是 CR 分數較懸 ê 兩個單音節詞 ê 語詞搭配 掠出來做人工檢查，考慮將遮 ê 語詞加入辭典；

(c) Uī-hōo Tâi-huā sū-tián, kā MI iah-sī CR hun-sòo khah kuân ê n̄ng-ê tan-im-tsiat-sū ê gí-sū tah-phuè liáh tshut-lâi tsò jîn-kang kiám-tsa, khó-lū tsiong tsia ê gí-sū ka-jíp sū-tián;

(d) 改進斷詞系統，利用 BMM 演算法結果做訓練資料重做，對有至少兩種斷法 ê 文句，用統計方法提懸正確率，猶閣有專有名詞、定量詞處理；

(d) Kái-tsìn tng-sû hē-thóng, lī-iōng BMM ián-suàn-huat kiāt-kó tsò hùn-liân tsu-liâu tīng tsò, tuì ū tsì-tsió nng-tsióng tng-huat ê bûn-kù, iōng thóng-kè hong-huat thê-kuân tsing-khak-lút, iah-koh ū tsuan-iú bing-sû, tīng-niû-sû tshú-lí;

(e) 整理台語詞類，總是這是大工程，嘛需要詳細 ê 規劃；

(e) Tsing-lí Tâi-gí sū-lūi, tsóng--sī tse-sī tuā kang-thīng, mā su-iàu siōng-sè ê kui-uē;

(f) 利用統計方法自動標記詞類，解決含糊性問題，嘛進一步處理較幼 ê 詞類；

(f) Lī-iōng thóng-kè hong-huat tsū-tōng piau-kì sū-lūi, kái-kuat hām-hô-sing bûn-tê, mā tsin-tsit-pōo tshú-lí khah iù ê sū-lūi;

(g) 做台語文句 ê 分拆(parsing)，利用語法樹(parsing tree)改進語詞搭配 ê 正確率。

(g) Tsò Tâi-gí-bûn bûn-kù ê hun-thiah(parsing), lī-iōng gí-huat-tshiū(parsing tree) kái-tsìn gí-sû tah-phuè ê tsing-khak-lút.