



ROCLING 2014

The 26th International Conference on
Computational Linguistics and Speech Processing



國立中央大學
National Central University

**Proceedings of the Twenty-Sixth Conference on
Computational Linguistics and Speech Processing
ROCLING XXVI (2014)"**

September 25-26, 2014"

National Central University, Zhongli, Taiwan"

Sponsored by:"

Cuuqekvqpp'hqt'Ego r wcvkqpcnNkpi wku'eu'cpf'Ej kpgug'Ncpi wci g"

Rtqeguukpi "

P cvkqpcnEgptcn'Wpkxgtukv{ "

"

Co- Sponsored by:

Academic Sponsor"

Kpu'kwg'qh'Kphqto cvkqp'Uekgpeg.'Cecf go lc'Ukplec"

"

Government Sponsors"

O kpkwt { 'qh'Gf wecvkqp"

O kpkwt { 'qh'Uekgpeg'cpf "Vgej pqmqi { "

"

Industry Sponsors

CUWUVgMEgo r wgt'Kpe0'

E { dgtqp'Eqtr qtcv'kqp"

Ej wpi j y c'Vgrgego 'Ncdqtcvqtkgu"

Kpf wutken'Vgej pqmqi { "T gugctej 'Kpu'kwg"

Hqtvgo gf lc"

Hktuv'Rwdrkuj gf 'Ugr vgo dgt'4236"

D{ 'Vj g'Cuuqekvqp'hqt'Ego r wcvkqpcn'Nkpi wkuu'cpf 'Ej kpgug'Ncpi wci g'Rtqeguulpi " *CENENR+ "

Eqr {tki j 4236'vj g'Cuuqekvqp'hqt'Ego r wcvkqpcn'Nkpi wkuu'cpf 'Ej kpgug' Ncpi wci g'Rtqeguulpi " *CENENR+ 'P cvkqpcn'Egptcn'Wpkxgtuk{ . 'Cwj qtu'qh'Rcr gtu"

Gcej 'qh'vj g'cwj qtu'i tcpvu'c'pqp/gzenwukxg'hlegpug'vq'vj g'CENENR'cpf 'P cvkqpcn' Egptcn'Wpkxgtuk{ 'vq'r wdrkuj 'vj g'r cr gt'lp'r tlpvgf 'hqt o 0Cp{ 'qyj gt 'wuci g'ku'r tqj kdkgf " y kj qw'vj g'gzt tgu'r gto kuukq'qh'vj g'cwj qt'y j q'o c{ "cnuq'tgvckp'vj g'qp/rkp'xgtukq" cv'c'mqecvqp'vq'dg'ugrgevgf "d{ 'j ko lj gt0'

Ej kc/J w'kEj cpi . 'J ukp/O kp'Y cpi . 'Lgp/V| wpi 'Ej kpg. 'J wpi /| w'Mcq. 'Uj kj /J wpi 'Y w' *gf u0'

Rtqeggf kpi u'qh'vj g'Vy gpv{/Ukzvj 'Eqphgtgpeg'qp'Ego r wcvkqpcn'Nkpi wkuu'cpf "

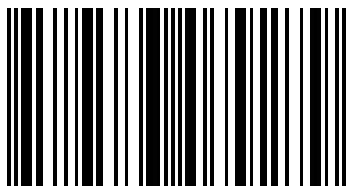
Ur ggej 'Rtqeggf kpi " *TQENR I "ZZXK"

4236/2; /474234/2; /48"

CENENR"

4236/2; "

KUDP <978-957-30792-7-9



Welcome Message from Chairs"

Qp" dgi cñi" qh" vj" g" r tqi tco "eqo o kvgg." k' ku" qwt" r rgcwtg" vq" y graqo g" { qw" vq" vj" g" P cwkqpcñi" Egpwtcñi" Wpłkgtukł. " Lj qpi nk" Vcły cp." hqt" vj" g" 48vj" "Eqphgtgpeg" qp" Eqo r wcvkqpcñi" Nłpi wkłkłeu" cpf "Ur ggej" " Rtqegułłpi " *TQENRPI +." vj" g" hci uj kr" eqphgtgpeg" qp" eqo r wcvkqpcñi" nłpi wkłkłeu." pcwtcñi" rpi vci g" r tqegułłpi ." cpf" ur ggej" r tqegułłpi " łp" Vcły cp' TQENRPI " ku" vj" g" cppwci" eqphgtgpeg" qh" vj" g" Eqo r wcvkqpcñi" Nłpi wkłkłeu" cpf "Ej łpgug" Nłpi vci g" Rtqegułłpi " *CENENR+ y j lej "ku" j grf "łp" cwwop "łp" f hgtgpv'ekłgu" cpf "wpłkgtukłgu" łp" Vcły cp' Vj ku" { gct" y g" tgegkłgf "47" uwdokłkłqu." gcej "qh" y j lej "y cu" tğłkly grf " d{ " y q" vq" hqt" tğłkly gtu" qp" vj" g" dcuku" qh" qtki łpckłł. " pąxgrł. " vgej płeci" uqwpf pguu. " cpf" tğrłcpeg" vq" vj" g" eqphgtgpeg' Vqcmł. " y g" ceegr v' 38" qtcñi" r cr gtu" cpf "7" r quvgt" r cr gtu" y j lej "eqxgt" vj" g" ctgcu" qh" ur ggej " cpf" ur gcngrt" tgeqi płkłqp. " vgzv' o łpłpi. " ur ggej " r tqegułłpi " cpf" u{ pıj guku. " cpf" pcwtcñi" rpi vci g" r tqegułłpi O' Y g" ctg" i tcvghwi" vq" vj" g" eqpłkdwłqp" qh" vj" g" tğłkly gtu" hqt" vj" g" gzłcqtł łpct { " ghqtłv" cpf "xcnwdrg" eqo o gpwł

TQENRPI " 4236" hgcwtgu" vq" f kłkłpi wkłj grf " rgewtgu" łtqo " vj" g" tgpqy pgrf " ur gcngrtu" łp" ur ggej" r tqegułłpi " cu" y gni" cu" pcwtcñi" rpi vci g" r tqegułłpi O' F t' Hłcpni" Uqppi " *Rłkłekr cñi" Tgugctej gt łłgugctej " O cpci gt "qh" O UTC +y kłł' rgewtg" qp " \$Ugctej " hqt" vj" g" Grgo gpwt { " Rctłkrgu" łp" J wo cp " Ur ggej " / " Enwgu" hqt" vj" g" Eqo o qp " Wpku" Cetqu " F hgtgpv' Ur gcngrtu" cpf " Nłpi vci gu" cpf " F t' O J cpi " Nk" *Ej kgh' Uekpłv" qh" J wcy gk' Vgej pąni łgu+ y kłł' ur gcnłqp " \$Ugo cpłk " O cvej łpi < Vj g" P gz v' Dki " Vj łpi " hqt" P cwtcñi" Nłpi vci g" Rtqegułłpi A' S' Vj ku" TQENRPI " cnuq" hgcwtgu" vq" F qevqtcñi" Eqpuqtłkwo u. " qpgr" łpł wut { " Vtceni" cpf " qpgr" Cef go k " F go q" Vtceni" y j lej " r tqxłf g" hqtwo u" cpf " uj qy / cpł / vgn " hqt" Rj F " uwf gpw. " łpł wutłeni" cpf " cef go k " tğugctej gtu" cpf " f gıgrı gtu

Hłpcmł. " y g" cr r tgelcvg" { qwt" gpıj wukłk " r ctłkłk cłkqp" cpf " uwr r qtł' Y kłj gu" c" uweeguhwi" cpf " łt włshwi" TQENRPI " 4236" łp' Lj qpi nk

Ej k/J w' Ej cpi "

J ulp/O łp" Y cpi "

I gpgtcñi" Ej cku"

Lgp/Vł wpi "Ej łgp"

J wpi / [w' Młq"

Rtqi tco "Eqo o kvgg" Ej cku"

Organizing Committee"

Honorary Chair"

Lipi /[cpi 'Iqw'P cvkqpcn'Egptcn'Wplxgtukʌ "

General Chairs"

Ej kc/J wk'Ej cpi . 'P cvkqpcn'Egptcn'Wplxgtukʌ "

J ukp/O kp'Y cpi . 'Cecf go kc'Ukplec"

Advisory Committee"

Lcuqp'UOEj cpi . 'P cvkqpcn'Vulpi 'J wc'Wplxgtukʌ "

J ukp/J uk'Ej gp. 'P cvkqpcn'Vcly cp'Wplxgtukʌ "

Mgj /Ikpp'Ej gp. 'Cecf go kc'Ukplec"

Ukp/J qtpi 'Ej gp. 'P cvkqpcn'Ej kcq'Vwpi 'Wplxgtukʌ "

Y gp/Nkcp'J uw. 'Cecf go kc'Ukplec"

Ej wTgp'J wcpi . 'J qpi 'Mqpi 'Rqn' vgej ple'Wplxgtukʌ "

Ej kp/J wk'Ngg. 'I gqti kc'Kpukswg'qh'Vgej pqrqi { "

Nkp/uj cp'Ngg. 'P cvkqpcn'Vcly cp'Wplxgtukʌ "

J ck| j qw'Nk'Kpukswg'ht'Kphqeo o "Tgugctej "

Ej kp/[gy 'Nkp. 'O letquqhw'Tgugctej 'Cukc"

J gngp'O gpi . 'Ej kpgug'Wplxgtukʌ "qh'J qpi 'Mqpi "

Ikcp'Uw'Kpukswg'ht'Kphqeo o "Tgugctej "

Mgj /[kj "Uw'Dgj cxlqt'F guki p'Eqtr qtcvqp"

J ukcq/Ej wcp'Y cpi . 'P cvkqpcn'Vulpi 'J wc'Wplxgtukʌ "

Lj kpi /Hc'Y cpi . 'P cvkqpcn'Ej gp'Mwpi 'Wplxgtukʌ "

Ej wpi /J ukgp'Y w'P cvkqpcn'Ej gp'Mwpi 'Wplxgtukʌ "

Steering Committee"

Ej kc/Rkpi 'Ej gp. 'P cvkqpcn'Uwp'[cv'Ugp'Wplxgtukʌ "

Lipi /Uj kp'Ej cpi . 'P cvkqpcn'Ej kp cp'Wplxgtukʌ "

Dgtrkp'Ej gp. 'P cvkqpcn'Vcly cp'P qto cn'Wplxgtukʌ "

Mwepi /J wc'Ej gp. 'P cvkqpcn'Vcly cp'Wplxgtukʌ "

J wpi /l cp'I w'P cwkqpcn'Vcly cp'Wpkxgtuk{ "qh'Uekpeg"cpf "Vgej pqrqi { "
 \ j cq/O kpi 'I cq.'P cwkqpcn'Vcly cp'Wpkxgtuk{ "
 Lj /Uj kpi 'Lcpi .'P cwkqpcn'Vcly cp'Wpkxgtuk{ "
 [wcp/Hw'Nkcq.'P cwkqpcn'Vclr gk'Wpkxgtuk{ "qh'Vgej pqrqi { "
 Ej cq/Nkp'Nkw'P cwkqpcn'Ej gpi ej k'Wpkxgtuk{ "
 Y gp/J ukpi 'Nw'P cwkqpcn'Ej gpi 'Mwpi 'Wpkxgtuk{ "
 Uj wEj wcp'Vugpi .'Ceef go kc'Ukplec"
 [wgp/J ukp'Vugpi .'P cwkqpcn'Vcly cp'P qto cn'Wpkxgtuk{ "
 [kj /Tw'Y cpi .'P cwkqpcn'Ej kcq'Vwpi 'Wpkxgtuk{ "

Program Committee Chairs"

Igp/Vl wpi 'Ej kgp.'P cwkqpcn'Ej kcq/Vwpi 'Wpkxgtuk{ "
 J wpi /l w'Mcq.'P cwkqpcn'Ej gpi /Mwpi 'Wpkxgtuk{ "

Publicity Chair"

Tlej ctf 'Vl qpi /J cp'Vuck'P cwkqpcn'Egpt'cn'Wpkxgtuk{ "

Local Arrangement Chair"

Ikc/Ej kpi 'Y cpi .'P cwkqpcn'Egpt'cn'Wpkxgtuk{ "

Publication Chair"

Uj kj /J wpi 'Y w'Ej cq{cpi 'Wpkxgtuk{ "

Industry Track Chair"

Y gp/J ukpi 'Nw'P cwkqpcn'Ej gpi /Mwpi 'Wpkxgtuk{ "

Academic Demo Track Chair"

[gqw'lkpp'Ej gp.'Uqwj gtp'Vcly cp'Wpkxgtuk{ "

Doctoral Consortium Chairs"

Tlej ctf 'Vl qpi /J cp'Vuck'P cwkqpcn'Egpt'cn'Wpkxgtuk{ "
 [w'Vucq.'Ceef go kc'Ukplec"

Keynote 1 - "

Search for the “Elementary Particles” in Human Speech – Clues for the Common Units Across Different Speakers and Languages"



Frank Soong"

Rtlpkr cniTgugctej gt"ITgugctej 'O cpci gt'qh'O UTC"

Vj wtuf c{ .Ugr vgo dgt'47"

32-22"6"33-22"

Nqecvqp-<Kp>gtpcvqpcniEqphgt gpeg'Egpgvt"

Biography

HicpmiUqqpi "ku" c"Rtlpkr cniTgugctej gt" cpf "Tgugctej "O cpci gt'qh' y g"Ur ggej "I tqwr ."y j gtg"ur ggej " o qf grlpi . "tgeqi pklqp. "u{pvj guku"tuguctej "ku"eqpf wevfg 0J g"tgegkxgf "j ku"DU."O U'cpf "Rj 0F . "cm'kp"GG" Htqo " y j g" P cvkqpcni' Vcly cp" Wplxgtukx{ . " y j g" Wplxgtukx{ " qh" Tj qf g" Kircpf" cpf " Ucphqtf " Wplxgtukx{ . " tgur gevkgxgf 0J g"lqkpgf "Dgm'Ncdu" Tgugctej . "O wtct{ "J km" P L" WUC" kp"3; : 4. "y qtngf "y j g" hqt"42" { gctu" cpf "tgvtgf "cu" c" F kvkpi wkuj gf "O go dgt'qh" Vgej plecni' Uclh' kp"42230' Kp" Dgm'Ncdu. "j g" j cf "y qtngf "qp" xctkqwu" cur gew" qh" ceqwvku" cpf "ur ggej " r tqeguipi . " kpenw' kpi <ur ggej " eqf kpi . " ur ggej " cpf " ur cngt" tgeqi pklqp. " uvcej cvke" o qf grlpi " qh" ur ggej " uki pcnu. " ghlekpv' ugctej " cni qtkj o u. " f kuetko kpcvkg" vclpki . " f gtgxtgdgtcvkqp" qh" cwf kq" cpf "ur ggej " uki pcnu. " o letqr j qpg" ctct{ " r tqeguipi . " ceqwvke" gej q" epegnrcvqp. " j cpf u'htgg" pqku{ " ur ggej " tgeqi pklqp 0' J g" y cu" cnuq" tgur qpukdng" hqt" vcpuhgttkpi " tgeqi pklqp" vgej pqmi { "Htqo "tuguctej "vq" CV(V"xqlkg/cevkcvgf "egni' r j qpgu' y j lej "y gtg"tcvfg" d{ "y j g" O qdkg" Qhleg" O ci c' lpg" cu' y j g' dguv' co qpi "eqo r gvkpi " r tqf wev' gxcnxcvgf 0J g' y cu' y j g' eq/ tgekr lgpv' qh' y j g" Dgm'Ncdu" Rtgu' kfpv' I qnf "Cy ctf "hqt" f gxgnr kpi "y j g" Dgm'Ncdu" Cwqo cvke" Ur ggej "Tgeqi pklqp" *DNCUT+ "uqhwy ctg" r cenci g 0J g" xkukgf "Icr cp" y leg" cu" c" xkuklpi "tuguctej gt-<htuv' Htqo "3; : 9" vq"3; : : . " vq" y j g" P VV" Grgewt/ Eqo o wplecvkqp" Ncdu. "O wcu' j kpg. " Vqnf{ q= y j g" Htqo "4224/4226. "vq" y j g" Ur qnpg" Ncpi wci g" Vtcpu' cvkqp" Ncdu. " CVT. " M{ qvq 0' Kp" 4226. " j g" lqkpgf " O letquqhw" Tgugctej " Cuk" *O UTC+ " Dgk' kpi . " Ej kpc" vq" rcf " y j g" Ur ggej " Tgugctej " I tqwr 0' J g" ku" c" xkuklpi " r tqhguuq" qh' y j g" Ej kpgug" Wplxgtukx{ "qh" J qpi "Mqpi " *EwJ M+ cpf "y j g' eq/ f tgevt' qh' EwJ M/ O UTC' lqkpv' Tgugctej " Ncd. " tgegpv{ " r tqo qvfg "vq" c" P cvkqpcni' Mg{ " Ncd' qh' O kpknt { "qh' Gf wecvkqp. "Ej kpc 0J g' y cu' y j g' eq/ ej ck" qh' y j g"3; ; 3" KGGG" Kp>gtpcvqpcni' Ctf gp" J qwug" Ur ggej "Tgeqi pklqp" Y qtnij qr 0J g" ku" c" eqo o kvgg" o go dgt'qh' y j g"

KGGG"Ur ggej "cpf "Ncpi wci g"Rtqeguulpi "Vgej plecni"Eqo o kwgg"qh"vj g"Uki pcni"Rtqeguulpi "Uqelgv"cpf " j cu"ugt xgf "cu"cp"cuuqelcvg"gf kqt"qh"vj g"Vtcpuvelqpu"qh"Ur ggej "cpf "Cwf kq"Rtqeguulpi 0J g"r wdrukj gf " gzvpuksgrf "cpf "eqcwj qtgf"o qtg"vj cp"422"vej plecni"r cr gtu"lp"vj g"ur ggej "cpf "uki pcni"rtqeguulpi "hgrf u0 J g"ku"cp"KGGG"Hgny 0

Abstract

Kp"vj ku"vcm"y g"y kni"tckug"cp"lpvgtgukpi "qt"gxgp"lpvki wpi "s wguvqp<"Ecp"y g"hkp"vj g"öngro gpvt {" r ctvenguo"qh" c"r gtuqpau"ur ggej "lp"qpg"ncpi wci g"cpf "wug"vj go "hqt"ur ggej lur gcngt"tgeqi plkqp"cpf " tgpfgt kpi "j kulj gt "xqlcg"lp" c" f khtgpv"ncpi wci g" A" C"r qukskxg" { guo" ecp"o cng"öngro gpvt {" r ctvenguo" wughwi" hqt" o cp {" cr r nlevlqpu." g0 0" o kzgf " eqf g" VVU." ugeqpf " ncpi wci g" rgtplkpi ." ur ggej /vq/ur ggej " vcpurvqp." gve0Y g"vt {" vq"cpuy gt "vj g"s wguvqp"d {" rko kkp " qwtugrk gu"htuv"vq"j qy "vq"vclp" c" VVU"lp" c" f khtgpv"ncpi wci g"y kj "ur ggej "eqmgev" f"tqo "c"o qpqtkpi wci"ur gcngt0Cf f kkpncm." c"ur ggej "eqtr wu" lp"vj g"vcti gvgf "pgy "ncpi wci g"ku"tgeqtf gf "d {" c"tghgtgpeg"ur gcngt0Y g"vj gp" wug"qwt"övtclgevqt {" vclpi " cni qt kjo .ö" kpxgpvgf "hqt"u{pvj guk kpi "j ki j "s wcrkx." wps"ugrgevqp"VVU."vq"övkrgö"vj g"vclgevqt kgu"qh"vj g" ugvpegu" lp" vj g" tghgtgpeg" ur gcngtau" eqtr wu" y kj " vj g" o quv" cr r tqr tkcv" ur ggej " ugi o gpwu" lp" vj g" o qpqtkpi wci" ur gcngtau" f cvc0 Vq" o cng" vj g" vclpi " r tqr gt" cetquu" vy q" f khtgpv" tghgtgpeg" cpf " o qpqtkpi wci"ur gcngtu." vj g" f khtgpeg" dgw ggp" vj go "pggf u"vq"dg"gs wcrk gf "y kj " cr r tqr tkcv" xqecni vcev"ncpi vj "pqto crk cvkqp." g0 0" c"dkp gct"y ctr kpi "hpevqp"qt" hqto cpv"o cr r kpi 0Cni"vkrf ." ugvpegu" ctg" vj gp" wugf "vq" vclp" c" pgy "J O O/dcugf " VVU" qh" vj g" o qpqtkpi wci"ur gcngt" dw" lp" vj g" tghgtgpeg" ur gcngtau"ncpi wci g0F khtgpv"ncpi vj "wpku"qh"vj g" -ngro gpvt {" r ctvenguo"j cxg" dggp" vclgf "cpf "c" rcdgn rguu" hco g" ncpi vj " *32" o u" ugi o gpwu" j cxg" dggp" hqwpf " vq" { kgrf " vj g" dguv" VVU" s wcrkx0 Uqo g" r tgrko kpct {" tguwu" cnuq" uj qy " vj cv" vclpki " c" ur ggej " tgeqi plk gt" y kj " ur ggej " f cvc" qh" f khtgpv" ncpi wci gu"vvpf u"vq"lo r tqxg"vj g" CUI"r gthqto cpeg"lp" gcej "lpf kxk wci"ncpi wci g0Cnuq."lp"cf f kkp"vq" vj g" hcev"vj cv"vclp"öngro gpvt {" r ctvenguo"qh"j wo cp"ur ggej "lp" f khtgpv"ncpi wci gu"ecp"dg" f lueqxtgf " cu" hco g/rgxgn"ur ggej " ugi o gpwu." vj g" o qwj "uj cr gu"qh" c"o qpqtkpi wci"ur gcngt" j cxg" cnuq" dggp" hqwpf " cf gs wcvg" hqt" tgpfgt kpi "vj g" rkr u"o qxgo gpv"qh"vcmkpi "j gcf u"lp" f khtgpv"ncpi wci gu0Xctkqu" f go qu"y kni dg"uj qy p"vq" kmuutcvg" qwt" hkp kpi u"lp" vcmkpi "j gcf " rkr u" tgpfgt kpi ." ur gcngt" tgeqi plkqp" lp" f khtgpv" ncpi wci gu." ur ggej " tgeqi plkqp" o qf gni"vclp gf "y kj " vj g"j gr "qh" f cvc" eqmgev" lp" f khtgpv"ncpi wci gu0

Keynote 2 - "

Semantic Matching: The Next Big Thing for Natural Language Processing?"



Hang Li

Ej lgh'Uelgpvku'qh'J wcy gk'Vgej pqm' lgu"

Hlf c{. 'Ugr vgo dgt"48"

2; <22'6'32-22"

Nqecvqp-<KpvtpcvqpcnEqphgtgpeg'Egpygt"

Biography

J cpi "Nk'ku"ej lgh'uekpvku'qh'v'j g"P qcj)r"Ctni'Ncd'cv'J wcy gk'J g"ku'cnuq"cf lwpv'r tqhguqt "qh'Rgnkpi " Wpkxgtuks{ "cpf "P cplkpi "Wpkxgtuks{ 0J ku'tgugctej "ctgcu"lpenmf g"lphqto cvkqp'tgtlgxcn'pcwtcn'rcpi wci g" r tqeguukpi . "ucvknlecl'no cej lpg"ngctplkpi . "cpf "f cv" o lpkpi 0J g"i tcf wcvf "Itqo "M{qvq"Wpkxgtuks{ "lp" 3; ; : "cpf "gctpgf "j ku'Rj F "Itqo "v'j g"Wpkxgtuks{ "qh'Vqn{ q"lp"3; ; : 0J g"y qtngf "cv'v'j g"P GE"ned"lp"Lcr cp" f wtkpi "3; ; 3"cpf "4223."cpf "O letquqhv" Tgugctej "Cuk" f wtkpi "4223" cpf "42340'J g"lqkpgf "J wcy gk' Vgej pqm' lgu" lp" 42340'J cpi "j cu" o qtg" v'j cp" 322" r wdrlecvkpu" cv' vqr "kpvgtpcvqpcn' lqwtpcnu" cpf " eqphgtgpegu. lpenmf lpi "UK KT."Y Y Y ."Y UFO ."CEN."GOP NR."EON."P RU."cpf "UK MFF 0J g"cpf " j ku' eqngci wgu" r cr gtu" tgegkxgf " v'j g" UK MFF -2: "dguv' cr r rlecvkqp" r cr gt" cy ctf . "v'j g" UK KT)2: "dguv' uwf gpv' r cr gt" cy ctf . "cpf "v'j g" CEN)34" dguv' uwf gpv' r cr gt" cy ctf 0J cpi "j cu'cnuq" dgpp" y qtnkpi "qp" v'j g" f gxgrq o gpv' qh' ugxgtcn' r tqf wevu' Vj gug" lpenmf g" O letquqhv" US N' Ugtxgt "4227." O letquqhv" QHleg"4229" cpf "QHleg"4232." O letquqhv" Nkxg" Ugtcej "422: ." O letquqhv" Dkpi "422; "cpf "Dkpi "42320J g" j cu'cnuq" dgpp" xgt{ "cevkg"lp" v'j g" t'gugctej "eqo o wpkkgu"cpf "ugt'xgf"qt "ku'ugt'xkpi "v'j g"vqr "eqphgtgpegu"cpf "lqwtpcnu"0Hqt" gzco r ng. "lp" 4234. "j g" ku' v'ceni' eq/ej ck" qh' v'j g" y gd" ugctej " v'ceni' qh' Y Y Y)34= ugplqt" r tqi tco " eqo o kvgg" o go dgtu"qt" ctgc"ej cku"qh' Y UF O)34." MF F)34." EMO)34." CEO N)34." CKTU)34=eq/ej ck" qh' MF F)34" uwo o gt" uej qqn" gve0" cpf " cp" gf kqtkcn' dqctf " o go dgt" qp" v'j g" lqwtpcn' qh' v'j g" Co gtlecp" Uqelgv{ "hqt "Kphqto cvkqp" Uelkpeg. "CEO "Vtcpucevqp" qp" Kpvnki gpv' U{ ugo u' cpf "Vgej pqm' { . "cpf "v'j g" lqwtpcn'qh'Ego r wgt' Uelkpeg{ "Vgej pqm' { 0'

Abstract

O quv'qh'pcwtcn'ncpi wci g'r tqeguakpi "P NR+"cumu."uwej "cu"lphqto cwkqp"tgtlqxcn"s wguakqp"cpuy gt kpi ." cpf "o cej kpg"tcurv'kqp."ctg"dcugf"qp"o cej kpi "dgy ggp"ncpi wci g"gzr tguakpu0'Vj ku"cr r tqcej "y qtmu" s wkg'y gni'kp"r tcevek="ku"rko kcvkqp"ku"cmq"qdxkqwu."j qy gxt0'Uqo gwko gu'o kuo cej "dgy ggp"ncpi wci g" gzr tguakpu"ecp"qeevt0'Y g"cti wg"vj cv"bugo cpvk"o cej kpi 0'ku"cp"ghgevkxg"cr r tqcej "vq"qxgteqo g"vj g" ej cmgpi g." vj cv' ku" vq" eqpf wev' o qtg" ugo cpvk" cpcn(uku) cpf " r gthqto " o cej kpi " dgy ggp" ncpi wci g" gzr tguakpu" cv' ugo cpvk" rxxgn0' kp" vj ku" vcm" Ky kni'htuv"r qlpv"qww"y j { "ugo cpvk" o cej kpi "ecp"j gr " uki pkk'ecpv("gpj cpeg"vj g"r gthqto cpeg"qh"PNR0'Ky kni'v gp"lwukh{ "o { "cti wo gpv'y kj "uqo g"gzco r ngu0' O qtg"ur gek'k'ecm{ ."Ky kni'kpv'qf weg"qwt"tgegpv'y qtnl'qp"vukpi "o cej kpg"ngctpkpi "vgej pls wgu"vq"eqputwev' o qf gni"ht"ugo cpvk"o cej kpi 0'Vj gug"lpenw'g"rcv'p'ur ceg"o qf gni'ht"s wgt { "f qewo gpv'o cej kpi "kp" ugctej ." utkpi "tg/y tkkpi "ngtpgn'ht"s wguakqp" cpuy gt kpi ." cpf "f ggr " o cej kpi " o qf gni'ht" uj qtv' vgzv' eqpxgtucv'kqp0'

Proceedings of the Twenty-Sixth Conference on Computational Linguistics and Speech Processing" ROCLING XXVI (2014)"

TABLE OF CONTENTS

Preface i

Oral Session 1: Speech and Speaker Recognition

運用概念模型化技術於中文大詞彙連續語音辨識之語言模型調適

Po-Han Hao, Su-Cheng Chen and Berlin Chen1

探究新穎語句模型化技術於節錄式語音摘要

Shih-Hung Liu, Kuan-Yu Chen, Yu-Lun Hsieh, Berlin Chen, Hsin-Min Wang, and Wen-Lian
Hsu3

台灣情緒語料庫建置與辨識

Bo Chang Chiou and Chia-Ping Chen21

Sparse Representation Based Speaker Identification

Kuang-Yao Wang, and Jia-Ching Wang31

Oral Session 2: Speech Processing and Synthesis

利用核依賴估計來進行多軌自動混音之研究

Tsung-Ting Wu and Chia-Hui Chang42

Research of Hakka Word Segmentation Processes on Chinese-to-Hakka's Text-to-Speech
System

Hsin-Wei Lin, Feng-Long Huang, Ming-Shing Yu, and Yih-Jeng Lin58

基於發音知識以建構頻譜 HMM 之國語語音合成方法

Hung-Yan Gu, Ming-Yen Lai, Wei-Siang Hong, and Yan-Hua Chen78

| | |
|---|----|
| Some Prosodic Characteristics of Taiwan English Accent | |
| Chao-yu Su, Chiu-Yu Tseng, and Jyh-Shing Roger Jang | 89 |

Oral Session 3:Text mining

| | |
|--|-----|
| 網頁商家名稱擷取與地址配對之研究 | |
| Yu-Yang Lin and Chia-Hui Chang | 91 |
| Public Opinion Toward CSSTA: A Text Mining Approach | |
| Yi-An Wu and Shu-kai Hsieh..... | 94 |
| Towards Automatic Enrichment of Standardized Electronic Dictionaries by Semantic Classes | |
| Bilel GARGOURI, Imen ELLEUCH, and Abdelmajid Ben Hamadou | 96 |
| Collaborative Ranking between Supervised and Unsupervised Approaches for Keyphrase Extractionen | |
| Gerardo Figueroa and Yi-Shin Chen..... | 110 |

Oral Session 4:Natural Language Processing

| | |
|--|-----|
| Semantic Representation of Ellipsis in the Prague Dependency Treebanks | |
| Marie Mikulová | 125 |
| Sketching the Dependency Relations of Words in Chinese | |
| Shih Meng-Hsien and Shu-kai Hsieh..... | 139 |
| 使用中文字筆畫構形資料庫校正字形相似之別字 | |
| Tao-Hsing Chang, Hsueh-Chih Chen, and Jian-Liang Zheng | 153 |
| 學術論文簡介的自動文步分析與寫作提示 | |
| 黃冠誠 吳鑑城 許湘翎 顏孜曦 張俊盛..... | 163 |

Poster Session:

| | |
|---|-----|
| 以二維共振峰分布建立語者音色模型及其在語者驗證上之應用 | |
| Jia-Guu Leu, Jyh-Bin Shiau, Chang En Pu, Ming-Ching Lee, and Chia-Long Wu | 165 |
| Unsupervised Approach for Automatic Keyword Extraction from Arabic Documents | |
| Arafat Atwi Awajan..... | 175 |
| Testing Distributional Hypothesis in Patent Translation | |

| | |
|--|-----|
| Hsin-Hung LIN and Yves Lepage | 185 |
| Spectrum Analysis of Cry Sounds in Preterm and Full-Term Infants | |
| Li-mei Chen, Yu-Hsuan Yang, Chyi-Her Lin, Yuh-Jyh Lin and Yung-Chieh Lin | 193 |
| Web-Based Recording and Visualization Framework for Moving Trajectories | |
| Po-An Yang, LiJung Chi, Kun-Ta Chuang, Seth Chen, Jonathan Tsai and Yung-Chung Ku | 204 |

運用概念模型化技術於中文大詞彙連續語音辨識之語言模型調適

Ngxgtci kpi 'Eqpegr v'O qf grkpi 'Vgej pks wgu'hqt 'Ncpi wci g'O qf gn'
Cf cr wvkqp'kp'O cpf ctkp'Ncti g'Xqecdwrct { 'Eqpvkpwqu'Ur ggej "
Tgeqi pkkqp"

郝柏翰 Rq/J cp'J cq"國立臺灣
師範大學資訊工程學系
822692: 4uB pvpwQf w0y "

陳思澄 Uuw/Ej gpi 'Ej gp"國立
臺灣師範大學資訊工程學系
82469293uB pvpwQf w0y "

陳柏琳 Dgtrkp'Ej gp"國立臺灣
師範大學資訊工程學系
dgtrkpB pvpwQf w0y "

摘要

語言模型已被廣泛地使用在各種自然語言相關的研究議題上；譬如，在語音辨識上，語言模型是關鍵的組成，其主要的功能通常是藉由已解碼的歷史詞序列資訊來預測下一個詞彙為何的可能性最大，以協助語音辨識系統從眾多混淆的候選詞序列假設中找出最有可能的結果[3, 4]。本論文旨在於發展新穎動態語言模型調適技術，用以輔助並彌補傳統 N 連 N -gram 語言模型不足之處。為此，我們本提出所謂的概念語言模型 *Eqpegr v'Ncpi wci g'O qf gn'ENO+，其主要假設是認為每一句的語句都是用來代表語者內心隱含而欲傳達的概念，並藉由語言*及語音+來具體表達相對應的語意內容。而概念模型最主要的目的則是希望能夠獲取使用者欲表達的概念，並假設在語音辨識過程中，同一概念之中歷史詞序列中所有詞彙以及待預測詞彙具有共同出現的關係，進而藉此關係達到預測待預測詞彙出現機率的目的。

在實作上，概念模型會使用*搜尋+與初步語音辨識結果相關的同領域文件*或調適語料+內表述的若干概念，用以近似語者內心欲傳達的真正含意，並基於此來建立概念語言模型。而概念語言模型的建立是分兩個面向來探討，它們分別是「詞彙」面向與「文件群聚」面向。首先，在實作上，概念模型會使用*搜尋+與初步語音辨識結果近似同領域文件*或調適語料+內表述的若干概念，用以近似語者內心欲傳達的真正含意，並基於此來建立概念語言模型。而概念語言模型的建立是分兩個面向來探討，它們分別是「詞彙」面向與「文件群聚」面向。首先，我們發展所謂的概念語言模型 *Y qtf/dcugf "Eqpegr v'Ncpi wci g'O qf gn+，並應用於語言模型調適。在建構詞概念語言模型時，我們期望能夠針對每一語句不同的語意內容*第一階段語音辨識結果，以詞圖[5]表示+，在調適語料的若干相關的文件中挑選一組具有代表性的概念關鍵詞組，藉以描述任一對歷史詞序列中所有詞

彙與待預測詞彙之間的相依關係。其次，我們亦發展所謂的群聚概念語言模型 *Enwugt/dcuqf 'Eqpegr v'Ncpi wci g'O qf gn，假設在調適語料的文件集內之文件可以由一組概念類別來獲得語句可能的概念分布，並做為語言模型預測的根據。實際上，概念類別的求取可透過一般分群演算法諸如 K/O gcpu"演算法[6]而求得；而根據每一語句的語意內容與各個概念類別的相關性，挑選一組具有代表性的概念類別組，藉以描述任一對歷史詞序列中所有詞彙與待預測詞彙之間之共同出現關係，並獲得對於待測詞彙之預測機率。再者，我們嘗試以不同方式來估測此種概念語言模型，並將不同程度的鄰近資訊 *Rtqzko k{ "Kqhto cvkq"融入概念語言模型以放寬其既有詞袋 *Dci /qh/Y qtf u"假設[5]的限制。

本論文是基於公視電視新聞語料庫來進行大詞彙連續語音辨識 *Ncti g" Xqecdwt { 'Eqpvkpwqu'Ur ggej "Tgeqi pklq. "NXEUT+實驗，以比較本論文所提出語言模型調適技術與其它當今常用技術之效能，包括了觸發對模型 *Vtki i gt/Rck" O qf gn、機率式潛藏語意分析 *Rtqcdklwle" Ncvpv' Ugo cpke" Cpcn{uku+以及狄利克里分配 *Ncvpv'F kkej rpv' Cmjecvqpp+ [7]等。實驗結果顯示我們的語言模型調適技在以字錯誤率 *Ej ctcevt "Gttqt "Tcvg. "EGT 評估標準之下，其它當今常用的語言模型調適技術相較，有相當具競爭性的表現，並且對於僅使用 N 連語言模型的基礎語音辨識系統能有明顯的效能提升。

關鍵詞：語音辨識、語言模型、概念模型、鄰近資訊

致謝：本論文之研究承蒙教育部 6'國立臺灣師範大學邁向頂尖大學計畫 *324L3C2: 22+與行政院科技部研究計畫 *O QUV" 325/4443/G/225/238/O [4." PUE" 325/4; 33/K225/523." PUE" 323/4443/G/225/246/O [5" "、 PUE" 323/4733/U/225/279/O [5" "、 PUE" "323/4733/U/225/269/O [5" "和 PUE" 324/4443/G/225/236/"O [5+之經費支持，謹此致謝。

參考文獻：

- [3]_"T0Tqughrf ."ōVy q" f gecf gu"qh"ucvkwkccn'ncpi wci g"o qf gkpi <"Y j gtg" f q" y g" i q" htqo "j gtg.ō"Proceedings of IEEE."xqr0': . 'pq0': .4222.'r r 034926349: .42220'
- [4]_"'L0' T0' Dgngi ctf c." ōUcvkwkccn' ncpi wci g" o qf gn' "cf cr vkwkccn' t gxlgy cpf 'r gtr gevkwgu.ō"Speech Communication." xqr0'64.'pq0'33.'r r 0'; 5632: .42260'
- [5]_" U0' Qtvo cppy. " J 0' P g{ " cpf " Z0' Cwdgtv." ōC" y qtf " i ter j " cri qtkj o " hqt" rcti g" xqecdwt { "eqpvkpwqu"ur ggej "tgeqi pklq.ō"Computer Speech and Language."xqr0' 33.'65/94.'3; ; 90'
- [6]_"T0Dcg| c/[cvgu'cpf "D0Tkdgtq/P gvq."Modern Information Retrieval: the Concepts and Technology behind Search."Cf f kuqpp/Y gurg { "Rtqhguwkqpcn"42330'
- [7]_"F 0'Drgk'cpf "L0'Nchgtv{ ."ōVqr ke"o qf gn.ō"lp" C0'Utkcucxc"cpf "O 0'Ucj co k" *gf u0:." Text Mining: Theory and Applications."Vc { rqt"cpf "Hcpeku."422; 0'

Vj g'4236'Eqphgtgpeg'qp'Ego r wcvkqpcrNlpi wku'cpf 'Ur ggej "
Rtqeguiki "TQENR I '4236.'r r 05/42"
Í "Vj g'Cuuqekvqp'hqt'Ego r wcvkqpcrNlpi wku'cpf 'Ej kpgug"
Npci wci g'Rtqeguiki "

探究新穎語句模型化技術於節錄式語音摘要

Kpxguiki cvkpi 'P qxgn'Ugpgpeg'O qf gkpi 'Vgej pks wgu'hqt "
Gzvtcevxg'Ur ggej "Uwo o ctł cvkqp"

劉士弘 Uj kJ /J wpi 'Nkw"陳冠宇 Mxcp/[w'Ej gp."

謝育倫 [wNwp'J ulgj ."王新民 J ulp/O kp'Y cpi ."許聞廉 Y gp/Nkcp'J uw"
中央研究院資訊科學研究所

}lqwtpg{.'n{ej gp."o qtr j g.'y j o ."j uw B kuOkpkecQf w0y "

陳柏琳 Dgtrkp'Ej gp"國立臺灣師
範大學資訊工程學系
dgtkpb pwpQf w0y "

摘要

近幾年來，基於語言模型化*Npci wci g"O qf gkpi ."NO +架構之摘要方法已初步在節錄式語音摘要任務上展現具競爭性的效能。在此架構下，對於被摘要文件每一句候選語句之語句模型的建立，可透過虛擬相關回饋*Rugwf q"Trgxcppeg"Hggf dcem"RTH策略來獲得較可靠的參數估測。一般來說，虛擬相關回饋在執行上可分為兩個階段：3+相關資訊*或者說虛擬相關文件+的選取；4+語句模型化與參數重新估測。首先，有別於現有基於語言模型化架構之摘要方法都聚焦在語句模型參數的重新估測，本論文深入探討與應用各種適合於節錄式語音文件摘要的虛擬相關文件選取技術，用以強化語句模型的參數估測。再者，本論文更進一步地考量使用每一語句的非相關性*P qp/trgxcppeg+資訊對於虛擬相關文件選取的影響。同時，我們亦額外嘗試基於重疊分群*Qxgtncr r gf "Enuvgtkpi +概念來有效地選取重要的虛擬相關文件。最後，本論文探索使用三混合模型*VtkO kzwtg"O qf gm"來表示每一語句，期盼其能更精確地表示一句語句之獨特詞彙使用和語意相關資訊。本論文的語音文件摘要實驗語料是採用公開的公視廣播新聞*O CVDP+；實驗結果顯示，相較於其它現有虛擬相關文件選取技術，我們所發展的虛擬相關文件選取技術能提供相當不錯的摘要效能改善。

關鍵詞：節錄式自動摘要、虛擬相關回饋、語句模型、非相關資訊、重疊分群

一、緒論

隨著網際網路的普及與蓬勃發展，大量含有語音資訊的多媒體內涵快速地傳遞並分享於全球各地，像是電視新聞、課程演講、會議錄音和語音郵件等。由於這些多媒體內涵透過自動語音辨識*Cwqo cvk"Ur ggej "Tgeqi pklqp."CUT+處理後都可以表示成語音文件及其對應的自動轉寫*Cwqo cvk" Vtcpuetr w+供瀏覽使用，節錄式語音摘要*Gzvtcevxg"Ur ggej "Uwo o ctł cvkqp+在過去十年逐漸成為熱門研究議題[39]47_。節錄式語音摘要目標在於根據一定的摘要比例，從語音文件中選取重要語句並組合成摘要，以期能夠扼要的表示語音文件主要的主题或語意資訊；藉此，使用者能迅速地瀏覽大量多媒體內涵並能充分理解其中主题或語意資訊。

當前主流的節錄式語音摘要方法大致可分為三類：*3+基於文件結構或語句位置資訊來選取重要語句；*4+基於特定詞彙或語意資訊之非監督式*Wpuwr gt xkugf +摘要方法；*5+需要使用人工摘要標註來訓練模型之監督式*Uwr gt xkugf +摘要方法。第一類摘要方法大都是簡單地從文件的緒論*Hvtqf wevklq+或結論*Eqpenwukq+所在段落擷取出若干語句來組成摘要[3]；此類方法僅適用在特定具有一致結構的文字或語音文件上，因此在實際應用上有其侷限性。另一方面，第二類摘要方法通常將自動摘要任務視為如何排序並挑選具代表性*或重要性+語句之問題，其方法通常是以非監督式方式產生出一種語句層次的摘要特徵或重要分數以供語句排序使用。而第二類摘要方法又可進一步地分成三小類：*K+以向量*Xgevt+為基礎的方法，包含有向量空間模型*Xgevt"Ur ceg"O qf gn"XUO +: _、潛藏語意分析*Ncvpv"Ugo cpve"Cpcn"ulu."NUC +: _[35]及最大邊際關聯*O czlo cn"O cti kpcnT grgxcpeg."O O T +[4]等；*KK+以圖*I tcr j +為基礎的方法，具代表性的有馬可夫隨機漫走*O ctmqx"Tcpf qo "Y cm"O TY +[52]、詞彙排序*NgzTcpm][9] [45]及最小支配集*O lpko cnF qo kpcvpi "Ugv+[4]：_等；*KKK+以組合最佳*Eqo dlpcvqtken"Qr vko k cvkq+為基礎的方法，包括有次模*Uwdo qf wrctk/[+[36]以及線性整數規劃*Nlpgct"Kvgi gt" Rtqi tco o lpi +[44]_等。最後，第三類監督式摘要方法通常將自動摘要之任務視為二元分類問題*Dkpc{ "Erucukhecvkq+，亦即將語句區分為摘要語句或非摘要語句。在訓練階段，必須事先準備好一些訓練文件以及其對應的人工標註過摘要資訊，然後結合各種詞彙、語意或音韻等特徵來表示語句，並且透過各種分類器的學習機制進行摘要*分類+模型的訓練；在測試階段，對於將被摘要之文件，此類方法將文件裡的每一句語句進行二元分類，即可依所設定摘要比例來取摘要語句以產生出摘要。在此類方法中，較著名的包括簡單貝氏分類器*P c'xg/Dc{ gu"Erucukhgt +[33]、高斯混合模型*I cwukcp"O kzwtg"O qf gn"l O O +[46]、隱藏式馬可夫模型*J kf gp"O ctmqx"O qf gn"J O O +[8]、支援向量機*Uwr r qtV"Xgevt"O cej kpgu."UXO +[32]與條件隨機場域*Eqpf kkpccn"Tcpf qo "Hgrf u."ETH+[4]：_等。由於監督式摘要方法所使用的模型在訓練時必須使用一定數量文件及其對應經人工標註過摘要資訊，所以當它們被應用到新的摘要任務或應用領域時，相較上述兩類摘要方法而言，是會耗費許多人力與時間的。值得一提的是，自動摘要研究也可從其它不同面相來進行探討，包括了來源、需求、方式、用途等，有興趣的讀者可參考相關文獻[39] [42] [43] [47] 進行更深入的瞭解。

有別於上述的摘要方法，近期有一些基於語言模型化*Ncpi wci g"O qf gkpi ."NO +架構之摘要方法被提出，並且初步在節錄式文字或語音摘要任務上展現不錯的效能。在此架構下，對於被摘要文件每一句候選語句之語句模型的建立，可透過虛擬相關回饋*Rugwf q" Tgrgxcpeg"Hggf dcem"RTH策略來獲得更加可靠的參數估測。一般來說，虛擬相關回饋在執行上可分為兩個階段：3+相關資訊*或者明確地說，虛擬相關文件+的選取；4+語句模型化與參數重新估測。本論文同樣也基於語言模型化架構來發展語音摘要方法，其貢獻主要有三方面。首先，有別於現有基於語言模型化架構之摘要方法都聚焦在語句模型參數的重新估測，本論文深入探討與應用各種適合於節錄式語音文件摘要的虛擬相關文件選取技術，用以強化語句模型的參數估測。其次，本論文更進一步地考量使用每一語句的非相關性*P qp/tgrgxcpeg+資訊對於虛擬相關文件選取的影響。同時，我們亦額外嘗試基於重疊分群*Qxgtcr r gf "Enwvgtkpi +概念來有效地選取重要的虛擬相關文件做為語句模型的參數估測之依據。最後，本論文探索使用三混合模型*Vtk/O kzwtg"O qf gn來表示每一語句，期盼其能更精確地表示一句語句之獨特詞彙使用和語意相關資訊。

本論文後續安排如下：第二節首先介紹使用語言模型於節錄式語音摘要任務之原理，然後闡述虛擬相關回饋的觀念及其現有虛擬相關文件選取技術；第三節將介紹本論文提出之新穎式虛擬相關文件選取技術；第四節則介紹現有各種關聯模型，並且說明如何結合語句關聯性資訊來改進語句模型之估測，使其得以更精準地代表語句的語意內容；第五節介紹實驗語料與設定以及摘要評估之方法；第六節說明實驗結果及其分析；

最後，第七節為結論與未來研究方向。

二、使用語言模型於語音文件摘要

在過去近二十年來，各種語言模型在資訊檢索任務中已被廣泛地應用，並且有不錯的實務成效[56]。近期在語音摘要領域，亦開始有一些基於語言模型化架構的非監督式摘要方法被提出。本節將先簡介兩種常見的、基於語言模型化架構的摘要方法：其一為使用語句語言模型生成文件的文件相似度量值 $P^*_{qewo\ gpv} Nkngrkj\ qqf\ "O\ gcuwg.\ "FNO\ +$ 的摘要方法[5]；另外為使用庫爾貝克/萊伯勒離散度量值 $Mwmdcem\ Ngkdngt\ "F\ kxgti\ gpeg\ "O\ gcuwg.\ "MN\ +$ [35][37] 來計算文件模型和語句模型間之距離的摘要方法。接著，我們將闡述如何利用虛擬相關回饋 $Rugwf\ q\ "T\ grgxcpeg\ "Hggf\ dcem\ +$ 概念來獲得更可靠的語句模型估測，並介紹數個在資訊檢索領域已被發展出的新穎虛擬相關文件選取技術。

4B、文件相似度量值

我們可以把語音文件摘要任務視為是資訊檢索的問題。一般來說，資訊檢索 $Kphqto\ cvkqp\ "Tgvtlgxcn\ "K\ +$ 旨在尋找相關文件 $Tgrgxcpv\ "F\ qewo\ gpw\ +$ 來回應使用者所送出的查詢 $S\ wgt\ \{ \ +$ 或資訊需求 $Kphqto\ cvkqp\ "P\ ggf\ +$ 。同樣地，在從事語音文件摘要時，我們可將每一篇被摘要文件視為是查詢，而文件中的每一句語句 $Ugpv\ gpeg\ +$ 視為候選資訊單元 $Ecpf\ kf\ cvg\ "Kphqto\ cvkqp\ "Wpk\ +$ ；據此，我們可以假設在被摘要文件中，與文件本身愈相關的語句就愈有可能是可用來代表文件主旨或主題之摘要語句。

當給定一篇被摘要文件 D 時，文件中每一句語句 S 的事後機率 $P^*_{S\ \sim\ D\ +}$ 可以用來表示語句 S 對於文件 D 的重要性。當使用語言模型來計算 $P^*_{S\ \sim\ D\ +}$ 時，我們透過貝氏定理 $Dc\ \{ \ gu\ vj\ gqtgo\ \ +$ 將 $P^*_{S\ \sim\ D\ +}$ 展開成[5]：

$$P^*_{S\ \sim\ D\ +} = \frac{P^*_{D\ \sim\ S\ +} P^*_{S\ +}}{P^*_{D\ +}} \quad *3+$$

其中 $P^*_{D\ +}$ 是文件 D 的事前機率，由於 $P^*_{D\ +}$ 不影響語句的排序結果，故可省略不討論；另一方面， $P^*_{S\ +}$ 是語句 S 的事前機率，可以使用各式非監督式方法或監督式方法來求得[5]。在此先假設語句的事前機率為一個均勻分布 $Wpkh\ qto\ "F\ kwt\ kd\ w\ kqp\ +$ ，所以 $P^*_{S\ +}$ 亦可省略。最後， $P^*_{D\ \sim\ S\ +}$ 是語句 S 所形成的語言模型生成文件 D 之機率*或稱作文件相似度*，可以用來表示文件 D 與語句 S 之間的相似關係，如果語句 S 生成文件 D 的機率值愈高，代表語句 S 與文件 D 愈為相似*語句愈能代表文件 D ，即愈有可能是摘要語句。我們可以更進一步地假設文件 D 中詞彙與詞彙之間是獨立的，並且不考慮每一個詞彙在文件 D 中發生的順序關係*即詞袋假設 $Dci\ /qh\ Y\ qtf\ "Cuuwo\ r\ v\ kqp\ +$ ，則語句 S 生成文件 D 的文件相似度量值 $F\ qewo\ gpv\ Nkngrkj\ qqf\ "O\ gcuwg.\ "FNO\ + P^*_{D\ \sim\ S\ +}$ 可拆解成文件 D 中每一個詞彙 w 個別發生的條件機率之連乘積：

$$P^*_{D\ \sim\ S\ +} = \prod_{w \in D} P^*_{w\ \sim\ S\ +}^{C^*_{w.\ D\ +}} \quad *4+$$

此種方法是為語句 S 建立一個語句模型 $Ugpv\ gpeg\ "O\ qf\ gnt\ P^*_{w\ \sim\ S\ +}$ ， w 是一個出現在文件 D 中的詞彙， $C^*_{w.\ D\ +}$ 是詞彙 w 出現在文件 D 中的次數。其中，我們可利用最大化相似度估測 $O\ czko\ wo\ "Nkngrkj\ qqf\ "Gurko\ cvkqp.\ "ONG\ +$ 的方式來建立每一個語句的語句模型：

$$P^*_{w\ \sim\ S\ +} = \frac{C^*_{w.\ S\ +}}{|S\ \sim|} \quad *5+$$

在*5+中， $C^*_{w.\ S\ +}$ 表示詞彙 w 在語句 S 中出現的次數， $|S\ \sim|$ 則表示語句 S 所含的總詞數。

值得注意的是，由於語句 S 通常僅由少數字詞所組成，因此容易遭遇資料稀疏的問題，這會使得語句模型使用最大化相似度估測時，不僅可能無法準確地估測每一個詞彙在語句中真正的機率分佈，也可能因為某些詞彙的條件機率值為零，導致語句 S 產生文件 D 的機率值為零。為了減輕上述的現象，本論文使用平滑化技術藉由使用以大量文字語料訓練而成的背景單連語言模型來調適語句模型^[55]，故 $P^*_{D|S}$ 可進一步地表示成：

$$P^*_{D|S} = \prod_{w \in D} [\lambda \cdot P^*_w|S + (1-\lambda) \cdot P^*_w|BG]^{C^*_{w,D}} \quad (6)$$

其中， $P^*_w|BG$ 是詞彙 w 由背景單連語言模型所產生的機率值。

4.4、庫爾貝克-萊伯勒離散度量值

語言模型使用於文件摘要的研究中，除了可被用於計算語句生成文件的可能性外，另一種方式為藉由庫爾貝克/萊伯勒離散度量值來評估文件中每一個語句的重要性。當使用庫爾貝克/萊伯勒離散度量值於摘要任務中，被摘要文件 D 和 D 中的每一個語句 S 都將分別被描述為一個單連語言模型；當相對於被摘要文件 D 的文件模型，語句模型的離散度量值愈小時，則代表語句與文件愈相關，亦即語句 S 愈重要。在此摘要架構下，排序語句重要性的公式如下^[37]：

$$KL^*_{D|S} = \sum_{w \in V} P^*_w|D \cdot \frac{P^*_w|D}{P^*_w|S} \quad (7)$$

其中， V 表示一個由語言裡所有可能的語彙所形成的集合。本論文的研究中，文件模型 $P^*_w|D$ 的建立方式與語句模型相同^{參照式(5)}。當我們更進一步地對⁽⁷⁾作分析時，可以發現當文件模型僅使用最大化相似度估測的前提下，採用庫爾貝克/萊伯勒離散度量值所得到的語句排序將與使用文件可能性⁽⁸⁾測量方式⁽⁹⁾所得到的結果是相同的^[6]。

由於使用庫爾貝克/萊伯勒離散度量值時，不僅每一句語句被表示成語句模型，每一篇被摘要文件 D 亦被視為一個文件⁽¹⁰⁾，而文件模型在經由各式語言模型調適與平滑化的技巧下，可以有系統地、適當地調適文件模型的機率分佈；因此相較於文件相似度值⁽¹¹⁾只能針對語句模型進行調適，庫爾貝克/萊伯勒離散度量值⁽¹²⁾能透過不同模型參數估測技術的使用來求取語句模型和文件模型，以獲得最佳的自動摘要效能。

4.5、虛擬相關回饋

通常，文件中的語句僅由少許的詞彙所組成，當語句模型使用最大化相似度估測時，容易遭遇資料稀疏的問題；再者，由這語句 S 中些許的表面詞彙是遠不夠正確估算語句 S 與被摘要文件 D 之間的相似度⁽¹³⁾，所以藉由背景語言模型進行語句模型之調適為最常見的方法之一^{參照式(6)}。

為了有效解決語句的資料稀疏及相似度被低估的問題，我們可利用在資訊檢索領域被廣泛應用的虛擬相關回饋技術來強化語句模型⁽¹⁴⁾。為此目的，當虛擬相關回饋運用於文件摘要領域中時，會將每一語句 S 當成是一個查詢⁽¹⁵⁾，然後輸入到一個資訊檢索系統中，找出一些與語句 S 最可能相關的文件，而這些文件就稱之為虛擬相關文件

*Rugwf q "Tgrgxcpv" F qewo gpu+；一個最簡單的方式即是選取排名最前面 *檢索分數最高+ 的幾篇文件 *Vqr /tcpngf "F qewo gpu+。有了這些虛擬相關文件後，就可以利用它們來增進語句模型以解決語句資料稀疏及其相似度低估之問題。

當使用虛擬相關回饋技術時，通常會遭遇到兩個挑戰。第一為如何選取、純化 *Rwtlh{ + 虛擬相關文件，意即如何去除虛擬相關文件的冗餘性 *Tgf wpf cpe{ + 和摒除非相關資訊 *P qp/tgrgxcpv"Kphto cvkqp+。第二則是如何有效地運用虛擬相關文件來重新估測語句模型或進一步調適。對於第二個挑戰，已經有許多學者提出各種不同的關聯模型方法，常見的有關聯模型 *Tgrgxcpeg" "O qf gn" "TO +j34_ 和簡單混合模型 *Uo r rg" "O kzwg" "O qf gn" UO O +j54_ 等 *我們將在第4節介紹上述各種關聯模型+。然而對於第一個挑戰的研究，雖亦有少數學者提出一些虛擬相關文件選取方法，但相較於第二個挑戰的研究，仍顯較匱乏。

關於虛擬相關文件的選取方式，過去在資訊檢索領域有學者陸續提出間隔式最高 K *I cr r gf "Vqr "K+ 選取法 [49]、群中心 *Enwugt "Egptqkf + 選取法 [49] 以及主動式/關聯多元密度 *Cevkxg/TFF + 選取法 [53]，其中 TFF " 三個英文字母分別代表關聯 *Tgrgxcpeg+、多元 *F kxgtuk{ + 以及密度 *F gpub{ +，下面將簡單介紹上述虛擬相關文件選取方法。

間隔式最高 K 選取法就是在最高排序文件中每間隔 J *例如兩個間隔， J 為 4+ 挑選出 K 個相關文件出來當作是最高排序文件，簡單的例子如下：假設最高排序文件有 32 篇，我們每間隔 4 篇要挑出最高 5 篇出來，則第一、第四及第七篇會被挑選出來當作是最高排序文件。間隔式最高 K 選取法的主要思想是要挑選出具有多元性的文件，但其實此選取方法也是相當不穩定的。群中心選取法則是先將最高排序文件作分群 *Enwugt kpi +，分群方法可以是任意的，常用的分群方法為 K 中心 *K/o gcpu+ 分群法，然後再從分出來的每一群中挑選出一篇最相關的文件，以此構成新的最高排序文件，由分群的觀念可知，群中心選取法旨在選取出具有多元性的文件，與間隔式最高 K 選取法相較，群中心選取法是一個比較穩定的選取方法。主動式/關聯多元密度選取法為同時考量最高排序文件中的關聯性、多元性以及密度性的一種貪婪 *I tggf { + 選取法 [53]。以上選取方法常見於資訊檢索領域中，有興趣的讀者可參考相關文獻 [7]，本論文是首次將上述方法運用在 *語音+ 文件摘要任務中並做深入探討。

三、新穎式虛擬相關文件選取

基於主動式/關聯多元密度選取法，除了考量到虛擬相關文件 *最高排序文件+ 中的關聯性、多元性以及密度性之外，我們認為非相關性 *P qp/tgrgxcpeg+ 資訊 *在這裡是指語句的非相關性資訊+ 也是相當重要的線索，可以用來幫助重新選取更好的虛擬相關文件，因此本論文提出額外考量非相關資訊以改進主動式/關聯多元密度選取法，稱之為主動式/關聯多元密度非相關 *Cevkxg/TFFP + 選取法。另一方面，本論文也提出使用重疊分群 *Qxgtm r r gf "Enwugt + 的概念來幫助重新選取更有效的虛擬相關文件，茲介紹如下：

5B、主動式-關聯多元密度非相關選取

假設語句 S 已經輸入到資訊檢索系統中並得到了最高排序文件 $D_{Top?} \{D_1, D_2, \dots, D_{M_i}\}$ ，那主動式/關聯多元密度非相關選取則是從最高排序文件中迭代地 *Kgtcvkxgn{ + 同時考慮四種重要因素 *關聯性、多元性、密度性以及非相關資訊+ 來重新選取更具代表性的文件集。具體地說，最高排序文件中的每個候選 *Ecpf kf cvg+ 文件 D_m 都會有著同時考量四種因素的一個線性結合的分數，其選取公式如下：

$$D_m = \text{cti o cz} \left[(3 - \alpha - \beta - \gamma) \cdot M_{Rel}(S, D_m) + \alpha \cdot M_{NR}(S, D_m) + \beta \cdot M_{Diversity}(D_m) + \gamma \cdot M_{Density}(D_m) \right] \quad (8)$$

其中 D_R 為已經選入的文件集， $M_{Rel}(S, D_m)$ 、 $M_{NR}(S, D_m)$ 、 $M_{Diversity}(D_m)$ 及 $M_{Density}(D_m)$ 分別代表候選文件 D_m 的關聯性量值、非相關資訊性量值、多元性量值、以及密度性量值，而 α 、 β 、 γ 為可調參數且其總和為 3 即 $\alpha + \beta + \gamma = 3$ 。值得一提的是，式 8 與早期用於資訊檢索及摘要領域的經典公式最大邊際關聯 $O\text{ czko crl} O\text{cti kpcnl} T\text{grgxcpeg}$ 相似 [4]。另外，關聯性量值 $M_{Rel}(S, D_m)$ 可定義為文件 D_m 語句 S 的負庫爾貝克/萊伯勒離散度量值 $1/KL(D_m \sim S)$ 。下面將分別介紹非相關資訊量值、多元性量值、以及密度性量值。

5.3.3、非相關性資訊量值

對於一個語句 S ，其非相關性資訊通常可以從第一次資訊檢索時排在最後面的一些文件最低排序文件 NR_S 來代表，那麼語句 S 的非相關模型 P_{NR_S} 便可由最低排序文件來估測，而非相關性資訊量值可由下面式子表示：

$$M_{NR}(S, D_m) = KL(NR_S \| D_m) \quad (9)$$

其中 NR_S 表示語句 S 的非相關性資訊，亦即最低排序文件。若非相關性資訊量值越小，表示候選文件 D_m 與對應於語句 S 的非相關性資訊很相像，則不應該被選取當成代表性文件；反之，則表示候選文件 D_m 與語句 S 的非相關資訊離較遠，應該有機會被選取起來當成代表性文件。實際上，虛擬相關文件最高排序文件相對於最低排序文件是相當少量的，所以實作上我們可利用全部的文件集來估測語句 S 的非相關模型，更明確的說，就是利用背景語言模型 P_{BG} 來當作是語句 S 的非相關模型 P_{NR_S} 。

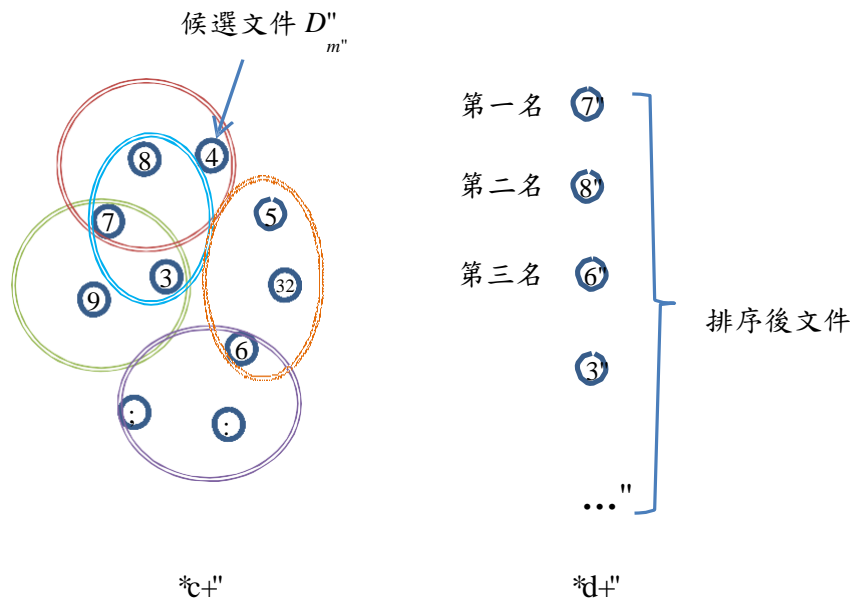
5.3.4、多元性量值

在資訊檢索領域中，多元性 $F\text{lxgtullec}v\text{qp}$ 的考量已備受矚目，相較於傳統只考量關聯性的查詢使得有太多重複及冗餘查詢結果，額外考慮多元性是希望能提供查詢者多元且多樣性的查詢結果，以滿足不同需求的使用者。而在虛擬相關回饋下，若最高排序文件的重複性文件太多，則會導致後面的關聯模型估測有太多的冗餘資訊，導致模型估測不良，所以多元性量值考量的點就是希望選取過的文件或與已選文件相近的文件不要再被選取進來，亦即去除最高排序文件的冗餘性，使得其後面的關聯模型估測更精準。多元性量值就是從已選的文件 D_P 中找出與候選文件 D_m 有最小對稱離散度量的值，可由下列表示：

$$M_{Diversity}(D_m) = \frac{3}{4} \cdot \left[KL(D_j \sim D_m) + KL(D_m \sim D_j) \right] \quad (10)$$

5.3.5、密度性量值

另一方面，最高排序文件中的結構 $U\text{twewt}c\text{ni}$ 資訊可以被當成是一個線索來幫助選取代表性文件，其主要目的是希望在考量多元性資訊的同時，也應避免選取到過度極端文件，因為過度極端文件很可能是錯誤的資訊。為了實現這個想法，我們可以利用計算最高排序文件中的候選文件 D_m 與其他候選文件 D_h 的負平均對稱 $P\text{gi}c\text{v}x\text{g} C\text{xgtci}g U\text{o}o\text{gtle}$ 離散度量值來達成，其公式如下所示：



圖一、使用重疊分群概念的候選文件選取示意圖，

*c+利用 k/PP 為每個候選文件 D_m 找出重疊分群，並計算其重疊分群個數。

*d+依據每一候選文件 D_m 的重疊分群個數做排序。

$$M_{Density}(D_m) = \frac{-3}{|\mathbf{D}_{Vqr}| - 3} \cdot \sum_{\substack{D_h \in \mathbf{D}_{Vqr} \\ D_h \neq D_m}} [KL(D_h \sim D_m) + KL(D_m \sim D_h)] \quad *+, "$$

其中 \mathbf{D}_{Top} 為最高排序文件之個數。若負平均對稱離散度量值越大，表示此候選文件 D_m 與最高排序文件中的其他候選文件 D_h 很接近，亦即可能是最高排序文件中心點*密度大的文件+，所以此候選文件 D_m 可能是個重要的文件，應該要被選入；反之，就會離中心點較遠*密度較低+，有可能會是不重要的文件，就不應該被選入。

504、重疊分群

本論文提出使用重疊分群* $Qxgtncr r gf "Enwugt$ +的概念來重新選取更好的虛擬相關文件*即 \mathbf{D}_p +以利接下來各種不同的關聯模型估測。其重疊分群選取方法為一個三步驟的演算法，示意圖如圖一所示，其演算法描述如下：

30' 第一步驟：計算最高排序文件中兩兩候選文件的相似度，在此是將候選文件表達成向量空間模型* $Xgevt "Ur ceg" O qf gm$ +並使用餘弦相似度* $Equlpg "Ulo krtkw$ +量值來做計算。

40' 第二步驟：利用 k -最近鄰居* k/PP +來為每個候選文件 D_m 找出 k 個最接近的相關文件，並形成一個群*每個候選文件都會形成一個分群，而此分群中會有 $k-3$ '個文件+。

50' 第三步驟：對於每個候選文件 D_m 都去計算重疊分群的個數*以圖一例子來說明，候選文件編號 7"被三個分群所包圍，所以其重疊分群個數為 5+，並且按照重疊分群個數來對每個候選文件作排序，排序後即可得到新的最高排序文件。

重疊分群的概念在資訊檢索領域中已有些許研究[3: _]，它是利用最高排序文件的結構資訊來幫助訓練語言模型，而本論文的思想是要利用重疊分群的概念來找出支配* $F qo lpcvg$ +文件，若一個候選文件的重疊分群個數很多的話，表示它是很重要的且能夠支配其他候選文件，則應該要被選為代表性文件。本論文是首次將重疊分群的概念用在*語音+文件摘要任務上。

四、各種關聯模型之簡介

當我們透過資訊檢索系統已取得虛擬相關文件*最高排序文件，即 D_{Top} ，或進一步地使用上一節提出之方法來改善虛擬相關文件後*即 D_P ，接下來就要做模型估測，底下介紹常見的模型包含有關聯模型* $Tgrxcpeg"Oqf gn"TO$ 、簡單混合模型* $Uko r ng"O k z w t g"O q f g n"UO O$ 以及三混合模型* $Vtk O k z w t g"O q f g n"Vtk O O$ 。

6B、關聯模型

關聯模型的基本假設是認為每一語句 S 皆是被用來描述一個概念、想法或主題，我們稱之為語句的關聯類別* $Tgrxcpeg"Ernu"R_S$ 。在本論文中，我們的目標是想進一步地模型化關聯類別所代表的資訊，藉此來豐富語句模型所能傳達的語意內容或主題特性。然而，實際上每一語句 S 的關聯類別 R_S 是非常難以求得的；為此，我們透過虛擬相關回饋* $Rugwf q" Tgrxcpv" Hggf dcem$ 來尋找與關聯類別可能相關的一些文件，並藉由這些文件來近似關聯類別 R_S 。更明確地，在實作上我們將虛擬相關文件*最高排序文件* D_{Top} 或透過上一節所介紹的選取方法來產生較佳的虛擬相關文件 D_P 用以代表關聯類別 R_S 。接著，透過檢視詞彙 w 與語句 S 在這些虛擬相關文件中同時出現之關係，可計算出詞彙與語句的聯合機率[34]：

$$P_{TO}(w, S) = \sum_{D_m \in D_R} P(w, S | D_m) \cdot P(D_m) \quad *32*$$

當我們進一步地假設在給定某一篇虛擬相關文件時，詞彙與語句是獨立的，並且語句內的詞彙也是獨立且不考慮其先後次序*即所謂的詞袋假設*，則透過虛擬相關回饋所估測的語句模型為：

$$P_{TO}(w, S) = \frac{\sum_{D_m \in D_R} \prod_{w \in S} P(w | D_m) \cdot P(D_m)}{\sum_{D_m \in D_R} \prod_{w \in S} P(w | D_m)} \quad *33*$$

我們稱之為關聯模型* $Tgrxcpeg"Oqf gn"TO$ 。關聯模型的優點在於藉由虛擬相關文件的資訊，可以更清楚地知道語句所蘊含的資訊、所欲表達的內涵，所以相較於傳統使用最大化相似度估測的語句模型，可更準確地表達語句的語意內容或主題特性，以提升摘要的成效。

6C、簡單混合模型

簡單混合模型的基本想法是假設由虛擬相關回饋技術所得到的虛擬相關文件是相關的且能從最高排序文件中估測比較好的簡單混合模型 $P_{UO O}(w, S)$ ，更明確地說，簡單混合模型是假設虛擬相關文件 D_P 裡的詞彙 w 是源自於二種成分混合模型* $Vy q/Eqo r qpgpv"O k z w t g"O q f g n$ ，其一為簡單混合模型 $P_{UO O}(w, S)$ ，另一為背景語言模型 $P(w, BG)$ 。簡單混合模型的估測是藉由期望值最大化* $Gzr gevckqp"O czko k cvkqp."GO$ 演算法來最大化虛擬相關文件的對數相似度* $Nqi / Nkgrkj qqf$ 以進行模型的估測，其虛擬相關文件的對數相似度的定義如下[54]：

$$LL_{D_R} = \sum_{D_m \in D_R} \sum_{w \in V} c(w, D_m) \cdot \ln [\frac{1}{\alpha} P_{UO O}(w, S) + \alpha \cdot P(w, BG)] \quad *34*$$

其中 α 為平衡參數，用來控制模型估測時是要比較偏好簡單混合模型或是背景語言模

型， $c_w D_m$ 為詞彙 w 在虛擬相關文件 D_m 的次數，式*34 的最大化可透過期望值最大化迭代更新式來達成：

期望值步驟：

$$\tau_w^{l+1} = \frac{\alpha \cdot P_{\text{UO}}^{l+1} \cdot c_w \cdot D_m}{\alpha \cdot P_{\text{UO}}^{l+1} \cdot c_w \cdot D_m + (3 - \alpha) \cdot P_w \cdot BG} \quad (35)$$

最大化步驟：

$$P_{\text{SMM}}^{l+3} \cdot c_w \cdot D_m = \frac{\sum_{D_m \in \mathbf{D}_P} c_w \cdot D_m \cdot \tau_w^{l+1}}{\sum_{w' \in V} \sum_{D'_m \in \mathbf{D}_P} c_{w'} \cdot D'_m \cdot \tau_{w'}^{l+1}} \quad (36)$$

其中 l 表示期望值最大化的第 l 次迭代。這個簡單混合模型的估測會加強具有獨特性 Ur geklek 的詞彙之機率，例如某詞彙沒有在背景語言模型中有好解釋 Y gm Gzr nlpf 則會被加強其機率，這樣使得此模型為更具有鑑別 F kuetko kpcpv 能力的語句模型；反之，若是沒有獨特性的詞彙，則其機率就會被背景語言模型所吸收。

6.5、三混合模型

另一方面，本論文嘗試將三混合模型 $\text{VtkO kwtg} \text{O qf gm}$ 用於語音摘要任務。三混合模型可視為是複雜化後的簡單混合模型；它更進一步的假設虛擬相關文件 \mathbf{D}_P 裡的詞彙 w 是源自於三種成分模型 $\text{Eqo r qpgpv} \text{O qf gm}$ ，其一為文件模型 P_w/D_m ，其二為三混合模型 $P_{\text{tkO}} \cdot w/S$ ，最後為背景語言模型 P_w/BG 。三混合模型的估測也是藉由期望值最大化演算法來最大化虛擬相關文件的對數相似度以進行模型的估測，其虛擬相關文件的對數相似度的定義如下：

$$LL_{\mathbf{D}_R} = \sum_{D_m \in \mathbf{D}_R} \sum_{w \in V} c_w \cdot D_m \cdot [\lambda - \mu \cdot P_{\text{tkO}} \cdot w/S + \lambda \cdot P_w \cdot D_m + \mu \cdot P_w \cdot BG] \quad (37)$$

其中 λ 和 μ 為平衡參數，用來控制模型估測時是要比較偏好三混合模型或文件模型亦或是背景語言模型， $c_w D_m$ 為詞彙 w 在虛擬相關文件 D_m 的次數，式*37 的最大化可透過期望值最大化迭代更新式來達成：

期望值步驟：

$$\begin{cases} r_{w, D_m} = \frac{c_w \cdot D_m \cdot [\lambda - \mu \cdot P_{\text{tkO}} \cdot w/S]}{[\lambda - \mu \cdot P_{\text{tkO}} \cdot w/S + \mu \cdot P_w \cdot BG + \lambda \cdot P_w \cdot D_m]} \\ e_{w, D_m} = \frac{c_w \cdot D_m \cdot \lambda \cdot P_w \cdot D_m}{[\lambda - \mu \cdot P_{\text{tkO}} \cdot w/S + \mu \cdot P_w \cdot BG + \lambda \cdot P_w \cdot D_m]} \end{cases} \quad (38)$$

最大化步驟：

$$\begin{cases} l_{\text{TriMM}} \cdot c_w \cdot S = \sum_{D_m \in \mathbf{D}_P} \frac{r_{w, D_m}}{\sum_w r_{w, D_m}} \\ \hat{P}_w \cdot D = \frac{e_{w, D_m}}{\sum_w e_{w, D_m}} \end{cases} \quad (39)$$

運用此三混合模型來調適語句模型時，庫爾貝克/萊伯勒離散度量值的公式*參照式*7 可進一步地表示成：

$$KL^*D \sim S^+ = \sum_{w \in V} P^*w \sim D^+ \frac{P^*w \sim D^+}{\gamma \cdot P^*w \sim S^+ + (3 - \gamma) R_{V|O}^*w \sim S^+} \quad *3: +$$

其中 $2 \leq \gamma < 3$ ，當 $\gamma = 2$ 代表使用三混合模型取代原本的語句模型。

關聯模型、簡單混合模型及三混合模型在資訊檢索領域中已被廣泛應用[34][54]；，但在摘要任務中卻是相對較少研究的，值得一提的是，雖然關聯模型、簡單混合模型已初步被應用在摘要任務上[6][3]；，但三混合模型卻是本論文首次引入到(語音)文字摘要任務中。

五、實驗語料及評估方法

7B、實驗語料

本論文實驗語料庫為公視新聞語料 *O cpf ctkp "Ej kpgug" Dtqcf ecuv" P gy u" Eqtr wu." O CVDP +，是由中央研究院資訊科學研究所耗時三年與公共電視台合作錄製並整理的中文新聞語料，其錄製內容為每天一個小時的公視晚間新聞深度報導。我們抽取其中由 4223 年 3 月到 4224 年 1 月總共 427 則新聞報導，區分成訓練集 *共 3: 7 則新聞+ 以及測試集 *共 42 則新聞+ 兩部分，其詳細的統計資訊如表一所示。全部 427 則語音文件長度約為 907 小時，我們先做人工切音，切出真正含有講話內容的音訊段落，再經由語音辨識器自動產生出的語音辨識結果稱之為語音文件 *Ur qngp" F qewo gpv." UF +，因此語音文件中只包含有語音辨識錯誤之雜訊；另一方面，我們將此 427 則語音文件藉由人工聽寫的方式，產生出沒有辨識錯誤的正確文字語料，我們稱之為文字文件 *Vgzv" F qewo gpv." VF +，每則文字文件再經由三位專家標記摘要語句，我們將此標記的人工摘要做為語音文件與文字文件的正確摘要答案。藉由比較語音文件和文字文件的摘要效能，我們可以觀察語音辨識錯誤對於各種摘要方法之影響。本研究的背景語言模型訓練語料取材自 4223 到 4224 年的中央社新聞文字語料 *Egptcn" P gy u" Ci gpe{." EP C +，並且以 UTK 語言模型工具訓練出經平滑化的單連語言模型，我們假設此單連語言模型為明確度中的非相關資訊之來源。另外，本論文蒐集 4224 年中央通訊社的約十萬則同時期新聞文字文件做為建立關聯模型時的檢索標的[6]，關於語句 S 的虛擬相關文件 *最高排序文件+ 篇數為 42 * 也就是 $D_{Top-2} 42 +$ ，而經由各種不同虛擬相關文件選取方法的篇數為 5 * 亦即 $D_p-2 5 +$ 。

7C、評估方法

自動摘要的評估方法主要有兩種，一為主觀人為評估，另一為客觀自動評估；前者為請幾位測試人員來為系統所產生的摘要做評估，給分的範圍為 3/7 分，後者則是預先請幾位測試者依據事先定義好的摘要比例挑選出適合的摘要語句，系統所產生的摘要句子將與測試者所挑選出的句子計算召回率導向的要點評估 *T gecm" Qtkpvgf " Wpf gtuwf {" hqt" I krlpi " Gxcnvcvkp." TQM G+ [38]。由於主觀人為評估非常耗時耗力，所以目前多數自動摘要方法皆採用召回率導向的要點評估做為文件摘要的評估方式，本論文亦採用此種評估方式。TQM G 方法是計算自動摘要結果與人工摘要之間的重疊單位元 *Wpku+ 數目占參考摘要 *T ghgtgpeg" Uwo o ct{ + 長度 * 單位元總個數+ 的比例。估計的單位元可以是 N 連詞 *N/i tco +、詞序列 *Y qtf " Ugs wgpegu+，如：最長相同詞序列或詞成對 *Y qtf " Rcku+。由於此方法是採用單位元比對的方式，不會產生語句邊界定義的問題，並且適合於多份人工摘要的評估。其評估的分數有三種，TQM G/3 *Wpli tco"，簡寫為 T/3+、

表一、實驗語料統計資訊

| | 訓練集 | 測試集 |
|--|----------------------------|----------------------------|
| 語料時間 | 4223 133 129/4224 123 144" | 4224 123 145/4224 12: 144" |
| 文件個數 | 3: 7" | 42" |
| 文件平均持續幾秒 | 34: 06" | 36304" |
| 文件平均詞個數 | 54802" | 4: 205" |
| 文件平均語句個數 | 4202" | 4505" |
| 文件平均字錯誤率 *Ej ctcevgT'Gttqt'Tcvg.'EGT+ " | 4: 0' " | 4: 0' " |
| 文件平均詞錯誤率 *Y qtf'Gttqt'Tcvg.'Y GT+ " | 5: 0' " | 5: 06' " |

TQM G/4*Di tco，簡寫為 T/4+和 TQM G/N*Nqi guv'Ego o qp'Uwdugs wpeg，簡寫為 T/N+分數，TQM G/3"是評估自動摘要的訊息量，TQM G/4"是評估自動摘要的流暢性，TQM G/N"是最長共同字串，本論文希望觀察摘要的流暢性，因此，實驗數據主要是以TQM G/4"分數為主。本論文所設定的摘要比例為 32'，其定義為摘要所含詞彙數占整篇文件詞彙數的比例，也就是以詞彙做為判斷摘要比例的单元。在挑選摘要語句過程中，若選到某語句中的某個詞彙時就已經剛好達到摘要比例，為了保持語句語意完整性，此語句剩下的詞彙也會被挑選成為摘要。

六、實驗結果

本論文主要著重在虛擬相關文件選取方法之發展與改進，是屬於非監督式摘要方法的範疇，因此所比較的對象以非監督式摘要方法為主；除此之外，本論文亦嘗試與現今最被廣為使用的監督式機器學習方法做比較，即支持向量機*UXO +32_。

8B、基礎實驗

首先，我們比較庫爾貝克/萊伯勒離散度*MN+與數個非監督式摘要方法之摘要成效，包含有最長語句摘要*Nqi guv'Ugpwpeg."NU+、首句摘要*NGCF +48_、向量空間模型*Xgevqt"Ur ceg"O qf gn"XUO +: _、潛藏語意分析*Ncvpv'Ugo cpve"Cpcnf uku."NUC +: _、最大邊際關聯*O czko cn' O cti kpcn'Tgrgxcpeg."O O T+4_以及整數線性規劃*Kvgi gt"Nlpgct" Rtqi tco o kpi ."KNR+44_。一般來說，文件中長句可能蘊含有較豐富的主題資訊，因此依據文件中語句長度做排序後，依序選取最長語句做為摘要結果是一種簡單的摘要方法。除此之外，也有學者研究發現，文件常以開門見山法的方式來提點出主題，因此文件開頭的前幾個語句經常是具代表性的語句，首句摘要即是以此概念出發，選取前幾句語句來形成整個文件的摘要。最長語句摘要*NU+及首句摘要*NGCF +都僅適用在一部分具有特殊結構的文件上，因此它們的缺點就是有其侷限性。另外，向量空間模型是把文件和語句分別視為一個向量，並使用詞頻/反文件頻*VH/KF H特徵來計算每一維度的權重值，文件與語句間的關聯性是藉由餘弦相似度量值來估測，當語句分數較高時，則越有機會成為此文件的摘要。潛藏語意分析是在向量空間的假設下更進一步地使用奇異值分解*Ukpi wrct"Xcnvg"F geqo r quklqp."UXF +來找到可能的潛藏語意空間，使之能在考量潛藏語意的情況下進行文件與語句的關聯性量測。最大邊際關聯可視為是向量空間模型的一個延伸，在做語句排序時考量了冗餘性以達到更好的摘要結果。整數線性規劃是一個全域*1 mqdcn"

表二、基礎實驗結果

| | | F/ueqtg" | | |
|------|-------------|---------------|---------------|---------------|
| | | TQM G/3" | TQM G/4" | TQM G/N" |
| VF " | NU" | 2047" | 202; ; " | 203: 5" |
| | NGCF " | 20532" | 203; 6" | 20498" |
| | XUO " | 20569" | 2044: " | 204; 2" |
| | NUC " | 20584" | 20455" | 20538" |
| | O O T" | 2058: " | 2046: " | 20544" |
| | MN" | 20633" | 204; ; " | 20583" |
| | KNR" | 0.442" | 0.337" | 0.401" |
| UF " | NU" | 203: 3" | 20266" | 2035: " |
| | NGCF " | 20477" | 20339" | 20443" |
| | XUO " | 20564" | 203; ; " | 204: 9" |
| | NUC " | 20567" | 20423" | 20523" |
| | O O T" | 0.366" | 0.215" | 0.315" |
| | MN" | 20586" | 20432" | 20529" |
| | KNR" | 2056: " | 2042; " | 20528" |

的限制性最佳化*Eqputclpv'Qr wo k cvkqp-的語句選取方法[44]。表二為本論文之基礎實驗結果。首先，在 VF "的實驗中，MN"的摘要效果比 NU、NGCF、XUO、NUC、O O T"等非監督式摘要方法來得好些；因 NU"與 NGCF "僅適用於特殊文章結構上，所以若被摘要文件不具有某種特殊的文章結構，其摘要效能就會有限。相較之下，MN"是較具一般性的摘要方法，因此比較不會受限於文章的結構之影響，故摘要效能比 NU"以及 NGCF "來得彰顯。MN"與 XUO "皆使用淺層的詞彙*詞頻-資訊，但由於 MN"是計算語句模型與文件模型之間的距離關係，對於代表語句與文件的語言模型，我們較容易透過各種技術來進行模型的估計與調適，進而獲得較好的摘要成果。O O T"在選取時多考慮了冗餘資訊，所以摘要效果也比 XUO "來得好些。NUC "在潛藏語意空間計算文件與語句的餘弦相似度量值，其結果顯示也會較 XUO "好。整數線性規劃是一個全域選擇方法，所以在 VF "上可以得到最好的摘要效能。另一方面，在 UF "的實驗中，MN"同樣較優於 NU、NGCF "之摘要方法，但 XUO "的結果則稍微較 MN"好一點，我們認為這可能是因為 XUO "比較不受到語音辨認錯誤的影響。原以為 KNR"也會在 UF "中得到最好的摘要效能，結果反而是 O O T"得到最好的摘要效能，可能的原因是 KNR"受到語音辨識錯誤的影響比較大，造成其摘要結果不彰。

通常語音文件主要會有語音辨識錯誤和語句邊界偵測錯誤的問題，但我們有先經人工切音，因此摒除了語句邊界偵測錯誤的問題，藉由比較 VF "與 UF "之實驗結果，我們可以觀察語音辨識錯誤率對摘要結果的影響性。比較各式方法，UF "比 VF "下降了 30' 0' 的 TQM G/4"摘要效能，由此可知語音辨識錯誤率對摘要效能是有顯著的影響性。為了減緩語音辨認錯誤的問題，在未來我們將嘗試使用音節*U{ mcdg-為單位來建立語句以及文件模型；或利用詞圖*Y qtf "I tcrj +-、混淆網路*Eqphwkqp" P gy qtm-來含括更多的可能正確候選詞彙以裨益模型估測；更可利用韻律資訊*Rtquqfle" kphqto cvkqp-等聲學線索來輔助減緩語音辨認錯誤對摘要效能的影響。

表三、關聯模型之實驗結果*使用最高排序文件前三篇*Vqr 5+*

| | | F/ueqtg" | | |
|------|----------|---------------|---------------|---------------|
| | | TQM G/3" | TQM G/4" | TQM G/N" |
| VF " | MN" | 20633" | 204; : " | 20583" |
| | TO " | 20672" | 20558" | 20622" |
| | UO O " | 20658" | 20547" | 205: 7" |
| | VtkO O " | 0.457" | 0.350" | 0.404" |
| UF " | MN" | 20586" | 20432" | 20529" |
| | TO " | 20596" | 20448" | 20543" |
| | UO O " | 20597" | 20443" | 20536" |
| | VtkO O " | 0.379" | 0.228" | 0.325" |

804、基礎關聯模型之實驗

使用關聯模型於語句模型之建立時，需要做一次的資訊檢索來為每個語句找出虛擬相關文件，由同時期的新聞文字文件*共 323.48: "篇+中為每一語句選取出 42"篇虛擬相關文件，但為了要與後續虛擬相關文件選取方法作公平的比較，因此此基礎關聯模型實驗是取前三篇*Vqr 5+來進行關聯模型之估測與相關實驗]6_。由於文件中的語句通常相對簡短，因此當使用最大化相似度估測建立語句模型時，容易遭遇資料稀疏的問題，不容易獲得精準的模型，故我們期望考慮額外的關聯資訊於語音文件摘要，亦即藉由虛擬相關文件來重新估測並建立語句的語言模型，能獲得進一步地摘要成效。重新估測後的關聯模型則可與原本的語句模型相結合或取代之，相結合的參數調整在本實驗中是採用經驗設定*Go r ktlecnUgwłpi +。實驗結果如表三所示，在 VF "與 UF "之摘要成效上，使用關聯模型*TO +、簡單混合模型*UO O +及三混合模型*VtkO O +皆能比基礎的 MN"實驗較好，尤其是三混合模型*VtkO O +相較於 MN"在 VF "及 UF "的 TQM G/4"結果上能有 704' 與 30' "的改進。接著，我們比較不同關聯模型的摘要成效，首先是關聯模型*TO +與簡單混合模型*UO O +的比較，從表三的實驗結果得知關聯模型在 VF "上表現比簡單混合模型來得好，但在 UF "似乎在 TQM G/3"就沒比簡單混合模型好，不過 UF "的 TQM G/4"跟 TQM G/N"都還是比簡單混合模型的效果好。關聯模型的假設是強調詞彙 w 與語句 S 在這些虛擬相關文件中同時出現之關係*參照式*32+來估測模型，而簡單混合模型是強調訓練好的模型能讓有獨特性的詞彙得到更多的機率值因而讓模型具有鑑別能力，兩者皆有其好處。最後，三混合模型*VtkO O +因複雜化了簡單混合模型*UO O +，額外多考量文件模型的影響力，因此相較於關聯模型及簡單混合模型能得到更佳的摘要效能，三混合模型相較於關聯模型在 VF "上有明顯的進步，於 TQM G/4"結果能有 306' 的改進，但在 UF "上，於 TQM G/4"結果只有微量的 204' 改善。

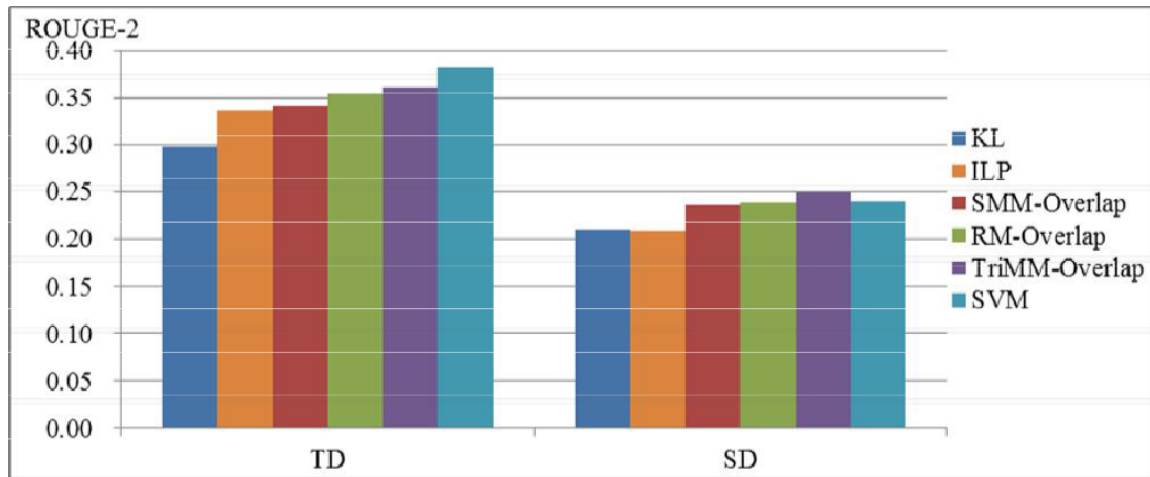
在關聯模型的相關實驗中，語音辨識錯誤也是影響摘要效能非常嚴重，在三混合模型的數據中，UF "比 VF "劇烈下降了 3404' 的 TQM G/4"摘要效能，在未來研究中，我們認為可以以次詞索引*Uwdy qtf "Kpf gzłpi +的方式來建立關聯模型以減緩語音辨識錯誤之影響。

表四、各種虛擬相關文件選取方法於關聯模型之實驗結果

| F/ueqtg" | | TO " | | | UO O " | | | VtkOO " | | |
|----------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | | T/3" | T/4" | T/N" | T/3" | T/4" | T/N" | T/3" | T/4" | T/N" |
| VF | Vqr 5" | 20672" | 20558" | 20622" | 20658" | 20547" | 205: 7" | 20679" | 20572" | 20626" |
| | I crr gf "K" | 20673" | 2055: " | 20623" | 20655" | 20539" | 205: 7" | 20676" | 20565" | 20628" |
| | Egptqkf " | 2066; " | 20556" | 20624" | 2065; " | 20553" | 205: ; " | 20678" | 20575" | 20629" |
| | Cevkxg/TFF" | 20682" | 20563" | 2062: " | 2066; " | 20564" | 20622" | 20685" | 20577" | 20636" |
| | Cevkxg/TFFP" | 20686" | 20574" | 20633" | 20677" | 0.346" | 20627" | 20688" | 0.367" | 20643" |
| | Qxgttr r gf " | 0.470" | 0.354" | 0.416" | 0.460" | 20563" | 0.410" | 0.471" | 20584" | 0.422" |
| UF | Vqr 5" | 20596" | 20448" | 20543" | 20597" | 20443" | 20536" | 2059; " | 2044: " | 20547" |
| | I crr gf "K" | 20596" | 2044: " | 20544" | 20593" | 2043: " | 20535" | 20598" | 20447" | 20537" |
| | Egptqkf " | 20596" | 20449" | 20536" | 20599" | 20449" | 20542" | 205: 2" | 20455" | 2054: " |
| | Cevkxg/TFF" | 2059; " | 2044: " | 20554" | 2059: " | 2044; " | 20543" | 205: : " | 20464" | 20557" |
| | Cevkxg/TFFP" | 205: 5" | 2045; " | 20552" | 205: 2" | 20448" | 20549" | 205: 3" | 20466" | 2055; " |
| | Qxgttr r gf " | 0.386" | 0.239" | 0.334" | 0.382" | 0.236" | 0.332" | 0.396" | 0.250" | 0.345" |

805、各種虛擬相關文件選取方法於關聯模型之實驗

本小節的實驗是由第一次虛擬相關文件*最高排序文件+"D_{Top} 中*D_{top}~242+再精鍊選取出較佳的虛擬相關文件 D_p *D_p~25+，我們比較所提出兩種新穎的選取方法*即主動式/關聯多元密度非相關*Cevkxg/TFFP+和重疊分群*Qxgttr r gf+與其他現有選取方法*即間隔式最高 K 選取法*I crr gf "K+、群中心選取法*Egptqkf+以及主動式/關聯多元密度*Cevkxg/TFF+於各種關聯模型*TO、UO O"及 VtkOO+之摘要效能比較，實驗結果如表四所示，與基礎關聯模型只使用前三篇*Vqr 5+虛擬相關文件的結果相較*參照表三+，大部分透過虛擬相關文件選取方法都會比只使用 Vqr 5"的摘要結果還要來得好，除了 I crr gf "K 無論在 VF"和 UF"中，使用 UO O"與 VtkOO"都會有比 Vqr 5"差的摘要效能，因為 I crr gf "K 是一個較不穩定的虛擬相關文件選取方法，在本實驗中有比較差的結果是可以預期的。群中心選取法*Egptqkf+表現尚可，在 VF"及 UF"中，於各種關聯模型*TO、UO O"及 VtkOO+下，比 Vqr 5"及 I crr gf "K 都要來得好。Cevkxg/TFF"因在選取虛擬相關文件時同時考量了關聯性*Tgrxcpeg+、多元性*F kgtuk{+以及密度性*F gpuk{+，用於不同的關聯模型訓練時，相對於 Vqr 5、I crr gf "K 以及 Egptqkf，無論在 VF"或 UF"中，都可以得到更好的摘要結果。Cevkxg/TFFP"在多考量了非相關*P qp/tgrxcpeg+資訊的情況下，其實驗結果都會比現有的選取方法較佳，相對於 Cevkxg/TFF、Egptqkf、I crr gf "K 以及 Vqr 5"無論在 VF"或 UF"中，各種關聯模型*TO、UO O"及 VtkOO+下都會得到比較好的摘要結果，所以證實非相關資訊確實一個有用的選取線索。最後本論文所提出的重疊分群*Qxgttr r gf+選取方法無論在 VF"或 UF"中，於各種關聯模型下*TO、UO O"及 VtkOO+都可以得到最好的摘要效果，驗證了重疊分群在利用虛擬相關文件中結構化資訊確實可以找到具代表性的文件以利各種關聯模型的模型訓練或參數估測。



圖二、UXO 與其他非監督式摘要方法之比較

806、與監督式模型之比較

除了各式非監督式摘要方法外，本論文亦嘗試比較支持向量機 (SVM) 於文件摘要之成效，比較的對象有基礎 (MN) 以及使用重疊分群選取方法於不同的關聯模型中 (TO/Qxgtrr、UO/Qxgtrr 和 VtkO/Qxgtrr)。支持向量機是現今常見的監督式機器學習方法之一，近年來已有學者將其運用至文件摘要領域之中 [32]。本論文使用訓練集的 3:7 篇文件進行支持向量機模型的訓練語料，我們為文件中的每一語句抽取 57 維特徵 [35]，包括有韻律特徵 (Rtquqf ke) Hgcwtgu、語彙特徵 (Ngzkecn) Hgcwtgu、結構特徵 (Utwewtcn) Hgcwtgu 以及基本的模型特徵 (O qf gn) Hgcwtgu 等資訊，其核心函數設定為半徑式函數 (Tcf kcnDcuku) Hwpevkqp，其中 UXO 的參數設定都是使用預設值。

實驗結果如圖二所示。一如預期地，UXO 與其他各式非監督式模型相比較，在 VF 實驗上 (其 TQWIG/4 為 20:5) 是表現最好的方法，這是由於監督式機器學習藉由使用人工標註的摘要句子進行模型之訓練，其使用的資訊較非監督式機器學習方法多且正確，因此其摘要的效果也較非監督式機器學習來的好。值得一提的是，使用重疊分群虛擬相關文件選取方法於三混合模型中 (VtkO/Qxgtrr) 摘要之成效在 UF 上可比監督式機器學習方法的 UXO 來的好一些，此一實驗結果令人感到驚訝，因為本論文所探討之各式摘要方法僅考慮了文件與語句中的單一種特徵值，即藉由詞彙分佈資訊來挑選語句，而支持向量機不僅使用了 57 種特徵值，更需要使用人工標註的正確答案進行模型的訓練。我們認為，此結果之原因可能是由於支持向量機之摘要技術在語音辨識錯誤的情況下 (在此實驗中，訓練集與測試集的詞錯誤率達 62%)，未必能真的有效學習分辨摘要與非摘要語句。

七、結論與未來方向

本論文基於語言模型化架構來發展語音摘要方法，其貢獻主要有三方面。第一，有別於現有基於語言模型化架構之摘要方法都聚焦在語句模型參數的重新估測，本論文首次深入探討與應用各種新穎虛擬相關文件選取技術於節錄式語音文件摘要任務，用以強化語句模型的參數估測。第二，本論文更進一步地考量使用每一語句的非相關性 (Pqp/tgrgxcpeg) 資訊對於虛擬相關文件選取的影響。同時，我們亦額外嘗試基於重疊分群 (Qxgtrr r gf) (Enwvgtlpi) 概念來有效地選取重要的虛擬相關文件。第三，本論文探索使用三混合模型 (VtkO lzwtg) (O qf gn) 來表示每一語句，期盼其能更精確地表示語句之詞彙使用和語意相關資訊。一系列的實驗結果顯示，本論文所提出之方法的確能較其它現有

的非監督式摘要方法有更加的摘要效能表現。未來，我們的研究將有三個主要的方向：

首先，本論文所提出之虛擬相關文件選取

方法是建構在向量空間或語言模型空間上，並沒有考慮到語意空間的相似度量，我們將進一步的研究是否可以在潛藏語意空間中來選取較好的虛擬相關文件，以期獲得更好的摘要成效；其次，目前所發展的關聯模型僅運用於重建語句的語言模型，我們將嘗試使用被摘要文件的關聯資訊來重新估測並建立文件的語言模型；最後，我們希望能將非監督式方法所形成的特徵結合於更加複雜且有效的監督式機器學習方法如 ETH 或深度類神經網絡 F ggr "P gwtcnlP gy qtm'Ngctplpi ."F P P 等+中，並融合其它語音文件所獨有之特徵諸如音韻與語者特徵等+，期望訓練後的模型能夠在語音文件摘要上獲得更好的表現。

致謝

本論文之研究承蒙教育部國立臺灣師範大學邁向頂尖大學計畫*324L3C2: 22+與行政院科技部研究計畫 *O QUV"325/4443/G/225/238/O [4."P UE"325/4; 33/K225/523."P UE"323/4443/G/225/246/O [5"、 P UE"323/4733/U/225/279/O [5"、 P UE"323/4733/U/225/269/O [5"和 P UE"324/4443/G/225/236/"O [5+之經費支持，謹此致謝。

參考文獻

- [3_] RØ Dczgpf crg." *Machine-made Index for Technical Literature – an Experiment*, IDO " Lqwtpcn'qh'Tgugctej "cpf "F gxgnr o gpv."Xqf04."P q06."r r 05766583."3; 7: "
- [4_] IØ Ectdqpgm' cpf " IØ I qf uvglp." *The Use of MMR Diversity-based Reranking for Reordering Documents and Producing Summaries*, Rtqeggf lpi u" qh" vj g" 43^y" Cppwcn' Kvgtpcvkpcn' CEO "UK KT" Eqphgtgpeg" qp" Tgugctej "cpf "F gxgnr o gpv" lp" Kphqto cvkqp" Tgtlgxcn" *UK KT+."r r 05576558."3; ;: "
- [5_] [ØVØ Ej gp." DØ Ej gp" cpf " J ØO Ø' Y cpi ." *A Probabilistic Generative Framework for Extractive Broadcast News Speech Summarization*, KGGG"Vtcpuvcvqpu"qp" Cwf lq. "Ur ggej " cpf "Ncpi wci g" Rtqeguulpi ."Xqf039."P q03."r r 0; 76328."422; "
- [6_] DØ Ej gp." J ØEØ Ej cpi ." MØ [Ø' Ej gp." *Sentence Modeling for Extractive Speech Summarization*." Rtqeggf lpi u" qh" vj g" Kvgtpcvkpcn' Eqphgtgpeg" qp" O wko gf kc" ("Gzr q" *KEO G+."4235"
- [7_] DØ Ej gp." [ØY Ø' Ej gp" cpf " MØ [" Ej gp." Gpj cpe lpi " S wgt { " Hqto wrcvqp" hqt" Ur qngp" F qewo gpv" Tgtlgxcn" Lqwtpcn'qh' Kphqto cvkqp" Uelgpeg" cpf "Gpi lpggt lpi ."Xqf052."P q05."r r 0 775678; ."4236"
- [8_] LØO Ø' Eqptq { " cpf " F ØRØ' QaNgct { ." *Text Summarization via Hidden Markov Models*, Rtqeggf lpi u" qh" vj g" 46^y" Cppwcn' Kvgtpcvkpcn' CEO "UK KT" Eqphgtgpeg" qp" Tgugctej "cpf " F gxgnr o gpv" lp" Kphqto cvkqp" Tgtlgxcn" *UK KT+."r r 06286629."4223"
- [9_] I ØGtnp" cpf " F ØTØ' Tcf gx." *LexRank: Graph-based Lexical Centrality as Saliency in Text Summarization*." Lqwtpcn'qh' Ct vhekn' Kvgnki gpv" Tgugctej ."Xqf044."P q03."r r 0679669; ." 4226"
- []: _ [ØI qpi " cpf " ZØNkw." *Generic Text Summarization using Relevance Measure and Latent Semantic Analysis*, Rtqeggf lpi u" qh" vj g" Kvgtpcvkpcn' CEO "UK KT" Eqphgtgpeg" qp" Tgugctej "cpf "F gxgnr o gpv" lp" Kphqto cvkqp" Tgtlgxcn" *UK KT+."r r 03; 647."4223"

-]3_ "F0'J kgo utc." U0' Tqdgtrvqp." cpf "J 0' \ ctcic q| c." *Parsimonious Language Models for Information Retrieval.* Rtqeggf kpi u" qh" yj g" kpvgtpcvkqpcn' CEO "UK KT" eqphgtgpeg" qp" Tgugctej "cpf "f gxgnr o gpv'kp' kphqto cvkqp" Tgtkxncn" *UK KT+ "r r 039: 63: 7."4226"
-]32_ "C0'Mqre| ." X0'Rtcdncnto wt vj k'cpf "L0'Mrks." *Summarization as Feature Selection for Text Categorization.* Rtqeggf kpi u" qh" yj g" kpvgtpcvkqpcn' Eqphgtgpeg" qp" kphqto cvkqp" cpf "Mpqy rfi g'O cpci go gpv" *E KMO + "r r 05876592."4223"
-]33_ "L0'Mx kge." *A Trainable Document Summarizer,* Rtqeggf kpi u" qh" yj g" Cppwcn' kpvgtpcvkqpcn' CEO "UK KT" Eqphgtgpeg" qp" Tgugctej " cpf " F gxgnr o gpv' kp' kphqto cvkqp" Tgtkxncn' *UK KT+ "r r 08: 695."3; ; 7"
-]34_ "X0'Ncxtgpmq" cpf "Y0'D0'Etqhv." *Relevance -based Language Models,* Rtqeggf kpi u" qh" yj g" 46^y " Cppwcn' kpvgtpcvkqpcn' CEO "UK KT" Eqphgtgpeg" qp" Tgugctej " cpf " F gxgnr o gpv' kp' kphqto cvkqp" Tgtkxncn' *UK KT+ "r r 03426349."4223"
-]35_ "U0'J 0' Nkp" cpf " D0' Ej gp." *Improved Speech Summarization with Multiple-hypothesis Representations and Kullback-Leibler Divergence Measures,* Rtqeggf kpi " qh" yj g" 32^y " Cppwcn' Eqphgtgpeg" qh" yj g" kpvgtpcvkqpcn' Ur ggej " Eqo o wplecvkqp" Cuuqekcvkqp" *kpvgtur ggej + "r r 03: 6963: 72."422; "
-]36_ "J 0' Nkp" cpf "L0' Dkro gu." *Multi-document Summarization via Budgeted Maximization of Submodular Functions.* Rtqeggf kpi "qh' P CCEN' J NV." r r 0; 346; 42."4232"
-]37_ "U0'J 0' Nkp." [00' 0' [gj " cpf "D0' Ej gp." *Leveraging Kullback-Leibler Divergence Measures and Information-rich Cues for Speech Summarization,* KGGG" Vtcpuvcvqpu" qp" Cwf kq." Ur ggej "cpf "Ncpi wci g' Rtqeguulpi 0Xq103; . 'P q06." r r 0: 936: : 4."4233"
-]38_ "E0' [0' Nkp." *ROUGE: Recall-oriented Understudy for Gisting Evaluation.* 4225"]Qprkpg_0' Cxckrcdng < j wr < l j c { f p 0 u l q f w T Q W L G I 0'
-]39_ "[0' Nkw" cpf "F 0' J cmrpkVwt." *Speech Summarization.* "kp" I 0' Vwtcpf "T0' F0' O qt' k']Gf _." Ur qngp "Ncpi wci g" Wpf gtucpf kpi < U{ ugo u" hqt" Gz vcevkpi "Ugo cpvle" kphqto cvkqp" htqo " Ur kgej . "Y kgej . "4233"
-]3_ "Z0' Nkw" cpf "Y 0'D0'Etqhv." *Cluster-based Retrieval Using Language Models.* "Rtqeggf kpi u" qh" yj g" kpvgtpcvkqpcn' CEO "UK KT" eqphgtgpeg" qp" Tgugctej " cpf " F gxgnr o gpv' kp' kphqto cvkqp" Tgtkxncn' *UK KT+ "r r 03: 863; 5."4226"
-]3;_ "U0'J 0' Nkw." MD[0' Ej gp. "[0' N0' J ulgj . "D0' Ej gp. "J 0' 0' 0' Y cpi . "J 0' E0' [gp. "Y0' N0' J uw." *Effective Pseudo-relevance Feedback for Language Modeling in Extractive Speech Summarization.* "Rtqeggf kpi u" qh" yj g" KGGG" kpvgtpcvkqpcn' Eqphgtgpeg" qp" Ceqwulku. "Ur ggej . " cpf "Uki pcn' Rtqeguulpi " *ECUUR+ "4236"
-]42_ "U0'J 0' Nkw." MD[0' Ej gp. "D0' Ej gp. "J 0' 0' 0' Y cpi . "Y0' N0' J uw." *Improving Sentence Modeling Techniques for Extractive Speech Summarization.* "TQENKI " ZZ X < Eqphgtgpeg" qp" Eqo r wcvkqpcn' Nkpi wkuu" cpf "Ur ggej "Rtqeguulpi . "4235"
-]43_ "K0' O cpk" cpf "O 0' V0' O c { dwt { ." *Advances in Automatic Text Summarization,* Eco dtki g < O K' Rtguu." 3; ; ; "
-]44_ "T0' O eF qpcrf ." *A Study of Global Inference Algorithms in Multi-document Summarization.* "Rtqeggf kpi u" qh" Gwtqr gcp" Eqphgtgpeg" qp" kphqto cvkqp" Tgtkxncn' *GEKI+ " r r 0' 7796786." 42290'

- 145_ "T0' O lj cægc" cpf "R0' Vctcw." *TextRank Bringing Order into Texts*, Rtqeggf lpi u" qh' Go r klcenO gj qf "lp' P cwtcnNcpi wci g'Rtqeguikpi "GO P NR+."r r 06266633."4226"
- 146_ "I 0'0 wttc{."U0Tgpcn."cpf "L0'Ectrgwc." *Extractive Summarization of Meeting Recordings*, Rtqeggf lpi u"qh"vj g"Eqphgtgpeg"qh"vj g"Kvgtpcvkqpcn'Ur ggej "Eqo o wplecvkqp"Cuqekcvkqp" *Kvgtur ggej +."r r 07; 567; 8."4227"
- 147_ "C0'P gpnqxc" cpf "M0'O eMggy p." *Automatic Summarization*."Hqwpf cvkqpu" cpf "Vtgpf u"lp" Kphqto cvkqp" Tgtlgxcn" Xqt07. "P q0465-3256455."4233"
- 148_ "I 0' Rgpp" cpf "Z0'\ j w" *A Critical Reassessment of Evaluation Baselines for Speech Summarization*, Rtqeggf lpi u"qh" Cppwcn' O ggkpi "qh"vj g"Cuqekcvkqp" hqt "Eqo r wcvkqpcn' Nkpi wkuu."r r 0692669: ."422: "
- 149_ "Z0'Uj gp" cpf "E0'\ j ck" *Active Feedback in Ad Hoc Information Retrieval*."Rtqeggf lpi u"qh' CEO "UK KT" Eqphgtgpeg" qp" Tgugctej " cpf " F gxgr o gpv' lp" Kphqto cvkqp" Tgtlgxcn' *UK KT+."r r 077-88."4227"
- 14_ "E0'Uj gp" cpf "V0'Nk" *Multi-document Summarization via the Minimum Dominating Set*." Rtqeggf lpi u"qh"vj g"Kvgtpcvkqpcn'Eqphgtgpeg"qp"Eqo r wcvkqpcn'Nkpi wkuu"*EQNKPI +."r r 0; : 66; 4."4232"
- 14; _ "F0'Uj gp." L0'V0' Uwp." J 0'Nk" S0' [cpi ." cpf "\ 0'Ej gp." *Document Summarization using Conditional Random Fields*, Rtqeggf lpi u"qh"Kvgtpcvkqpcn'Lqkv'Eqphgtgpeg"qp"Ct vlcen' Kvgnti gpeg"*KECK."r r 04: 8464: 89."4229"
- 152_ "Z0'Y cp" cpf "L0' [cpi ." *Multi-document Summarization using Cluster-based Link Analysis*, Rtqeggf lpi u"qh" vj g" Cppwcn' Kvgtpcvkqpcn' CEO "UK KT" Eqphgtgpeg" qp" Tgugctej " cpf " F gxgr o gpv'lp" Kphqto cvkqp" Tgtlgxcn" *UK KT+."r r 04; ; 6528."422: "
- 153_ "\ 0'Zw" T0'Cngm" cpf "[0' j cpi ." *Incorporating Diversity and Density in Active Learning for Relevance Feedback*."Rtqeggf lpi u"qh' Gwtqr gcp "Eqphgtgpeg"qp" Kphqto cvkqp" Tgtlgxcn' *GEKT+."r r 04676479."4229"
- 154_ "E0'Z0'\ j ck" cpf "L0'Nchgtv\ ." *Model-based feedback in the language modeling approach to information retrieval*." Rtqeggf lpi "qh"vj g" Kvgtpcvkqpcn'Eqphgtgpeg"qp" Kphqto cvkqp" cpf " Mpqy rfi g'O cpci go gpv" *EKMO+."r r 06256632."4223"
- 155_ "E0'Z0'\ j ck" cpf "L0'Nchgtv\ ." *A Study of Smoothing Methods for Language Models Applied to Ad Hoc Information Retrieval*, Rtqeggf lpi u"qh"vj g" Cppwcn' Kvgtpcvkqpcn' CEO " UK KT" Eqphgtgpeg"qp" Tgugctej " cpf " F gxgr o gpv'lp" Kphqto cvkqp" Tgtlgxcn" *UK KT+."r r 0' 556/564."4233"
- 156_ "E0'Z0'\ j ck" *Statistical Language Models for Information Retrieval: A Critical Review*," Hqwpf cvkqpu" cpf "Vtgpf u"lp" Kphqto cvkqp" Tgtlgxcn" Xqt04. "P q06."r r 0596435."422: "
- 157_ "L0'\ j cpi " cpf "R0' Hwpi ." *Speech Summarization without Lexical Features for Mandarin Broadcast News*, Rtqeggf lpi u"qh' P CCEN' J NV, Eqo r cplqp" Xqnwo g."r r 04356438."4229"

Vj g"4236"Eqphgtgpeg"qp"Eqo r wcvkqpcrNkpi wku"cpf "Ur ggej "Rtqeguulpi ""
TQENRPI "4236."rr043/52" ""
©"Vj g"Cuuqekvq"lqt"Eqo r wcvkqpcrNkpi wku"cpf "Ej kpgug"Ncpi wci g"Rtqeguulpi "

台灣情緒語料庫建置與辨識

Bo-Chang Chiou and Chia-Ping Chen

摘要

現在已有許多公開情緒語料庫被實驗於語音情緒辨識的研究上，但是並沒有一個語料庫以台灣常用的語言所錄製。語音情緒辨識會因為語音中不同的文本和語言等資訊影響最後的辨識結果，因此公開的外國語言語料庫不一定適用於台灣語言的情緒辨識研究中。為了解決語料庫的問題我們自行錄製的台灣語言情緒語料庫，採用了最普遍的國語、台語以及客家語三種語言。我們的語料庫仿照德語公開語料庫 GO Q/FD 錄製，對每種語言採用十位語者、十句文本以及七種情緒，每個語言收錄七百句，並且在錄製完後進行人工辨識，以一定的人工辨識率做為篩選條件，以確保一定的辨別度。錄製完後以一個龐大的聲學特徵集搭配支持向量機作為後端分類器，以此實驗做為三種語言的基準辨識率。"

關鍵字：語音情緒辨識、情緒語料庫建置

一、緒論

在人類社會中言語溝通扮演著非常重要的角色，人們透過語言來交流和傳遞訊息，基於此理念下發展出許多人機介面科技* wo cp/eqo r wgt "kpvgtcevkq."J EK，如蘋果公司的 \$UkIS 便是此語音領域的代表作之一，將人類從鍵盤滑鼠輸入限制到能不需手動 *Vqwej nguu 的聲控操作。而人類的交流不只單單透過語言文字，還有情感的交流，情感表達能包含人們所處的狀態，因此有了情感運算的研究。情感運算橫跨許多領域如電腦科學、心理學、認知哲學與工程學等，在語音領域相關為語音情緒辨識及情緒語音合成。基本的情感辨識透過分析人的臉部表情以及聲音來做辨識，更進一步則同時分析對話語意作為分辨依據。"

13_這本小書開始了情感運算的時代，此書定義了情感運算，從一個資訊工程研究者的角度來描述與說明情感運算的應用及重要性，解釋基本的訊號處理與機器學習分類的觀念。情感運算為一個龐大的研究議題，其中橫跨多個領域如資訊科學、心理學、感知科學和工程學等等。14_對語音情緒辨識做了大略的概論，首先整理了研究中普遍常見的情緒語料庫，探討基本特徵參數如音高 *Rkej +、能量 *Gpgti {+、基頻 *Hvpf co gpvcnHtgs wgpe {.H2+ 和梅爾倒頻譜係數 *O gnHtgs wgpe {"Egr utcrnEqghkkgpw."O HEE+等，後端則實驗高斯混合模型 *I cwulcp "O kzwtg "O qf gn "I O O +、支持向量機 *Uwr r qtv "Xgevqt "O cej kpg."UXO +、隱藏式馬可夫模型 *J kf f gp "O ctmqx "O qf gn "J O O +和類神經網絡 *Ct vkkkcrn "P gwtcrn "P gvy qtmu." CPP +。"

研究語音情緒辨識就必須具有情緒語音語料庫用於實驗，而語料庫的收錄會直接影響研究的實驗結果，目前常用的公開情緒語料庫如德語的 Dgtrkp"Fcwcdug"qh"Ur ggej "Go qvkqp" *GO Q/FD+5_、HCWCkldq]6和 Xgtc/Co /O kwci *XCO +7_。]8提及在收集情感資料時，人們處於實驗室等非平常習慣的場合所收集到的情感資料實際上會與真實世界表達出的有所差異，因此情緒語音語料庫的錄製大致能分為兩類；第一類為引導性錄音，此類多在實驗室或者是錄音室中錄製，如德語的 GO Q/FD，其透過在錄音室錄製並且有三位語言學專家在旁監督與指導，如此可以確保錄製出乾淨並且情緒表達度和差異性高的情緒語料。此種錄製方式較需要語者的表達能力，在語者表達能力較差的情況下無法在錄音室錄製出含有足夠情緒表達的語句-相對的第二種則是自發性的情緒表達語句，如 HCW' Cldq 和 XCO，HCW' Cldq 透過錄製小孩子與機器人的互動，便可以直接收集到小孩子自然表達的情緒語句。而 XCO 則是從談話性節目中節錄片段再從中挑選出情感較為活躍的句子，如此收錄的語料庫也較自然呈現。"

" "" 在語音情緒辨識中提出作為基準實驗的有兩者。由於情感辨識相對於語音辨識及語者辨識為較新的領域，因此較早並沒有一個公認的基準可供參考比較，所以 K̂ VGTURGGEJ " 422; " Go qvkqp" Ej cmgpi g]9便採用了具有明確訓練集及測試集的 HCW' Cldq 情緒語料庫，在前端使用過零率*] gtq"etqukpi "tcvg." ET+、能量、基頻、泛音噪音比 *] cto qpleu/vq/pqlug." J PT +和梅爾倒頻譜系數。為了解決 HCW' Cldq 五類資料量中有兩類佔了近 : 2' 的資料數量嚴重不平均的狀況，使用 UO QVG*U{pyj gvk"O k̂qtkv "Qxgtuco r ikpi " Vgej pqmqi {+; 方法來增加數量較少的類別的樣本數，後端則以隱藏式馬可夫模型以及支持向量機做為辨識器，以此做為 HCW' Cldq 的一個實驗基準。]; 實驗；個情緒語料庫分別為 Cktr n̂pg" Dgj cxkqwt" Eqtr wu" *CDE+、 Cwf k̂xk̂wcn' k̂vgtguv" Eqtr wu" *CXK̂+、 Fcpluj " Go qvkqpcr"Ur ggej " *FGU+、 GO Q/FD、 gP VGTHCEG、 Ugpuk̂xg" Ct v̂k̂ekcn' Nk̂vpggt " *UCN+、 Uo ct v̂k̂qo、 Ur ggej " Wpf gt " Ulo ŵv̂gf " cpf " Cewcn' Ut̂guu" *UWUCU+、 和 XCO 等資料庫，這些語料庫分別有不同的錄音條件，包含引導式錄音、自發性語音與各種的錄音環境等等。實驗分為兩種，第一種為擷取一般音框特徵並搭配隱藏式馬可夫模型+高斯混合模型做為後端辨識器，另一種為截取大量的聲學資訊並經由泛函將音框特徵轉換為句子特徵，如此以方便對應至後端支持向量機的訓練與辨識，此篇可視為各情緒語料庫採用聲學模型時的基準實驗參考。"

" "" 目前常見的公開語音情緒辨識的語料庫並沒有以台灣本土語言錄置而成，而語音中包含了語言、語者、文本、腔調等等的資訊，這些皆會影響語音情緒辨識的準確度。因此為了提供一個較適合用於台灣語言的情緒語料庫，我們自行錄製了台灣最普遍的國語、台語及客家語三種語言的情緒語料庫。我們仿造 GO Q/FD 錄製方式錄製，其中相同的七種情緒、五位男性五位女性共十位語者以及十句文本，錄製完後進行人工辨識測驗。最後再採用與]; 相同的基準實驗設定，作為我們語料庫的基準辨識率。"

" "" 以下為整篇論文的架構，第一節為緒論，第二節詳細介紹我們仿造的德語語料庫 GO Q/FD，第三節詳列我們自行錄製的台灣語料庫資訊，第四節為台灣語料庫的基準實驗，最後為我們的結論。"

二、現有公開語料庫

我們錄製的語料庫參考已廣泛實驗於語音情緒研究的德語語料庫 Dgtrkp"Fcwcdug"qh' Go qvkkp"Ur ggej *GO Q/FD+。GO Q/FD 由 32 位德語母語語者所錄製而成，經由麥克風錄製 6: mJ | 後降為 38mJ |。錄音環境在隔音良好及具有高品質錄音設備的錄音室，因此表達的情緒並非自然呈現。語者由 5 位語言學專家從 62 位新聞播報員中挑選，在面試時每人對每種情緒透過麥克風錄製一句測試句，再將這些句子交予三位語言專家，專家根據自然度及可辨識度挑選出男女各 7 位作為錄音語者。語料庫共採用 9 種情緒分別為中性、生氣、害怕、開心、傷心、噁心和無聊，文本設計為五句長句及五句短句，文本內容為日常生活中的句子並且不帶有任何情緒詞彙，如此使語者在錄音時能較自然的表達語句，且在情緒表達過程中不會受到文本的內容而影響。錄製時語者對相同一句分別以七種情緒表達，因此每位語者至少收錄 92 句，整個語料庫至少包含 922 句，最後全部共錄製約 : 22 句，其中有部分句子保留二至三個版本。為了確保語句的品質，請 42 位測試者以隨機輪播的方式將每一句做七選一的人工辨識，並且對句子的自然度以 3 至 322 評分，測驗結束後捨棄辨識率低於 : 2' 和自然度低於 82' 的句子，篩選後七種情緒的句數分別為生氣 349 句、無聊 : 3 句、噁心 68 句、害怕 8; 句、開心 93 句、傷心 84 句和中性 9; 句共 757 句。"

三、台灣情緒語料庫

為了研究跨語言的情緒語音辨識實驗，我們仿照了 GO Q/FD 的設計方式錄製了自己的情緒語料庫，分別用相同的七種情緒表達五句長句以及五句短句，每種語言各錄製 922 句，表一列出我們的十句文本，表二與表三分別為台語與客家語的發音文本。我們選擇台灣常見的三種語言/國語、台語及客家語/作為我們的台灣錄製語言，由於客家語有不同腔調，因此我們挑選語者時選擇最為普遍的四縣腔與海陸腔作為我們的錄製腔調。錄音環境為中山大學電資大樓 H7/239D 實驗室，錄音麥克風使用 CVJ /CV; ; 64，錄音介面卡為 Uqwpf" Drcugt"Z/HR'Uwttqwpf"708，取樣率為 38MJ |，錄音時語者與麥克風間距約 42 公分。檔案命名格式為 9 碼，分別為第一碼語言 KF、第二三碼語者 KF，第四碼情緒 KF，第五六七碼文本 KF，如 J 25Cv23 音檔則為客家語、語者 25、情緒為生氣和 v23 文本。接著分別對三種語言各找 7 位男性以及 7 位女性作為錄音者。相同的我們在錄製完後也請了 32 位測試者做人工辨識，但我們的測驗中並沒有自然度的評分。我們從中保留辨識率 82' 以上的句子，表四列出三種語言挑選後各情緒所剩餘的句數，圖一表示三種情緒各個人工辨識率的句數，由圖中可看出多數錄製的人工準確度在 82' 以上，所以我們的情緒語料庫具有一定的情緒辨別度。所有的錄音者和測試者皆為大學生或是研究生，並且非語言學或者音樂學系學生。"

"

表一、台灣情緒語料庫 10 句文本

| 文本 K" | 文本" |
|-------|--------------------|
| ∅3" | 你的早餐放在桌上" |
| ∅4" | 晚上他有一個約會" |
| ∅5" | 最近常常睡不飽" |
| ∅6" | 等下一一起去餐廳吃飯" |
| ∅7" | 過年要買一件新衣服" |
| ∅8" | 這禮拜放假，要跟朋友一起出去玩" |
| ∅9" | 昨天早上出門的時候，外套被鉤子勾到" |
| ∅: | 今天一整天都沒吃東西好餓喔" |
| ∅; | 這班火車人很多，很多人都沒位子坐" |
| ∅2" | 早上去騎腳踏車，下午散步去買東西" |

表二、台語對應文本

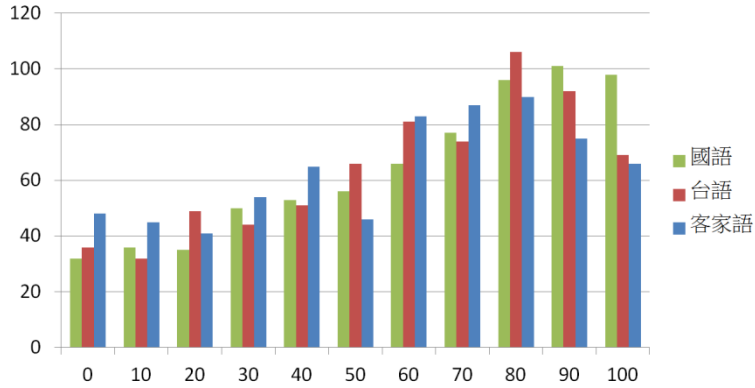
| 文本 K" | 文本" |
|-------|--------------------|
| ∅3" | 你的早頓放在桌頂" |
| ∅4" | 暗時依有一個約會" |
| ∅5" | 最近定定暍不飽" |
| ∅6" | 等咧做伙去餐廳食飯" |
| ∅7" | 過年愛買一件新衫" |
| ∅8" | 這禮拜放假，欲佻朋友做伙出去企桃" |
| ∅9" | 昨透早欲出門的時陣，外套被鉤仔鉤到" |
| ∅: | 今仔日規工攏無食物件就夭" |
| ∅; | 這班火車人足濟，真濟人攏無位湯坐" |
| ∅2" | 透早去騎腳踏車，下晡散步去買物件" |

表三、客家語對應文本

| | |
|-------|--|
| 文本 K" | 文本" |
| ㄩ3" | *ㄅ一ㄩ>+*ㄍㄨㄣˊ九*ㄅㄛㄌㄨㄣˊㄩㄣˊ擗擗轟" |
| ㄩ4" | *ㄉ一補壓*ㄍㄨㄣˊ一>+武依炸右*ㄉ一+" |
| ㄩ5" | 追*ㄅ一ㄨㄣˊㄌㄨㄣˊ擗擗雖某報" |
| ㄩ6" | *ㄉㄍㄨㄣˊ哈*ㄅ一ㄨㄣˊㄌㄨㄣˊ哈*ㄉ一+蠶堂*ㄉ一+翻" |
| ㄩ7" | 勾輾*ㄉ一+埋億良心賞夫" |
| ㄩ8" | *ㄅ一ㄩ>+禮辦*ㄅ一ㄨㄣˊㄌㄨㄣˊ*ㄅ一ㄨㄣˊㄌㄨㄣˊ+; *ㄉ一+老品" *ㄉ一ㄌㄨㄣˊ+*ㄅ一ㄨㄣˊㄌㄨㄣˊ哈粗*ㄉ一+*ㄅ一ㄨㄣˊㄌㄨㄣˊ+" |
| ㄩ9" | 醜比則損*ㄉ一+處*ㄉ一ㄌㄨㄣˊㄌㄨㄣˊ>+*ㄅ一ㄨㄣˊㄌㄨㄣˊ+死*ㄉ一+又+; 歪掏*ㄅ一ㄨㄣˊㄌㄨㄣˊ>+*ㄅ一ㄨㄣˊㄌㄨㄣˊ+又+; *ㄅ一ㄨㄣˊㄌㄨㄣˊ+又+; 豆" |
| ㄩ:" | *ㄍㄨㄣˊㄌㄨㄣˊ>+*ㄅ一ㄨㄣˊㄌㄨㄣˊ+逆*ㄍㄨㄣˊㄌㄨㄣˊ+逆*ㄅ一ㄨㄣˊㄌㄨㄣˊ+朱某奢懂席喊 窠喔" |
| ㄩ;" | *ㄅ一ㄩ>+*ㄅ一ㄨㄣˊㄌㄨㄣˊ+*ㄉ一+又+查擗懂*ㄅ一ㄨㄣˊㄌㄨㄣˊ+又+; " 懂*ㄅ一ㄨㄣˊㄌㄨㄣˊ+擗朱某*ㄅ一ㄨㄣˊㄌㄨㄣˊ+*ㄅ一ㄨㄣˊㄌㄨㄣˊ+後仇" |
| ㄩ32" | 則損*ㄉ一+*ㄅ一ㄨㄣˊ>+師朗查; *ㄉ一ㄌㄨㄣˊ>+租三*ㄅ一ㄨㄣˊㄌㄨㄣˊ+ *ㄉ一+擗懂席" |

表四、台灣情緒語料庫各語言篩選後句數

| 情緒" | 國語" | 台語" | 客家語" | 跨語言總數" |
|-----|-------|------|------|--------|
| 生氣" | 7: " | 95" | 84" | 3; 5" |
| 無聊" | 88" | 77" | 7; " | 3: 2" |
| 噁心" | 68" | 77" | 76" | 377" |
| 害怕" | 79" | 83" | 86" | 3: 4" |
| 開心" | 89" | 7: " | 6: " | 395" |
| 傷心" | 79" | 78" | 74" | 387" |
| 中性" | : 9" | 86" | 84" | 435" |
| 總數" | 65: " | 644" | 623" | 3483" |



圖一、各人工辨識率的句數

四、基準實驗

4.1 實驗設定

我們實驗中採用 $qr\ gpGCT"2\ \emptyset]32$ 與 $Y\ gnc"5\ \emptyset\]33$ 兩套工具用於特徵參數擷取和訓練分類器。輸入的音訊處理採用 47 毫秒的漢明窗 $*J\ co\ o\ kpi\ "y\ kpf\ qy\ +$ 並且處理過程中每次偏移 32 毫秒。在訓練分類器之前使用極值正規法 $*O\ kp/O\ cz"pqto\ crk\ cvkqp+$ 將特徵參數縮放為 2 至 3，公式定義為算式 $*3+$ ， X 與 $X\ \emptyset$ 分別為正規化前與正規化後的值。後端以支持向量機做為分類器，支持向量機的核心函式 $*ngt\ pgn'hwpe\ vkp+$ 為基於序列最小化演算法 $*Ugs\ wgp\ vcn'O\ kpo\ cn'Qr\ vo\ k\ cvkqp+$ 的多項式核心 $*Rqn\ pqo\ kn'Mgt\ pgn+$ ， $Mgt\ pgn$ 函式為算式 $*4+$ ， r 設為 3。實驗方式以每次從語料庫中挑出一位語者作為測試語者，而其他語者當作訓練語者 $*Ngcxg/Qpg/Ur\ gcngt/Qw.'NQUQ+$ ，經過 32 次訓練與辨識後統計結果便為此語料庫的辨識率。"

"

$$V' = \frac{Max_V - V}{Max_V - Min_V} \dots \dots \dots *3+$$

$$K(x, y) = \langle x, y \rangle^p \dots \dots \dots *4+$$

4.2 基準特徵集

本論文使用聲學特徵作為實驗的特徵參數，表五和表六分別列出大型特徵集所包含的低階參數 $*ny / r\ xgn'f\ guet\ lr\ vqt.'NNF+$ 以及泛函 $*hwpe\ vkp\ cn+$ 。首先從音訊中擷取每個音框為 78 維的特徵，再透過泛函將這些音框特徵轉換為一個句子一組的特徵向量，78 個低階參數和 5 個泛函再計算一階與二階動態特徵後共得到 8774 個特徵參數，此便為我們的基準特徵集。"

表五、56 個低階參數

| | | |
|-------------------|---|----------|
| Hgcwtg'I tqwr " | Hgcwtg'kp'I tqwr " | %qh'NNE' |
| Tcy 'uki pcn' | \ gtq/etqukpi /tcvg" | 3" |
| Uki pcn'gpgti { " | Nqi ctkj o " | 3" |
| Rkej " | Hwpc co gpvci'ht gs wgspe { 'H2'kp'J 'xlc'Egr ut wo "cpf " C wqeqttrgvkqp"*CEH0Gzr qp'gpcm{ 'uo qqyj gf "H2" gpxgqr gO' | 4" |
| Xqleg'S wcrkv{ " | Rtqdcdrkv{ "qh'xqlekp'i " | 3" |
| Ur gevten' | Gpgti { 'kp'dcpf u'2/472'J . '2/872'J . '472/872'J . '3/6" m' '47' . '72' . '97' . ; 2' 'tqm'qh'r qkpv.'egpvtqf . 'hwz." cpf 'tgr0r qu0qh'ur gev wo 'o cz0'cpf 'o kp0("34" O gr'ur gev wo "('Dcpf '3/48'('48'Egr utcrn('O HEE '2/34" ("35" | 34" |
| O gr'ur gev wo " | Dcpf "3/48" | 48" |
| Egr utcrn' | O HEE "2/34" | 35" |

表六、39 個泛函

| | |
|---|------------|
| Hwpcvkpcnu. "gve0 | %qh'hwpc0' |
| Tgur gev'xg'tgr0r quk'kqp'qh'o cz0lo kp0'xcnwg" | 4" |
| Tcpi g"*o cz0'o kp0" | 3" |
| O cz0'cpf "o kp0'xcnwg"/"ctkj o gvle"o gcp" | 4" |
| Ctkj o gvle"o gcp. 'S wcf tvle"o gcp" | 4" |
| P wo dgt'qh'pqp/ gtq'xcnwgu" | 3" |
| I gqo gvle. 'cpf "s wcf tvle"o gcp'qh'pqp/ gtq'xcnwgu" | 4" |
| O gcp'qh'cduq'nwg'xcnwgu. 'O gcp'qh'pqp/ gtq'cdu0'xcnwgu" | 4" |
| S wct'vkgu'cpf "kpvt/s wct'vkg'tcpi gu" | 8" |
| ; 7' "cpf "; : ' "r gtegv'kqg" | 4" |
| Uf 0f gxk'v'kqp. 'xctk'peg. 'mxt'v'uku. 'ungy'pguu" | 6" |
| Eg'pvtqf " | 3" |
| \ gtq/etqukpi 'tcvg" | 3" |
| %qh'r gcmu. "o gcp'f'k'v0'dy p0'r gcmu. "ct'v'0'o gcp'qh'r gcmu. "ct'v'0'o gcp" qh'r gcmu"/"qxg'cm'ct'v'0'o gcp" | 6" |
| Nk'pgct't'gi t'gu'kqp'eq'gh'le'k'p'u'cpf "eq'tt'gur 0'cr r tqz'ko cv'kqp"gtt'qt" | 6" |
| S wcf tvle't'gi t'gu'kqp'eq'gh'le'k'p'u'cpf "eq'tt'gur 0'cr r tqz'ko cv'kqp"gtt'qt" | 7" |

動態特徵計算為 $f_{gnc} \cdot t_{gi} \cdot t_{guikqp}$ ，係數差公式採用 $J_{kf} \cdot g_{p'} \cdot O_{ctnqx'} \cdot O_{qf} \cdot g_{ri} \cdot V_{qri} \cdot M_{kv}$ 定義如算式 5+，其中 Y 設為 4。

$$\Delta x_t = \frac{\sum_{i=1}^W (x_{t+i} - x_{t-i}) \cdot i}{2 \sum_{i=1}^W i^2} \quad (5+)$$

4.3 實驗

根據表五與表六的 8774 維大型特徵集搭配支持向量機作為國語、台語和客家語三種語言的基準實驗，對每個語言的 32 位語者進行 NQUQ 實驗後分別得到國語 8:07'、台語 720' 和客家語 7:07' 辨識率。表七至表九分別列出國語、台語和客家語三種語言的詳細辨識結果。

從表七至表九中可看出雖然在人工辨識上可得到至少 82' 的辨識率，但在基準實驗中許多情緒的辨識率不及 82'。此實驗結果也反映出錄音狀況，七種情緒中以生氣、無聊和害怕三種情緒的辨識率較佳，此三種情緒對錄音者來說較好表達；而錄音者普遍認為噁心只單純利用語調來表達並不夠直覺，因此在錄音的情緒表達上會相對較差，三種語言的最高辨識率僅 72'；另在錄音過程中傷心和無聊兩種情緒對語者表達來說會有相似的狀況，實驗也可看出傷心的辨識結果有不少部分被分類為無聊，因此造成傷心的辨識率普遍較低。在台語的中性僅 390' 的準確率，推測因各個語者在表達方式上的差異性過大造成分類器無法順利辨識，主要為語速與聲調的差異，而人工辨識時人們可以由語者的特性去推敲出其情緒，但辨識器上卻無法做出此判斷，所以造成辨識率不理想的狀況。

表七、國語基準實驗混淆矩陣

| 結果 \ 答案 | 生氣 | 無聊 | 噁心 | 害怕 | 開心 | 中性 | 傷心 | 準確率 |
|---------|----|----|----|----|----|----|----|------|
| 生氣 | 66 | 3 | 4 | 2 | 2 | 2 | 5 | 970 |
| 無聊 | 2 | 7 | 2 | 2 | 2 | 6 | 5 | 706 |
| 噁心 | 8 | 6 | 43 | 2 | 9 | 6 | 6 | 670 |
| 害怕 | 4 | 2 | 2 | 67 | 8 | 4 | 4 | 900 |
| 開心 | 9 | 2 | 6 | 8 | 62 | 3 | 2 | 709 |
| 傷心 | 3 | 32 | 9 | 4 | 5 | 53 | 5 | 7606 |
| 中性 | 7 | 6 | 6 | 5 | 32 | 3 | 82 | 802 |

表八、台語基準實驗混淆矩陣

| 結果\答案 | 生氣 | 無聊 | 噁心 | 害怕 | 開心 | 中性 | 傷心 | 準確率 |
|-------|----|----|----|----|----|----|----|------|
| 生氣 | 67 | 3 | : | 2 | 36 | 2 | 7 | 8308 |
| 無聊 | 2 | 65 | 2 | 2 | 2 | 9 | 7 | 9:04 |
| 噁心 | 8 | 6 | 49 | 4 | 6 | : | 6 | 6:08 |
| 害怕 | 9 | 3 | 6 | 5; | 8 | 4 | 4 | 850 |
| 開心 | 34 | 7 | 6 | 5 | 46 | 3 | ; | 6306 |
| 傷心 | 2 | 38 | : | 5 | 3 | 47 | 5 | 6608 |
| 中性 | 38 | 35 | 8 | 8 | 32 | 4 | 33 | 3904 |

表九、客家語基準實驗混淆矩陣

| 結果\答案 | 生氣 | 無聊 | 噁心 | 害怕 | 開心 | 中性 | 傷心 | 準確率 |
|-------|----|----|----|----|----|----|----|------|
| 生氣 | 66 | 3 | 4 | 2 | 8 | 3 | 6 | 9408 |
| 無聊 | 2 | 7; | 2 | 2 | 2 | 8 | 3 | :80 |
| 噁心 | 8 | 6 | 43 | 2 | 4 | 5 | 7 | 7202 |
| 害怕 | 4 | 2 | 2 | 67 | 8 | 5 | 8 | 8608 |
| 開心 | 9 | 2 | 6 | 8 | 38 | 4 | : | 5505 |
| 傷心 | 3 | 32 | 9 | 4 | 3 | 37 | 4 | 4:0 |
| 中性 | 7 | 6 | 6 | 5 | 7 | 2 | 5; | 840 |

五、結論

在此篇論文中我們建立了一個台灣語言的情緒語料庫，其中以台灣常見的國語、台語和客家語三種語言錄製而成。我們仿照公開語料庫 GO Q/FD 錄製，共包含了生氣、無聊、噁心、害怕、開心、中性和傷心共七種，每種語言的語者為十位，分別為五位男性以及五位女性，文本共十句，每位語者以七種情緒來表達同一句文本。在錄製後進行人工辨識的檢驗，以人工辨識率 82% 作為篩選基準，以此確保我們語料具有一定的辨別度，篩選後三種語言的所剩句數為國語 65 句、台語 644 句和客家語 623 句。另我們使用了一個龐大的聲學特徵集並搭配支持向量機做為我們的基準實驗，其中分別在三種語言可以得到國語 8:0%、台語 720% 和客家語 7:0% 的辨識率，便為三個語言的基準辨識率。此資料庫及實驗數據可做為未來台灣語音情緒辨識使用及參考。在未來若可以我們希望能夠更提升語料庫的品質，如改善台語的中性辨識率過低的狀況，對於篩選後句數剩餘過少的語者進行重錄或以新的語者資料代替，提升品質的同時也增加語料庫的總句數。

參考文獻

- 13_ " T0Y 0Rlectf ."Chgevxg"eqo r wkpi 0'O K'Rtguu."3; ; 90'
- 14_ " O 0'Gr'C {cf k "O 0'U'Mco gn" cpf "H0'Mctt{ ."δUwtxg{ "qp" ur ggej "go qvqp" tgeqi pkkqp< hgcwt gu."ercuukhcevkqp" uej go gu."cpf "f cvdcugu.δ" kp"Rcwgtp" Tgeqi pkkqp."xqr0'66."r r 0' 79467: 9.'O ct042330'
- 15_ " H0'Dwtj ctf v."C0'Rcgej ng."O 0'Tqrhgu."Y 0'H0'Ugpf m gkgt."cpf "D0'Y gluu."δC" f cvdcug"qh' i gto cp"go qvqpcn'ur ggej .δ"kp"Rtqeggf kpi u"qh'K'P VGTURGGEJ ."r r 0'373963742."42270'
- 16_ " U' Uxgf n" *Automatic classification of emotion related user states in spontaneous children's speech* ORj F "vj guku."Wpkxgtukv{ "qh'Gtrpi gp/P wtgo dgti ."422; 0'
- 17_ " O 0'I tko o ."M0'Mtquej gn"cpf "U0'P ctc{cpcp."δVj g"Xgtc"co "O kwei "i gto cp"cwf kq/xkawn' go qvqpcn' ur ggej " f cvdcug.δ" kp" Rtqeggf kpi u" qh' KGGG" K'vgtpcvqpcn' Eqphgtgpeg" qp" O wko gf k'cpf "Gzr q."r r 0': 876: 8: ."422: 0'
- 18_ " F 0'O eF wih."T0'Mcrkwd{ ."cpf "T0'Rlectf ."δEtqy f uqwtkpi "hcekn'tgur qpugu" vq" qprkpg" xkf gqu.δ"kp"KGGG"Vtcpucevqpu"qp"Chgevxg"Eqo r wkpi ."xqr0'5."r r 0'678668: ."42340'
- 19_ " D0'Uej wrngt."U0'Uxgf n"cpf "C0'Dcvkpgt."δVj g"K'P VGTURGGEJ "422; "go qvqp"ej cmgpi g.δ" kp"Rtqeggf kpi u"qh'K'P VGTURGGEJ ."r r 0'5346537."422; 0'
-]: _ " I 0'Y gluu"cpf "H0'Rtqxquv."δVj g"ghge'v"qh'ercuu'f kwtkdwkqp"qp"ercuukhgt"rgctpkpi <"Cp" go r k'lecni'uwf { .δ"gej 0tgr 0'F gr ctwo gpv'qh'Eqo r wgt "Uelgpeg."Twi gtu"Wpkxgtukv{ ."42230'
-]; _ " D0'Uej wrngt."D0'Xrcugpnq."H0'G{dgp."I 0'Tki qm"cpf "C0'Y gpf go wj ."δCeqwvke"go qvqp" tgeqi pkkqp<"C" dgpej o ctn' eqo r ctkuqp" qh' r gthqto cpegu.δ" kp" Rtqeggf kpi u" qh' KGGG" Y qtmij qr "qp"qh'Cwqo cvk"Ur ggej "Tgeqi pkkqp"Wpf gtucpf kpi "CUTW:"r r 0'7746779." 422; 0'
- 132_ " H0'G{dgp."O 0'Y qmō gt."cpf "D0'Uej wrngt."δQr gpGCT"kpvtqf welpi "vj g"o vplej "qr gp/uqwtg" go qvqp" cpf " chge'v" tgeqi pkkqp" vqmqk.δ" kp" Rtqeggf kpi u" qh' CEKK' K'vgtpcvqpcn' Eqphgtgpeg"qp"Chgevxg"Eqo r wkpi "cpf "K'vgnki gpv'K'vgtcevqpcn'cpf "Y qtmij qr u."422; ." r r 0'3680'
- 133_ " K0'J 0'Y kvgp."G0'Hicpm"N0'Vtki i ."O 0'J cm"I 0'J qmō gu."cpf "U0'U'Evppkpi j co ."δY gnt< Rtcvkecn'o cej kpg'rgctpkpi "vqnu"cpf "gej pks wgu'y kj "lxc'ko r ngo gpvcvqpu.δ'3; ; ; 0'

基於稀疏表示之語者識別

Sparse Representation Based Speaker Identification

王光耀 "Mwcpil / cq'Y cpi "
國立中央大學資訊工程學系
F gr ctwo gpv'qh'Ego r wgt'Uekpeg'cpf 'lphqto cvkq'Gpi kpggtkpi "
P cvkqpcn'Egptcn'Wpkxgtuk{"

王家慶 "Lc/Ej kpi "Y cpi "
國立中央大學資訊工程學系
F gr ctwo gpv'qh'Ego r wgt'Uekpeg'cpf 'lphqto cvkq'Gpi kpggtkpi "
P cvkqpcn'Egptcn'Wpkxgtuk{"
[ley B eukQewQf w0y "](#)

摘要

稀疏表示分類器 $U^T E$ 是一種基於影像稀疏表示的機器學習方法。在影像以及人臉辨識上的研究上，稀疏表示分類器具有非常好的辨識效果以及強健性。有鑑於 $U^T E$ 在影像辨識上的高鑑別能力，近幾年已有許多基於稀疏表示的語者識別方法相繼被提出。本論文提出一套基於稀疏表示的辨識系統，我們提出以機率型主成份分析 $R^T k$ 並加入檢定的方式調整特徵值選取，使語者高斯混合模型 I 中每個高斯的維度可以針對資料的不同作調整。接著，我們在稀疏字典上加強，透過低秩矩陣還原 N 以及核化降維 M 對說話內容 U 以及通道 E 變異補償，使字典增加鑑別性。最後利用稀疏表示分類器進行分類。根據實驗結果顯示，不論是參數的改進、字典的處理，對辨識率都有一定程度的提升。此外，與傳統的 i/x 語者識別系統相比，提出的系統則具有更佳的辨識率表現。

關鍵詞：稀疏表示、語者識別、核化稀疏表示、超級向量

一、簡介

以生物特徵作為辨識基礎的研究已經有長達數十年的歷史，包含人臉、語音和指紋。其中聲音因具有取得容易、非侵入性、運算量少、輸入具有便利性等優點，語者識別一直是近年來熱門的主題。語者識別利用語者語音的特性來識別使用者身份，大致上的研究方向分成語音特徵擷取以及分類演算法兩部分，其中以分類演算法的研究最多，

包含的層面也最廣，包括語者模型的建立、分類機制、因為周邊環境的影響所需的補償辦法等研究。目前來說高斯混和模型 [3] [4] 和支持向量機 [Uw r qt v] Xgevt " O cej kpg." UXO [5] 是經典的分類方法。特徵參數方面，會將一段語音切割成數個音框，音框代表此語者在短時間內的發聲特徵，幾種常見的有梅爾倒頻譜係數 [O gn/uecrg] Hgs wgepe {" Egr utcn' Eqghlekpw." OHEE [7]、感知線性預測係數 [Rgtegr wcn' Nkpgct" Rt gf le vkp 'Eqghlekpw.' RNRE [8] 等。"

" "" 基於 I O O 超級向量 [I O O /Uw r gtxgevt [9] [32] 的辨識方法是一種綜合 I O O 模型表示以及 UXO 分類方法優點的語者識別演算法。I O O 超級向量是一種語者模型的向量表示法，通用背景模型 [Wpkxgtucn' Dceni tqwpf " O qf gn" WDO 經過輸入語音調適後成為 I O O 語者模型，其中各個高斯成份 [Eqo r qpgpvu+] 的期望值被串在一起形成一個超級向量。此向量表示可以用來代表整個語音檔的特徵向量，最後，再以 UXO 作為分類方法來達到辨識語者的目的。然而語者識別系統會因為錄製工具和背景不同，造成通道差異、說話內容差異和環境差異的干擾，導致辨識系統的辨識率低落。擾動屬性投影 [P wkucpeg' Cwtkdwg' Rtqlgevkp. 'PCR+ 是一種有效補償通道干擾的方式 [9] [33]，藉由設計一個最佳化問題使干擾最小，配合著前面提到的方法，更有效地提高辨識系統的強健性。在 PCR 之後陸續有聯合因素分析 [Lkpv' Hcevt " Cpcn {uku". 'IHC [34] [35]、i/xgevt [36] [38] 等技術被用來補償上述提到的差異，透過對語音參數的拆解，將一些語者不相關的資訊刪除，得到去除干擾後最能代表語者的特徵。"

" "" 近年來對於複雜的資訊 [例如訊號、圖像]，往往希望可以用較簡化的方式呈現，特別在訊號處理的部分，分析時經常先將資料轉換至不同的定義域，並且假設其在轉換後，會呈現稀疏分布 [46]。在近期稀疏表示的發展中，UTE 在 [47] [48] 中被提出，UTE 是一個 P qpr ctco gtle 學習方法，不需訓練過程但是需要訓練資料，以及可以直接參考訓練資料對策是資料進行分類的動作。實驗結果顯示出在人臉辨識的應用上，UTE 有著比 MP P [MP gctguv' P gki j dqtu [49] 以及 P gctguv' Uwdur ceg [P U [4] [4] 更好的辨識率。近幾年也有少數研究將 UTE 應用於語者驗證 [38] [3] 的問題上，然而此方面的研究仍屬於剛起步的階段，仍有許多問題值得我們探討。"

" "" 本論文提出一套基於核化稀疏表示的語者識別系統，文章的編排共分成七個部分 < 第一部分為簡介，第二部分為參數擷取，我們以 RREC 建構超級向量 [RREC 超級向量+ 取代 I O O 超級向量] [42] [53]，並且以巴雷特檢定 [Dct vgw' Vguv+] 的方式調整特徵值的選取，第三部分則是描述如何利用特徵參數建立稀疏表示分類器。第四部分描述我們提出的兩種變異補償方法，低秩矩陣還原以及核化稀疏表示分類器 [Mgtpgn' Urctug" Tgr tgugpv vkp 'Ercuukhgt. "MUTE+。第五部分為實驗部分，展示提出之改良方法是否具有其必要性。最後在第六部分則是結論。"

二、參數擷取"

" "" 在本論文中，我們利用機率型主成分分析建構超級向量 [42] [55]，並以提出巴雷特檢定 [Dct vgw' Vguv+] 主成分的個數。傳統上，超級向量是由高斯混和模型建構而成，這裡

加入主成分分析的概念，並希望能以機率分布模型的形式與高斯模型對應，使得資料點由原本高斯混和模型轉成 RREC 混和模型，再透過 $\mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 的轉換，形成新的超級向量，稱作 RREC 超級向量 [42] [43] [55]，其中，主軸的挑選，我們引入巴雷特檢定 [44] [45] 的概念，建立假說，找到臨界的特徵值。而目前語者識別問題中， i/x_{gevt} 是表現最好的參數之一。藉由訓練出總體變異矩陣，將原本的超級向量轉到更低維的空間，使 i/x_{gevt} 更加表現出語者及通道的資訊，因此，我們將 RREC 超級向量轉換到總體變異空間上，希望得到更具鑑別力的 i/x_{gevt} 使辨識效能提升。

4.3 基於機率型主成分分析之因素分析模型

傳統的 I/O 超級向量，並沒有考量到其聲學特徵參數有高度的冗餘性 [43] [53]，因此應該採用更低的子空間來表示。在 [43] [53] 中比較成份分析 [46] 與主成分分析 [47] 的關聯性，觀察成份分析的數學式：

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon} \quad (3)$$

其中 \mathbf{x} 代表高維的資料， \mathbf{z} 代表低維變數，又稱 $\mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 。回想主成分分析的數學式，找出主成分 \mathbf{V} 使資料降維：

$$\mathbf{x} = \boldsymbol{\mu} + \mathbf{V}\mathbf{z} \quad (4)$$

我們可將成份分析視成對資料 \mathbf{x} 做 REC 加上一個噪音項，並且導入機率的觀念解釋。在式 (3) 中， \mathbf{x} 是 $K \times 3$ 的原始資料， \mathbf{W} 為 $K \times J$ 的轉換矩陣，其中 $J < K$ ，而 $\mathbf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 假設為一高斯分布 $N(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ ，噪音 $\boldsymbol{\varepsilon} \in N(0, \sigma^2 \mathbf{I})$ ，根據以上的假設可以知道原始資料 \mathbf{x} 能建模成 $N(\boldsymbol{\mu}, \sigma^2 \mathbf{I} + \mathbf{W}\mathbf{W}^T)$ ，稱作 RREC 模型。當遇到更複雜的資料時，希望藉由混合多組 RREC 模型，如下式：

$$p(\mathbf{x}) = \sum_{c=1}^M w_c p(\mathbf{x} | c) \quad (5)$$

$$p(\mathbf{x} | c) = N(\boldsymbol{\mu}_c, \sigma_c^2 \mathbf{I} + \mathbf{W}_c \mathbf{W}_c^T) \quad (6)$$

其中 M 是高斯成份的個數，這可以跟高斯混合模型做對應，一般表示資料以多個高斯模型描述，這裡的基本單位以 RREC 模型取代。

因此，我們希望找出 \mathbf{W}_c 及 σ_c 來推算 $\mathbf{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$ ，方法是利用最大概似估計 [48] [49]，因為已經知道 \mathbf{x} 的分布，所以藉由 O/N/G 得到的 \mathbf{W}_c 及 σ_c 可進而推出 $\mathbf{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$ \mathbf{z}_c

$$\mathbf{z}_c = (\mathbf{W}_c^T \mathbf{W}_c + \sigma_c^2 \mathbf{I})^{-1} \mathbf{W}_c^T (\mathbf{x} - \boldsymbol{\mu}_c) \quad (7)$$

當原始資料 \mathbf{x} 進入系統後，會先以 RREC 混合模型，並且透過式 (7) 的轉換得到屬於

每個高斯成份的 $N(\boldsymbol{\mu}_{z,c}, \boldsymbol{\Sigma}_{z,c})$ 。當所有的參數都做了 RREC 後，則原始以 OHEE 為基礎的 WDO 也必須調整，形成以 $N(\boldsymbol{\mu}_{z,c}, \boldsymbol{\Sigma}_{z,c})$ 為參數下構建的新 WDO：

$$p(\mathbf{z}) = \prod_{c=3}^M w_c N(\boldsymbol{\mu}_{z,c}, \boldsymbol{\Sigma}_{z,c}) = \prod_{c=3}^M w_c N(\mathbf{z}; \boldsymbol{\mu}_{z,c}, \boldsymbol{\Sigma}_{z,c}) \quad (8)$$

其中 $\boldsymbol{\mu}_{z,c} = E\{\mathbf{W}_c^T \mathbf{W}_c + \sigma_c^4 \mathbf{I} + \mathbf{W}_c^T \mathbf{x} - \boldsymbol{\mu}_c\}$

$$\boldsymbol{\Sigma}_{z,c} = E\{\mathbf{z}_c \mathbf{z}_c^T - \boldsymbol{\mu}_{z,c} \boldsymbol{\mu}_{z,c}^T$$

$$= \mathbf{W}_c^T \mathbf{W}_c + \sigma_c^4 \mathbf{I} + \mathbf{W}_c^T E\{\mathbf{x} - \boldsymbol{\mu}_c\} + (\mathbf{x} - \boldsymbol{\mu}_c)^T \mathbf{W}_c + \mathbf{W}_c^T \mathbf{W}_c + \sigma_c^4 \mathbf{I} \quad (9)$$

$$= \boldsymbol{\Lambda}_c^{-3} \boldsymbol{\Lambda}_c - \sigma_c^4 \mathbf{I} + \boldsymbol{\Lambda}_c \boldsymbol{\Lambda}_c - \sigma_c^4 \mathbf{I} = \mathbf{I} - \sigma_c^4 \boldsymbol{\Lambda}_c^{-3}$$

4.4 巴雷特檢定

在上個小節中提到的 RREC 中，有一個問題是需要討論的，那就是主軸個數的挑選。我們假設後面不要的特徵值，其數值小到足以視為相同，因此，定義假說 $H_2: \lambda_{k+3} = \lambda_{k+4} = \dots = \lambda_n$ ，希望迴圈由 $n/3$ 往前找到門檻值 k 。而如果我們將特徵值的大小視為高斯分布，則巴雷特檢定 [44, 45] 可以整理成相似度的檢驗，如下式：

$$T = \frac{\prod_{q=k+3}^n \lambda_q}{\left(\frac{\sum_{q=k+3}^n \lambda_q}{n-k} \right)^{n-k}} \quad (10)$$

當觀測資料數 m 夠多的話，可將 T 視為一個 χ^2 分布，且將上式近似於下式：

$$T \approx (m - n + 3) \left(n - k + \sum_{q=k+3}^n \lambda_q \bar{\lambda} - \sum_{q=k+3}^n \lambda_q \right) \quad (11)$$

其中 $\bar{\lambda}$ 是後 n/k 個特徵值的平均，在檢定過程中，當 $T > \chi^2_{\alpha, n}$ 時，則否決假說 H_2 ，檢定終止， q 即為我們挑選的主軸個數，其中 α 為 χ^2 的顯著水平。

4.5 $i/xgevt$

啟發於早期 IHC 在語者識別上的應用，Fujimori et al. 提出的一個新的分析方法 $i/xgevt$ [36/38]，不像 IHC 將語者和通道分開， $i/xgevt$ 僅用一個總體變異性空間，他發現 IHC 中的通道部分仍包含了能用來識別語者的資訊，所以將 IHC 中分開的變異部分，合併成一個單一的總體變異性空間超級向量，藉由合併錄音方式的變異性，提升其辨識性，表示如下：

$$\mu = m + Tw \quad (33)$$

m 是 WDO 超級向量，與 IHC 中使用的相同； T 代表的是所有變異性的矩陣； $i/xgevt$ 代表的是總體變異性元素 w 。

最後，總結參數擷取的整體架構，參數擷取包含三個部分，第一，由 Wpikgtucrl Dcemi tqwpf Fcvc 首先訓練出 WDO，並透過 RREC 將 WDO 轉換成以 NcvgpvHcevt 為參數的 WDO，第二，當輸入語音進入系統後，擷取其 NcvgpvHcevt，接著，對新的 WDO 調適產生 RREC 超級向量。最後，在基於 RREC 超級向量參數下訓練總變異矩陣，並將輸入語音轉換成 $i/xgevt$ 。

三、稀疏表示分類器

我們利用稀疏表示具鑑別力的特性，應用於語者識別問題上，利用 $i/xgevt$ 作為特徵參數，以固定維度的向量表示語音訊號。假設有 C 個不同語者類別的訓練資料，每一筆訓練資料為一個 $i/xgevt$ ，我們需要建構一個跨類別的字典 $D \in \mathbb{R}^{P \times Q}$ ，作法是將所有語者類別的字典組在一起

$$D = [D_1 \ D_2 \ \dots \ D_C] \quad (34)$$

其中， D_j 表示第 j 個類別的字典，由類別 j 的 $i/xgevt$ 串接而成。此外， P 代表參數維度， Q 是訓練資料個數。在測試資料 y 與字典 D 已知的情況下，我們希望求出稀疏係數 $x = [x_1 \ x_2 \ \dots \ x_k]$ ，使原訊號與重建訊號之間的誤差能越小越好，且 x 要符合稀疏特性，如下式：

$$\min_x \|y - Dx\|_4 + \lambda \|x\|_1 \quad (35)$$

在解出稀疏係數 x 後，決策上，將 k 類別各自的字典 D_j 及係數 x_j 還原 y ，並與 y 計算誤差，獲得最小誤差者即為所屬類別：

$$j^* = \arg \min_j \|y - D_j x_j\|_4 \quad (36)$$

四、字典處理及變異補償

在在以 i/x_{gevqt} 為參數的架構下， i/x_{gevqt} 的總變異矩陣是假設訓練資料標籤不同所訓練而成，因此，它存在與語者相關的變異，這是我們想要的，同時它存在一部分錄音通道及說話內容造成的變異，因此，以往的研究都會加入 $N_{kpgct}'F_{kuetko} \ kpcpv' Cpcn\{uku'NF C+$ 及 $Y_{kij} \ kp/Ercuu'Eqxctkpeg'P_{qto} \ crk \ c'kqp''Y_{EEP}$ 去補償錄音通道的變異[55]，希望即使是在錄音通道不同的狀態下，找出類別間的變異最大化，類別內的變異最小化的空間，使變異排除，另外說話內容的變異也可以由 NFC 的處理解釋，因為它將語者的說話內容做變異最小的假設，消除說話內容變異造成的問題。

本節提出對 i/x_{gevqt} 字典的處理及補償辦法，包括以低秩矩陣還原以及核化稀疏表示分類器 $M_{gtpgn}'UTE$ 分別對說話內容及錄音通道變異做補償，並增加字典的鑑別性。

"

6B 低秩矩陣還原字典處理及變異補償

低秩矩陣還原[56]提供一種訊號拆解的方式，假設訊號可以被拆解成一個低秩 N_{qy} T_{cpm} 矩陣 A 及稀疏誤差 $U_{rctug}'G_{ttqt}'+E$ 的和，如下式：

$$D = A + E \quad \text{*****} \quad *37+$$

在 A 與 E 皆未知的情況下，我們希望讓 A 能以最少的 $tcpm$ 還原原訊號，並且 E 能符合 U_{rctug} 的特性，因此整理成以下的最佳化式：

$$\min_{A,E} \text{rank}(A) + \gamma \|E\|_2 \quad \text{*****} \quad *38+$$

進而，整理成一個凸最佳化的問題，並用 $C_{wi} \ o \ gpvgf'N_{ci} \ t_{cpi} \ g'O_{wnk} \ r_{kgt}'CNO$ 求解。

透過低秩矩陣還原，我們將每個語者類別的字典分別作低秩矩陣還原分解，得到與原來字典大小相同的低秩矩陣 A ，取代原來的字典，讓字典中每個語者的特徵更為凸顯。

"

6C 核化稀疏表示分類器

藉由核化方法 $M_{gtpgn}'V_{tlem}$ ，稀疏表示分類器可進一步非線性化，稱作核化稀疏表示分類器[54]。輸入空間的資料經過非線性核化映射 $M_{gtpgn}'O_{cr} \ r_{kpi}$ 投射至高維參數空間，讓原本在輸入空間混淆的參數資料在高維空間變成可分離的。透過更進一步的核化降維方法，我們可得到測試資料的稀疏組合係數 $U_{rctug}'E_{qo} \ dlpc'kqp'E_{qgthkclgpv}$ 來進行分類的動作。

我們利用稀疏表示具鑑別力的特性，應用於語者識別問題上，利用 i/x_{gevqt} 作為特

"

徵參數，以固定維度的向量表示語音訊號。假設有 c 個不同語者類別的訓練資料，每一筆訓練資料為一個 $\mathbf{x}_i \in \mathbb{R}^D$ ，我們需要建構一個跨類別的字典 \mathbf{A} ，作法是將所有語者類別的字典組在一起 $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_c]$ ，其中 $\mathbf{a}_i \in \mathbb{R}^D$ 。一筆測試資料 $\mathbf{x} \in \mathbb{R}^D$ 令 Φ 表示 $\mathbb{R}^D \rightarrow \mathbb{R}^F$ 的非線性映射，可將輸入空間 \mathbb{R}^D 的資料投射到高維空間 F 。

$$\Phi: \mathbf{a} \in \mathbb{R}^D \rightarrow \mathbf{a}' \in \mathbb{R}^F \quad (39)$$

與稀疏表示分類器相同，參數空間的測試資料可表示為訓練資料的線性組合：

$$\mathbf{x}' = \sum_{i=1}^c x_i \Phi(\mathbf{a}_i) = \Phi_{\mathbf{A}} \mathbf{x} \quad (40)$$

其中 $\Phi_{\mathbf{A}} = [\Phi(\mathbf{a}_1), \dots, \Phi(\mathbf{a}_c)]$ ，且 $\mathbf{x}' = [x_1, x_2, \dots, x_c]^T$ 。

根據稀疏表示的概念，稀疏係數向量 \mathbf{x} 可以從下式最佳化問題解出：

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\text{arg min}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathbf{x}' - \Phi_{\mathbf{A}} \mathbf{x}\|_2 \leq \epsilon \quad (41)$$

當核化空間的維度是未知且遠高於訓練資料的個數時，會導致式子 (41) 的解會不夠稀疏，因此會需要在參數空間進行降低維度的動作。令 \mathbf{P} 表示投射矩陣 $\mathbf{P} \in \mathbb{R}^{F \times F}$ 。則 (41) 可改成下式：

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\text{arg min}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathbf{P}(\mathbf{x}' - \Phi_{\mathbf{A}} \mathbf{x})\|_2 \leq \epsilon \quad (42)$$

具有最小重建誤差 $\|\mathbf{P}(\mathbf{x}' - \Phi_{\mathbf{A}} \mathbf{x})\|_2$ 的類別則為核化稀疏表示分類器的分類結果：

$$\hat{i} = \underset{i}{\text{arg min}} \|\mathbf{P}(\mathbf{x}' - \Phi_{\mathbf{A}} \delta \mathbf{x}_i)\|_2 \quad (43)$$

其中 $\delta \mathbf{x}_i$ 是 $\hat{\mathbf{x}}$ 的第 i 個語者的非零稀疏係數。

MUTE 即結合 MFC 及 UTE 兩個方法，將原始特徵向量由 MFC 投影至高維後，再由 MFC 降維，在先前的文獻有提到， \mathbf{x} 存在變異，因此，我們透過 MFC 補償，希望錄製方式不同產生的干擾能藉由最大化變異及最小化變異得到補償。

五、實驗結果

我們以 P KUV4227 做為語者資料庫，P KUV 每年都會錄製語者資料庫，許多產業界、學術界都會以此作為評估效能的標準，而 4227 年發布的資料庫以電話對話語音數據為主，同時收集一些輔助麥克風接收的數據，這些數據主要來自英語演講，並包括四種額外語言。我們挑選其中 0.3eqp 的 Eqpf kkkp 做為測試資料庫，其中 :eqp 為訓練資料，而 3eqp 是測試資料。

為了驗證巴雷特檢定是否能進一步改善辨識系統，我們定義另一套主軸個數選取方式（各個高斯成份在統一主軸個數下與透過檢定方式挑選適合主軸個數下的差別，最後

再經由轉換到 i -xgevqt 上。藉由比較巴雷特檢定與我們定義的主軸選取方式來判斷巴雷特檢定的必要性。在我們的實驗中，如表一所示，固定每個高斯成份的主軸個數為 49 時來到最高 98Q 3'，不過隨著主軸個數的減少辨識率也降低。而加入巴雷特檢定後，我們設定巴雷特檢定的 $\alpha = 2Q27$ 。各個高斯成份降維的數目不固定，有效的提升辨識率到 99Q23'。在接下來的實驗中，我們會以這個辨識率為 99Q23' 的系統作為後續實驗的基準 *Dcugrlpg+，測試提出採用的改進方法是否有效。"

《表一》不同主軸個數與巴雷特檢定選取之辨識率比較。"

| 方法" | | 辨識率" |
|---|---------|----------|
| RREC/UX" | 52個主軸" | 98Q2' " |
| | 49個主軸" | 98Q 3' " |
| | 46個主軸" | 98Q62' " |
| | 43個主軸" | 98Q42' " |
| | 3: 個主軸" | 97Q ;' " |
| | 37個主軸" | 94Q45' " |
| RREC/UX"- "Dctvrgw"Vguv"- "i/xgevqt"- "UTE" | | 99Q23' " |

在以下實驗中，我們比較了四種不同語者識別系統。Dcugrlpg 我們使用目前 ucvg/ql/vj g/ctv 辨識方法，i/xgevqt "dcugf "equlpq" f kucpeg，即以 i -xgevqt 為參數下，將輸入特徵與每個 Encuu 字典的：筆 i -xgevqt 做內積求平均，挑選內積最大的語者類別，表示相似度為最大，通常在參數後面會加上 NFC 與 Y EEP 對錄音通道變異補償。其餘三個系統分別為基於稀疏表示器的系統、基於核化稀疏表示器作字典補償之識別系統以及基於低秩矩陣還原之識別系統。實驗數據如表二所示，透過低秩矩陣還原方法建構的字典具有最高的辨識率：2Q79'，比 Dcugrlpg 以及沒有變異補償的系統分別多了 35Q85' 以及 5Q78' 的辨識率。與沒有變異補償的系統相比，Mgtprn'NFC 達到變異性補償的效果，增加了字典的鑑別性，辨識率提升了 4Q46'。"

《表二》不同語者識別系統之辨識率比較。"

| 方法" | 辨識率" |
|--|-----------|
| I O O /UX"- "i/xgevqt"- "NFC"- "Y EEP"- "EF "*Dcugrlpg+" | 88Q 6' " |
| RREC/UX"- "Dctvrgw"Vguv"- "i/xgevqt"- "UTE" | 99Q23' " |
| RREC/UX"- "Dctvrgw"Vguv"- "i/xgevqt"- "Mgtprn'UTE" | 9; Q47' " |
| RREC/UX"- "Dctvrgw"Vguv"- "i/xgevqt"- "UTE "*Nqy /Tcprn' O cvtkz "Tgeqxgt { "Dcugf "F kvlpqct { +" | : 2Q79' " |

六、結論

這篇論文提出一套基於稀疏表示分類器為基礎的辨識系統，在前端以 RREC/Uwr gtxgevqt 為參數，加入巴雷特檢定作為準則，決定每個高斯 Eqo r qpgpv 挑選主軸的辦法，使每個高斯 Eqo r qpgpv 的維度可以針對資料的不同，決定適當的維度，接著，訓練出總變異矩陣，將 Uwr gtxgevqt 投映至總變異空間上，以 i-xgevqt 作為辨識參數。在字典的建構上，我們提出以低秩矩陣還原以及 Mgtpgn'UTE 進行變異補償，去除說話內容以及通道變異造成的干擾。從實驗結果看來，RREC/Uwr gtxgevqt 在未做巴雷特檢定前就有比 I O O/Uwr gtxgevqt 好的效果，而再加入巴雷特檢定後，效果更加提高，與傳統基於 i-xgevqt 的識別系統相比，辨識率提升了 32.09%。此外，再加入兩種變異補償方法後，低秩矩陣還原以及 Mgtpgn'UTE 分別可以再提升系統的辨識率 5.78% 以及 4.46%。

參考文獻

- [3] F O C O T g { p q r f u " c p f " T O E O T q u g . " o T q d w u v " v z v l p f g r g p f g p v " u r g c n g t " k f g p w h l e c v k p " w u l p i " I c w u l c p " o k z w t g " o q f g n . o " IEEE Trans. Speech Audio Process. " x q r 0 5 . " p q 0 3 . " r r 0 9 4 6 : 5 . " l c p 0 3 ; ; 7 0
- [4] D O N O R g m q o " c p f " L O J O N O J c p u g p . " o C p " g h l e k p v " u e q t k p i " c n i q t k j o " h q t " I c w u l c p " o k z w t g " o q f g n " d c u g f " u r g c n g t " k f g p w h l e c v k p . o " IEEE Signal Process. Lett. " x q r 0 7 . " p q 0 3 3 . " r r 0 4 : 3 6 4 : 6 . " P q x 0 3 ; ; : 0
- [5] Y O O O E c o r d g m " L O R O E c o r d g m " F O C O T g { p q r f u . " G O U k p i g t . " c p f " R O C O V q t t g u / E c t t c u s w k m q . " o U w r r q t v " x g e v q t " o c e j k p g u " h q t " u r g c n g t " c p f " n p i w c i g " t g e q i p k k p . o " Comput. Speech Lang. " x q r 0 4 2 . " r r 0 4 3 2 6 4 4 ; . " 4 2 2 8 0
- [6] L O N O I c w x c k p " c p f " E O J O N g g . " o C z k o w o " a p o s t e r i o r i g u w k o c v k p " h q t " o w n k x c t k c v g " I c w u l c p " o k z w t g " q d u g t x c v k p u " q h " O c t n q x " e j c k p u . o " IEEE Trans. Speech Audio Process. " x q r 0 4 . " p q 0 4 . " r r 0 4 : 3 6 4 ; . : 3 ; ; 6 0
- [7] O O T O J c u c p . " O O L c o k n " O O I O T c d d c p k " c p f " O O U O T c j o c p . " o U r g c n g t " k f g p w h l e c v k p " w u l p i " O g r i t g s w p e { " e g r u t c n i e q g h l e k p v . o " 3 r d i n t e r n a t i o n a l C o n f e r e n c e o n E l e c t r i c a l & C o m p u t e r E n g i n e e r i n g I C E C E 2 0 0 4 . " 4 : / 5 2 " F g e g o d g t " 4 2 2 6 . " F j c n e . " D c p i n f g u j 0
- [8] J O J g t o c p u m { o o R g t e g r w c n i " N l p g c t " R t g l e v x g * " R N R + " C p c n { u k i " q h " U r g g e j . o " Journal of the Acoust. Society of Amer. , : 9 < 3 9 5 : / " 3 9 7 4 . " C r t k n " 3 ; ; 2 0
- [9] Y O E c o r d g m " F O U w t k o . " F O T g { p q r f u . " c p f " C O U q n q o q p q h h " o U X O " d c u g f " u r g c n g t " x g t h l e c v k p " w u l p i " c " I O O " u w r g t x g e v q t " n g t p g n " c p f " p c r " x c t k c d k k v { " e q o r g p u c v k p . o " k p " Proc. ICASSP . " V q w r q w u g . " H i c p e g . " 4 2 2 8 . " r r 0 ; 9 6 3 2 2 0
- [] _ Y O E c o r d g m " F O U w t k o . " c p f " F O T g { p q r f u . " o U w r r q t v " x g e v q t " o c e j k p g u " w u l p i " I O O " u w r g t x g e v q t u " h q t " u r g c n g t " x g t h l e c v k p . o " IEEE Signal Process. Lett. " x q r 0 3 5 . " p q 0 7 . " r r 0 5 2 : 6 5 3 3 . " O c { " 4 2 2 8 0
- [] _ V O M l p p w p g p " c p f " J O N k " o C p " q x g t x l g y " q h " v z v l p f g r g p f g p v " u r g c n g t " t g e q i p k k p < " H i q o " h g c w t g u " v q " u w r g t x g e v q t u . o " Speech Commun 0 " x q r 0 7 4 . " p q 0 3 . " r r 0 3 4 6 6 2 . " 4 2 3 2 0
- [32] D O I O D O H c w x g . " F O O c v t q w h " P O U e j g h t g t . " L O H O D q p c u t g . " c p f " L O U O F O O c u q p . " o U c y g / q h / v j g / c t v " r g t h q t o c p e g " k p " v z v l p f g r g p f g p v " u r g c n g t " x g t h l e c v k p " j t q w i j " q r g p / u q w t e g " u q h y c t g . o " IEEE Trans. Audio, Speech, Lang. Process. " x q r 0 3 7 . " p q 0 9 . " r r 0 3 ; 8 2 6 3 ; 8 : . " 4 2 2 9 0

- 133_ C0Uqno qpqhh "Y 0O 0Eco r dgm"cpf "E0S wknpp."öEj cpggnleqo r gpucvqp"ht "UXO "ur gcngt "tgeqi plkqp.ö
in Proc. Odyssey04."4226."r r 0796840
- 134_ Q0I ngo dgm"N0Dwti gv."P 0Dtwo o gt."cpf "R0Mgpp{."öEqo r ctluqp"qh'ueqtkpi "o gj qf u"wgf "lp"ur gcngt "
tgeqi plkqp"y kj "lqpv'hcevt"cpn{uku.ö"lp"Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.. "Vckr gk"
Vcky cp."Cr t0422; .r r 06279662820'
- 135_ R0Mgpp{."R0Qwngv."P 0F gj cm"X0I w vc."cpf "R0F wo qwej gn"öC"uwf {"qh'lpvgtur gcngt"xctkcdkkl{ "lp"
ur gcngt"xgthhecvtqp.ö"IEEE Trans. Audio, Speech, Lang. Process.. "xqr038."pq07."r r 0; : 26; : . : "Lwr0422: 0
- 136_ C0Mpcpi cuwpcctco . "T0Xqi v"FOF gcp."U0Ukf j ctcp."cpf "O 0O cuqp."ök/xgevt"dcugf "Ur gcngt "Tgeqi plkqp"
qp"Uj qt v"Wwgtcpegu.ö"lp"Interspeech."42330'
- 137_ R0O cvglnc."Q0I ngo dgm"H0Ecuwrf q."Q0Rrej qv."R0Mgpp{."N0Dwti gv."cpf "L0Egtpqen{."öHwm/eqxctkpeg"
wdo "cpf "j gcx{/vckgf "r f c"lp'k/xgevt"ur gcngt"xgthhecvtqp.ö"Proc. ICASSP '11."r r 06: 4: 66: 53."42330'
- 138_ I0O 0M0Mwv."I0'Gr r u."G0'Co dknktclcj . "ök/xgevt"y kj "ur ctug"tgr tguvpcvqp"enrukhhecvtqp"ht"ur gcngt "
xgthhecvtqp."ö"Speech Commun."42350'
- 139_ M0'J wpi " cpf " U0' Cxkl gpv." öUr ctug" Tgr tguvpcvqp" ht" Uki pcr' Enrukhhecvtqp.ö" Neural Information
Processing Systems."42280'
- 13: _ I0O 0M0Mwv."G0'Co dknktclcj . "I0'Gr r u."cpf "T0Vqi pgtk"öUr gcngt"xgthhecvtqp"wupi "ur ctug"tgr tguvpcvqp"
enrukhhecvtqp.ö"lp"Proc. ICASSP."Oc{"4233."r r 0676: 667730'
- 13: _ T0Ucglf k"C0J wto cncpqp."V0Xktvcpqp."cpf "F 0C0'xcp"Nggwy gp."öGz go r rct/dcugf "Ur ctug"tgr tguvpcvqp"
cpf "Ur ctug"F kuetlo lpcvqp"ht"P qlug"tqdwu"Ur gcngt "K gpvkhhecvtqp.ö"lp"Odyssey speaker and language
recognition workshop."Ukpi cr qtg."42340'
- 142_ V0'J cuqp"cpf "I0'J 0'N0'J cpugp."öHcevt"cpn{uku"qh"ceqwute"hgcvwtgu" wupi "c"o lzwg"qh'r tqdcdkklwte"
r tlpekr cnleqo r qpgpv'cpn{| gt u'ht"tqdwu"ur gcngt"xgthhecvtqp.ö"in Proc. Odyssey."Ukpi cr qtg."Lwp042340'
- 143_ O 0' Vkr lpi " cpf " E0' Dkuj qr . " öO lzwg" qh" r tqdcdkklwte" r tlpekr cnl' eqo r qpgpv' cpn{| gt u.ö" Neural
Computation."xqr033."pq04."r r 066566: 4."3; ; ; 0'
- 144_ O 0' [0'Ngg."Rgtegr wcn'Hcevt"Cpcn{uku"ht"Ur ggej "Gpj cpego gpv."O cuvt "Vj guku."Kpukwg"qh'Ego r wgt "
Uelgpeg"cpf "Kphqto cvqp"Gpi lpggtkpi . "P cvkpcnEj gpi "Mwpi "Wplxgtuk{"*4226-0'
- 145_ O U0'Dctvngw."öVguu"qh'uki plcepeg"lp"hcevt"cpn{uku.ö"British Journal of Psychology."Ucvkulecn"Ugevtqp"
5.'996: 7.'3; 72"
- 146_ O 0'Cj ctqp."O 0'Grf . "cpf "C0'Dtvenunglp."öM/UXF <"cp"eni qt kj o "ht"f guki plpi "qxgteqo r ngv"flvqpctkgu"
ht"Ur ctug"tgr tguvpcvqp.ö"IEEE Trans. Signal Process.. "xqr076."pq033."r r 0653366544."P qx042280'
- 147_ C0' [0' [cpi . " I0' Y tki j v" [0' O c . " cpf " U0' U0' Ucut { . " öHgcwtg" ugrgevtqp" lp" hceg" tgeqi plkqp<" C" ur ctug"
tgr tguvpcvqp" r gtur gevkg.ö" GGEU" F gr v0" Wplx0' Ecrkhtplc." Dgtngg{." EC." Vgej 0' Tgr qt v"
WEDIGGEU/4229/; ; ."42290'
- 148_ I0' Y tki j v." C0' [0' [cpi . " C0' I cpugj . " U0' U0' Ucut { . " cpf " [0' O c . " öTqdwu" hceg" tgeqi plkqp" xlc" ur ctug"
tgr tguvpcvqp.ö"IEEE Transaction Pattern Analysis and Machine Intelligence."xqr053."pq04."r r 04326448."
422; 0'
- 149_ T0F wf c."R0J ctv"cpf "F 0Uqtm"Pattern Classification."4pf "gf klqp"P gy 'l qtn"Y kg{ ."42220'

- 14: _ U\ 0Nk"öHceg"tgeqi pkkqp"dcugf "qp"pgctgu"npogct"eqo dlpcvkuu.ö"IEEE Computer Society Conference on Computer Vision and Pattern Recognition."r r 0: 5; ó: 66."3; ; : 0'
- 14; _ M'E0'Ngg."L0'J q."cpf "F0'Mlki o cp."öCes wklpi "npogct"uwdur cegu"lqt"lceg"tgeqi pkkqp"wpf gt"xctkcdrg" rki j vpi .ö"IEEE Transaction Pattern Analysis and Machine Intelligence."xqr049."pq07."r r 08: 668; : ."42270'
- 152_ PUV" 4227" Ur gcngt" Tgeqi pkkqp" Gxcnvcvqp" r rcp." j wr <ly y y pkuö qx lur ggej hguvulur ml4227 lut g/27 agxcn rcp/x80 f h0"
- 153_ V0'J cucp"cpf "L0'J 0'N0'J cpugp."öCequwle"lcevqt"cpn{uku"lqt"tqdwuv"ur gcngt"xgtkLecvkuu.ö"IEEE Trans. Audio, Speech, and Lang. Process.."xqr043."pq06."r r 0 64/: 75."Cr t042350'
- 154_ N0' \ j cpi ." Y 0' F 0' \ j qw" R0' E0' Ej cpi ." L0' Nkw" \ 0' [cp." V0' Y cpi ." cpf " H0' \ 0' Nk" öMgtpgn' ur ctug" tgr tgugpvcvqp/dcugf "ercuuklgt.ö"IEEE Trans. Signal Processing."xqr082."pq06."r r 038: 6638; 7."Cr t042340'
- 155_ C0' Mpcu cuwpf ctco ." F0' F gcp." T0' Xqi v" O0' O eNctgp." U0' Ukf j ctcp." O0' O cuqp." " öY gli j vgf "NF C" vgej pls wgu"lqt"lxgevt"dcugf "ur gcngt"xgtkLecvkuu.ö" "IEEE International Conference on Acoustics, Speech and Signal Processing."r r 069: 3669: 6."42340'
- 156_ L0'Y tki j v"C0'I cpguj ."U0'Tcq."cpf "[0'O c."öTqdwuv'Rtlpek rri'Eqo r qpgpv"Cpcn{uku<Gzcev"Teqxgt {"qh" Eqttw vgf "Nqy /TcpiO cvtegu'xlc"Eqpxgz"Qr vlo k cvkqp.ö"Uwdo kwgf "vq'yj g"Lqwtprl'qh'yj g"CEO ."422; 0'

Vj g'4236'Eqphgt gpeg'qp'Eqo r wcvkqpcnNkpi wknku'cpf "
Ur ggej 'Rtqeguukpi "'TQENR I '4236.'r r 064/79"
Í 'Vj g'Cuuqekvqp'hqt'Eqo r wcvkqpcnNkpi wknku'cpf "
Ej kpgug'Ncpi wci g'Rtqeguukpi "

利用核依賴估計來進行多軌自動混音

Automatic Multi-track Mixing by Kernel Dependency Estimation"

吳宗庭 Vuwpi 'Vlpi 'Y w'

國立中央大學資訊工程研究所

F gr ctvo gpv'qh'Eqo r wgt'Uelgpeg'cpf "Kphqto cvkqp'Gpi kpggtkpi "

P cvkqpcn'Egptcn'Wpkxgtukv{ "

101522015@cc.ncu.edu.tw"

張嘉惠 Ej kc/J wk'Ej cpi "

國立中央大學資訊工程研究所

F gr ctvo gpv'qh'Eqo r wgt'Uelgpeg'cpf "Kphqto cvkqp'Gpi kpggtkpi "

P cvkqpcn'Egptcn'Wpkxgtukv{ "

ej.kcB@ewu.edu.tw "

摘要

近年來由於數位音樂的蓬勃發展，錄音器材越來越普及。使得非混音專業人士也能利用錄音界面 *Cwf kq" Kpvgt hceg+錄製出不錯的成品。但是一旦錄製了多軌 *O wnk/Vtcem' Tgeqtf kpi +就會面臨到混音 *O kz kpi +的問題，即需要把多軌的聲音混合在同一個軌中。混音牽扯到許多音響及聲學心理學的相關技術與知識，非專業人士要混出尚可的成品有一定的難度，所以我們提出了自動多軌混音系統 *Cwqo cvk' O wnk/vtcem' O kz kpi "U{urgo +，希望藉由監督式學習的方式學習各軌間混音參數的調配，產生每首的基礎混音 *Dcukle" o kz/f qy p+來幫助非混音專業人士也能混出不錯的成品 *O kz/f qy p+。由於混音參數取得不易，我們會先藉由分軌及混音好的關係估計出各個混音參數，接著利用其參數進行混音模型 *O qf gm+的建立。在參數學習 *Rctco gvt' Ngctpkpi +方面由於每軌的混音參數是有依賴關係的 *F gr gpf gpe{ +，我們採用了核依賴估計 *Mgt pgn' F gr gpf gpe{ 'Guko cvkqp+3的參數學習 *Rctco gvt' Ngctpkpi +方式來預測每軌的混音參數。

Abstract"

F wg"vq"vj g'tgxqmwkqp"qh'f ki kcn'o wuke."r gqr ng"ecp"etgcvg"tgeqtf kpi u"kp" c"j qo g"uwf kq" y kj "ej gcr gt"i gct'J qy gxgt"o wnk/vtcem'tgeqtf kpi u'pggf "vq"dg"o kzgf "vq"eqo dkpg"vj go "kpvq" qp"qt"o qtg"ej cppgn'Vj g's wgvkqp"ku"vj cv'o kz kpi "tgs wktgu'dceni tqwpf "npqy rgi g"kp"uqwpf " gpi kpggtkpi " cpf "r u{ej qceqwku'K'ku'f k hlew"vq"i gv'i qqf "o kz f qy p" hqt "pqp/ur gekrku'kp" uqwpf "gpi kpggt'K"vj ku'r cr gt."y g"wug"uwr gtxkugf "ngctpkpi "o gjv qf "hqt"cwqo cvkcm{ "o kz kpi " o wnk/vtcem' tgeqtf kpi " kvq" eqj gtgpv' cpf " y gm/dcrpegf " r kgeg' F wg" vq" ncem' qh" o kz kpi " r ctco gvtu."htuv'y g"guko cvg"vj g'y gki j v'qh'o kz kpi "r ctco gvtu"d{ "wukpi "vj g'tgrv'kqp"dgvy ggp" tcy "o wnk/vtcem'cpf "o kz f qy p'0I kxgp"vj g"o kz kpi "r ctco gvtu" hqt"cp{ "o wuke"i gptg."y g" wug" ngtpgn'f ggepe{ "guko cvkqp"o gjv qf "vq"etgcvg"qtw"o kz kpi "o qf gr'Vj g"gzr gtko gpv'uj qy "MF G'ku"

cdng"q"o cng"o qtg"ucvkucvqt {"guko cvkqp"j cp"tgcvkpi "gcej "r ctco gvg"lpf gr gpf gpn{ 0'

關鍵詞：核依賴估計，音樂資訊檢索，音樂製作，混音

Mg{y qtf u<MgtprnF gr gpf gpe {"Guko cvkqp."O wuke"KT".O wuke"Rtqf wevkqp."O kzłpi 0'

"

一、緒論

在音樂的製作*O wuke" Rtqf wevkqp+上大致分為三個階段，創作編曲*Rtg/Rtqf wevkqp+、聲音錄製*Rtqf wevkqp+、後製*Rquv/Rtqf wevkqp+。其中後製又分為混音*O kzłpi +及母帶後製*O cuvgtłpi +兩部分。混音在音樂製作上是一個非常重要的過程。其主要的工作是要把先前錄製好的多軌*O wnk/Vtcem+的聲音，如人聲、吉他、爵士鼓等聲軌混合進同一個立體聲軌*uvgtgq"ej cppgn+或單聲軌*O qpq"Ej cppgn+中。

近年來由於數位音樂的蓬勃發展，錄音器材越來越普及。使得非混音專業人士也能利用錄音界面*Cwf lq" Kpvt hceg+錄製出不錯的成品。但是一旦錄製了多軌*O wnk/Vtcem' Tgeqtf łpi +就會面臨到混音*O kzłpi +的問題，即需要把多軌的聲音混合在同一個軌中。混音牽扯到許多音響及聲學心理學的相關技術與知識，非專業人士要混出尚可的成品有一定的難度，混出來的結果往往會照成整首歌聆聽的清晰度降低、聲音不扎實、音量落差太大、空間感不夠、聲音雜亂等問題。而且混音的處理方式基本上會隨著樂器、音樂類型而有所不同，不同的音樂類型會有不同的混音風格*O kzłpi "U{rg+，這更加增加了一般非專業人士學習混音的難度。所以要如何藉由電腦來幫助混音便是本篇論文的目標。接下來將會詳細介紹混音的相關背景知識。

1.1 混音(Mixing)"

混音在音樂製作上是非常重要的過程，不同的混音方式在最後的成品上會有截然不同的模樣。混音的好壞會影響整首歌的表現，好的混音可以掩蓋瑕疵、放大優點，提升整體的質感。在混音的過程中，混音師*O kzłpi "Gpi łpggt+會依照各音軌樂器間的頻率*Hgs wpe{+、響度*ıqwf pguu+、音色、音場定位*Rcpqtco le"Rqukkqp+、空間感等聲音元素加以調配，以讓每個音軌*vtcem+樂器最佳化，讓每個音軌在最後混在一起時一樣能保持清晰，保有層次，使得音樂呈現更生動、更動聽。

在開始混音前，混音師會先作混音規劃*O kzłpi "F guki p+，規劃整首歌的音像*Uqwpf "Kı ci g+，決定每個音軌在整首歌裡的定位以及其重要性。如下圖 3"為一首爵士樂的混音規劃示意圖，我們可以發現在這首歌中人聲*Xqecn+設計在音像的正中間，吉他*I vkct+分別落在人聲的左右。音量方面主吉他*Ngcf "I vkct+略大於人聲等等，混音的過程中主要就是調配這些音量*Xqıno g+、等化*Gs wırkı cvkqp+、擺位*Rcp+來達成我們的混音規劃讓整體更加和諧，讓個樂器融入其中。

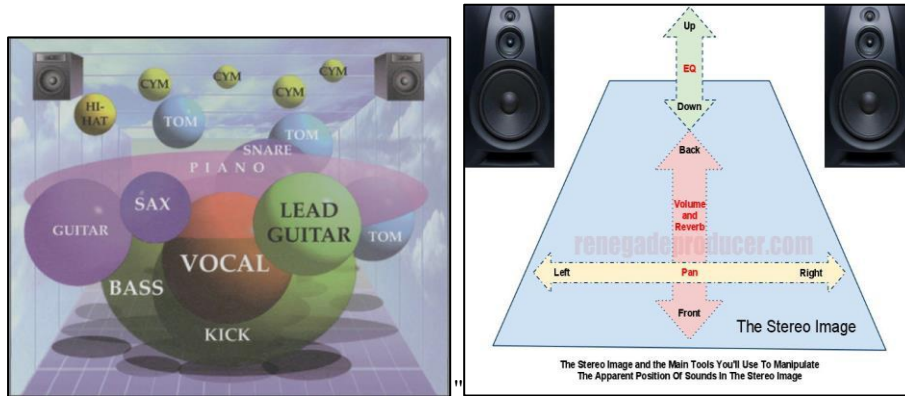


圖 3"混音設計示意圖：人聲居中 圖片來源：
The Art of Mixing – David Gibson

1.2 研究動機(Motivation)"

在自動混音的研究中，大多都是利用聲音特徵間的關係估計以及預測其混音的參數，所使用的聲音長度大多為一首歌中 52"秒的片段，最後所混音出來的成品當然也會相似於這 52"秒的片段。但是混音實際上是會因為橋段的不同而有不同的混音方式的，例如在主歌與副歌的混音方式會是不一樣的，後者其橋段通常為歌曲中激昂的部份，配器使用會前者多，在音量或頻率上比例會有所不同。

再者歌曲音樂類型也會影響混音的方式，同樣是爵士鼓在流行歌與爵士樂中音色與其占有的比例也會有所不同，同樣的樂器在不同的音樂類型中會有不同的音色以及角色，例如，在爵士樂中鼓手也可以是樂曲中的主角，在流行歌中主角往往是人聲或是電吉他。所以若利用 52"秒片段建立的模型來套用在整首歌曲的混音將會造成整首歌較平淡無味。在 [4] 的訪問中也有提到，混音其實是個別的，混音師基本上都會依照各音軌的聲響、音樂類型不同而有不同的處理方式，再者音樂感受是非常主觀的，不同的人會有不同的偏好，所以比較難去訂出一個通則來進行混音。

一個混音模型包括音量、頻率、樂器擺位等多個參數，大多數的研究都是獨立的去預測各軌的混音參數，如多線性迴歸，利用各軌特徵間的關係去建立各軌的迴歸模型。但實際上每軌間的混音參數是有依賴關係的，舉例來說有一軌的音量上升勢必就會有一軌的音量下降，這樣整首歌的音量才不會忽大忽小，若是對各軌獨立去建立其模型的話，最後的成品將會喪失其依賴關係。所以在參數預測的方法選用上我們認為需要考慮到依賴性的問題。

由以上三點，本篇論文採用對不同的音樂類型不同的橋段，利用和依賴估計的方法來建立最後的模型。首先我們會先需要使用者先提供一些該歌曲的訊息，例如音樂形態、橋段、分軌的樂器標籤等等，接著在對個別的橋段套用該音樂類型的核依賴模型預測出混音參數後即完成混音。圖 4"

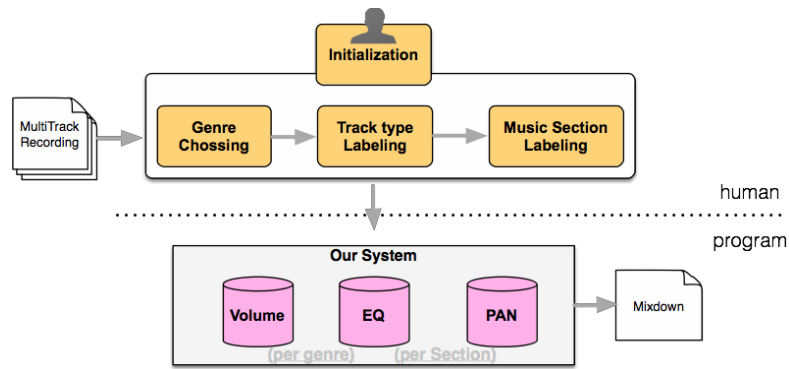


圖 4" 系統使用概念圖

本篇論文的架構如下，第二章將會討論自動混音的相關研究，第三章則會介紹本篇論文的研究方法及所使用的資料集，第四章為實驗，我們會對我們的混音模型去做交叉驗證來評估模型的正確性以及依賴性的功用，第五章為結論以及未來工作。

二、 相關研究

在介紹混音參數預測模型之前，我們必須設法取得一些歌曲的混音參數以作為訓練資料，不過混音參數在實際上是難以取得的，由於混音師使用的軟體不同，不同的器材軟體會有不同的設定、基準及刻度，這讓得到混音參數這一類的資訊變得非常困難，且混音師在混音時也鮮少會把參數記錄下來。所以在多軌混音的相關研究中大部份的研究著重於如何估計出混音參數，包括：音量、頻率、動態等等。另一部份的研究著重在如何建立混音參數模型。

2.1. 混音基本元素(Basic Factor of Mixing)"

在混音的過程中音量平衡是件很重要的事，混音師會調整每個音軌間彼此之間的音量，決定各軌在這首歌中的音量比例，即是決定各軌在音像的前後順序。若其中一音軌的音量比其他音軌還要大很多的話，整首歌將會聽起來頭重腳輕。

在音色修正的過程中，混音師會用到等化器。的工具來去對每個音軌的頻率做修正調整。例如我們發現吉他的音色跟其他的樂器相比太亮太尖銳，以至於無法融入這首歌中，我們就可以用等化器去對吉他的高頻部分做衰減。頻率間的平衡在混音過程中也是重要的過程之一。

樂器擺位是將錄製好的聲音訊號放置於新的雙聲道或多聲道的聲場。由於我們一般音響設備的環境基本上是以雙聲道為主，雙聲道的混音可以在聆聽上增加平面的聽感，而不是只有一點。所以我們在混音的時候會決定各樂器的擺位，看是要擺在中間還是擺在靠近左聲道左喇叭還是右聲道右喇叭等等來增加整體的空間感，而非全部的樂器都擠在一起。在實作上即是分別調整左聲道與右聲道的音量，如下式

$$Left_output = \cos(p) * input$$

$$Right_output = \sin(p) * input \quad (1)$$

其中 p 為偏離中央點的角度， $Left_output$, $Right_output$ 分別為左聲道右聲道的輸出。

2.2. 混音參數估計(Mixing Parameter Estimation)

如前一節研究動機所提到，由於混音參數難以取得，我們需要利用原始分軌和混音成品之間的關係來估計出每首歌的混音參數。在 [5] 的研究中，他們採用了訪問的方式，訪問了線上的混音師，了解他們如何混音、如何處理聲音，利用訪問後得到的一些通則來當作他們最後混音模型建立的依據。在大部分自動混音的研究中，在混音參數估計上大多是假設其分軌和混音成品是聲音特徵的線性組合關係，如圖 5 所示 Z_k 為第 k 聲軌的特徵向量， β_k 為第 k 軌之混音參數權重， \hat{y} 為混音成品的特徵向量。利用最小平方法的方式最小化 $\|y - X\hat{\beta}\|$ 至 $\text{col}(X)$ 平面的距離來做各音軌混音參數權重估計。

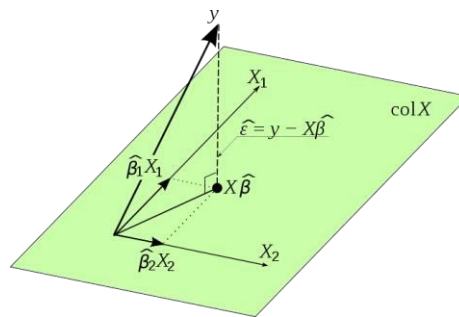


圖 5 最小平方法

由於採用線性組合的假設，每軌間彼此要是線性獨立，這會影響最後聲音特徵的選用以及音軌的選擇，選擇含有較多串音的音軌，例如錄製時會連大鼓小鼓等其他鼓組的聲音一併錄進，在實務上因為包含多個樂器的聲音，使得音軌間彼此並非獨立會造成其估計結果會有誤差。

2.3. 混音參數預測

參數預測大多是採用機器學習的方法，如 [5] 採用了多線性迴歸的方法，利用大量的數據建立最後的回歸模型，利用最佳化的方法去最小化分軌及成品間的尤拉距離。在 [5] 中同時利用了的方法，將混音參數視為一個潛藏的狀態，利用聲音特徵當作動態系統的輸出。另外在 [8] 採用了不一樣的方式，利用例子混音。其概念類似理髮師的概念，混音前提供例子供使用者選擇，最後藉由複製該例子的混音設計的方式來達成混音。

三、 研究方法

本系統主要流程如下圖 6 所示，我們會先把原始的分軌錄音檔依照使用者提供該音樂的資訊做前處理，接下來做聲音特徵的擷取，以便之後模型的訓練及測試。由於不同的類型、不同橋段的音樂會有不同的混音方式，在模型建立時我們會特別依照不同的音樂類型建立個別的混音參數模型，在依不同的橋段去建立模型，如圖 7 我們會對 TQEM 的音樂類型建立 Intro、Verse、Chorus，RQR 音樂類型也建立其三個橋段的模型，以此類推。訓練的部分有兩大步驟，第一步驟是混音參數估計，由於混音參數難以取得，原始訓練資料中也無此資訊，我們會先利用原始分軌錄音和混音成品做最小平方估計，藉此來估計出訓練資料中每首歌的混音參數的權重，第二步驟是核依賴估計的模型建立，我們會利用每首歌的混音參數權重及特徵向量來當作訓練核依賴估計的依據，訓練好的模型將會用來預測各個混音參數的權重。最後依照模型預測出的權重進行混音。

在接下來的章節我們會詳細介紹各步驟的做法，在章節 5.1 會先對我們所使用的資料集做介紹以及前處理的部分。章節 5.2 會介紹我們如何利用最小平方來估計各個混音參數及為何需要做混音參數估計。章節 5.3 我們會介紹核依賴估計的核心概念以及在本篇論文的問題中該如何設計。

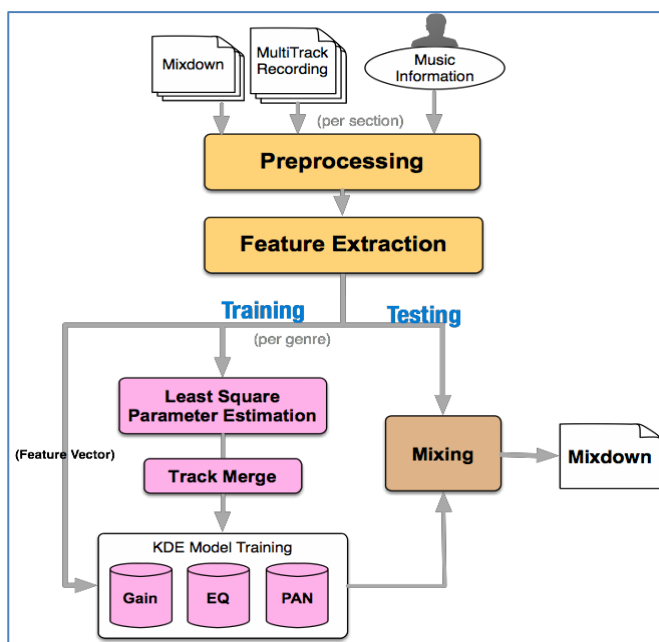


圖 6 系統架構圖

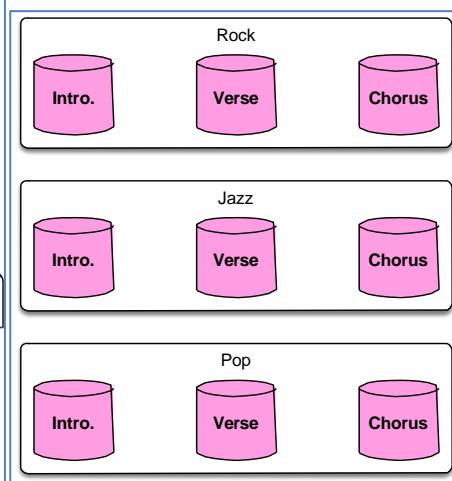


圖 7 混音參數模型

3.1 資料集 (DATASET)

我們使用的資料集 *F cvc"Ugv* 是來自國外一本關於混音的專門書籍 *öO kz lpi "Ugetgw* *hqt" Ub cm'Uwf kqö]9_*，此書有提供多首原始分軌檔案給讀者用於混音練習用，其含括的音樂類型搖滾、爵士、鄉村等多種音樂類型，如下表 3 所示，此書將相似的音樂類型分成四大類。此資料集較特別的點在於提供的聲音檔案長度是整首歌 **Hvni'O wnkctcem*，一般以往音樂資訊探勘 **O wule" kphqto cvkqp" Tgtkgxcn* 研究所使用的資料集大多是 42 秒 52 的長度，鮮少有提供整首歌的資料集。此特點有利於幫助我們對於不同的音樂橋段 **O wule" Ugevkqp* 去建立不同的模型，來讓最後的模型能更適用於實際的情況，此資料集所提供的混音成品 **O kz f qy p+* 一樣也是整首歌的長度。檔案格式為無損 YCX 檔 **vpeqo r tguugf "Y CX' hkgu. "46dk' cpf "66Ø' hJ | "uco r ng' tcvg+*。

表 3 音樂類型統計表

| Genre | # of Song |
|---|-----------|
| Alt Rock / Blues / Country Rock / Indie / Funk / Reggae | 9 |
| Rock / Punk / Metal | 39 |
| Pop / Singer-Songwriter | 32 |
| Acoustic / Jazz / Country / Orchestral | : |
| Total | 64 |

資料集前處理的部分，由於此資料集涵蓋的音樂類型包括 *Cn" TqemDwgu*、*TqemRvpm*、*Rqr lUkpi gt/Uqpi y tkgt*、*Ceqwule llc| lEqvpt* 等，每個音樂類型所使用的配器會略有不同，例如在爵士樂中管樂器的使用比例會比搖滾樂來的高。這樣會造成之後無法將此資料集套用至我們的模型中。所以為了解決此問題我們統計了各音樂類型所使用的配器，決定出每個類型基本音軌 **Dcuke" Vtcm* 如表 4，我們定義了 34 種音軌的形態。前處理的部分會先對分軌做音樂格式上的轉換，待每首歌的混音參數權重估計出來後會做基本軌的合併 **Vtcm O gti g+*，將同一類型的音軌依照其權重先合併成該類型基本音軌，將同一類型的檔案先合併成一個-例如在某一首歌中吉他錄了兩把，我們會先將這兩把的錄音合併成一個，以方便之後訓練及測試該類型的音樂。

表 4 基本軌表

| 12 Basic Track Type | | | | | |
|---------------------|----------|-------------------------|------------------------|----------------|---------------|
| *3-Mlen | *4-Upctg | *5-J kj cv' | *6-VQO " | *7-Ftwo Tqgo " | *8-QXGTJ GCF |
| *9-RGTEWUUKP | *: -DCUU | *; +Grgettle I vkct" | *32-Ceqwule" I vkct | *33-+Ngcf XqzC | *34-+DcenXqz' |

3.2 混音參數估計 (PARAMETER ESTIMATION)

我們的系統最終的目的是要希望藉由訓練資料 **Vtcklpi "F cvc+* 來對每一種混音參數 建立一個模型，藉由每一軌的特徵來預測其參數的值。所以訓練的過程中勢必需要原始

的分軌檔案及最後各混音的參數來進行監督式學習 [16]。但是實際上混音參數的資訊是非常難取得的。由於混音師使用的軟體不同，不同的器材軟體會有不同的設定、基準及刻度，而且混音師在混音時也鮮少會把參數記錄下來，這讓得到混音參數這一類的資訊的取得變得非常困難。為了之後的監督式學習，我們必須先估計出資料集中每首歌的混音參數，利用原始分軌檔案 [17] 及最後混音成品 [18] 來估計出每首歌的混音參數，針對每一首歌求得其混音參數當作之後監督式學習 [16] 的依據。

為了從原始分軌檔案估計出其混音參數，我們假設原始分軌與最後混音的成品 [19] 的關係是一個線性的組合 [20]，如下 [21] 式為一首歌的線性組合關係。

$$\alpha_1 U_1 + \alpha_2 U_2 + \dots + \alpha_k U_k = V \quad (5)$$

α_i 為第 i 軌的混音參數權重， $U_i = [u_{1i}, u_{2i}, \dots, u_{Ni}]^T$ 為第 i 軌特徵向量， V 為最後混音結果的特徵向量 [22]，每一軌抽取 N 個 [23] 做為其代表。其中在不同的混音參數會用不同的聲音特徵，例如在音量參數方面會採用聲音的方均根來當作衡量的依據，頻率參數方面會採用聲音的頻譜 [24] 等等。利用此線性組合的關係，我們可以利用最小平方方法 [25] 來估計出混音參數 [26] 的數值，最小平方方法 [25] 是以觀測值 W 與預測值 \hat{W} 之差的平方和作為最佳化的目標函數 [27]。以音量參數為例，令 u_{Nk} 為第 k 軌第 n 個音框 [28] 的方均根值 [29]，每個音框長度約為 42 毫秒，則 [30] 即可表示為

$$\begin{bmatrix} u_{11} & u_{12} & u_{13} & \dots & u_{1k} \\ u_{21} & u_{22} & \dots & \dots & u_{2k} \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ u_{N1} & \dots & \dots & \dots & u_{Nk} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \vdots \\ \alpha_k \end{bmatrix} \approx \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ \vdots \\ v_N \end{bmatrix} \quad (6)$$

在混音時由於混音的參數眾多，本篇論文只討論了其中三個最重要參數：[31]、[32]、[33]，在估計不同的參數時我們會使用不同的聲音特徵。接下來會對各參數所用的特徵來做介紹

音量 [34]

音量參數也被稱作增益 [35]，主要在控制每個音軌間的音量使得整體音量達成平衡。本篇論文使用了方均根 [36] 的方式去測量同一個音框 [37] 中各音軌的聲壓 [38]，以此作為音量參數的特徵向量，寫成矩陣形式如 [39] 式。 I_m 為第 m 軌的權重， P 為音框的長度。

$$\begin{bmatrix} RMS_{11} & RMS_{12} & RMS_{13} & \dots & RMS_{1k} \\ RMS_{21} & RMS_{22} & \dots & \dots & RMS_{2k} \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ RMS_{N1} & \dots & \dots & \dots & RMS_{Nk} \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ \vdots \\ g_k \end{bmatrix} \approx \begin{bmatrix} V_{RMS_1} \\ V_{RMS_2} \\ \vdots \\ \vdots \\ V_{RMS_N} \end{bmatrix} \quad (7)$$

頻率參數估計

頻率參數即是控制整首歌中各音軌在頻率上的平衡。在混音過程中為了修正音軌的音色或頻率時會用到等化器來幫助我們對聲音的頻率作調整，也就是說頻率參數即是等化器參數。在設計等化器時會先將整個頻譜切成多段，如切成三塊的話即是高頻、中頻、低頻。接著選出各頻段的中心頻率及頻寬，後即完成設計。本篇論文在頻率參數方面一樣採用多頻段等化方式來模擬實際等化器的操作，即是把頻率參數估計的問題切成了多個聲音參數的子問題。如下圖 8 所示。

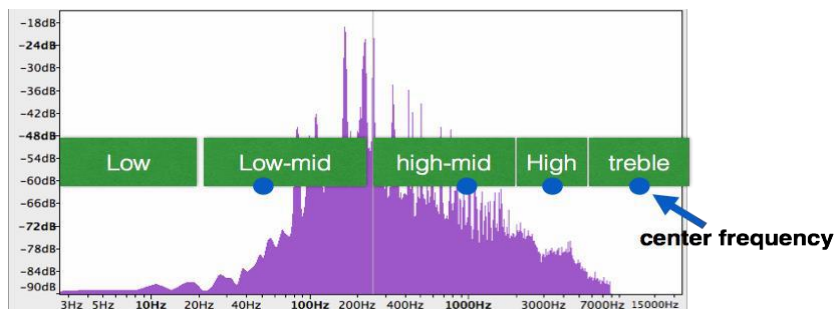


圖 8 多頻段頻譜

在進行頻率參數估計時，我們會先用快速傅立葉轉換先得到各音軌的頻譜，接著把各音軌的頻譜依照我們預先分段的頻率分別去解迴歸問題，以估計出各頻段在各軌間的平衡參數。如下式

$$\begin{aligned}
 \alpha_{Treble_1} U_{Treble_1} + \alpha_{Treble_2} U_{Treble_2} + \dots + \alpha_{Treble_k} U_{Treble_k} &= V_{Treble} \\
 \alpha_{High_1} U_{High_1} + \alpha_{High_2} U_{High_2} + \dots + \alpha_{High_k} U_{High_k} &= V_{High} \\
 \alpha_{H-mid_1} U_{H-mid_1} + \alpha_{H-mid_2} U_{H-mid_2} + \dots + \alpha_{H-mid_k} U_{H-mid_k} &= V_{H-mid} \\
 \alpha_{L-mid_1} U_{L-mid_1} + \alpha_{L-mid_2} U_{L-mid_2} + \dots + \alpha_{L-mid_k} U_{L-mid_k} &= V_{L-mid} \\
 \alpha_{Low_1} U_{Low_1} + \alpha_{Low_2} U_{Low_2} + \dots + \alpha_{Low_k} U_{Low_k} &= V_{Low}
 \end{aligned} \tag{7}$$

樂器擺位

樂器擺位即是決定該音軌在左聲道及右聲道之音量比例，如式 3。在本篇論文我們將擺位參數的問題轉化成前一章節音量參數估計的問題，即左聲道做一次音量參數估計，右聲道做一次音量參數估計，這樣即可決定該音軌在左右聲道之比例，如下式 8、9 所示。

$$\alpha_1 u_{L1} + \alpha_2 u_{L2} + \dots + \alpha_k u_{Lk} = V \tag{8}$$

$$\alpha_1 u_{R1} + \alpha_2 u_{R2} + \dots + \alpha_k u_{Rk} = V \tag{9}$$

3.3 核依賴估計模型建立(Kernel Dependency Estimation Model)

估計完資料集中每首歌的三種混音參數（ G 、 S 、 R ）後，我們有了每首歌的特徵向量 Z ，及每首歌的混音參數 μ ，即可藉由學習出 Z 間的關係建立我們最後的混音模型。在模型建立的部分我們會依照不同的音樂類型，不同的音樂橋段建立一組模型，其中一組模型包括了音量模型兩個（左聲道、右聲道），頻率模型五個

利用此模型建立的方法讓最後的成品更能應用在實際的歌曲中，訓練模型的樣本數我們自各訓練歌曲中隨機抽取 5222 個樣本來當作我們的訓練資料。如圖 9， w_{pm} 為該模型的第 p 個樣本的第 m 軌的聲音特徵值， α_{pm} 為第 p 個樣本之第 k 軌的混音參數權重值 M 為分軌的個數，在本篇論文為 34。

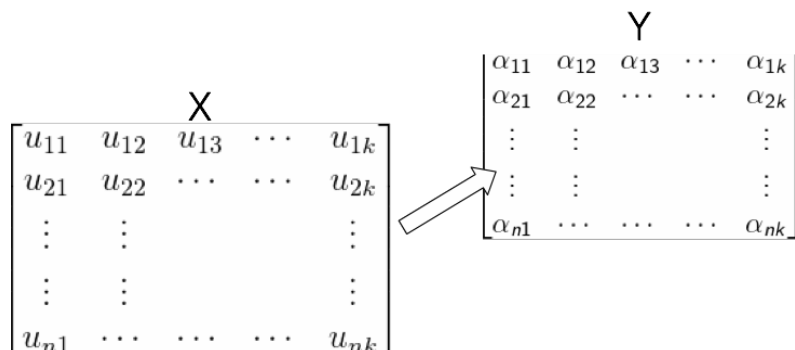


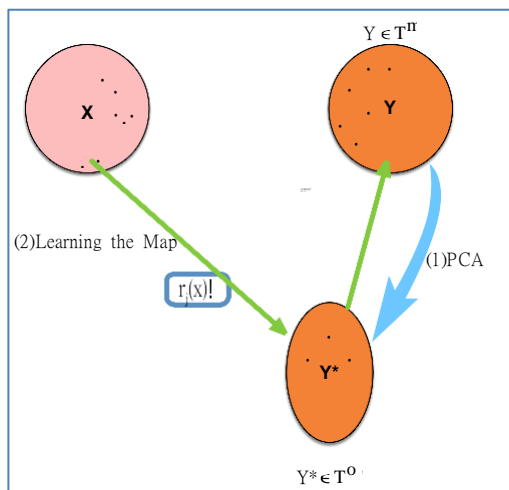
圖 9 混音參數預測示意圖

這個問題基本上是可以利用常見的參數預測的方法分別去求解，例如多線性迴歸、動態線性系統。但是由於在混音時每一軌間的參數間是有互相影響的，其中一軌的音量變大後勢必其他軌會變小聲些，彼此是依賴關係的。若用各別訓練模型來預測我們的混音參數的話，我們將會得到 m 個獨立的迴歸模型，這對於混音結果可能會產生較不和諧的影響。所以本篇論文採用了核依賴估計的方式來解決上述的問題，用核依賴估計訓練出我們最後的模型。接下來將會先介紹核依賴估計的基本精神。

3.3.1 核依賴估計(Kernel Dependency Estimation)

核依賴估計以下簡稱 MFG 是一種用於尋找輸入 Z 與輸出 Y 間依賴關係的學習架構， MFG 的流程如下圖所示，步驟主要分為三個步驟：

1. 投影：對輸出 Y 做主成份分析，將輸出 Y 也就是參數向量 $\{\alpha_k \in T^m\}$ 投影至 o 個主成份所形成的空間上，即對輸出 Y 做降維的動作，投影至較低的維度下形成 Z 。
4. 學習：對每個基礎原件，我們會學習一個對應函數 t ，對應函數 t 會將會把 Z 對應至 Y 第 1 個基礎原件，即是在較低的空間上去解 o 個迴歸問題。
5. 預測：對於一個新的輸入 z ，我們可以利用先前學習好的對應函數求出 $\{t_i\}$ ，最後再將 $\{t_i\}$ 投影回原本的空間上，即求出 Z 所對應的 Y 值。



圖：MF G 示意圖

用於本篇論文的主題來說，我們的 X 就是我們每首分軌的聲音特徵向量， λ 就是我們原先估計出來的混音參數 $\lambda \in T^m$ ，在 MF G 的過程中我們會先將 λ 降維至 o 維的 REC 空間上，接著在 o 維的空間上我們去學習 o 個基礎原件的對應函數。在對應函數方面我們使用了 $t_k f_i g^t g_i t g u k q p$ 的方法來解決 o 個迴歸問題。在做預測時，新的一首歌分軌特徵向量 z_θ ，會先利用先前學習的 o 個對應函數 t_i 求出 $\{\hat{\alpha}_i\}$ ， $\{\hat{\alpha}_i\} \in T^o$ ，接著再將 $\{\hat{\alpha}_i\}$ 投影回原本的 m 維空間上即求出各軌混音參數權重。

四、 實驗

實驗分為三個部分作討論，第一部分的實驗是評估由核依賴估計所建立的模型的正確性。第二部分是評估 MF G 依賴性 $F g r g p f g p e \{ \}$ 的效果，比較不同的 o 值對模型正確性的影響程度。第三部分是評估不同的音樂類型混音方式的差異性。實驗評估的方式是利用一次挑一個交叉驗證對同一音樂類型的歌做均方誤差 $O g c p U s w c t g G t t q t$ 評估。均方誤差的計算方式如下：

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{\alpha}_i - \alpha_i)^2 \quad * \cdot +$$

$\hat{\alpha}_i$ 代表由 MF G 模型所預測的混音參數權重， α_i 是由參數估計得出實際的混音參數權重值。

4.1 一次一個交叉測試(Leave-one-out Cross Validation)

表 5 為每個音樂類型副歌音量參數的一次一個交叉驗證 $N g c x g / q p g / q w E t q u i X c r k f c v k q p$ 的結果，其值為 MF G 模型所預測的值與實際混音參數權重的均方誤差值 $O g c p U s w c t g G t t q t$ ，測試時會先對測試歌曲做隨機抽樣 5222 個樣本來做測試， o 值皆為 1。由表可知由 MF G 所訓練出的模型均方誤差平均大約落在 204 左右的數值，預測出的權重有著還不錯的準確度，但由於是採用隨機抽樣的方式，若抽到的聲音樣本為較

安靜的樣本即是說在該樣本的當時有較多軌音量是趨近於 2 時會導致均方誤差值飆高，形成 Q_{wkt} ，如搖滾類中的第 39 首的結果。

表 5 副歌音量模型交叉驗證結果

| | TqemO gvcn" | RQR" | Lc Eqpvt { " | Cn'TqemHxpm" |
|--------------|---------------|-------|-----------------|--------------|
| uqpi 3" | 2059" | 2065" | 2067" | 2049" |
| uqpi 4" | 206: " | 2046" | 2027" | 20: 9" |
| uqpi 5" | 20; 4" | 2079" | 203: " | 2062" |
| uqpi 6" | 2084" | 2052" | 203: " | 2075" |
| uqpi 7" | 20427" | 2078" | 205: " | 205: " |
| uqpi 8" | 2057" | 2092" | 2059" | 20648" |
| uqpi 9" | 2088" | 2089" | 2039" | 203: " |
| uqpi : " | 2057" | 2077" | | |
| uqpi ; " | 2074" | | | |
| uqpi 32" | 2042" | | | |
| uqpi 33" | 20; ; " | | | |
| uqpi 34" | 2062" | | | |
| uqpi 35" | 20; 4" | | | |
| uqpi 36" | 206: " | | | |
| uqpi 37" | 2048; " | | | |
| uqpi 38" | 2096: " | | | |
| uqpi 39" | 2.175" | | | |
| Mean" | 204; 2" | 2097" | 2076" | 20; ; " |

圖 31 是其中一種音樂類型中副歌交叉驗證的詳細圖表，Z 軸代表其隨機抽樣測試的聲音樣本 $U_{qwpf} U_{co r ng}$ ，[軸為其所對應均方誤差，圖表中同樣顏色的線條代表同一首歌的聲音樣本。可以發現 MF G 所訓練出的模型在對同一首歌的聲音樣本時會有一致效果。

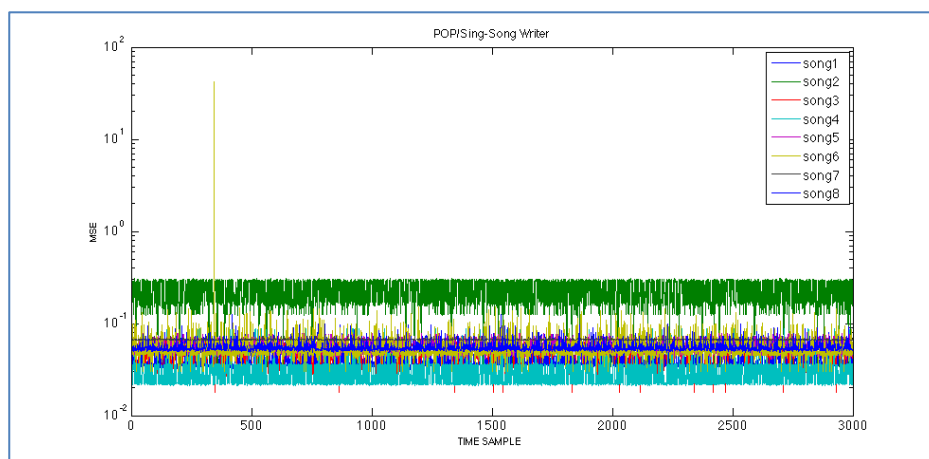


圖 31 $RQR U_{pi} / U_{qpi} 'y tkgt'$

"

在另一方面由音軌的觀點來看，如圖 32，我們可以發現在 Mem、Ftwo Tqgo、Qxgtj gcf 等音軌上各類型會有較大的誤差出現，其原因要歸咎於在實際上多軌同步錄音 $O_{wnk} t_{cen} T_{geqt} f_{kpi}$ 時，單一軌也會參雜著其他軌的聲音 \ast 串音 \ast ，如 Qxgtj gcf、Ftwo Tqgo 等軌會包含 n_{lem} 、 u_{pctg} 等其他鼓組的樂器。此原因違反了當初研究方法一開始的假設 \llcorner 我們假設分軌即混音成品間是個線性組合的關係，分軌間必須要是線性獨立 \llcorner 但由於串音的因素會導致估計及預測有效果不佳的情形。

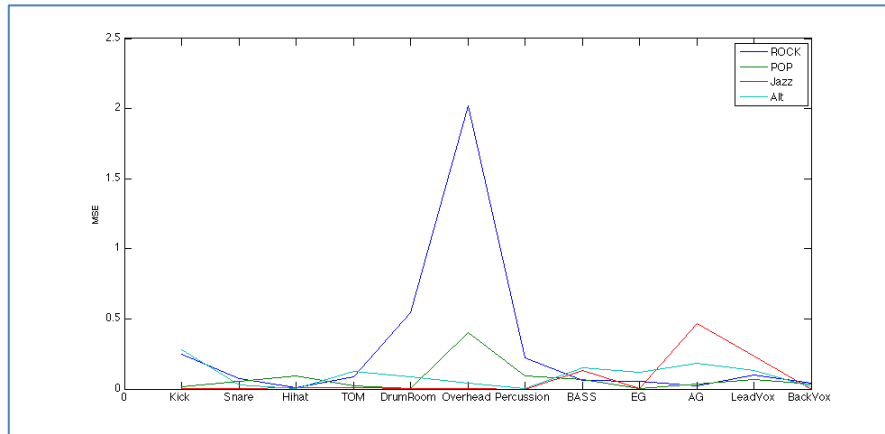


圖 32"音軌比較圖

在頻率*Gs wrk c vlp+模型方面我們同樣也作了交叉驗證，由表中可知大致上各分軌的準確率約為 2037，顯示在頻率參數預測上有著較好的表現。

表 6"頻率模型交叉驗證

| | TqemO gcn" | RQR" | Lc Eqpwt {" " | Cn'TqemHwpm" |
|--------------|---------------|---------------|------------------|---------------|
| Mleni" | 2033" | 203; " | 2025" | 202: " |
| Upctg" | 2065" | 2037" | 2033" | 203: " |
| J kj cv" | 2025" | 2023" | 2042" | 2043" |
| VQO " | 202: " | 2039" | 2035" | 2029" |
| FTWO TQOQO " | 2035" | 2026" | 2058" | 2023" |
| QXGTJ GCF " | 2039" | 204: " | 2025" | 2068" |
| RGTEWUQP " | 202: " | 2028" | 2047" | 2027" |
| DCUU" | 2033" | 2029" | 2026" | 2039" |
| GI " | 2028" | 203; " | 2046" | 2027" |
| CI " | 2022" | 2026" | 2037" | 2028" |
| NGCF XQZ " | 2043" | 2035" | 2027" | 2047" |
| DCEMXQZ " | 2024" | 2025" | 2054" | 2039" |
| Mean" | 0.012" | 0.011" | 0.016" | 0.015" |

4.2 KDE 方法的效果(Effect of KDE Method)"

第二部分的實驗是要來評估 MF G"即其依賴性的成效，首先我們討論了不同的 o "值所帶來的影響，結果如下圖 33，為 TQEM"在副歌時不同 o "值的均方差，Z"軸為不同的 o "值，["為所對應其混音。我們可以發現當 o "值越低 MF G"模型會有較好的結果，較低的 o "值也可加速 MF G"的計算*gi (REC"ur ceg"從 T³⁴"降低為 T: +。這結果也顯示了依賴性對於混音參數預測的幫助。

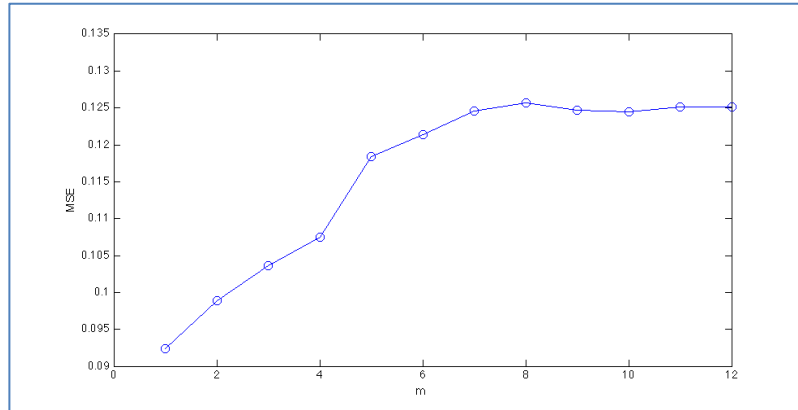


圖 33"不同 m 值比較

接著我們實作了相關研究中 [5] 所使用的多線性迴歸*無考慮依賴性+與本篇論文的 MF G*考慮依賴性+的方法做比較，結果如下圖 34，可發現在大多數的軌上 MF G比多線性迴歸有較好的表現。顯示其依賴性估計方式較能考量各軌之間的平衡。

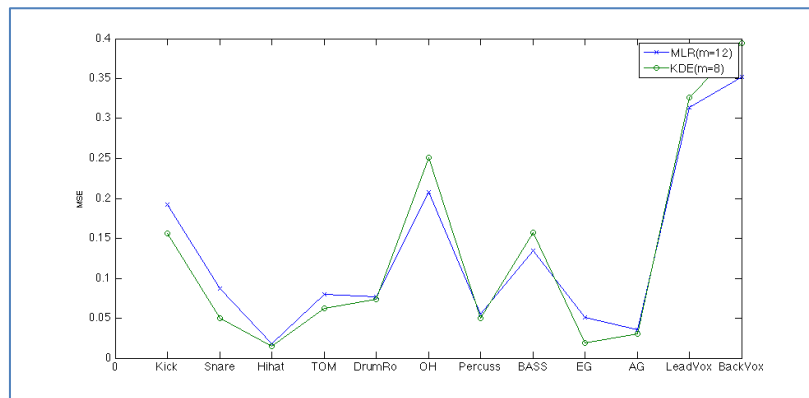


圖 34'MFG'與 ONT'比較圖

4.3 跨類型測試(Cross Genre Testing)"

在第一部分的實驗模型的建立以及測試都侷限在同一個音樂類型中，在第三部分的實驗，我們想要評估不同的音樂風格的音樂是否有其混音特色，我們將會套用由不同類型所訓練出的模型來看看其效果是否有差別。結果如下圖 13，Z'軸分別是先前定義的基本軌，['軸是其對應的均方誤差值，當中測試資料為 TQEM'這一類的歌，一共有 39"首，圖中的線條為分別用 RQR、lc||、Cn'Tqem所訓練出的模型套用至 TQEM'類別的測試結果。我們可以發現將別的類別的模型套用在不同類型的歌時會導致模型的正確度下降、誤差增大，如圖中的 RQR"與 IC\ \ "的結果，兩類別的模型套用在 TQEM'音樂上其結果顯示不太適合。我們也發現音樂類型相似的歌其混音方式會較相近如圖中 TQEM'與 Cn' TQEM'的均方誤差值較為接近。經由此實驗我們可以得知不同的音樂類型有其不同的混音方式，所訓練出來的模型有其獨特性。

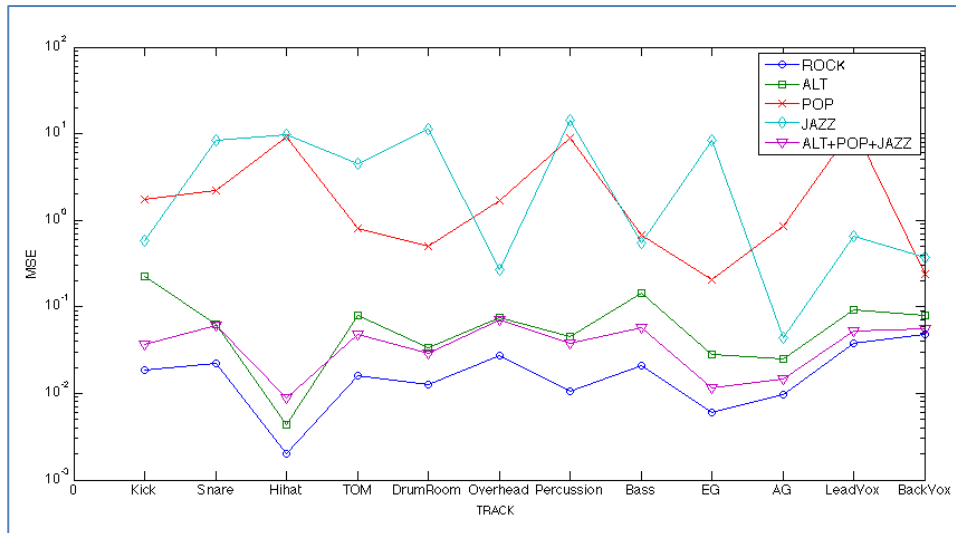


圖 13 類別交叉驗證

五、 結論與未來工作

在音樂製作上混音是非常重要的過程，音樂成品品質好壞，取決於混音是否混得好而且混音牽涉到許多音響及聲學心理學的相關技術與知識，非專業人士要混出尚可的成品有一定的難度，所以本篇論文提出一個利用監督式學習來進行自動多軌混音的系統。與其他篇相關研究不同的是其主要的核心方法是核依賴估計，利用混音參數間的依賴性來做混音參數的預測。另一個與其他相關論文不同的是訓練的單位不同，由於混音其實是非常個別的，混音師基本上都會依照各音軌的聲響、音樂類型不同而有不同的處理方式，所以在本篇論文的模型建立的過程我們會依照不同的音樂類型以及橋段建立不同的模型，以方便最後實際上的利用。由實驗結果得可知，不同的音樂類型其混音方式是有其獨特性的。

未來工作方面，由於本篇論文有分成四大音樂類型去做模型建立，導致各類別訓練資料有偏少的傾向，未來希望合併其他資料集以補足其數量。再者由於混音參數是非常難取得的，本篇論文是採用估計的方式估計其資料集中的混音參數權重，未來可以改用相關研究的估計方式或是實際收集混音參數以讓最後的混音模型能更貼近實際上的情況，例如建構一個線上混音系統的方式讓使用者實際混音記錄其混音參數。另外在 MF G 方法的部分由於本篇未套用其核函數的部分，未來可以在做投影前先套用核函數來提升混音模型的效果。實驗部分評估的對象也可以新增跟資料集的混音成品做比較試著看看目前的混音模型與實際上的差距如何。

參考文獻

- [1] J. Weston, O. Chapelle, A. Elisseeff, B. Scholkopf, and V. Vapnik, "Kernel Dependency Estimation," *Neural Information Processing Systems*, 2002.
- [2] J. Scott and Y. E. Kim, "Instrument Identification Informed Multi-track Mixing,"

- International Society for Music Information Retrieval*, 2013.
- [3] J. Scott and Y. E. Kim, "Analysis of Acoustic Features of Automated Multi-Track Mixing," *International Society for Music Information Retrieval*, 2011.
 - [4] D. BARCHIESI and J. REISS, "Reverse Engineering of a Mix," *Audio Engineering Society*, 2009.
 - [5] D. Barchiesi and J. Reiss, "Automatic Target Mixing Using Least-squares Optimization of Gains And Equalization Settings," *Digital Audio Effects(DAFx-09)*, 2009.
 - [6] H. Katayose, A. Yatsui, and M. Goto, "A Mix-Down Assistant Interface with Reuse of Examples," *Automated Production of Cross Media Content for Multi-Channel Distribution*, 2005.
 - [7] Mike.Senior, "Mixing Secrets for the Small Studio," 2011.
 - [8] D. Ward, J. D. Reiss, and C. Athwal, "Multi-track mixing using a model of loudness and partial loudness," 2012.
 - [9] Kolasinski and Bennett, "A Framework for Automatic Mixing Using Timbral Similarity Measures and Genetic Optimizatio," *Audio Engineering Society*, 2008.
 - [10] S. H. Nielsen and E. Skovenborg, "Evaluation of Different Loudness Models with Music and Speech Material," *Audio Engineering Society*, 2004.
 - [11] B. C. J. Moore, B. R. Glasberg, and M. A. Stone, "Why Are Commercials so Loud? ' Perception and Modeling of the Loudness of Amplitude-Compressed Speech," *J. Audio Eng. Soc*, 2003.
 - [12] Balster and Alex, "Audio Control Facilities in Modern Recording Studios," *Audio Engineering Society*, 1972.
 - [13] N. Montecchio and A. Cont, "Accelerating The Mixing Phase In Studio Recording productions By Automatic Audio Alignment," *International Society for Music Information Retrieval*, 2011.
 - [14] E. R. R. 128, "Algorithms To Measure Audio Programme Loudness And True-peak Audio Level," *EBU*, 2010.

Vj g'4236'Eqphgtgpeg"qp'Eqo r wcvkqpcn'Nkpi wku'eu'cpf "
Ur ggej 'Rtqeguulpi "'TQENRPI "4236.'r r 07: /99"
Í "Vj g'Cuuqekvqp'ht'Eqo r wcvkqpcn'Nkpi wku'eu'cpf "
Ej kpgug'Ncpi wci g'Rtqeguulpi "

中文轉客文文轉音系統中的客語斷詞處理之研究

Research on Hakka Word Segmentation Processes in Chinese-to-Hakka Text-to-Speech System

黃豐隆¹、余明興²、林昕緯²、林義証³

¹國立聯合大學 資訊工程所, flhuang@nuu.edu.tw

²國立中興大學 資訊科學與工程所, msyu@dragon.nchu.edu.tw; lin@sinwei.tw

³建國科技大學 資訊管理學系, yclin@ctu.edu.tw

摘要

語言(Language)是文化傳承與推廣的首要工具,尤其是少數族群的語言,如:台灣的客語或原住民語言。臺灣的客家族群約佔總人口七分之一,為閩南語語系外之第二大族群。根據近年來相關臺灣客語使用狀況調查報告指出,阻礙客語傳承之主因是:不太會講。由於台灣學習環境使然,導致連客籍家庭的學童亦少能以客語說話、交談,具有聽、說客語能力者逐年下降,能說客語的人口大量減少,台灣出現客語失聲、客家文化失傳之危機。

我們為了建置線上客語的數位學習系統,已開發出以大量合成單元為基礎的客語四縣腔及海陸腔的中文轉客文的文轉音系統(Hakka Text-to-Speech, HTTS),以及相關的應用系統,如:線上國客雙語有聲詞典 [13]、國客雙語有聲地圖社群系統 [14] 等。

我們的系統,主要是提供不太會講客語或不會講客語的使用者來使用、學習客語。因此系統的輸入為「中文文句」,輸出為「客語語音」。這樣的設計,學習者或使用能不需額外再學習客語輸入法、客語拼音,只需使用最熟悉的中文,即可透過本系統來學習客語。

為了更進一步改善與提升文轉音的效果,本論著重在改善系統中的客語文句分析模組的客語斷詞處理。在系統中,使用者輸入中文文句後,透過我們提出的客語斷詞方法,能將「中文文句」轉換為「客語文句及斷詞和詞性標記結果」。透過這個提升後的斷詞與詞性標記結果,來得到更佳的文句分析結果、提升文轉音中的文意正確性,如:韻律階層的求取、停頓類型的求取及讀音的求取。

本論文提出混合型的 N-Gram 序列分數算法,搭配中文斷詞模組及動態規劃演算法的客語斷詞方法。在嚴重資料稀疏的客語語料下,對中文轉客語斷詞結果的精確率有 80.78%。相較於傳統中文詞直翻客語詞的方法,已提升不少。

Abstract

Language is a major tool for cultural inheritance especially for the minority nationality, for example Hakka and aborigine language in Taiwan. As second ethnic besides Minnan dialect, the population of Hakka in Taiwan is one seventh. According to the recently reports of Hakka usage survey in Taiwan, the difficulties to inherit the culture of Hakka is missed in spoken Hakka language, the reason is the environments for learning and has led to the results of descending population for communicating by Hakka. It will become crucial for the cultural inheritance of Hakka.

Therefore, we has developed the Text-to-Speech method and system for Hakka language, and our goal is building environments for leaning the Hakka language, our some applied system such as:

“Web Hakka Phonetic Dictionary” [13] and “Blogging System of Bilingual Language by Integrating Mobile Cells and Google Map” [14], etc.

Our system is provided for users who interested in Hakka language, who can input the Chinese texts and system will output the speech of Hakka, users need not to learn the typing and phonetic writing of Hakka, and can take the advantage to learning Hakka with familiar language.

For the advanced improvements of Hakka Text-to-Speech, this article will emphasis on the word segmentation processing of Hakka text. In our system, when user enter the Chinese text, our proposed methods can convert the Chinese text to Hakka text and assign the part-of-speech for each Hakka text segments. By the better performance of text segments and part-of-speech in Hakka, We can improvements the Hakka text analysis module.

We proposed an hybrid N-gram sequence score, and Chinese word segmentation module developed by the dynamic programming algorithm, in the data-sparseness of Hakka corpus, the accuracy of Chinese to Hakka word segmentation is 80.78%.

Keywords: Hakka Text-to-Speech, Hakka Word Segmentation, Dynamic Programming, Hakka Text Analysis.

1. 緒論

一個斷詞系統的效果，通常跟語料的大小有關。但目前客語語料的收集非常困難，現有的電子資料，如：客委會初級、中高級的認證教材、教育部編著的國小客語教材等，對於自然語言處理來說，資料規模仍然屬極少量語料。因此，想要建置出更多的客語語料，幾乎都需要從客語書籍、文章中，透過人工輸入、建置成電子檔的方式來取得。但有了這些文本資料只是第一步，後續仍有許多的處理工作，如：斷詞、詞性標記的處理，擷取出這些語言特徵後，才能做更進一步的分析與應用。

客語詞的判定是一件嚴謹的事情，理論上我們必須遵照詞的定義¹來標記，但有極少數的情況下，我們仍會將詞組標記成一個客語詞，如：滑溜溜，在中文斷詞被斷為：滑/溜溜，我們視它是一個詞。而對於非客語語言專家的標記人員來說，其最有效率的方法，是透過具有平行資訊²的語料，先將中文語料輸入至中文的文句處理系統，取得中文的斷詞、詞性標記的特徵後，再對其對應的客語文章，以人工方式去判斷客語詞的邊界與詞性的標記。這個方法普遍被使用於同類型³的平行語料標工作記上，如 Tsai 的碩士論文 [1]也是用此方法。因為中文文句處理系統中，在文句斷詞資訊標記的技術方面已經相當成熟，而客語與中文的文法結構也相近，實際上中文的斷詞、詞性特徵，幾乎都能直接對應於客語詞。

目前客語語料的收集與建置，在學界有許多學者專家已積極的在做努力，建置出客語研究的相關基礎資料庫。如屏東教育大學的「學術研究基礎建置暨客家文化研究計畫 [2]」，他們歷時了至少三年的時間，在收集、建置客語語料及詞頻庫。這項創舉能有助於客語文句處理的發展，如：客語斷詞系統、客語文句分析系統、客語文句剖析系統、客語語音合成系統、智慧型的客語輸入法等，都非常需要足夠的客語語料來支持其發展。

¹ 詞(Word)，是指最小、有完整明確意義且可以自由使用的中文語言單位。

² 具有中文文句和客文文句 1 對 1 對應的平行資訊的語料。

³ 客語與中文都屬於漢語，文法結構幾乎相同，僅有少數俚語、特殊的客語構詞不同。

本研究旨在提出一個中文轉客文斷詞資訊的方法，針對嚴重資料稀疏的情況下，提出搭配中文斷詞模組與客語語言模型的混合式 N-Gram 序列分數的算法。透過兩階段方式，將中文文句以中文斷詞模組得到第一階段的斷詞及詞性標記結果後，再以國客語對照辭典找出所有可能被轉換的客語詞序列，以少量客語 Bi-gram、Uni-gram 語言模型為基底，搭配混合式 N-gram 序列分數的算法，找出分數最高的客語詞轉換序列，來得到第二階段轉換後的結果。

因礙於人力有限、語料收集的困難，仍有許多無法突破之處，如：語料的規模、人工標記資料的正確性。但使用本論文方法的客語斷詞法，以內部測試結果可知，若在訓練語料充足的情況下，能得到一個不錯的斷詞效能，其內部測試的精確度達 94.46%。相信在未來持續增加客語語料的規模後，本論文所提出的方法，效能會有更顯著的提升。

2. 文獻探討

2.1 中文斷詞

中文斷詞法每年都有新的研究與技術，甚至每年都會有舉辦斷詞比賽，知名的比賽如：SIGHAN 所舉辦的國際中文分詞競賽。第一屆競賽起始於 2003 年在日本札幌舉行。而之後每年都有相當多高手共襄盛舉。

在這麼多琳瑯滿目的斷詞技術中，常見的中文斷詞技術可分為三大類：(一)統計式斷詞法、(二)法則式斷詞法和、(三)混合式斷詞法。

(A) 統計式斷詞法

統計式斷詞法藉由收集詞彙資料，如詞彙長度和詞彙出現的頻率或次數等統計上的資訊。然後系統運用此訓練資料經由演算法分析來取得斷詞序列。常見的演算法如 Xue 使用的 Maximum Entropy [15]，最後實驗得到最好的 F 分數是 94.98%，或是 Lo 使用的 Conditional Random Field [4]，最後實驗得到最好的 F 分數是 96.40%。這些演算法都是利用字元之間的資訊當作特徵。然後把斷詞問題轉換成為字元之間的分類問題。

而過去常被使用的演算法是 Hidden Markov Model。如 Fu 和 Luck [5]從訓練語料中統計詞頻、字元在詞中出現位置的次數等資訊，組合過後做實驗。得到 F 分數最好可達 93.7%。而 Lin 和 Chang [4]使用兩階段特製化的方式，藉著擴充觀測符號及狀態符號來改善隱藏式馬可夫模型的斷詞效果。得到 F 分數最好可達 96.3%。

(B) 法則式斷詞法

法則式斷詞法主要是根據一些經驗法則做為斷詞的標準，藉以達到較好的斷詞序列。常見的法則如「長詞優於短詞」、「與左邊詞的結合優於與右邊詞的結合」。而這類型的斷詞法常被參考的是 Chen 和 Liu [17]，他們在該篇論文中提出了六條法則(heuristic rules)，並且根據這些法則解決了歧義性的問題及剔除一些較不可能的詞彙組合，以完成中文斷詞的工作。在實驗上效果相當不錯。

不過，此種斷詞法容易受到詞典的好壞而影響效能。若是句子中出現未知的新詞彙時，則正確率就可能下滑。

(C) 混合式斷詞法

每種方式的斷詞法都有好壞、優缺點，因此後來學者們才會嘗試去混合兩種斷詞方式。早期的 Nie 等人[18]提出結合詞典、經驗法則及統計資訊來對中文斷詞。而令人印象深刻的是 Wu 和 Jiang [19]提出結合剖析器和斷詞器的方法，在斷詞時先使用查詞典來產生所有可能的斷詞組

合，再使用經驗法則剔除不可能的詞彙組合，然後再利用剖析器解決剩下的歧義問題，最後得到斷詞序列。在實驗上 F 分數高達 99%。

而比較不同的做法是 Gao 等人 [20]提出的專有名詞法則，這些法則幫助系統抓取未知的專有名詞如：人名、地名、組織名、外來語音譯人名。在最後實驗上 F 分數約為 96%左右。

2.2 客語斷詞

因客語語料稀疏的緣故，目前針對中文轉客文的相關研究非常少，客語斷詞系統的實做，先以(一)輸入為客語、(二)輸入為中文，分為兩大種類。第一類是直接針對客語文句做斷詞，第二類是針對中文文句翻譯成客語斷詞結果。

而目前仍沒有任何一篇是針對客語斷詞做深入研究的論文，其中對客語斷詞有做效能評估的論文，也僅有 Tsai 的碩士論文 [1]—基於隱藏式馬可夫模型之客語文句轉音系統。顯見目前客語斷詞的研究，不管是語料的建置，還是斷詞的方法，仍非常多待探討與解決的問題。

(A) 輸入為客語

這一類的系統，適合具備客語輸入能力及熟悉客語的使用者，對於一般不熟悉客語的使用者而言，較不方便。這類系統常見的做法，是直接使用中文斷詞系統，對客文做斷詞。當然，這樣會有一些客語造字或客語用詞無法辨別的問題，針對這部份，是使用國客語對照辭典，來解決客語未知詞(Out of Vocabulary, OOV)問題。

如 Tsai 的論文 [1]，他們透過 Conditional Random Field 方法實做中文斷詞系統，並加入國客語對照外部辭典，配合客語構詞規則，實做出客語斷詞模組。最後的實驗效能，客語斷詞的 F 分數為 82.87%，客語詞性標記的 F 分數為 77.14%。

(B) 輸入為中文

這一類的系統，使用者不需使用客語輸入法，也不需熟悉客語，很適合客語初學者使用。這類系統常見的做法，是使用中文斷詞系統，先將輸入的中文文句斷詞，找出詞與詞性後，再將詞透過國客語的平行對照辭典，翻譯成客語詞。如本實驗室的線上客語語音合成系統，Wu [5]、Lo [6]的斷詞方法皆相同，都是使用 Jiang [7]所提出的中文斷詞系統，將中文文句斷詞後，再透過國客語對照辭典，將中文詞翻譯成客語詞。經測試後，其不含詞性標記的客語斷詞效能的 F 分數分別為 69.82%及 66.72%。

另一種是僅透過國客語對照辭典，將中文文句直翻成客語。如 Lee [8]，他們建置出一套國客語對照辭典，將輸入的中文文句字串切割成 1 到 4 字詞，並查找對照辭典、翻譯成客語詞。而他們沒有針對中文翻客語詞做效能評估，因此無法得知效果如何。

3. 準備工具及語料

3.1 中文斷詞工具

本論文的中文斷詞系統，是使用 Lai 於 2011 提出的「應用多詞及多詞性語言模型的中文斷詞及詞性標記方法 [9]」。此斷詞方法採用兩階斷式，第一階斷是斷詞，第二階斷是詞性標記。其斷詞的 F 分數有 96.69%，詞性標記的 F 分數 92.04%。

3.2 國客語對照辭典

我們所建置的國客語對照辭典，主要來源有：(一)客委會初級、中級暨中高級認證語料 [10, 11]、(二)台北市客委會-現代客語詞彙彙編。對於斷詞系統而言，辭典是決定正確率的重要因素，理論上辭典越大，斷詞效能也就越好。因此除了現有的辭典來源外，我們也利用標記客語

文句斷詞答案的同時，找出一些尚未被收錄在辭典中的國客語對照之詞目。最後，也針對每個詞目，進行人工校正工作，去除重複或不合理的詞目以及標記拼音。

表一、國客語對照辭典資料樣貌

| 欄位 | 內容 | 說明 |
|--------------------|------------------|----------|
| ID | 13267 | 資料庫中 ID |
| Chinese | 年輕人 | 中文詞用詞 |
| Hakka | 後生人 | 客語詞用詞 |
| Pinyin | heu1 sang2 ngin3 | 客語拼音 |
| Pos | Na | 客語詞性 |
| Pos_pattern | Na | 中文詞性組 |
| Hakka_pos_feq | 25 | 客語含詞性詞頻 |
| Hakka_nonpos_feq | 25 | 客語不含詞性詞頻 |
| Chinese_pos_feq | 21138 | 中文含詞性詞頻 |
| Chinese_nonpos_feq | 21138 | 中文不含詞性詞頻 |

表二、2014 版，興大國客語對照辭典分佈統計及比較

| 字詞 | 2014 興大客語辭典 | Wu[5] | Lo[6] | 交大客語辭典[1] |
|-----|-------------|-------|-------|-----------|
| 一字詞 | 3457 | 780 | 838 | 6747 |
| 二字詞 | 20800 | 16362 | 16540 | 18078 |
| 三字詞 | 7309 | 5654 | 5826 | 5095 |
| 四字詞 | 4093 | 3769 | 3861 | 4217 |
| 五字詞 | 312 | 273 | 283 | 250 |
| 六字詞 | 84 | 80 | 84 | 80 |
| 七字詞 | 65 | 63 | 68 | 60 |
| 八字詞 | 9 | 12 | 14 | 14 |
| 總計 | 36129 | 26993 | 27514 | 34541 |

3.3 客語四縣腔語言模型的建置

語言模型是斷詞系統中，用來選擇斷詞結果的重要元件。而語料經過統計詞頻後，就能得到該份語料的機率分佈模型，即是語言模型。因此，語料越大，能得到的統計資訊越多，語言模型能包含的情況越多，效能也會越好。

但是，現今電腦上的客語語料仍非常匱乏，能用來建置語言模型的語料非常有限。因此，我們開發了一個半自動式的客語語料建置工具來建置客語語料，透過這個工具，可以將具平行資訊⁴的客語語料，標記出客語斷詞資訊，再透過該語料統計出客語 Uni-gram 及 Bi-gram 語言模型。

⁴ 具有中文文句和客文文句 1 對 1 對應的平行資訊的語料。

3.3.1 語料來源

我們用來建置客語語言模型的客語語料，主要來源是客委會四縣腔初級、中高級客語認證教材 [10, 11]。這份語料，句子數分別有初級 1678 句、中高級 4962 句，共 6640 句，每一句都有中文、客語的詞目、拼音及例句，如表：

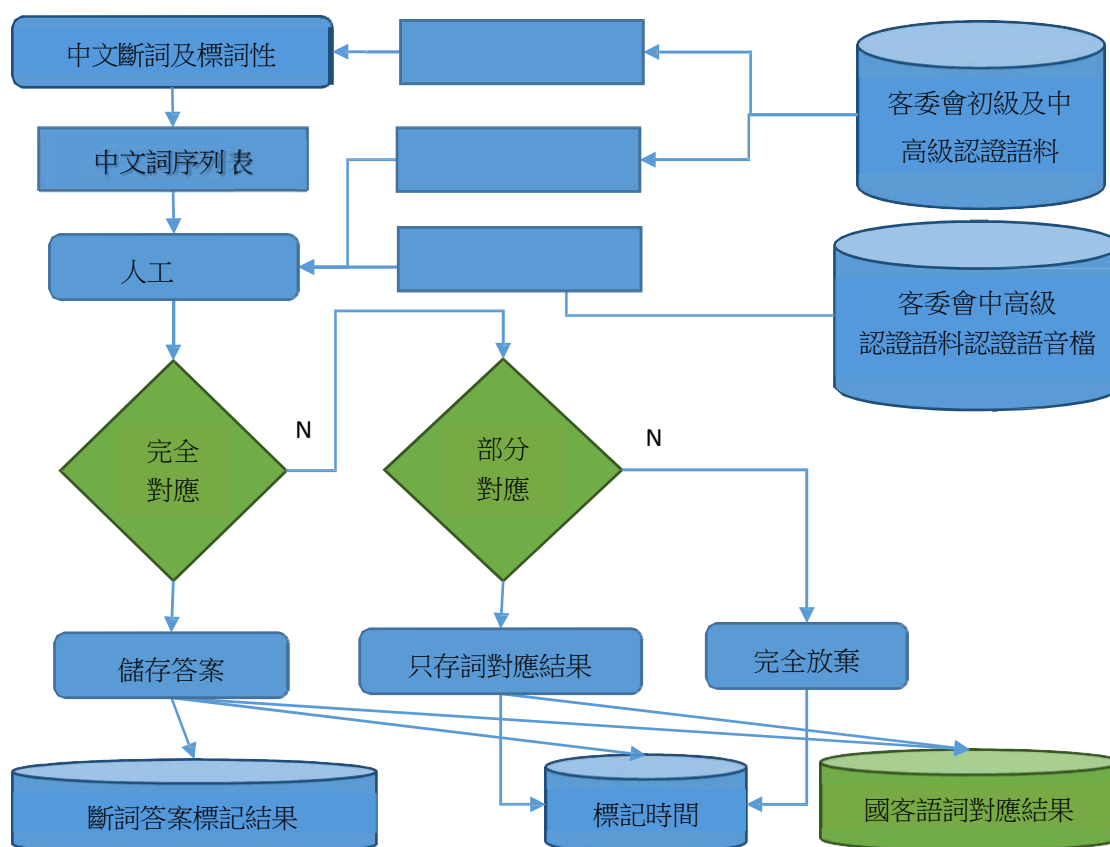
表三、客委會語料的資料樣貌

| | |
|------|-------------------|
| 客語詞 | 光線 |
| 中文詞 | 光線 |
| 客語拼音 | gong ` sien ˇ |
| 中文例句 | 這個房間的光線不好，不適合做書房。 |
| 客語例句 | 這隻房間个光線毋好，毋適合做書房。 |

3.3.2 客語斷詞標記工具及人工標記原則

客語與中文的文法結構相近，因此中文斷詞和詞性的標記，大部分都能與客語完全對應，僅有少部分的客語俚語或特殊用詞例外。而中文語料的處理，因目前中文斷詞系統的發展已相當成熟，因為中文語料的龐大，以規則法配合機率模型的混合式斷詞法所發展出來的中文斷詞系統，斷詞效能及詞性標記的 F 分數已達到 96.69%及 92.04%。因此都能直接以中文斷詞系統來得到可靠的斷詞及詞性特徵標記的結果。但客語文章的處理，因目前市面上及學界的客語斷詞系統仍處於發展中的階段，還沒辦法依賴任何客語斷詞系統來自動處理。因此我們開發出一套工具，以半自動的方法，快速的針對客語句子，做客語斷詞的標記。

3.3.2.1 客語斷詞標記工具介紹與操作



圖一、客語斷詞答案標記工具架構圖

我們將客委會四縣腔認證教材的客語語料 6640 句句子，依照句子的編號順序，由小到大分為語料 A 及語料 B，語料 A 有 4196 句，語料 B 有 2444 句，並再將 B 語料依編號順位分為 4 份，B1-B4，每份 611 句。

將 A、B 語料，分別找(標記者 1)一位碩士生、(標記者 2~5)四位大學在學生，使用本論文開發的客語斷詞答案標記工具，以半自動人工判斷方式，標記出客語斷詞的答案。其中每人的標記速度如表四：

表四、客語斷詞標記工具的標記時間統計，時間單位：秒。

| 標記者 | 1 | 2 | 3 | 4 | 5 | 總 平 均 |
|-------------------|-------|-------|-------|-------|-------|-------------|
| 語料編號 | A | B1 | B2 | B3 | B4 | |
| 處理句數 ⁵ | 4500 | 615 | 622 | 611 | 646 | |
| 儲存句數 | 4018 | 346 | 451 | 504 | 419 | |
| 總時間 ⁶ | 91798 | 28698 | 25297 | 18864 | 17219 | |
| 平均每句 | 20.39 | 46.66 | 40.67 | 30.87 | 26.65 | |

表中可看出，標記速度平均 33.04 秒能完成一句。而標記完成的資料，會儲存其 1.中文句子、2.中文斷詞及詞性標記結果、3.客語句子、4.客語斷詞及詞性標記結果、5.客委會語料中的句子編號，等五個欄位資料，儲存結果如下表：

表六、客語斷詞標記結果的資料樣貌

| | |
|--------|---|
| 中文句子 | 這個房間的光線不好，不適合做書房。 |
| 中文斷詞結果 | 這(Nep) 個(Nf) 房間(Nc) 的(DE) 光線(Na) 不(D) 好(VH) ， (COMMACATEGORY) 不(D) 適合(VH) 做(VC) 書房(Nc) 。 (PERIODCATEGORY) |
| 客語句子 | 這隻房間个光線毋好，毋適合做書房。 |
| 客語斷詞結果 | 這(Nep) 隻(Nf) 房間(Nc) 个(DE) 光線(Na) 毋(D) 好(VH) ， (COMMACATEGORY) 毋(D) 適合(VH) 做(VC) 書房(Nc) 。 (PERIODCATEGORY) |
| 句子編號 | 01-001 |

將這些標記完成的資料再經過一次經人工篩選後，最後確認有效句數為：A 語料 4018 句，B1-B4 語料共 1282 句，我們使用語料的分佈如表所示：

表七、客語語料的使用分佈

| | 訓練 | 測試 |
|----|-------|-------|
| 句數 | 4018 | 1282 |
| 詞數 | 45304 | 17646 |
| 字數 | 65572 | 25478 |

⁵ 處理句數的統計，包含重新處理曾經放棄的句子，因此實際處理可能會比分配到的筆數多。

⁶ 時間單位為秒。



圖二、標記工具操作畫面

3.3.2.2 標記原則

因為這份客語認證教材語料，專家們編審的主要目的是為了客語教學用，並不是為了要建立「國語/客語斷詞對應」的語料，所以在撰寫例句時並不會特別求強或注意到國客語的完全對應。

因此，我們在進行人工標記時也發現，其實大部分出現無法對應的情況，都可以以人工修改中文句子或用詞的方式，做適當的修飾與調整，來達到不影響文意、又能與客語句子對應完全的目的。但某些句子仍無法確認如何標記時，標記者也可選擇放棄該句的標記，或只存能對應到的詞，而不將這些未完成對應的句子視為斷詞答案。以下為標記時的 6 大標記原則：

原則 1：標記時，都以不修改客語句子為原則，但可跳過不影響文意的字串。

表八、標記原則 1 範例

| | |
|-----|--------------|
| 中文 | 這泉水的水質很甜很清澈。 |
| 客語 | 這窟泉水个水質當甜當清。 |
| 中文改 | 這泉水的水質很甜很清澈。 |
| 客語改 | 這泉水个水質當甜當清。 |

此例子中，窟這個字只是指一窟井或一窟水池的意思，省略後也不至於影響整句文意。

表九、範例 1 標記後樣貌

| | |
|-----|---|
| 中文改 | 這(Nep) 泉水(Na) 的(DE) 水質(Na) 很(Dfa) 甜(VH) 很(Dfa) 清澈(VH)。(P) |
| 客語改 | 這(Nep) 泉水(Na) 个(DE) 水質(Na) 當(Dfa) 甜(VH) 當(Dfa) 清(VH)。(P) |

原則 2：如果有很明顯且離譜的斷詞錯誤，要人工介入修正。

表十、標記原則 2 範例

| | |
|------|---|
| 原斷詞 | 這(Nep) 群(Nf) 小孩子(Na) ，(COMMACATEGORY) 每(Nes) 天都(Na) 在(P) 沙洲(Na) 上(Nes) 玩(VC) 摔跤(VA) 。(P) |
| 人工修正 | 這(Nep) 群(Nf) 小孩子(Na) ，(COMMACATEGORY) 每(Nes) 天(Nf) 都(Da) 在(P) 沙洲(Na) 上(Nes) 玩(VC) 摔跤(VA) 。(P) |

此例子中，「天都」一詞，被誤斷成一個普通名詞，這跟「天、都」意思不同，「天都」指的是地方名，而其正確斷詞應該斷成「天(Nf) 都(Da)」。

原則 3：可微調中文句子用詞及詞的順序以求對應到客語，但修改後的文意不能改變。

表十一、標記原則 3 範例 1

| | |
|-----|----------------------------|
| 中文 | 今天的天氣很好，太陽下山以後就可以看得到滿天的星星。 |
| 客語 | 今晡日个天時當好，日頭落山以後就看得著滿天个星仔。 |
| 中文改 | 今天的天氣很好，太陽下山以後就看得滿天的星星。 |
| 客語改 | 今晡日个天時當好，日頭落山以後就看得著滿天个星仔。 |

此例子中，若直接省略「可以」這個詞，也不會影響到原句的文意。

表十二、標記原則 3 範例 2

| | |
|-----|---------------------------|
| 中文 | 古時候的人會觀察天上的星宿變化，來判斷人間的吉凶。 |
| 客語 | 上早个人會觀察天頂星宿个變化，來判斷人間个吉凶。 |
| 中文改 | 古時候的人會觀察天上星宿的變化，來判斷人間的吉凶。 |
| 客語改 | 上早个人會觀察天頂星宿个變化，來判斷人間个吉凶。 |

此例子中，客語「天頂星宿个變化」與中文「天上的星宿變化」雖然詞的詞順序不同，但中文句子中的「的」調換後，也不影響文意。而修改句子的動作，我們只建議修改中文，客文 若非必要儘量保持原句。原因是較能保持客語句子原來的特性、結構、用語 等資訊。

原則 4：以客語句子為主體，發現中文詞和客語詞的對應有爭議時或太過模糊時，要找過一個最佳選擇的詞替換。

表十三、標記原則 4 範例

| | |
|-----|-----------------------|
| 中文 | 古早的人若是看到流星雨，會想到一些壞兆頭。 |
| 客語 | 上早个人係看著星仔瀉屎，會想著麼个壞兆頭。 |
| 中文改 | 古早的人若是看到流星雨，會想到什麼壞兆頭。 |
| 客語改 | 上早个人係看著星仔瀉屎，會想著麼个壞兆頭。 |

在此例子中，客語的「麼个」翻成中文「一些」，依客語文意來看過於模糊。因此，可透標記工具中的詞典查詢功能，找到客語詞「麼个」能翻成的國語詞有哪些，發現到「什麼」這個中文詞最貼近文意。因此，手動將中文句子的「一些」改為「什麼」，重新與「麼个」配對。

原則 5：若沒辦法標記完一整句，但部分詞能對應，則選擇「放棄，儲存詞配對結果」。

表十四、標記原則 5 範例

| | |
|----|--------------------|
| 中文 | 一身乾驚簡單粗陋的衣服都沒有能力買。 |
| 客語 | 一身腊食皮無才調買。 |

此例子中，我們發現這句除了「一/一、身/身、無/沒有、才調/能力、買/買」這些詞能對應外，其他詞皆無法非對應。因此此次的不納入正確答案範本內，但有部分詞能對應，也將這些詞的國客語對應結果儲存起來，待後續的工作中，仍可應用。

原則 6：配對時以「詞」為單位，不要以片語為單位。通常，非成語的片語，用法可能只是偶然出現，因此儲存這類資料沒意義。

表十五、標記原則 6 範例

| | |
|----|--------------------|
| 中文 | 現在外面既颶風又下雨，要怎麼樣回家？ |
| 客語 | 這下外背風合雨，愛仰般形轉屋下？ |

此例子中，發現到「既颶風又下雨」和「風合雨」看似能對應，但實際上可能只有在這句曾出現這種情況。因此，這種非成語的片語，也許只是偶然出現，我們不儲存此類的答案和詞配對。但依照原則 5，我們發現「要怎麼樣回家/愛仰般形轉屋下」可配對，成「要/愛、怎麼樣/仰般形、回家/轉屋下」。因此這句可儲存其「詞配對」成功的部份，但放棄儲存為斷詞答案。

4. 研究方法

本系統實驗有包括兩大種基底，(一)中文斷詞邊界優先、(二)客語詞邊界優先，經實驗發現第一種方法的外部測試正確率較高。但這可能跟標記語料時採用的：先中文斷詞、再人工對應客語詞的方法有關。但文章篇幅有限，本篇論文寫的實驗數據都以第一種方法為基底。

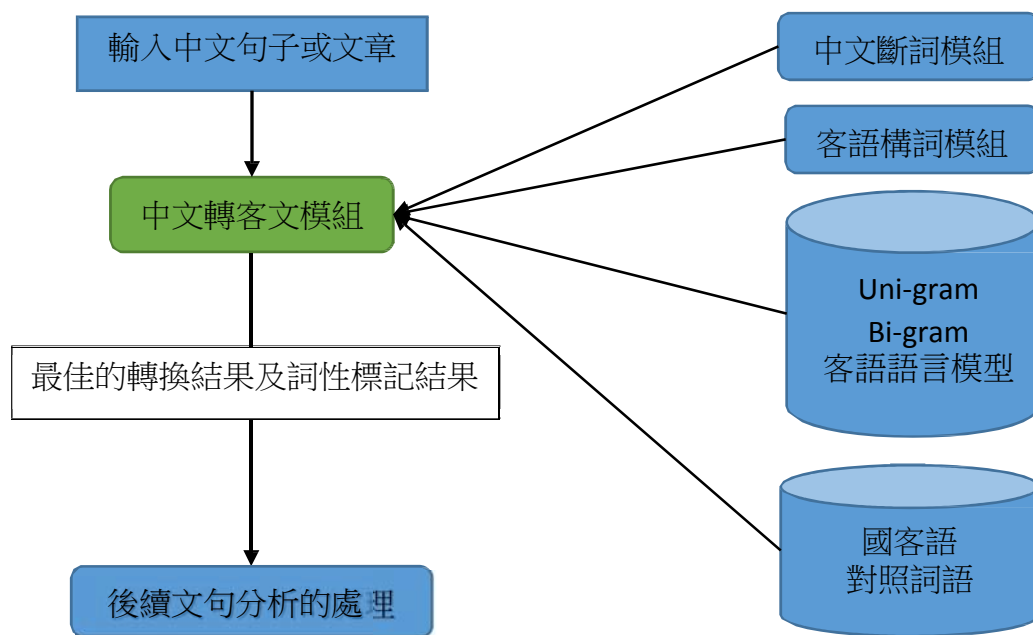
所謂(一)中文斷詞邊界優先，是先將輸入中文文句做斷詞，得到確定的詞邊界後，再以這些中文詞去查找國客語對照辭典，找出可能的客語詞轉換候選。而(二)客語詞優先，是先從國客語對照辭典中，找出所有可能被轉換的客語詞，最後再與中文斷詞邊界的詞一起做為候選詞。

兩種方法的不同之處：方法(一)的詞邊界已被固定，相對能找出的候選客語詞也較少、較侷限，方法(二)的詞邊界是自由的，可從國客語對照辭典及原中文斷詞斷的中文詞中，找出一條最佳的斷詞路徑。兩種方法其實各有好處，第一種方法在外部測試有較佳的正確率，第二種方法在內部測試有較佳的正確率。其意義是：如果在訓練語料充足的情況下，第二種方法會更好。而第一種方法適合用在訓練語料稀疏的情況。

4.1 系統架構

在本論文的客語斷詞系統中，輸入是「中文文句」，輸出是「客語文句」的斷詞及詞性標記結果。簡而言之，以中文文句輸入後，經此客語斷詞模組處理，產生具有斷詞與詞性的客語文句的輸出。標記斷詞及詞性兩種特徵後的結果，可再使用文法剖析器分析、得到文法結構樹，做更進一步的文句分析處理。如：用於 Hakka Text to Speech (HTTS) 中的停頓預估模型中，預測出句子中的 no break、minor break 及 major break 三種停頓類型，讓合出的語音可辨度更高、讓使用者能更輕易的聽懂句子內容。

因此，一個良好的客語斷詞系統是客語語言處理所不可或缺的，應是提升客語語音合成效果的重要因素。本系統中客語斷詞的架構，參見圖三。



圖三、客語斷詞模組架構圖

4.2 修改中文斷詞辭典

前文有提到，本論文採用第一種方法「中文斷詞邊界優先」為基底，因此為了讓一些客語用詞有機會成為被轉換的候選詞，我們將國客語對照辭典中，所有中文詞欄位的詞資料，都加入到中文斷詞辭典裡。如此做法能提高原本完全不會被找到的客語詞，有機會成為被轉換的候選詞。以下是一個例子：

表十六、中文斷詞邊界限制，範例 1

| | |
|-------|------------|
| 中文句子 | 泥鰍滑溜溜， |
| 中文斷詞 | 泥鰍/滑/溜溜/， |
| 中文翻客語 | 𩺰鰍仔/滑/溜溜/， |
| 正確答案 | 𩺰鰍仔/滑溜溜仔/， |

透過中文斷詞得到「泥鰍/滑/溜溜/，」的斷詞邊界，再透過以上方法找出最佳詞頻的客語詞並轉換為「𩺰鰍仔/滑/溜溜/，」。但其實我們的國客語對照詞典中，有收錄「滑溜溜/滑溜溜

仔」這筆個國客語對照資料，但礙於中文斷詞邊界之限制，只能由一個詞「滑溜溜仔」轉成「滑/溜溜」兩個詞。此情況會直接的影響到轉換的正確率。

表十七、中文斷詞邊界限制，範例 2

| | |
|-------|---------------------------|
| 中文句子 | 大家來去客家庄走一走。 |
| 中文斷詞 | 大家/來去/客家庄/ 走/一/走 。 |
| 中文翻客語 | 大家/來去/客家庄/ 行/一/行 。 |
| 正確答案 | 大家/來去/客家庄/ 遶遶啊 。 |

上面的例子「走一走」被斷成「走/一/走」，但實際上我們詞典有收錄「走一走/遶遶啊」這個對照詞組，但礙於中文斷詞邊界，我們沒辦法選到該詞組。

基於以上的原因，我們決定將國客語對照辭典中的中文詞欄位，加入到中文斷詞辭典。而這些新加入的中文詞，要決定出它們的詞頻大小，因為這兩個語料規模相差非常龐大，我們的中文斷詞辭典總詞頻數高達 462729801 個詞，而客語訓練語料僅有 47079 個詞。因此我們將中文斷詞辭典的平均詞頻取 Log 以 2 為底乘上 15(經實驗得到，15 有最佳的正確率)，再與國客語對照辭典的詞頻相乘，得到該詞新的詞頻。步驟如下：

1. 統計出中文斷詞辭典原始的分佈，我們得到其平均詞頻為 464 次。

表十八、中文斷詞辭典資料分佈

| | |
|------|-----------|
| 總詞數 | 995642 |
| 總詞頻 | 462729801 |
| 平均詞頻 | 464 |

2. 我們以下列公式，計算出每個要加入中文斷詞辭典中的中文詞，其新詞頻 $C(W_i)^*$ ，其中 $W_i \in (Woyd Length > 1)$ ：

$$C(W_i)^* = \log_2(464) * 15 * \lceil C(W_i) + 1 \rceil \quad (1)$$

一個國客對照詞「植物/植物」，其詞頻轉換的例子：

表十九、國客語對照辭典的中文詞詞頻換算

| | |
|----------------|--|
| 中文 | 植物 |
| 客語 | 植物 |
| 未含詞性特徵的詞頻 | 3 |
| 新詞頻 $C(W_i)^*$ | $Ceil(\log_2(464) * 15 * \lceil 3 + 1 \rceil) = 532$ |

依照上述方法，產生以下新增候選列表：

表二十、客語詞新詞頻候選表

| 中文詞 | 詞頻 |
|-----|-----|
| 植物 | 532 |
| 雕刻 | 266 |
| 信仰 | 443 |
| 八仙 | 355 |
| 筵席 | 178 |
| | |

3. 將步驟 2 產生的新增候選列表的結果，加入中文斷詞辭典，若有相同的中文詞，則其詞頻相加。

最後我們評估修改前與修改後的中文斷詞系統，其正確率的差異。測試為外部測試，使用「中研院中文平衡語料庫 3.0」。如表二十一：

表二十一、中文斷詞辭典修改前後比較

| | Precision | Recall | F-Measure |
|-----|-----------|--------|-----------|
| 修改前 | 97.16% | 96.21% | 96.69% |
| 修改後 | 97.15% | 96.16% | 96.65% |

可看到 F-Measure 略低 0.04%，但改善了以下問題：

表二十二、加入客語詞後的中文斷詞辭典的改善

| | |
|----------|------------|
| 中文句子 | 泥鯪滑溜溜， |
| 原中文斷詞 | 泥鯪/滑/溜溜/， |
| 改善後中文斷詞 | 泥鯪/滑溜溜/， |
| 原中文翻客語 | 鯪鯪仔/滑/溜溜/， |
| 改善後中文翻客語 | 鯪鯪仔/滑溜溜仔/， |
| 正確答案 | 鯪鯪仔/滑溜溜仔/， |

原本的斷詞系統無法將「滑溜溜」判斷出來，修正辭典後已能正確斷出，並轉為客語詞「滑溜溜仔」。

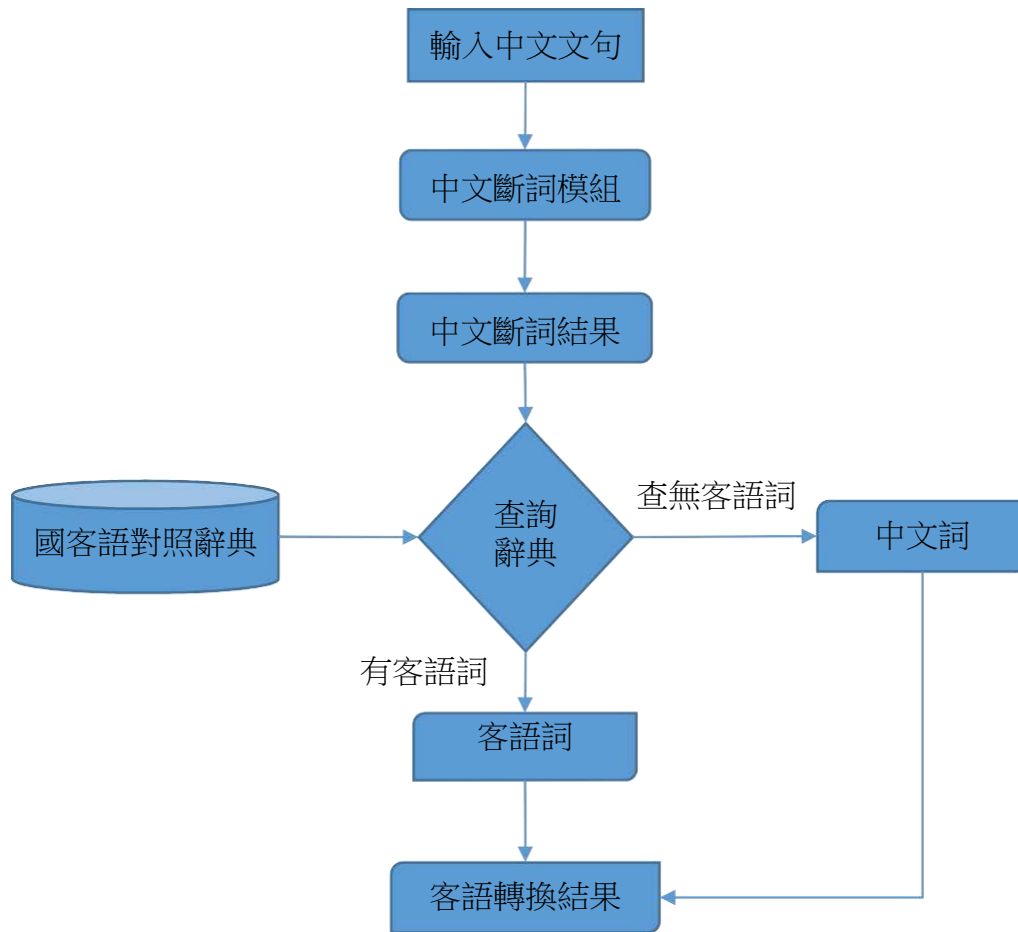
4.3 客語斷詞方法

4.3.1 中文斷詞搭配國客語對照辭典的直翻法

本方法的流程：

1. 透過中文斷詞系統得到中文斷詞及標詞性結果。
2. 查找國客語對照辭典，找出對應的客語詞，並選擇資料 ID 排序第一位者，並將中文詞轉換成該詞。
3. 若步驟 2 時查不到對應的客語詞，則沿用中文斷詞的結果。

圖四為此斷詞法的流程：



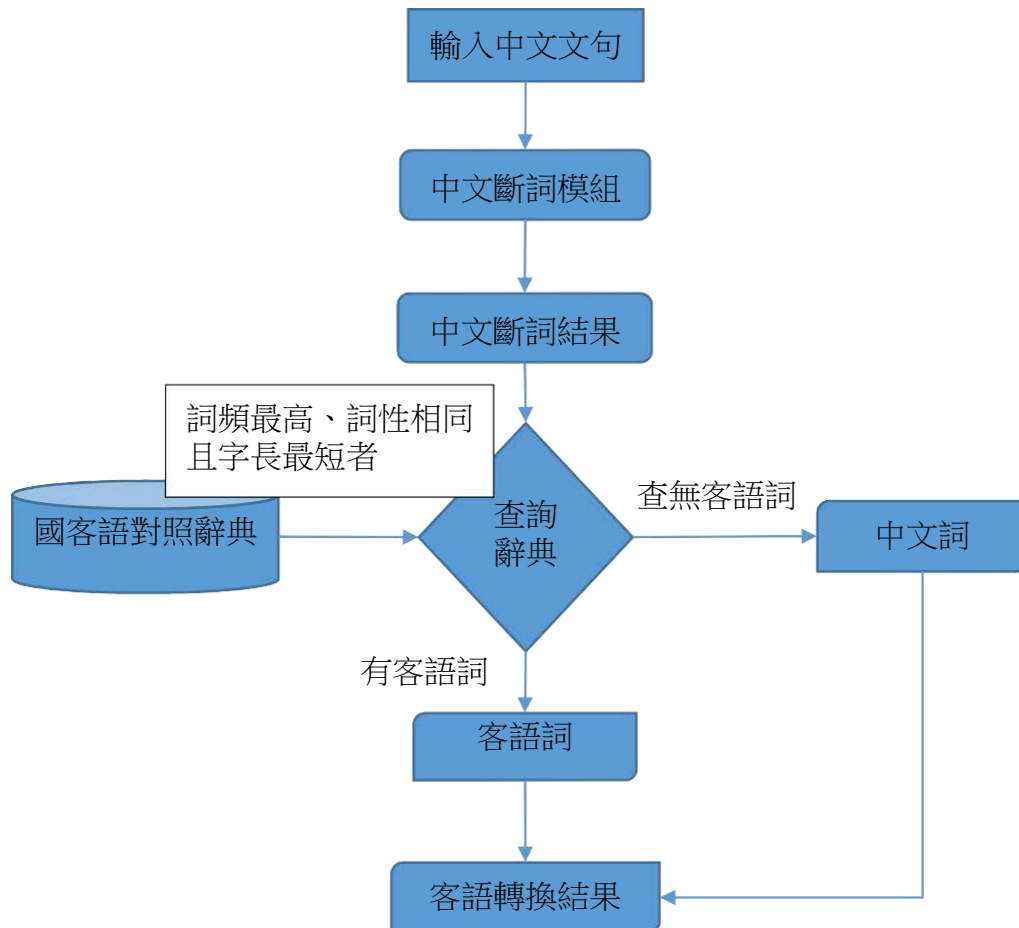
圖四、中文斷詞搭配國客語對照辭典的直翻法流程圖

4.3.2 中文斷詞搭配客語詞頻的直接翻譯法

以下是本方法的程式流程：

1. 透過中文斷詞系統得到斷詞及標詞性的結果。
2. 查找國客語對照詞典，找出對應的客語詞。
3. 挑選詞性與詞性標記結果相同者且詞頻最高、字長最短者為結果，若都找不到則以中文詞為結果。

圖五為本方法流程：



圖五、搭配客語詞頻的直接翻譯斷詞法

4.3.3 中文斷詞搭配客語 Uni-gram 及 Bi-gram 語言模型的混合式分數算法

本斷詞法，是先將輸入的中文句子，以中文斷詞模組做斷詞，得到一個確定的中文詞邊界後，再將這些中文詞查詢客語辭典，找出所有的可能候選詞。若找不到客語詞，則以僅有的中文詞當候選詞。最後將這些候選詞建立成一個所有斷詞路徑的有向圖，再以最短路徑演算法配合 Uni-gram 及 Bi-gram 混合式的分數算法，找出一條分數最佳的斷詞序列。

分數的計算方法：

本斷詞方法的分數計算方式，是混合式 Mix-gram(Uni-gram+Bi-gram)分數算法，如下列式子：

$$Seoye(\langle S \rangle, W_1, W_2, W_3, \dots, W_n) = \alpha \gamma \min - \{ \log_e [P(W_1 | \langle S \rangle) * P(W_1)] + \sum_{i=2}^n \log_e [P(W_i | W_{i-1}) * P(W_i)] \} \quad (2)$$

其中 $P(W_i | W_{i-1})$ 可利用 maximum likelihood estimation(MLE)來計算：

$$P(W_i | W_{i-1}) \approx \frac{C(W_{i-1}, W_i)}{C(W_{i-1})} \quad (3)$$

若遇到訓練資料 $C(W_{i-1}, W_i) = 0$ 時，我們將 Bi-gram 機率以為 α 值代替，轉換如下：

$$P(W_i | W_{i-1}) = \frac{C(W_{i-1}, W_i)}{\sum_{W \in V} C(W_{i-1}, W)}, \text{ if } C(W_{i-1}, W_i) > 0 \quad (4)$$

$$\alpha, \text{ if } C(W_{i-1}, W_i) = 0$$

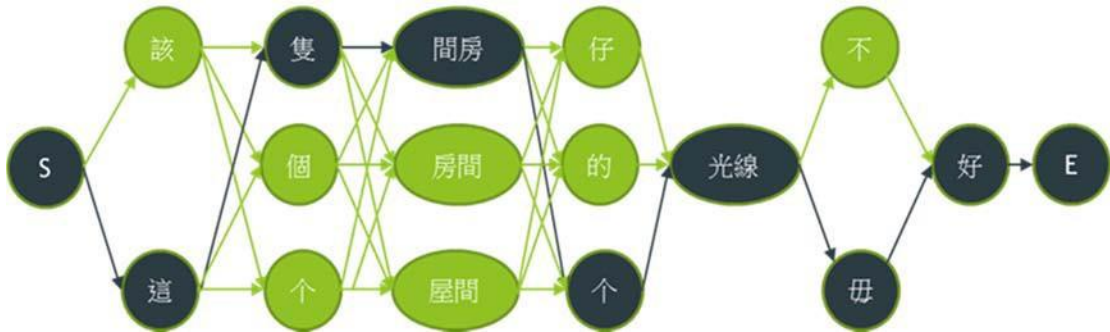
其中 $P(W_i)$ ：

$$P(W_i) = \frac{1 + C(W_i)}{v + (\sum_{j=1}^v C(W_j))} \quad (5)$$

$\sum_{j=1}^v C(W_j) = 47079$ ，為客語 Uni-gram 語言模型中的總詞頻數。

$V = 8931$ ，為 Uni-gram 語言模型的總 type 數(詞數)。

$\alpha = 10^{-3}$ ，為代替當 $C(W_{i-1}, W_i) = 0$ 時，用以替代 $\frac{C(W_{i-1}, W_i)}{\sum_{W \in V} C(W_{i-1}, W)}$ 的值，這是經由實驗得到，我們測試了 $10^{-1}, 10^{-2}, 10^{-3}, \dots, 10^{-8}$ 等值。



圖六、中文斷詞搭配客語 Uni-gram 及 Bi-gram 語言模型的混合式分數算法，範例。

以上圖為例，如要計算上條斷詞序列 $Seoye(S, \text{這隻間房個, 光線毋好})$ ，其 Mix-gram 的分數計算公式如下：

$$Seoye(S, \text{這隻間房個, 光線毋好}) = [P(\text{這}/S) * P(\text{這})] * [P(\text{隻}/\text{這}) * P(\text{隻})] * [P(\text{間房}/\text{隻}) * P(\text{間房})] * [P(\text{個}/\text{間房}) * P(\text{個})] * [P(\text{光線}/\text{個}) * P(\text{光線})] * [P(\text{毋}/\text{光線}) * P(\text{毋})] * [P(\text{好}/\text{毋}) * P(\text{好})]$$

5. 系統效果評估

5.1 評估方法

本測試所輸出的句子，是中文未斷詞的句子，輸出客語斷詞及標詞性的結果。目前只測試單純客語分詞的效能。效能評估的計算方法，我們使用精確率(Precision)、召回率(Recall)、以及 F-分數(F-score)來評估系統的效能，這三種方法的定義如下所示：

$$\text{率} = \frac{\text{系統正確斷出的詞數}}{\text{系統斷出的總詞數}} \quad (6)$$

$$\text{召回率} = \frac{\text{系統正確斷出的詞數}}{\text{標準答案的總詞數}} \quad (7)$$

$$F\text{分數} = \frac{4 \times \text{精確率} \times \text{召回率}}{\text{精確率} + \text{召回率}} \quad (8)$$

基於下列原因，除了上述的評估方法外，我們進一步運用編輯距離演算法(Levenshtein Distance)，評估轉換後的客語句子相似度(Similarity)。

原因 1：

中文翻客文的處理，是一個中文詞對多種可能客語字詞的問題，且客語常有一意多詞的問題，如：中文「收到」客語可翻成「收」或「收着」，若正確答案標記為「收着」，但系統輸出為「收」，在斷詞評估角度看，兩者詞不同，正確率為 0，但以相似度來看，兩者僅差一個字，相似度仍有 50%。

因此，除了斷詞邊界的評估標準外，其字串轉換後的相似度也可以是一個評估效能的方法。且客語文句的文法結構與中文幾近相同，沒有翻譯後文法結構對齊的問題，因此可用此方法計算其相似度。

原因 2：

本論文的斷詞系統，是用於客語語音合成系統中，因此，翻譯後的字串，距離標準答案的相似度越高，能將字念對的可能性也越高。

此演算法計算兩個字串 A、B 間，由字串 A 轉換成字串 B 的最小編輯距離(Insertions, Deletions 或 Substitutions)，計算方式如下：

$$D(i, j) = \min \begin{cases} D(i-1, j) + \text{InsertCost}(target_i) \\ D(i-1, j-1) + \text{SubstituteCost}(source_i, target_i) \\ D(i, j-1) + \text{DeleteCost}(source_i) \end{cases} \quad (9)$$

其中：

$$\text{SubstituteCost} = \begin{cases} 0 & \text{if } target[i] = source[j] \\ 1 & \text{otherwise} \end{cases}$$

$$\text{InsertCost} = 1$$

$$\text{DeleteCost} = 1$$

並由以下公式，將距離轉換成 0 到 1 之間的值，即為 A、B 字串間的相似度：

$$\text{Similarity}(A, B) = 1 - \frac{D(A, B)}{\text{Max Length}(A, B)} \quad (10)$$

5.2 測試語料

本實驗所使用的客語斷詞訓練與測試語料，是採用客委會四縣腔初級及中高級語料中的例句，共 6640 句。我們將這份語料，分為語料 A 及語料 B，語料 A 有 4196 句，語料 B 有 2444 句，並再將 B 語料依編號順位分為 4 份，B1-B4，每份 611 句。

而這些標記完成的資料，再經過一次經人工篩選後，確認有效句數為：訓練語料 4018

句，測試(B1-B4)語料共 1282 句，我們使用語料的分佈如表所示：

表二十三、客語語料的使用分佈

| | 訓練 | 測試 |
|----|-------|-------|
| 句數 | 4018 | 1282 |
| 詞數 | 45304 | 17646 |
| 字數 | 65572 | 25478 |

5.3 評估結果與討論

5.3.1 中文斷詞搭配國客語對照辭典的直翻法

實驗 A：使用 Lo[6]的國客語對照辭典。 **實**

驗 B：使用 Wu[5]的國客語對照辭典。

實驗 C：使用本論文所建置的國客語對照詞典。

表二十四、直接翻譯斷詞法效能評估-內部測試

| | Precision | Recall | F- Measure | 字串相似度 |
|-------------|-----------|--------|------------|--------|
| 實驗 A | 68.41% | 69.12% | 68.76% | 73.29% |
| 實驗 B | 72.66% | 73.42% | 73.04% | 75.68% |
| 實驗 C | 75.02% | 75.80% | 75.41% | 78.36% |

由此實驗可得知，詞典的校正與詞目的增加，能顯著的改善客語斷詞系統的效能。

5.3.2 中文斷詞搭配客語詞頻的直接翻譯法

實驗 A：輸入中文文句，評估其中文轉客語的斷詞效能。 **實**

驗 B：輸入中文文句，評估其中文轉客語及詞性標記效能。

表二十五、中文斷詞搭配客語詞頻的直接翻譯法-斷詞及詞性標記內外部測試結果

| | | Precision | Recall | F- Measure | 字串相似度 |
|-------------|-----------|-----------|--------|------------|--------|
| 實驗 A | 訓練 | 88.16% | 87.48% | 87.82% | 91.08% |
| | 測試 | 80.32% | 79.08% | 79.69% | 83.68% |
| 實驗 B | 訓練 | 87.46% | 86.78% | 87.12% | - |
| | 測試 | 79.90% | 78.66% | 79.27% | - |

以實驗結果來看，顯示客語詞頻的應用能顯著的提升選詞的正確率。但外部測試的正確率仍偏低，已接近未使用詞頻特徵的結果。此情況的原因，是因為用來統計客語詞頻的語料仍不足，因此造成統計資料稀疏的問題。但從實驗結果的特性觀察到，若能持續的增加客語語料、建置出更多的客語詞頻，能提升中文轉客文系統的效能。

而就目前的窘況而言，提升正確率的方法，僅能靠規則法(Rules-base)，如找出客語的構詞規則和文法規則，來提升客語斷詞的正確率。客語構詞部份，是下一階段即將進行的工作。

5.3.3 中文斷詞搭配客語 Uni-gram 及 Bi-gram 語言模型的混合式分數算法

實驗 A：輸入中文文句，評估其中文轉客語詞的斷詞效能。

實驗 B：輸入中文文句，評估其斷詞及詞性標記效能。

表二十六、中文斷詞搭配客語 Uni-gram 及 Bi-gram 語言模型的混合式分數算法

| | | Precision | Recall | F- Measure | 字串相似度 |
|-------------|-----------|-----------|--------|------------|--------|
| 實驗 A | 訓練 | 94.46% | 93.73% | 94.10% | 96.17% |

| | | | | | |
|------|----|--------|--------|--------|--------|
| | 測試 | 80.78% | 79.53% | 80.15% | 84.04% |
| 實驗 B | 訓練 | 93.96% | 93.24% | 93.60% | - |
| | 測試 | 80.37% | 79.11% | 79.74% | - |

加入 Bi-gram 後，內部測試的效能有顯著的提升，但外部測試僅略升 0.46%。原因是因為客語語言模型資料稀疏的關係，許多 Bi-gram pattern 未出現在客語 Bi-gram 語言模型中，或就算出現在 Bi-gram 語言模型中，也不符合句子實際的狀況，這是資料稀疏的問題。

而在下一階段的工作中，我們將要增加客語構詞規則以及持續增加客語語料，客語構詞包括了：重複、附加、附合、合併 等構詞規則。

6. 結論及未來工作

本論文針對中文轉客文轉音系統(Hakka Text-to-Speech System, HTTS)中的客語斷詞處理，已提出一個基礎的研究架構。但因為目前客語電子語料有嚴重不足的問題，對於本論文所探討的主題而言，是一項非常艱困的挑戰。研究之初，我們並沒有客語斷詞語料可使用，僅有少量一句句未處理的國客語對照文句。因此，我們投入了大量時間在客語語料的標記、建置及辭典的校正、標音。我們也持續的從國小客語教材、客語朗讀比賽文章 等電子文本中，人工建置出更多的客語斷詞語料及客語新詞。而目前用來測試的語料，皆是我們自行建置的，所以我們也非常需要更多不同來源的客語斷詞語料，做更公證客觀的效能評估、比較。

本論文在客語斷詞方法方面，不同於過去的架構，我們提出了使用客語 Uni-gram 及 Bi-gram 語言模型的混合式斷詞序列分數算法，可看見在內部測試的評估上，顯示在語料充足的情況下，將會有不錯的斷詞表現。但詞性標記的正確率僅能做為參考，因為我們自己標記產生的標準答案，其詞性大部份都是依照中文斷詞系統給出的詞性為主，除離譜的錯誤外，很少當下進行修正。因此，語料的斷詞、詞性標記資訊，仍需要經專家再一次做更嚴謹的校正。

本論文提出混合型的 N-Gram 序列分數算法，搭配中文斷詞模組及動態規劃演算法的客語斷詞方法。在嚴重資料稀疏的客語語料下，對中文轉客文的外部測試精確率有 80.78%，內部測試有 94.46%。相較於傳統中文詞直翻客語詞的方法，已獲得提升。相信未來持續增加客語語料的規模後，使用本論文所提出的方法，效能會有更顯著的提升。

客語斷詞的應用層面極廣，不僅止使用於我們中文轉客文轉音系統中的文句分析模組，還可獨立用於客語的數位學習、客語文句處理、客語語音辨識 等領域。本論文提出的研究方法，應能提供未來客語斷詞相關研究做為基礎與參考。

我們接下來要進行的工作有：

1. 持續擴充國客語對照辭典。
2. 加入客語構詞規則。
3. 最佳化語言模型平滑化問題，如：Good-Turing Katz、Kneser-Ney。
4. 持續標記、建置客語語料。

致謝

客家委員會提供給本論文實驗用之客語認證語料以及獎助部份經費，特此致謝。

參考文獻

- [1] 蔡依玲，基於隱藏式馬可夫模型之客語文句轉語音系統，國立交通大學電信工程所碩士論文，2009。
- [2] 鍾屏蘭、江俊龍，學術研究基礎建置暨客家文化研究計畫，屏東教育大學客家文化所計畫成果報告書，2009。
- [3] 羅永聖，結合多類型字典與條件隨機域之中文斷詞與詞性標記系統研究，國立台灣大學資訊工程所碩士論文，2008。
- [4] 林千翔，*Chinese Word Segmentation using Specialized HMM*，國立中央大學資訊工程所碩士論文，2005。
- [5] 吳俊毅，線上客語語音合成系統中產生韻律訊息之研究，國立中興大學資訊科學與工程所碩士論文，2010。
- [6] 羅丞邑，以資料探勘之技術解決線上客語語音合成系統中多音字發音歧義之研究，國立中興大學資訊科學與工程所碩士論文，2011。
- [7] 江昶毅，應用多種特徵的中文斷詞及詞性標記方法，國立中興大學資訊科學與工程所碩士論文，2010。
- [8] 李雪貞，客語語音合成之初步研究”，國立臺灣科技大學資訊工程所碩士論文，2002。
- [9] 賴亦傑，“應用多詞及多詞性語言模型的中文斷詞及詞性標記方法”，國立中興大學資訊科學與工程所碩士論文，2011。
- [10] 客委會出版，客語能力認證基本詞彙-中級、中高級暨語料選粹四縣版上冊。
- [11] 客委會出版，客語能力認證基本詞彙-中級、中高級暨語料選粹四縣版下冊。
- [12] 林東毅，客語文句翻語音系統之實作，國立交通大學電信工程所碩士論文，2007。
- [13] 黃豐隆，線上國客雙語有聲詞典建置之研究，全國計算機會議(NCS-2009)，台灣，2009。
- [14] 黃豐隆，國客雙語有聲地圖社群系統，聯合大學資工所客委會計畫成果報告書，2013。
- [15] Nianwen Xue, *Chinese Word Segmentation as Character Tagging*, Computational Linguistics 2003 February, Vol. 8, No. 1, pp.29-48.
- [16] Guhong Fu and K.K. Luke., *A Two-Stage Statistical Word Segmentation System for Chinese*, Proceeding of The Second SIGHAN Workshop on Chinese Language Processing 2003, Vol. 17, pp.156-159.
- [17] Keh-Jiann Chen and Shing-Huan Liu, *Word Identification For Mandarin Chinese Sentences*, Proceedings of COLING 1992, pp.101-107.
- [18] Jian-Yun Nie, Marie-Louise Hannan, and Wanying Jin, *Unknown Word Detection and Segmentation of Chinese Using Statistical and Heuristic Knowledge*, Communications of COLIPS 1995, Vol. 5, pp.47-57.
- [19] Andi Wu, Zixin Jiang, *Word Segmentation In Sentence Analysis*, International Conference on Chinese Information Processing in Beijing China 1998, pp.169-180.
- [20] Jian-feng Gao, Mu Li, and Chang-Ning Huang, *Improved Source-Channel Models for Chinese Word Segmentation*, the 41st Annual Meeting on Association for Computational Linguistics 2003, Vol. 1, pp.272-279.

Vj g'4236"Eqphgtpeg"qp"Eqo r wcvkqpcnNkpi wku'cpf "Ur ggej "Rtqegu'kpi "
 TQENRPI "4236."rr09:/: : "
 © Vj g'CuuekcvkqpfqtEqo r wcvkqpcnNkpi wku'cpf Ej kpgug'Nkpi wci g'Rtqegu'kpi "

基於發音知識以建構頻譜 HMM 之國語語音合成方法

A Mandarin Speech Synthesis Method Using Articulation-knowledge Based Spectral HMM Structure"

古鴻炎*、賴名彥*、洪尉翔*、陳彥樺*

Hung-Yan Gu, Ming-Yen Lai, Wei-Siang Hong, and Yan-Hua Chen

摘要

在有限語料的情況下，本論文提出一種 J O O "的結構設計，來掌握各個語音單元之文脈相依的頻譜特性，以便改進合成語音的流暢度。此外，在決策樹之文脈分群方法之外，我們依據音素的發音知識，來作文脈分群而大幅降低文脈組合數量。為了評估所提出的 J O O "結構，我們使用三種不同的 J O O "結構方式去建造對應的國語語音合成系統，以作相互的比較。在這些系統裡，使用的韻律參數值是一樣的，都是使用之前研究的 CPP "模組來產生；但是頻譜係數則是使用各自的 J O O "模型來產生；至於信號波形的合成，則都是使用之前研究的基於諧波加雜音模型*J P O -的信號合成模組。聽測實驗的結果顯示，使用本論文提出的 J O O "結構所合成出的語音，比用其它 J O O "結構所合成的明顯地更為流暢；此外，依據錄音語句與合成語句之間的平均頻譜距離的量測結果，也顯示本論文的 J O O "結構，比其它 J O O "結構更能夠降低頻譜距離。

關鍵詞：語音合成、J O O "結構、發音知識、頻譜流暢度、離散倒頻譜係數

Abstract"

Kp"vj ku"r cr gt. "c"pgy "J O O "utwewtg"ku"rtqr qugf "v"y qtm'y kj "c"rko kgf "v'cl'k'kpi " eqtr wu"kp"qtf gt"v"q"qdv'cl'k"ko r tqxgf "u{p'vj g'v'e/ur ggej "h'w'g'p'e{ 'O'Ur g'v'e'c'n' h'w'g'p'e{ "ku" ko r tqxgf "d'g'c'w'g"vj ku"J O O "utwewtg"ecp"o qf gn'vj g"eqp'v'z'v'f gr g'p'f'g'p'v'ur g'v'e'c'n' ej ctce'v'g't'k'u'k'u"qh'c"ur ggej "w'p'k'o' Kp" c'f'f'k'k'q'p." k'p'u'v'g'c'f"qh' w'k'p'i "c" f'g'ek'k'q'p"t'g'g"v"q" e'n'w'g't"eqp'v'z'w."vj g"hp'q'y r'g'f' i'g"qh'r'j q'p'g'o'g"ct'v'e'w'e'v'k'q'p"ku"dc'ug'f"v"q"e'n'w'g't"eqp'v'z'w" c'p'f" t'g'f'w'eg"vj g" g'p'q't'o'q'w'u" s'w'c'p'v'k'v' "qh' eqp'v'z'v' eqo d'k'p'c'v'k'p'u'o' V'q" g'x'c'n'w'c'v'g"vj g" r' t'q'r'q'u'g'f"J O O "utwewtg."y g"eqp'ut'w'e'v'vj t'g'g"O'c'p'f'c't'k'p"ur ggej "u{p'vj g'k'u"u{u'v'g'o' u" g'c'ej "w'g'u"q'p'g" f'k'h'g't'g'p'v'J O O "utwewtg" h'q't" eqo r c't'k'q'p'u'o' Kp" vj g'g' u{u'v'g'o' u."vj g" r' t'q'u'q'f'k'e"r'c't'c'o' g'v'g't'u'c't'g"cm'i' g'p'g't'c'v'g'f"y'k'j"u'c'o'g"CPP"o'q'f'w'g'u"u'w'f'k'g'f"r' t'g'x'k'q'w'u'k' "

*國立臺灣科技大學資訊工程系 Fgr ctvo gpv'qhEqo r wgt'Uelgpeg'cpf 'Kphqto cvkq'Gpi kpggtkpi ."

PcvkqpcnVcky ep'Wpkxgtuks{ 'qh'Uelgpeg'cpf 'Vgej pqrni { "

G'o ckn' }i vj {.'O: 837296.'O 32337257.'O 324372227; B o c'k'p'w'u'v'g'f' w'q'y "

"
"

dw' 'y g' ur gevtn'eqghlekpwu'ctg' i gpgtcv'gf' y kj' 'f khtg'gpv' J O O 'cf qr v'gf' d{ 'ku'
eqttgur qpf kpi 'u{ ugo O'Cu'v' 'y g' u{ p'v' guku'qh'uki p'cn'y cxghqto . 'y g'uki p'cn'o qf gn'
j cto qple'r n'u'p'q'lug'o qf gn' *J P O + 'u'w'f k'gf' 'r t'g'x'k'q'w'u' 'ku'eqo o q'pn' 'cf qr v'gf' 'k'p' 'y g'
y' tgg' u{ ugo u' O' C'eeq'tf k'pi 'v' 'y g' t'g'u'w'u' q'h' r'k'ug'p'k'pi 'v'g'u'u. 'y g' ur g'gej' 'u{ p'v' g'uk' gf' d{ '
y' g' u{ ugo ' w'ul'pi 'y g' r' t'q'r' q'ug'f' 'J O O ' u't'w'e'w't'g' 'ku' k'p'f' g'g'f' 'o q't'g' 'h'w'g'p'v' 'y' c'p' 'y' g'
ur g'gej' gu' u{ p'v' g'uk' gf' d{ 'y g' q'y g' t' 'y q' u{ ugo u' O' k'p' 'c'f' f'k'k'q'p. 'cx'g't'c'i' g' ur gevtn'
f' k'nc'p'eg'u'ctg'o g'c'u'w't'g'f' 'd'g'y' g'g'p' 't'g'e'q't'f' g'f' 'u'g'p'v'g'p'eg'u' c'p'f' 'u{ p'v' g'v'le' 'u'g'p'v'g'p'eg'u' V'j' g'
t'g'u'w'u' u'j' q'y 'y' c'v' 'y' g' 'J O O ' u't'w'e'w't'g' r' t'q'r' q'ug'f' 'j' g't'g' 'c'n'u'q' 'q'd'v'c'k'p'u' u'o c'm'g't' 'c'x'g't'c'i' g'
ur gevtn' f' k'nc'p'eg' 'y' c'p' 'y' g' q'y g' t' 'y q' 'J O O ' u't'w'e'w't'g'u' O'

Keywords: Ur g'gej' 'U{ p'v' g'uk' 'J O O ' U't'w'e'w't'g. 'C't'v'k'w'r'v'k'q'p' 'M'p'q'y' r'g'f' i' g. 'Ur gevtn'
H'w'g'p'e' { . 'F' k'u'e't'g'v' 'E'g'r' u't'c'n' 'E'q'g'h'le'k'p'w'u' O'

1. 緒論

近年來許多研究者早已利用隱藏式馬可夫模型 *k'f' f' g'p' 'O' c't'n'q'x' o' q'f' gn' 'J O O + 來建造語音單元 *如音素、音節等+之頻譜演進 *ur gevtn' 'r' t'q'i' t'g'u'k'q'p' +模型 *q'uj' k'o' w't'c' 'et al.' '3; ; ; : 'g'p' 'et al.' '4229+ [c'p' 'et al.' '422; ; : 'I' w' 'et al.' '4232+ , 之後在合成一個語句時, 就會使用訓練得到的 J O O "模型來產生一序列的頻譜特徵向量, 然後使用所產生的頻譜特徵向量序列去合成出語音信號。使用 J O O "來作語音信號的合成, 通常能夠獲得增進的可理解性 *k'p'v'g'n'k'i' k'd'k'k'v' +與流暢性 *h'w'g'p'e' { +。更好的是, V'q'm'f' c' "等人基於 J V M' *J O O " 'v'q'q'n' 'n'k'u' +所發展的 J V U "語音合成軟體 *g'p' 'et al.' '4229+ , 提供公開的程式原始碼、並且可供下載, 所以研究語音合成時, 使用 J V U "能夠減少許多時間與氣力。不過, 當未使用全域變異數 *i' m'q'd'c'n' 'x'c't'l'c'p'eg. 'I' X +匹配時, J V U "軟體所產生的頻譜包絡會發生過於平滑的現象, 使得合成出的語音變得悶悶的 *o' w'h'g'f' + *V'q'f' c' ('V'q'm'f' c. '4227+。

在本論文中, 我們並不打算沿用、修改 J V U "的程式碼, 因此必須自行發展 J O O "模型訓練的程式、及頻譜特徵向量序列之產生程式, 這樣的決定是因為, 我們想要研發一個具有彈性 *h'g'z' k'd'k'k'v' +的語音合成系統, 能夠容易地擴增額外的功能。例如音色轉換 *k'o' d't'g' 't'c'p'u'h'q'to' c'v'k'p' +之功能, 能夠把合成語音的音色從成年女性轉變為男孩 *I' w' ('V'uck' '4235+ ; 另一項預計擴增的功能則是, 同步地播放合成出的語音信號及其對應的拼音符號, 此功能可應用於人形機器人上, 讓語音發聲和嘴型同步。除了擴增額外功能的原因之外, 我們與其他研究者使用 J V U "的經驗 *J' u'k' 'et al.' '4232+ , 發現 J V U "所產生的國語語音的基週軌跡 *r' k'e'j' " 'e'q'p'v'q'w't'u' +聽覺上並不令人滿意, 因此本論文建造的國語語音合成系統, 就決定使用不同的方法來產生各個音節的基週軌跡。

在先前的一次研究中 *I' w' 'et al.' '4232+ , 我們曾嘗試去建立音節單位的 J O O "模型, 以掌握音節內的頻譜演進 *ur gevtn' 'r' t'q'i' t'g'u'k'q'p' +方式, 但是此 J O O "模型所合成出的語音信號並不够流暢, 在音節邊界處的頻譜不連續 *ur gevtn' f' k'u'e'q'p'v'k'p'w'k'g'u' +情形經常會被聽出來, 我們認為發生頻譜不連續的原因是, 設定的語音單位 *音節 *太大, 使得一個音節 Z "和前後音節所組合出的不同文脈數量非常龐大, 以至於無法為該音節 Z "建立良好的文脈

相依 J O O "模型。因此，本論文決定使用較小的語音單位，即聲母和韻母，接著根據標記檔裡記載的拼音符號序列，去建構文脈相依的 J O O "模型，並且依據發音知識來對 J O O "模型作分組訓練。J O O "結構的設計方式將會在第二節詳細說明。

我們建造的國語語音合成系統，其合成階段之處理流程如圖 3"所畫，區塊*c+進行輸入文句之分析；區塊*d+使用類神經網路**Ct v h e l c n " P g w t c n " P g y q t m " C P P*+模組來產生各個音節的基週軌跡參數和時長**f w t c v k p*+值，詳細作法可參考**1 w (" Y w " 4 2 2 ; +*；區塊*e+依據第二節介紹的挑選方法來為各個聲韻母選出對應的 J O O "模型；區塊*f+採用 *V q m w f c*"等人提出的方法**1 q u j k o w t c " e t a l . " 3 ; ; +*，去計算 J O O "各狀態分配到的時長音框數；區塊*g+採用先前研究提出的加權線性內插法**1 w " e t a l . " 4 2 3 2 +*，去產生各音框的頻譜係數；區塊*h+依據基週軌跡參數去計算出各個有聲音框的音高頻率**r k e j " h t g s w g p e { +*值；最後，區塊*i+使用諧波加雜音模型**J c t o q p l e " r n w u " p q k u g " o q f g n " J P O +*，去合成出信號波形，詳細作法可參考**1 w (" V u c k " 4 2 3 5 +*。

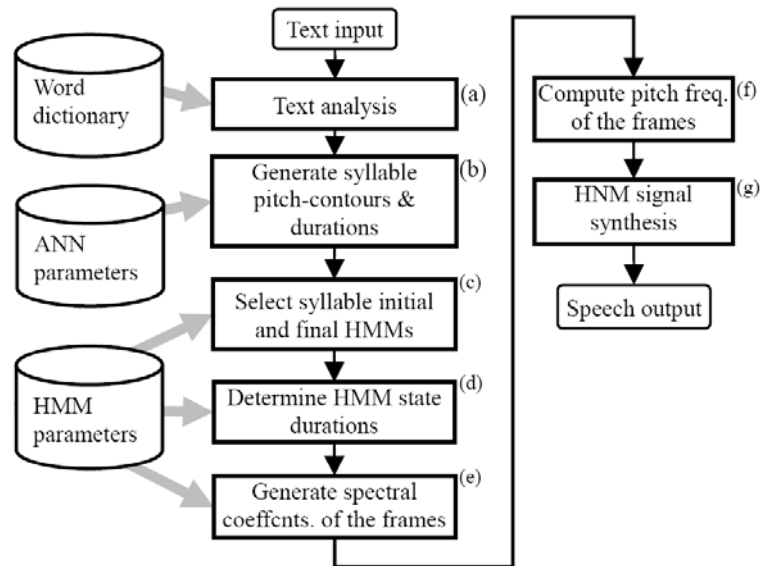


圖 1. 合成階段之主要處理流程

2. 聲韻母分類及文脈組合

在語音辨識領域，基本的 J O O "結構就如圖 4"所示之左至右**r g h v " q " t k i j v*結構，然而此結構並未處理左右兩邊界的文脈相依關係，因此其效能*辨識率+會顯現衰退的情形。

在傳統上，一個國語音節被分割成聲母*開頭子音+和韻母*結尾母音群+兩個部分，若採取聲、韻母作為語音單位，則一個語音單位的左右文脈之組合量可以比音節單位時的文脈組合量大幅減少。但是採用聲、韻母作為語音單位時，可能組合出的文脈量仍是非常龐大，如韻母的前後文脈組合量約有十二萬四千個，詳細數目是**4 3 - 5 9 + " 5 9 " " 4 3 - 5 9 +*

"
"

?"346.68: , 其中 43"表示 43"種聲母, 59"表示 59"種韻母, 因此有需要再對聲、韻母作進一步的分類。

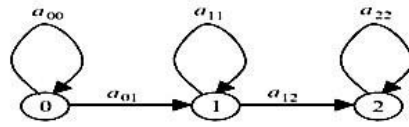


圖2. 左至右HMM結構

在一個聲母的左邊, 可能遇到前一音節的韻母或靜音, 在此我們把韻母尾端之可能發音口形分類成 33"類, 詳細分類方式如表 3"所列*含靜音öukrö+。相對地, 在一個聲母的右邊只會遇見韻母, 在此我們把韻母開頭之可能發音口形分類成:"類, 詳細分類方式如表 4"所列。當把聲母前後可能遇到的韻母發音口形作了分類之後, 再考慮國語共有 43"個聲母, 由此可推算出聲母可能組合出的文脈數量是*33- 3+ 43 :?4.238"個。

表1. 韻母結尾之發音口形分類

| | | | | | | | | | | | | |
|----------|----|----|----|----|----|----|-----|-----|-----|-----|------|------|
| ƙf gz" | 2" | 3" | 4" | 5" | 6" | 7" | 8" | 9" | : " | ; " | 32" | 33" |
| I guwtg" | c" | q" | ə | g" | k' | w" | {w" | kk' | gt" | p" | pi " | ukl' |
| ercuu" | | | | | | | | | | | | |

表2. 韻母開頭之發音口形分類

| | | | | | | | | |
|----------|----|----|----|----|----|----|-----|-----|
| ƙf gz" | 2" | 3" | 4" | 5" | 6" | 7" | 8" | 9" |
| I guwtg" | c" | q" | ə | g" | k' | w" | {w" | kk' |
| Ercuugu" | | | | | | | | |

類似於前一段的敘述, 如果目前音節沒有聲母時, 則在此音節韻母的左邊, 將會遇到前一音節的韻母或靜音, 在此也把前一音節韻母尾端之可能發音口形分成 33"類, 如表 3"所示。另一種情況, 當韻母左邊遇到的是本音節的聲母時, 在此我們根據聲母的發音位置把可能遇到的聲母分成 8"類, 詳細分類方式如表 5"所列。舉例來說, ld1"r l"b l"hl"的發音位置皆在嘴唇, 所以它們都被分類至ödü類別; 同樣道理, fl l"ml"pl"hl"的發音位置皆在齒槽, 所以把它們都分類至öfö類別。

表3. 聲母依發音位置之分類

| | | | | | | |
|------------|----|----|----|----|----|-----|
| ƙf gz" | 2" | 3" | 4" | 5" | 6" | 7" |
| Eqpuqpcpv" | d" | f" | " | j" | l" | i " |
| Ercuugu" | ㄉ | ㄈ | ㄌ | ㄌㄐ | ㄌ | ㄨ |

考慮目前音節韻母的右邊可能遇到的語音單元, 第一種情況是, 後接的音節沒有聲母, 在此我們把後接音節的韻母開頭之可能發音口形分成;"類, 詳細分類方式如表 3"中的;"個類別, 但不包含lp和lp1; 第二種情況是, 後接的音節具有聲母, 在此我們把後

接音節聲母的可能發音口形分成 8 類，詳細分類方式如表 5 所列；第三種情況是，後接音節不存在（即目前音節是語句的最後音節），相當於後面銜接的是靜音。當把韻母前後可能遇到的聲、韻母發音口形作了分類之後，再考慮國語共有 59 個韻母，由此可推算出國語韻母可能組合出的文脈數量是 $3^3 \times 8 \times 3 \times 59 \times 8 \times 3 \times 32.878$ 個。如果我們依據前述

的文脈組合方式去建造 J O O 模型，則需要訓練的 J O O 模型會有 4.238×32.878 個，這意味著在準備訓練語料時，必須錄製數倍於 32.878 個音節數量的音節發音。然而準備大量的訓練語料需要付出昂貴的費用，這暗示使用前述之文脈相依 J O O 模型是不切實際的。對於此問題的解決辦法，先前研究者提出使用決策樹 *f gekukqp vtgg 來對 J O O 模型作分群，再對各群分別去訓練一個共用的 J O O，如此就可大幅降低所需的訓練語料數量，如 J VU 軟體就是採取此種作法。不過，本論文採取另外一種研究方向 *crrrtqcej，就是先依據發音知識來對聲、韻母作分類（這相當於對 J O O 作分群），再研究新的 J O O 結構之設計，以便解決前述的需求大量訓練語料的問題。

3. 建構新的 HMM 結構

如圖 5 所畫的是擁有 8 個狀態之韻母 J O O 的原本結構，此結構表示該韻母 J O O 是左右文脈相依的，符號 H [aZ] 表示此 J O O 是韻母 [] 的模型 *H 表示左右文脈相依，並且 EZ 表示韻母 [] 前接的文脈樣式 *eqpvz v\ r g，E \ 則表示 [] 後接的文脈樣式。根據表 3 和表 5 得知文脈樣式 EZ 共有 $34 - 8 \times 3$ 種；相對地，從表 3 除了 p 與 pi l 之外，表 5 得知文脈樣式 E \ 共有 $32 - 8 \times 38$ 種。此外，國語有 59 種韻母 []，所以國語韻母的文脈相依之 J O O 模型總計多達 32.878 個。

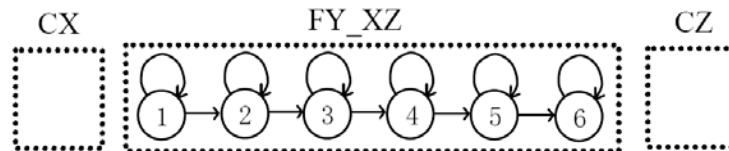


圖3. 左右文脈相依之HMM 模型結構

為了減少實作上需要投入的費用與人力（例如準備大量的訓練語料），因此我們嘗試對圖 5 的 J O O 模型 H [aZ] 去重作結構安排。我們的解決方法是，假設圖 5 中前面半數的狀態（即狀態 3、4、5）會相依於前接文脈 EZ，但是和後接文脈 E \ 不相關；類似地，我們也假設圖 5 中後面半數之狀態（即狀態 6、7、8）只相依於後接文脈 E \，但是和前接文脈 EZ 不相關。根據前述的兩項假設，圖 5 裡左右文脈相依之 J O O 模型 H [aZ]，就可以被分解成圖 6 與圖 7 之半段式 J O O 模型 I [aZ] 和 J [a \]，亦即我們要以半段式 J O O 模型 I [aZ] 和 J [a \] 之串接來取代 J O O 模型 H [aZ]。

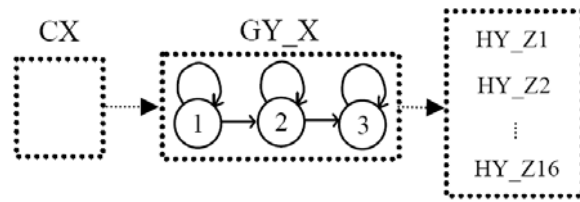


圖4. 韻母Y前半段之左文脈相依HMM結構

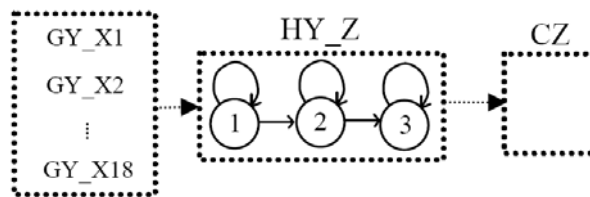


圖5. 韻母Y後半段之右文脈相依HMM結構

由於一個韻母前接的文脈被分類成 38 種文脈樣式，所以我們需要建立 38 個半段式 J O O 模型 I [a Z3 . I [a Z4 . ĩ . I [a Z3 : 如圖 7 裡列出，來掌握韻母 [的前半段部分，符號 I [之 I 表示半段式 J O O 之前半段。類似地，一個韻母後接的文脈被分類成 38 種文脈樣式，所以我們需要建立 38 個半段式 J O O 模型 J [a \ 3 . J [a \ 4 . ĩ . J [a \ 38 如圖 6 裡列出，來掌握韻母 [的後半段部分，符號 J [之 J 表示半段式 J O O 之後半段。如此，一個韻母需要建立的半段式 J O O 模型數量是 38 × 38 = 1444 個，而國語有 59 種韻母，所以總共需要建立的半段式 J O O 的數量是，1444 × 59 = 85196 個，3.47 萬個，比起 32.878 萬個左右文脈相依之韻母 J O O 少了許多。

關於國語 43 個聲母的 J O O 模型的建立，我們把前述之韻母 J O O 模型的假設套用進來。由於一個聲母前接的文脈被分類成 34 種文脈樣式，所以我們需要為一個聲母建立 34 個半段式 J O O 模型，來掌握該聲母的前半段部分；此外，一個聲母後接的文脈被分類成 38 種文脈樣式，所以我們需要為一個聲母建立 38 個半段式 J O O 模型，來掌握該聲母的後半段部分。如此，一個聲母需要建立的半段式 J O O 模型數量是 34 × 38 = 1292 個，而國語有 43 種聲母，所以總共需要建立的半段式 J O O 的數量是，1292 × 43 = 55576 個，5.5576 萬個，6.42 萬個，4.238 萬個左右文脈相依之聲母 J O O 少了許多。

4. 實驗評估

4.1 訓練階段

我們邀請了一位成年男性至隔音錄音室錄製 3.42 個語句的發音，這些語句的腳本 *uetkr v 是隨機地從無關的文章中挑選出，總計有 32.395 個音節，而錄音的取樣率為 44.272 J |。為了標記這些語句中各音節的拼音符號，我們首先使用 J VM 套件來作 hqtegf " crki po gpv" 處理，而得到初步的音節邊界之標記，然後以人工操作 Y cxgUwthgt" 軟體去更正錯誤的音節邊界標記。

我們將各個音節的音檔切割成一序列的音框，音框長度設為 734 個樣本點，而音框位移設為 34 個樣本點。每個音框經分析計算後擷取出 5 個頻譜參數 c_2, c_3, \dots, c_5 ，實際上是離散倒頻譜係數 *f kuetgyg" egr utcrn' eqgthlekgpw. " FEE+ " Ecrr 2" (" Oqwlpgu " 3; ; 8+, FEE" 係數的擷取方法，請參考我們先前的研究論文 *l w" (" Vuck" 422; +。此外，一個音框的週期性資訊也會被記錄在另一個維度中，即存入 c_5 ，如果一個音框被偵測出是週期性的，則設定 c_5 的值為 3，反之則設定 c_5 的值為 2。一個 J OO 經過訓練之後，我們就可以依據平均向量的 c_5 的值，來判斷各個 J OO 狀態是否為有聲 *xqlgef + 或無聲 *wpqxqlgef + 之狀態，然後就可為有聲的 J OO 狀態去產生基週軌跡。另外，頻譜參數的差分值也是有用的，所以我們把頻譜特徵向量增加 62 維，以儲存 FEE 係數的一階差分值。在訓練完一個 J OO 之後，除了記錄 J OO 的模型參數之外，也要記錄各個 J OO 狀態被駐留的音框個數之平均值與變異數。

在本論文中，一個韻母的前半與後半部分之半段式 J OO 模型皆是由 5 個狀態建造而成，並且狀態移轉方式都是由左至右，就如圖 6 和圖 7 所示。不過，對於一個聲母的半段式 J OO 模型，我們僅使用 4 個狀態去建造。關於高斯混合組件 *l cwulcp" o lzwtg" eqo r qpgpv 的數量，在每一個 J OO 狀態上，我們只設置一個高斯混合。對於半段式 J OO 之訓練，我們依據分段式 M' 中心 *ugi o gpvcn' M' o gcpu 演算法 *Tcdlpgt" (" Lwcpj . 3; ; 5+，去發展了訓練程式。

4.2 語音合成處理

本論文製作的國語語音合成系統，其合成階段之處理流程如圖 3 所示，區塊 *c 作 文句分析 *vgz v" cpcn(uku+，每次會從輸入的檔案讀取一個文句進來，然後以查詞典及檢查數個構詞規則的方式去作剖析，把讀入的文句切割成一序列的詞語 *y qtf u+，並且每一個詞語經由查詞典也可得知它的拼音符號。接著，在區塊 *d 產生各音節之音高軌跡 *r ke j " eqpvqt+ 及時長 *f wcvkqp+ 值，對於每一個音節，先準備好它的文脈資料項，再將文脈資料輸入至兩個類神經網路 *CPP+，以分別預測出音高軌跡參數和時長參數的值，關於 CPP 的輸出輸入資料項、及結構的細節，請參考我們先前的研究論文 *l w" (" Y w" 422; +。

在區塊 *e 挑選聲母、韻母之 J OO 模型，首先依據各個音節的拼音符號去查詢出對應的聲母和韻母之拼音符號與編號，並且決定聲、韻母在表 3、表 4、與表 5 的分類編號，然後依據各個單元 *聲母或韻母+ 的編號和其前後文脈的分類編號，從訓練出的半段式 J OO 模型中，找出一個單元 *聲母或韻母+ 對應的前、後兩個半段式 J OO，接著把聲、韻母的四個半段式 J OO 依序串接成一個音節的完整 J OO 模型。在區塊 *f 決定各個

JOO"狀態的駐留音框數，我們採用 Vqmfc"等人提出的方法*Vqmfc"et al.."4226+，依據 CPP"產生的音節時長值，去計算一個音節 JOO"之各個狀態所應分配到的時長音框數。接著，在區塊*g+產生各音框之頻譜係數*即 FEE"係數+，一個常被使用的方法是最大似然估計法*Oczko wo "rkngrkj qqf "Gurko cvg."ONG+*Vqmfc"et al.."4226+，不過，本論文使用的是先前我們提出的加權式線性內插法*y gli j vgf "rpgct"kpgrtr qrvqrp."Y NK" *I w"et al.."4232+，來計算各個音框的 FEE"係數。

在區塊*h+計算各個音框的音高值，一個有聲單元*例如聲母*lo l+和韻母*kl+的各個音框都必須被指派一個音高頻率值*單位 J |+，在此我們拿 CPP"產生的基週軌跡參數去作拉格蘭日內插*Nci tcpi g"kpgrtr qrvqrp+，以求得各個音框的音高頻率值。接著，在區塊*i+採用 JPO"信號模型作語音信號合成，我們把一個單元各音框的 FEE"係數與音高頻率值，按照音框次序逐個音框送給 JPO"語音信號合成模組，去合成出語音信號，關於 JPO"信號合成之細部處理方法，可參考我們先前的研究論文*I w("Vuck"4235=I w("Vuck"422; +。

4.3 頻譜距離量測

在此以代號 U E"與 U F"表示本研究所建造的兩個國語語音合成系統，U E"表示完整的系統，對於一個音節的前接與後接文脈都納入考慮，這意味聲母的前半段 JOO"的左文脈有作區分*例如圖 6"之 I [aZ+，並且韻母的後半段 JOO"的右文脈也有作區分*如圖 7"之 J [a\ +，所以對於各個不同的文脈樣式分類*如 EZ+，就需要去建造一個對應的半段式 JOO" *如 I [aZ+。相反地，在簡化的系統 U F"裡，對於一個音節的前接與後接文脈就不去作區分了，亦即聲母的前半段 JOO"不去區分它的左文脈，如此一個聲母就僅需訓練出一個前半段之 JOO，然而聲母後半段之 JOO，則仍然需訓練出數個半段式 JOO，以區分右文脈；類似道理，韻母的後半段 JOO"就不去區分它的右文脈，如此一個韻母就只需訓練出一個後半段之 JOO，然而韻母前半段之 JOO，則仍然需訓練出數個半段式 JOO，以區分左文脈。此外，以 U R"表示先前研究裡所建造的國語語音合成系統*I w"et al.."4232+，在 U R"系統裡，我們對於每一種國語音節都建造了一個或數個音節寬度的 JOO，至於一種音節所建造的 JOO"個數，則和該音節在訓練語句中的發音次數有關。在本研裡，我們使用相同的訓練語句，來訓練這三個系統*U E、U F"和 U R。

藉由這三個系統，我們就可以把錄音語句與合成語句相對應的音框拿去作頻譜距離的計算。詳細情形是，把 72"句測試語句的標記檔逐一輸入給 U E、U F"和 U R"系統去處理，以取得三個系統對應於各句測試語句的的合成音檔，然後分別拿各句的錄音語句去和合成語句去作 FEE"分析，以計算出兩個 FEE"頻譜特徵向量的序列。接著，偵測兩 FEE"序列中各組對應音框是否都為有聲，若對應的音框都偵測為有聲，就拿該組音框去計算音框之間的 FEE"向量幾何距離，然後我們拿 72"句測試語句的所有有聲音框算出的幾何距離，計算出一個跨語句的平均距離。

表 6"所列出的數值，就是這三個系統的平均頻譜距離，從表 6"可知 U E"系統所產生的音框 FEE"向量，最靠近於錄音語句分析出來的 FEE"向量，而 U R"系統所產生的音

框 FEE"向量，最遠離錄音語句分析出來的 FEE"向量。此外，只要一個音節中的聲母和韻母之間的文脈相依性有被掌握，即 U F"系統的情況，則量測出的 FEE"頻譜距離，就會比 U R"系統的好很多，這表示圖 6"和圖 7"所列出的半段式 J O O"結構，的確可幫忙掌握兩個相鄰語音單元之間的文脈相依之頻譜特性。

表4. 平均DCC 頻譜距離

| U{ungo " | U E" | U F" | U R" |
|-------------|---------------|---------------|---------------|
| Cxi 0f ku0' | 0.633" | 0.640" | 0.732" |

4.4 主觀聽測實驗

聽測實驗使用一篇訓練語句沒有用到的短文，該短文包括 92"個音節，我們將它分別輸入到三個系統去合成出語音音檔，在此分別以 Y E、Y F"和 Y R"來表示 U E、U F"和 U R"這三個系統所合成的音檔，這些音檔可到如下的網址去下載與試聽 < j wr <li vj {0ulg0pwu0gf w0y lj o o j crhl 。

透過 Y E、Y F"和 Y R"三個音檔，我們進行流暢度比較的聽測實驗，一共邀請了 34"位受測者，在第一次聽測實驗裡，受測者以隨機次序聆聽 Y E"和 Y F"音檔；在第二次聽測實驗裡，受測者以隨機次序聆聽 Y E"和 Y R"音檔；在第三次聽測實驗裡，受測者則以隨機次序聆聽 Y F"和 Y R"音檔。在各次聽測實驗中，當一位受測者聽完兩個音檔後，我們要求他給一個分數來顯示流暢度的聽測結果，評分的範圍為/4"到 4"分，4*/4+分表示後者*前者+比前者*後者+流暢很多，3*/3+分表示後者*前者+比前者*後者+稍微流暢，2"分表示分辨不出來。

三次聽測實驗之後，我們依音檔播放次序來調整分數的正負號，然後計算出各次實驗的平均分數和標準差，結果得到如表 7"所示的數值。從表 7"可知第一次和第二次實驗的平均分數為/20 55"和/20639，負的分數表示音檔 Y E"比 Y F"和 Y R"較為流暢，如果全部受測者都具有語音合成研究的背景，我們認為兩負分的絕對值應會更大，所以半段式 J O O"結構的確可以有效的提升合成語音的流暢度。至於第三次實驗的平均分數 20472，該分數的絕對值最小，我們覺得這表示 Y F"和 Y R"音檔的流暢度應無明顯的差異。

表5. 聽測實驗之平均評分

| | Y E'xu0Y F" | Y E'xu0Y R" | Y F'xu0Y R" |
|-------|----------------|----------------|---------------|
| CXI " | -0.833" | -0.417" | 0.250" |
| UVF " | 0.718" | 0.900" | 0.866" |

5. 結語

在本論文中，我們應用音素的發音知識於取代決策樹，來對一個語音單元*聲母或韻

母之左右文脈作分類，以降低文脈組合之數量；此外，更進一步研究提出文脈相依之半段式 J O O "結構，以便在有限語料的情況下，掌握一個語音單元的文脈相依頻譜特性。如此結合兩者，用以改進合成語音的流暢度。

為了評估本論文所提出的半段式 J O O "結構，我們進行了兩種實驗，即頻譜距離量測和流暢度聽測，量測出的平均頻譜距離顯示，使用半段式 J O O "結構所合成出的語句，和原始錄音語句之間的頻譜距離，可從 2054 減少到 2055；此外，聽測實驗的結果顯示，使用半段式 J O O "所合成出的語音，比使用另外兩種 J O O "結構的較為流暢，所以在訓練語句不充足的情況下，半段式 J O O "結構確實可改進合成語音的流暢度。

未來我們可在相同訓練語料的情況下，比較我們系統的合成語音與 J V U 軟體的合成語音，觀察兩者在客觀頻譜距離和主觀聽測上的差異。另外，本論文著重於改進語音單元之間頻譜銜接上的流暢度，未來可再考慮去改進韻律方面的流暢度，以更為提升系統整體的流暢度。

致謝

感謝國科會計畫之經費支援，國科會計畫編號：PUE 324/4443/G/233/34；。

參考文獻

- Ecr r². "Q0" ("G0O qwrpgru" *3; ; 8-0T gi wrctk cvkqp "gej pls wgu" hqt "f kuetgv" egr utwo "guko cvkqp" *IEEE Signal Processing Letters.* "5*6+ "322/3240"
- I w. "J 0[0" ("U0H0 Vuck" *422; +0C "f kuetgv/egr utwo "dcugf "ur getwo "gpxgnr g" guko cvkqp" uej go g" cpf "ku" gzco r ng" cr r rlec vqp" qh' xql eg" vcpuhqto cvkqp" *International Journal of Computational Linguistics and Chinese Language Processing.* "36*6+ "585/5: 40"
- I w. "J 0[0" ("E0[0Y w" *422; +0O qf gn'ur getwo /r tqi tguokp" y kj "F VY "cpf "CPP "hqt"ur ggej " u{p vj guku" *Proc. ECTI-CON.* "Rcwc{c. "Vj ckcpf. "3232/32350"
- I w. "J 0[0" O 0[0'Nck" ("U0H0 Vuck" *4232-0'Eqo dlpki "J O O "ur getwo "o qf gm" cpf "CPP" r tquf { "o qf gm" hqt"ur ggej "u{p vj guku" qh' u{ mcdng" r tqo kpgpv' rpi wci gu" *Proc. ISCSLP.* "Vclp cp. "Ur gekn' Uguokp" 30"
- I w. "J 0[0" ("E0'N0' Vuck" *4235-0' kpgi tcvpi "ur gcngt/pqpur gekhe" vko dtg" vcpuhqto cvkqp" vq" cp" J P O "dcugf "ur ggej "u{p vj guku" uej go g' *Journal of the Chinese Institute of Engineers.* "58*5+ "593/5: 30"
- J ulc. "E0'E0'E0J 0Y w" ("L0[0Y w" *4232-0'Gzr nklpi "r tquf { "j lgtctej { "cpf "f { pco le" hgcwt gu" hqt" r ke j "o qf gtrpi "cpf "i gpgtcvqp" lp" J O O /dcugf "ur ggej "u{p vj guku" *IEEE trans. Audio, Speech, and Language Processing.* "3: *+ "3; ; 6/42250"
- Tcdlpgt. "N0("DOJ 0Lxpi " *3; ; 5-0Fundamentals of Speech Recognition. "Rtgpvleg" J cno'

"
"

Vqf c."V0"("M0'Vqmf c"*4227-0'Ur ggej "r ctco gvgt"i gpgtcvqp"cri qtkj o "eqpukf gtlpi "i mdcn' xctkpeg" hqt" J O O/dcugf "ur ggej "u{pyj guku0' k" Proc. Eurospeech." Nkudqp." Rqtwi cn" 4: 23/4: 260'

Vqmf c."M0'J 0\ gp."("C0'Y 0'Drcem"*4226-0'Cp"J O O/dcugf "crr tqcej "v"o wnkpi wcn'ur ggej " u{pyj guku0' k" Text to Speech Synthesis: New Paradigms and Advances." Gf kqtu<" U0 P ctc{cpcp"cpf "C0'Cy cp."Rtgvleg"J cm"357/3750'

[cp." \ 0'LO"[0'S kp."("H0'M0'Uqpi "*"422; +0'Tlej "eqpvzv'o qf grkpi "hqt"j ki j "s wcrk{ "J O O/ dcugf "VVU0'k"Proc. INTERSPEECH."Dt ki j vqp."WM"3977/397: 0'

[quj ko wtc."V0" M0'Vqmf c."V0'O cuvnuq."V0'Mjdc{cuj k"("V0'Mkco wtc"*3; ; ; +0'Uko wncpgqwu" o qf grkpi "qh"ur gextwo ."r kej "cpf "f wtcvqp" k"J O O "dcugf "ur ggej "u{pyj guku0' k" Proc. Eurospeech."Dwf cr guv"J wpi ct{ ."4569/45720'

\ gp."J 0"V0'P qug."L0[co ci kuj k"U0'Ucnuq."V0'O cuvnuq."C0'Y 0'Drcem"("M0'Vqmf c"*4229-0'Vj g" J O O/dcugf "ur ggej "u{pyj guku"u{urgo "*"J VU"xgtukp"4020'k"Proc. 6th ISCA Workshop on Speech Synthesis."Dqpp."I gto cp{ ."4; 6/4; ; 0'

vj g'r qukkqp'qh'r tlo ct { "utguu0Rquv'r tlo ct { "u{ mcdngu"vqpf "vq"dg"tgf wegf "vq"pgct/vgt vct { "utguu" y j kg'r tg/r tlo ct { "u{ mcdngu"eqwf "dg"grxcvqf "vq"pgct/r tlo ct { "o ci pkswf g"lp"H20Vj wu"vj g"5/" y c { "r tlo ct { lueqpf ct { hgtvct { "eqpvcu'ku"o gti gf "kpq" c"dlpct { "utguu1pq/utguu"eqpvcu'y kj " tqdww'r tquf le"eqpvcu'dgy ggp"vj g'r tlo ct { "utguu"cpf "ku"hmjy lpi "u{ mcdng*+0Cu"t guwuu." vj g'r qukkqp/tgrvqf "o gti g'qh'vj g"ugeqpf ct { "y qtf "utguu"wtpu"qw"vq"dg" f khlewu"lqt"VY "N4" ur gcngtu"vq"ectt { "qw"cpf "k'o ki j v'dg"qpg'qh'uqtegu"qh"N4"Gpi rkuj "ceegp0Y g'cnuq"eqo r ctgf " VY "N4"Gpi rkuj "ceegp"cpf "VY "O cpf ctlp."vj g"vcti gv"N4"ur gcngtu'o qj gt "vqi wg."lp"qtf gt "vq" wpeqxt "lp"y j cv'y c { u"VY "N4"ceegp"eqwf "dg"cwtkdwgf "Itqo "vj gk"N3"O cpf ctlp"hgwtgu0 Hqmjy lpi "vj ku"rkp"qh'tgugctej ."vj g'r tgugpv"uwf { "lpeqtr qtcvq"r tquf le"hgwtgu"lqwpf "vq" eqpvlkdwg"vq"VY "N4"ceegp"lp"r tgxlqwu"uwf lgu"vq"lpxgugi cvg<3+"j qy "VY "ugeqpf /rcpi wci g" *N4+"Gpi rkuj "ku" f khgtgpv"lqo "N3"Gpi nkuj "d { "lpvgti tcvgf "r tquf le"hgwtgu."4+"k"cp { "tcpuht" ghhev" lqo " N4u" o qj gt " vqi wg" eqpvlkdwgu" vq" N4" ceegp" cpf " 5+" y j cv' ctg" vj g" ulo kctkkgulf khgtgpegu'dgy ggp"N3"cpf "N4"d { "r tquf le"vgo r rvgu"qh'vcti gv'y qtf ulugpvgegu0

Vq'vqu'j qy "VY "N4"Gpi rkuj "ku" f khgtgpv"lqo "N3"Gpi nkuj "d { "lpvgti tcvgf "r tuqf le"hgwtgu."N3" Gpi rkuj ."VY "N4"Gpi rkuj "cpf "VY "N3"O cpf ctlp"ctg"lf gpvkhgf "lqo "gcej "qj gt"d { "4"ercuukhgtu." UXO "cpf "MPPE0Tguwuu"uj qy "vj g'r tquf { "qh"VY "N4"Gpi rkuj "ku" f klpv"lqo "N3"Gpi rkuj = j qy gxg."VY "N4"Gpi rkuj "cpf "VY "O cpf ctlp"uj ctg"eqo o qp"r tquf le"ej ctcevgtku"y j lej " eqwf "dg" f khgtgpvcvqf "lqo "N3"Gpi rkuj 0Hvtj gt "cpcn { ugu"d { "lpf kxf wcnr" tquf le"hgwtg"uj qy " f klpv"N4"hgwtgu"qh"VY "Gpi rkuj "y j lej "o ki j v'cwtkdwg"vq"r tquf le"tcpuht"lqo "O cpf ctlp0 Cpqj gt "hgwtg"ku"vj g"rguu"vgo r q"eqpvcu'lp"ugpvpeg"vj cv'eqpvlkdwgu"vq" f khgtgpv'tj { vj o = y j kg"pcttqy gt "lqwf pguu"tcpi g"qh"y qtf "utguu"ku" { gv'cpqj gt "hgwtg"vj cv'eqpvlkdwgu"vq"rguu" utqpi ly gcnf klpvqpp"lp"VY "N4"Gpi rkuj 0

Rtquqf le"vgo r rvgu"qh'vcti gv'y qtf luvpvgeg"ctg"lmtj gt "eqo r ctgf "dgy ggp"cp { "y q"qh"N3 IN4" ur gcngtu0Y g"cuwo g" dgy ggp/ur gcngt"ulo kctk { "y qwf " dg"i tgcvt"y j gp"vj g"ur gcngt/r ct" ugrgegf " dngupi u"vq"vj g"uco g"ur gcngt"i tqwr "N3 IN4+"vj cp"y j gp"vj g"ur gcngt/r ct"ku"lqo " f khgtgpv'ur gcngt"i tqwr u"*N3 IN4+0Vj g'dgy ggp/ur gcngt"ulo kctk { "ku" f ghpgf "cu"cxgtci g'equkpg" o gcuwtg"dgy ggp"y q"ur gcngtu"d { "r tquf le"hgwtg"xgevtu0Vj gtguwv"fgqu"uj qy "N3"cpf "N4" ur gcngtu"r tqf weg"r tquf le"vgo r rvgu"y kj "i tgcvt"y kj lp/i tqwr "eqpukngpe { "tgr gevkgu" "dw" vj gk"y kj lp/i tqwr "r cwgtpu"ctg" f klpv"lqo "vj gk"eqwvgr ctv"i tqwr 0Qpg" f klpv"r cwgtp"y g" hqwpf "ku"lqwf pguu"qh'ugpvpeg"cpf "cpqj gt"r cwgtp"vj g"vko lpi lr kej "r cwgtpu"qh'y qtf 0Vj g'cdqyg" r tquf le"tcpuht"ghhev"cpf "f klpv"VY "N4"r cwgtpu"qh'r tquf { "ctg"lqwpf "lp"tgrvqpp"vq"u { pvcz/" lpf wegf "pcttqy "lqewu"cpf "rgzleqp/f ghpgf"y qtf "utguu"y j lej "gej q"qwt"r tgxlqwu"lqpf lpi u'y kj " tgi ctf "vq"r tquf le"tgerk cvqpu"VY "N4"Gpi rkuj 0

Y g'dgrkxg"vj g'cdqyg"cpn { ugu"y kj "lpeqtr qtcvq"r lpi wku"npqy rgi g'pqv'qpn { "uj gf "rki j v'qp" dgwt"wpf gtucpf lpi "qh"VY "N4"Gpi rkuj ."dw"ecp"cuq"dg"lo r rgo gpvqf "lp"lmtj gt"tghpo gpu"qh' vj g'ewttgpv"ECNN"cpf "ECRV"u { ugo u"0

網頁商家名稱擷取與地址配對之研究

Uqtg'P co g'Gz tcevkqp"cpf 'P co g/Cff tguu'O cvej kpi "qp"yj g'Y gd"

林育暘 Nkp' l w[cpi ."張嘉惠 Ej cpi 'Ej kc/J wk'
國立中央大學資訊工程學系

F gr ctvo gp'qh'Ego r wgt "Uekpeg"cpf "kphqto cvkqp"Gpi kpggtkpi "
P cvkqpcn'Egptcn'Wpkxgtukv{"
323744274B@ewqf.wly .'ej kcB eukg@ewqf.wly "

摘要

行動裝置的普及造就大量地域性查詢的需求，其中最常見的一種查詢，就是尋找附近的餐廳或加油站。然而當使用者在電子地圖上搜尋這些地點名稱 *RQK Rqkp' qh' kpgtgum時，經常無法找到，因為電子地圖上雖有地點名稱標註，但是相關資訊不足，而這些資訊其實大多可以在網頁中找到。因此使用者大多必須開啟瀏覽器搜尋商家名稱找出地址，並把地址輸入至電子地圖查詢路線。但行動裝置螢幕小，且輸入文字不便利，如果要反覆查詢將是一件耗時耗力的工作。如果這時候有一個商家地理資訊系統能事先將網路上的商家資訊進行整合，最後提供一個 CRR"直接讓使用者查詢，將可以大幅度減少使用者與裝置間的互動次數，有效的提供便利性。

為建構商家地理資料庫，Ej wcpj "等人]; 對於包含地址網頁提出三種抓取程式，並利用 Ej cpi "等人]4_的地址擷取程式，取得大量中文地址。NK'與 Ej cpi]3_並定義地址相關資訊擷取問題，希望藉此豐富每個 RQK'的相關資訊，提高地理資訊檢索'I gqi tcr j lecn'kphqto cvkqp'Tgtkxcn'I K'的召回率。然而不論是NK'與 Ej cpi]3_或Ej cpi "等人]4_或Ej wcpj "等人]; 的相關資訊擷取方法都僅能從多筆地址網頁擷取資訊，所得資訊有限，對於單筆地址網頁的相關資訊擷取仍尚無研究。

本研究從地址的角度出發，利用已抓取大量包含地址的網頁，先找出網頁中的地址，再藉由地址找出對應的商家名稱進行配對。換言之，給定一個已知地址，我們希望能透過網路資料擷取出該地點的名稱（如：商家名稱、政府單位…等）。舉例而言：當我們已有地址「新北市板橋區中山路二段 :: "號 5H」，我們希望能知道這個地址對應的名稱「大鈞醫學美容診所」，如此即可進一步藉由地址、名稱並利用搜尋引擎收集更多額外商家資訊。這些額外資訊不僅可以有效提昇 I K'檢索系統的召回率，也可提昇商家分類的準確率]32_。

在辨識商家名稱的部分，本篇論文使用了條件隨機域 *Eqpf kkpccn'Tcpf qo " Hkgf +]5_當作學習演算法。目前有許多關於中文組織名稱辨識的研究]6_"]7_"]8_"]9_，可以從新聞或一些較正式的文章中萃取出組織名稱，但是並沒有嘗試以一個 ETH'O qf gn'直接對各種網站中的整個網頁內容進行中文組織名稱辨識。這兩者之間不同處在於新聞類文章屬於較正式的文章體裁，因此容易出現行政機關與正式的組織名稱，例如：行政院和維德食品有限公司，但是整個網路上商家組織

名稱的命名方式傾向則不同，例如：吼牛排、努哇克咖啡、阿嬤祖傳菜包肉粽仙草...等，都是商家組織名稱。另外，一個完整的網頁內容有結構與非結構化的資訊交錯呈現，雖然結構化資訊會造成自然語言文字內容的破碎，但這些結構也隱含有可利用的資訊。

為了使商家辨識能以最少人力進行自動化學習，本研究使用自動標記方式建立訓練資料，我們先針對部份的黃頁網站撰寫 Rctugt"取得大量商家名稱與地址的組合，並以這些已經取得的商家名稱對網頁語料進行自動標記，再利用自動標記後的語料訓練 ETH"序列標記模型。然而一個地址可能出現在多個網頁之中，僅只仰賴其中一個網頁也有失之偏頗之慮，因此我們也收集了 I qqi rg"Upkr r gwu"當作訓練資料進行比較。本篇論文的第二個主題則是商家地址的配對，由於一個網頁可能包含多個商家名稱，我們對網頁以簡單的規則進行分類後，使用了啟發式 (j gwtkule) 的配對規則，利用各類型的網站所具有的表達特性，對地址與商家名稱進行配對。

本研究承續]: _]; 之研究，經由爬取網頁上包含地址的大量網頁（包括 [gmqy "Rci g"與 Uwthceg"Y gd) 進行商家名稱擷取。其中 [gmqy "Rci g"提供了大量商家名稱以及地址與商家的配對資料，而 Uwthceg"Y gd"則利用]4_之地址擷取模型擷取出了可能含有台灣地址的網頁與地址清單。本篇論文以已知可能含有台灣地址的中文網頁、每筆網頁的地址清單、大量商家名稱清單以及已知的地址與商家名稱配對資料為基礎，提出了一個商家名稱擷取系統，方法分為四大步驟：地址網頁的前處理、商家名稱命名實體辨認、及地址－商家名稱匹配。本研究在三個模型聯合標記商家名稱的方式下，地址與商家名稱的平均配對正確率為 2079。

關鍵詞：商家地理資訊檢索、商家名稱擷取、商家名稱與地址配對、序列標記、條件隨機域

Mg[y qtf u<"RQK"uqtg"pco g"gzvtcevkqp."pco g/cf f tguu"o cvej kpi ."ugs wpeg"ncdgrkpi ."eqpf kkpnci' tcpf qo "hgrf "

"
"

參考文獻

]3_ "U0[0'Nk"Cr r rkecvkqp"cpf "Gzvtcevkqp"qh"Rquvni' Cf f tguugu"cpf "Tgrcvgf "kphqto cvkqp." P cvkqpcni' Egpwtcni' Wpkxgtuk{ ."422; 0'

]4_ "E0J 0'Ej cpi ."E0[0J wcpj "cpf "[0U0Uw."\$Ej kpgug"Rquvni' Cf f tguu"cpf "Cuqekcvgf " kphqto cvkqp" Gzvtcevkqp.\$"Vj g"48vj "Cpwwni' Eqphtgpeg"qh'vj g" Lcr cpgug"Uqekgv{ "hqt" Ctvkhekni' kpgvni gpeg."42340'

]5_ "N0' F 0' Lqj p." O 0' Cpf tgy " cpf " P 0' E0' Hgtpcpf q." \$Eqpf kkpnci' Tcpf qo " Hgrf u<" Rtqdcdkkule" O qf gni" hqt" Ugi o gpvki " cpf " Ncdgrkpi " Ugs wpeg" F cv.\$" K EON" Rtqeggf kpi u"qh'vj g" Gki j vggpj " kpgvtcvkqpcni' Eqphtgpeg" qp" O cej kpg" Ngctkpi ."r r 0' 4: 4/4: ; ."42230'

]6_ "\ 0' Uwzkpi ." \ 0' Uwzcp" cpf " Y 0' Zkcqlk." \$Cwqo cvk" Tgeqi pkkqp" qh" Ej kpgug" Qti cpk cvkqp" P co g" Dcuqf " qp" Eqpf kkpnci' Tcpf qo " Hgrf u.\$" P cwtcni' Ncpi wci g" Rtqeguuki "cpf " Mpqy ngf i g" Gpi kpggtkpi ."r r 044; /455."42290'

]7_ "[0Zk kpi ." \$C" O gy qf "qh" Ej kpgug" Qti cpk cvkqp" P co gf "Gpvkku" "Tgeqi pkkqp"
Dcugf "Qp" Ucvkuecn" Y qtf "Hgs wgpe{ ." Rctv" Qh" Ur ggej "Cpf " Ngpi vj . \$" Dtqcf dcpf "
P gy qtnicpf "O wko gf k" Vgej pqmi { " *E/DP O V+ : 'r r 0859/863. "42330'

]8_ "N0[clwep." [0Lkpi "cpf "J 0Nkpi . "\$Ej kpgug" Qti cpk cvkqp" P co g" Tgeqi pkkqp" Dcugf "
qp" O wkr rg" Hgcwtgu. \$" Rckle" Cuk" eqphgtgpeg" qp" Kvgmki gpeg" cpf " Ugewtk{ "
Kphqto cvku. 'r r 0358/366. "42340'

]9_ "E0Y 0Y w." T0V0J 0Vuck" cpf "Y 0N0J uw." \$Ugo k/lqkv' rcdgnkpi "hqt" ej kpgug" pco gf "
gpvk{ "tgeqi pkkqp. \$" Rtqeggf kpi u" qh' vj g" 6vj "Cuk" kphqto cvkqp" tgvkxgn' eqphgtgpeg. "
r r 0329/338. "422: 0'

] : _ [0U0 Uw" Cuqekcvgf " Kphqto cvkqp" Gztcvkvqp" hqt" Gpcdnkpi " Gpvk{ " Ugctej " qp"
Ggewtqple" O cr . "P cvkqpcn' Egpvtcn' Wpkxgtuk{ . "42340'

] ; _ "J 0O 0Ej wpi . "E0J 0Ej cpi "cpf "V0[0Mcq. "\$Ghgevkg" Y gd" Etcy nkpi "hqt" Ej kpgug"
Cf ftguugu" cpf "Cuqekcvgf " Kphqto cvkqp. \$" k' GE/Y gd. "O wplej . "I gto cp{ . "42360'

Rwdrke'Qr kpkqp'Vqy ctf 'EUUVC<C'Vgzv'O kpkpi 'Cr r tqcej "

[kCp'Y w.'Uj wMckJ ukj ' "

"
"

Cdurtcev'

Tgegpw{\."y g'r wdrke"qh'Vcky cp"j cu"j cf "c"j gcvgf "f gdcvg"qp"y g"kuwv"qh'Etquu/
Utckv' Ugtxleg" Vtcf g" Ci tggo gpv' *EUUVC-0' Chgt" o qpvy u" qh" uko o gtkpi " vgpukapu"
dgy ggp" twkpi " r ctv{\ " cpf " qr r qukkqp" r ctv{\ " utqpi n{\ " dcengf " d{\ " yj g" uwf gpv/ngf "
Uwphqy gt"O qxgo gpv"y g'f gdcvg"j cu"hpckm{\ "tgcej gf "c"dtgcnkpi "r qkv"qp"O ctej "3: ." 4236."
cv'y j lej "uwf gpw" qeewr kgf " yj g"Ngk kurv\kg" [wcp0'F wtkpi " yj ku"r gkqf ."pqxgr'
eqo o wplecvkqp"uwej "cu'Hcegdqnm'uj ctkpi ."kpuvcpv'o guuci kpi ."cpf "f kuewukapu"qp"RVV"
j cxg"tguj cr gf " yj g"uqekn'o qxgo gpv' ukpeg" yj g{\ " ctg" gcuk{\ " ceeguukdrg" cpf " kpuvcpv{\ "
tgur qpf gf 0'Vj g'uqekn'o gf k"j cu'dgeqo g'y g'f qo kpcpv'uqwtg'kp"qr kpkqp"uj cr kpi "cpf "
yj g'ceeqo r cp{\ kpi 'ugpvko gpv'ur tgcf 0'

Vj g"gzvcev\kp" cpf "vcenkpi "qh"wr tkkpi "r qrkkecn'qr kpkapu"cpf "gxgpv"uwej "cu"
EUUVC"j cu'dgeqo g"qpg"qh" yj g"o quv'ko r qtwcpv' vqr keu" yj cv'tgegkxg" o wej "cwgpvkqp0'
Y kj "yj g"j wi g"co qpwu"qh'vgzu."k'ku'pqv'r quukdrg"vq"cpck{\ | g"cpf "kpvgr tgkpi "yj g'uqekn'
cpf "r qrkkecn'vgzu"o cpwcm{\ 0' kpuvgef ."y g"r tqr qug"vq" wug" yj g"vgz'v'o kpkpi "cr r tqcej ."
y j lej " cwqo cvkcm{\ " gzvcev' qr kpkqp" cpf " kphqto cvkqp" r tqhkgu" Itqo " yj g" vgzv0'
O qtqxtg."yj ku"cr r tqcej "cnuq'utgpi yj g'p'v'g'qdlgevxk{\ ."hqt"yj g'pqto u'ctg'ugv'a priori."
cpf "yj wuj wo cp'dkcugu'ctg'tgf wegf 0'

Cu"cr kqpggtkpi "y qtmkpi"yj g'eqpvz'v'qh'Vcky cp"uqekv{\ ."yj ku'tgugctej "cko u"v"vcev"
yj g'r wdrke"qr kpkqp"vqy ctf 'EUUVC"htqo "yj g'r gtur gev\kg"qh'vgz'v'o kpkpi 0'Vj g"cr r tqcej "
kpxqrxgu" yj g"o cpwcm{\ "gzvcev\kpi "qh"r qrkkecn'ucpeg" tgrcvgf "ng{\ y qtf u"cpf "r j tcugu."
uwr gtxkugf "o cej kpkpi "ngctpkpi ."cpf "c"ucv\kuecn'o qf gn'qh' yj g"v'gpf 0'Y g'hqewu"qp"yj g"
kpf kxk wcn'r quvu"qp"RVV"tcy gt" yj cp" pgy u" ukpeg" yj g{\ " ctg"o qtg" tgr tgugpv\kxg0' Vj g"
r qv'p'ken' r qrkkecn' qt " eqo o gtekn' cr r rkecvkpu" ctg" xcwcdrg0' Qpg" ecp" f kueqxtg" yj g"
r wdrke"qr kpkqp"cpf "tgur qpug'kp" c'uj qtv\ko g0'

Vj g"o cvgtken'y g" wugf "kp" yj ku'tgugctej "kpenf gu" c"rkuv"qh"o cpwcm{\ "etgcvgf "uggf "
y qtf u"cpf "r j tcugu" tgr tgugpv\kpi " yj g"r tq/cpf /eqp"r qrkkecn'r qmktk{\ ."tgur gev\kgn{\ 0'Vj gug"
y qtf u"ctg" v'ugvf "d{\ " yj g" vgzv"qh" yj g"y gdukng" 〇服貿東西軍〇^Ä." y j lej "ercuukhku" yj g"
uwr r qtv\kpi "cpf "qr r qukpi "vgzu0' Cpvy gt" tguqwtg"y g" wugf "kp" yj ku"y qtm'ku" yj g"RVV"
eqtr wu." y j lej "ku" c' r qr wrct "qprkpg" dwngv\p" dqctf "hcxqtgf "d{\ " b cp{\ "qh' yj g" {\ qwj 0'

· "I tcf wcv'k'pukwv"qh'Nkpi wku'ku.'P cvkqpcn'Vcky cp'Wpkxgtuk{\ "
Ä" j wr <lgehc0ur gcnkpi 0y ko j q0' j r "

Qwt "r tqegf wtgu" hqmjy "y g" eqo o qp "vgz v" o kpkpi "vej pls wgu" hgcwt gu" gz tcevkqpu." Ej kpgug" y qtf "ugi o gpvcvkqp" y kj "ewuqo "f levkqpc{. "guvdrkuj "y g" o qf gn' hqt "y g" UXO "ercuukhgt. "cpf "wukpi "y g" P / hqrf " "etquu" xcrkf cvkqp "hqt " gxcnvcvkqpu0" Y g" "ej qqug" "y g" ng{ y qtf u" cu" y g" hktuv" ugr "ukpeg" o cp{ "vgto u" ecp" r qvgrvcm{ "t gxcn' qp g' u" cvkwf g0' Hqt " kpuvcpeg. "y g" uwr r qt vgt "hqt " EUUVC" y qwf "ecm' uwf gpw" 霸占 "occupy. "y g" r ctrko gpv." y j kg" "y g" qr r qpgpv" y qwf "wug" 留守 "stay. "k" "y g" r ctrko gpv0' Vj gp" y g" wug" "y g" ng{ y qtf u" cu" hgcwt gu" vq" tckp" y g" UXO "ercuukhgt0' Vj g" i qrf "ucpf ctf u" qh' y g" vgz w" ctg" ej qugp' hqo "y g" 服貿東西軍 "y g" gdukg0' Vj g' tguwmu" ctg' uq y p' cu' hqmjy u<

| Ceewtce{ " | Rtgekukqp" | Tgecm" | H'ueqtg" | Uf 0F gx0" |
|------------|------------|---------|----------|------------|
| 20 72" | 20 72" | 20 7; " | 20 77" | 20 62" |

"

Y g" hmt vj gt "gz vpf gf " qwt " tguwmu" vq" f q" y g" vtgpf " cpcn{ uku0' Hktuv. " y g" cr r n{ " y g" kphqto cvkqp" i ckp" ecrewrcvfg " hqo " y g" r t gxlqwu" ercuukhgt. "cpf " y g" gp" y g" uwo " ng{ y qtf u" qh' gcej " r quv" cpf " uwo " qxgt " y g" r quu" qh' y g" uco g' f c{ 0' k' qvj gt' y qtf u. " y g" ueqtg' qh' gcej " f cvg' ku' ecrewrcvfg " cu' y g' hqmjy kpi " gs wcvkqp<

$$\text{Score} = \sum_i \sum_w IG(w) * C(w), \quad i = \text{post index}, w = \text{word}$$

Vj g' hki wtg' f go qputcvgu" y g' r qr wrctkv{ "qh' y ku" vqr le" qh' gcej " f c{ 0' Ugeqpf . " y g" ecrewrcvfg" y g' uwr r qt vki " kphqto cvkqp" i ckp" qxgt " y g" vqcn' kphqto cvkqp" i ckp. "cpf "cnuq" uwo " qxgt " y g" r quu" k' ppg' f c{ 0' Vj ku" hki wtg' uq y u" y g" tcvkq' qh' uwr r qt vki " EUUVC" hqo " y g" cpcn{ uku" qh' r quu0'

O kpkpi "cpf " tcentkpi " r qrkkecn' qr kpkqpu" hqo " vgz w" k' y g" uqekcn' o gf k' " ku" c" { qwpi " { gv" ko r qt vcp' t gugctej " ctgc" y kj " dqj " uelgpvkle" uki pkkcepeg" cpf " uqekcn' ko r cev0' Vj g" i qcn' qh' y ku" r cr gt " ku" vq" o qxg" qpg" ugr " hqy ctf " k' y ku" ctgc" k' " Ej kpgug" eqpvz0' Y g" uctvfg " hqo " y g" o cpwcm{ " etgcvgf " ng{ y qtf u" cpf " ng{ " r j tcugu" qh' EUUVC. " wugf " y g" " vq" dwnf " c" ercuukhgt " cpf " ecrewrcvfg " y g' k' kphqto cvkqp" i ckp. "cpf " y g" gp" f k' " y g" vtgpf " cpcn{ uku" qh' y g" RVV" eqtr wu0' Vj ku" cr r tqcej " kpxqrgu" kvgtf kuekr r pct{ " hkrf u" kpenf kpi " kphqto cvkqp" tgvtkxcn" f cvc" o kpkpi . " ucwku" ku. " o cej kpg" ngctpkpi . " cpf " eqo r wcvkqpcn' r kpi wku" ku0' Y g' j qr g' y cv' y ku" vgz v" o kpkpi " cr r tqcej " eqwf " f kueqxtg" y g" r wdike" qr kpkqp" vqy ctf " EUUVC. " cpf " hmt vj gt " t gxcn' r qrkkecn' ucpegu0' Hwmtg" y qtmu" kpenf g" o qtg" uqr j kulecvfg " rpi wci g" r tqeguukpi " vej pls wgu" cr r rkgf " vq" o qtg" dtqcf " f qo ckp" qh' r qrkkecn' vqr ku. " cu" y gm' cu" f g xgmqr kpi " f { pco le" tcentkpi " u{ ugo " i gctkpi " wr " hqt " { gct/ gpf " grgewkqp" 42360

Towards automatic enrichment of standardized electronic dictionaries by semantic classes

Elleuch Imen*, Gargouri Bilel* and Ben Hamadou Abdelmajid¹

Abstract

Kp"vj ku'r cr gt"y g'r tqr qug"cp"cr r tqcej "hqt"vj g"cwqo cvk" gptlej o gpv'qh"ucpf ctf k gf "grgevtqple" f levkqpcctkgu'd{ "vj g"ugo cpvk"ercuugu'Vj ku'cr r tqcej "eqpuku'qh'vj tgg'r j cugu'Vj g'htuv'r j cug'v'gcv' vj g"ugo cpvk"ercuukhlevkqpc"r tqegu"hwf gf "qp"vj g"uwf kgu'qh'I cuvq"i tqu'Vj g"ugeqpf "r j cug" r tqhkgu'htqo "vj g"gzkugf "uwdlgev'hkgrf u'kp"vj g"levkqpcct{u'ngzlecn'gpv'kgu'kp"qtf gt "v'cwtkdwg"vj g" uwkcdng"ugo cpvk"ercuugu'Vj g"hkpcn'r j cug"tgerk gu'u'pvcvk"cpnc'ugu'qh'vj g"vgz wcn'eqpv'gpv'qh' o gcpkpi u'ngzlecn'gpv'kgu'Vj ku'r j cug."clo u'v'ghkpg'vj g"uwdlgev'hkgrf "dcugf "gptlej o gpv'cpf "cuq" v'gcv'vj g'pqp"gptlej gf "o gcpkpi u'kp"vj g"ugeqpf "r j cug'Kp"cf f kkp."k'cwtkdwg'vj g'uco g'ugo cpvk"ercuugu'htqo "vj g"u'pqp{o "o gcpkpi u'Y g'wugf "cp"cxckcdng"ucpf ctf k gf "Ctdle" f levkqpcct{ "v'v'gugf " vj g'r gthqto cpeg'qh'vj g'r tqr qug" cr r tqcej 0

Keywords: "Cwqo cvk" gptlej o gpv." ucpf ctf k gf " grgevtqple" f levkqpcctkgu." ugo cpvk"ercuugu." Ctdle'ncpi wci g0

1. Introduction

Ugo cpvk"ncpy ngf i g."gur gekm{ "ugo cpvk"ercuugu"y j lej "clo "v'ej ctcevtk g"o gcpkpi u'qh'ngzlecn'wpku'kp" f levkqpcctkgu." j cxg" cwtcevgf " eqpukf gtdcng" kpwgtgu' kp" dqvj " rki wku" *Uvgf g." 3; ; : +." *Fqt." 3; ; 9+" cpf " eqo r wcvkqpcn'ncpi wku'eu"*Mkr r gt"gv'cn'4222+0Uwej "ugo cpvk"ercuugu'ecp"dg" f ghkpg'cu" c"ugo cpvk"ncpi wku" r tqr tkgv{ "ercuukh{ kpi "o gcpkpi u'cpf "ecp"vj g'ghqg'dg'wugf "cu"cxkcdng"o gcpu'qh'eqo r tgi gpf kpi "vj g'ur gekh" o gcpkpi "qh'r qn'uko qwu'ngzlecn'wpku'Vj wu"vj g"pggf "qh' f levkqpcctkgu"y kj "ugo cpvk"ercuugu"j cu" dgego g" c" pgeguukv{ "hqt" P cwtcn'ncpi wci g'Rtqeguukpi "P NR+" cr r rkev'kqu0

Hqt" xctkqwa" ncpi wci gu." xctkqwa" ugo cpvk"ercuukhlevkqpc" ctg" pqy " cxckcdng0 Y g" ecp" rku' vj g" xgtdu"ercuukhlevkqpc" *Rkpmgt." 3; ; : = Lcengpf qh" 3; ; 2= Ngxkp." 3; ; 5." Fwdqku" cpf " Fwdqku/Ej ctnkt." 3; ; 9+" vj cv' tgi tqw u'v'qj gj gt"xgtdu"vj cv'uj ctg"dqj "c"eqo o qp"ugo cpvk"ucpf "c"ugv'qh'u{pvcvk"cnv'gpc'v'kqu'Cuq."y g" pqv'eg"Y qtf P gv"*Hgmcdw ."3; ; : +vj cv'r tqxkf gu'ugo cpvk"qp'v'qmi lecn'ercuukhlevkqpc"cpf "Hico gP gv"*Hkmo qt g." 3; ; 7+" vj cv'j lgtc'ej lecn{ "ercuukh{ "ngzlecn'wpku" wukpi "xctkqwa" tgrv'kquj kr "cu"u'pqp{o {."cpv'q{o "cpf "ku/c" tgrv'kqu'J gy qxgt." vj g" tghgtgv'kcn'ercuukhlevkqpc" ku' dcugf " qp" ugo cpvk" hgcw'gu' rknq"]- 1/" j wo cp."]- 1/" eqpetgvg." gve'ej ctcevtk kpi " ugo cpvkcm{ " gcej " ngzlecn' wkv' qwukf g" qh' vj g" o gcpkpi u' eqpv'g'u' Qdlgev'ercuugu"*i tqu."3; ; 6+" f ghkpgu" c" ugo cpvk"ercuukhlevkqpc"dcugf "qp" uwt'ceg" tgerk' cvkqpc"qh'r tgf kecvg"cti wo gpv' uwt'wewt'g0C" ugo cpvk"ercuugu" tqr u'v'qj gj gt" r tgf kecvgu'cu'cti wo gpw'j cxkpi "vj g'uco g'u'pvcvk"eqpv'kqu' Tgn{ "qp" c" ugo cpvk"ercuukhlevkqpc"=v'q" o gj qf u'qh'gptlej o gpv'ngzlecn'tguq'wtegu"d{ "ugo cpvk"ercuugu"gzkuv' Vj g'htuv'qpg"ku"o cpwcn'K'ku'ej ctcevtk gf "d{ "vj g'nci g"pwo dgt"qh'ngzlecn'wpku"v'q"dg"ercuukhgf "Hico gP gv'

* HUGI U."DR032: . : "523: "Ucz."Vvpkuk"
G'o ckn' } "ko gp'ngwej ."dkgrf cti qwtk B hugi ut'pwp0

** KJO U."DR0464."5243"Ucn'kvG | k'Ucz."Vvpkuk"
G'o ckn'cdf gn' clkf @l'gpj co cf qwB kjo ut'pwp0

*Hkm qtg."3; : 7+"j ku"ku"y j { "k'ku"e"equwł "cpf "ko g/equwo kpi "o gvj qf 0'Vj g"ugeqpf "o gvj qf "ku"cwqo cve0'K' ecp"wg"eqtr qtc" *Hwej u(" J cdgtv."4226+." *Eqpf co kpgu."4227+"qt "kp"uqo g"ecugu."vgzu"qh"vj g"tgcvgf "ngzlecn' tguqtegu" *Tcułkt." 4223+" cpf " *Xcrgwg" gv" cn" 4228+0' Vj g" cwqo cve" o gvj qf " f qgu" pqv" pgegułkcvg" vj g" kpvtxgpkqp"qh"vj g"j wo cp"gzr gtv'f wtkpi "vj g"gtplej o gpv'r tqeguu" *Y knqp"gv'crl"4226+0'Dqj "o cpwcn'cpf " cwqo cve" o gvj qf " qh' gptlej o gpv' ngzlecn' tguqtegu" y kj " ugo cpve" ercuugu" tgs wktgu" vj g" kpuvkwkqp" qh' vj g" ugo cpve" ercuukhlecwkqp0' K' cf f kłqp." vj g" cdłkłv { "qh"vj g" ut wewt gđ" f lewkqct { "vq" tgegkxg" ugo cpve" ercuugu" ku" ko r qtvcv0' K' hcev" uqo g" o qf gnu"qh' ngzlecn' tguqtegu" f q" pqv' uwr r n { "vj g" chgevcwkqp" qh' vj g" ugo cpve" ercuugu" vq" ngzlecn' wpku0' ""

Kp" qtf gt" vq" r tqxkf g" c" wpkłkf " Htco gy qtn' hqt" o qf grkpi " ngzlecn' tguqtegu." kp" i gpgtcn" cpf " vq" hcełkvcg" vj g" gzej cpi g" cpf " kpgi tcvkqp" kpvq" P NR" cr r nlecwkpu." vj g" NO H' *Ngzlecn' O ctmw" Htco gy qtn" ucpf ctf " *Htceqr qwł (" I gqti g." 422: + " KQ" 46835" ku" r wdrukj gf 0' Vj ku" ucpf ctf " cmqj u" vj g" o qf grkł cvkqp" qh' cmi' r kpi wkułeu' hcxgn' uwej " cu" vj g" o qtr j qmı kcn" vj g" u { pvcve. " vj g" ugo cpve" cpf " vj g" u { pvcveq/ ugo cpve" qpgu0' "

Eqpukf gt kpi " vj g" ko r qtvcpeg" qh' vj g" ugo cpve" ercuugu" vq" ej ctcevgtk g" vj g" o gcpkpi " qh' ngzlecn' wpku." cpf " r tqłkłkpi " Htco " vj g" hkp" o qf gn' qh' NO H' ngzlecn' tguqtegu" vq" tgegkxg" ugo cpve" ercuugu" y g' r tqr qug" kp" vj ku' r cr gt" cp" cwqo cve" cr r tqcej " hqt" vj g" gptlej o gpv' qh' ucpf ctf kł gf " NO H' grgevtqpkle" f lewkqct kgu" d { " ugo cpve" ercuugu0' K' hcev." vj g" NO H' ucpf ctf " qh' gtu" r ct vewct " hgrf u" *KQ" Uwdłgev Hgrf + " vj cv' ecp" cuuku" vj g" kf gpv' hlecwkqp" qh' vj g" tgrxcpv' ugo cpve" ercuu" cpf " r tqxkf gu" u { pqp { o { " tgrv' kpuj kr u" vj cv' ecp" dg" wugf " vq" ko r tqxg" vj g" gptlej o gpv' r tqeguu0' Cnuq." kp" cp" NO H' f lewkqct { " vj g" o gcpkpi " qh' ngzlecn' gpv' kgu" ku" ceeqo r cpkfg" y kj " c" tlej " vgzwcn' eqpv' gp0' Vj g" r tqr qugf " cr r tqcej " ku" hqwpf gf " qp" c" ugo cpve" ercuukhlecwkqp" kpkłvcgf " d { " vj g" I cuvq" I tqwı" uwf kgu0' Cp" gzi gt kł gpvcwkqp" qh' vj ku" cr r tqcej " ku" ectłkf " qw" qp" cp" cxckrdng" ucpf ctf kł gf " NO H' Ctdle" f lewkqct { 0' ""

Vj g" pgzv' r ctv' qh' vj ku' r cr gt " ku" qti cpkł gf " cu" hqmıj u" Y g" y km' uctv' y kj " c" r tgu' gpvcwkqp" qh' uqo g" tgrv' gf " y qtmı" tgrv' gf " vq" ugo cpve" ercuukhlecwkqp" cpf " gptlej o gpv' o gvj qf u0' Vj gp." y g" y kmı' r tgu' gpv' vj g" NO H' ucpf ctf 0' Vj gtgchgt. " y g" y kmı' f gwckı' vj g" r tqr qugf " cr r tqcej " hqt" vj g" gptlej o gpv' qh' NO H' ucpf ctf kł gf " f lewkqct kgu" y kj " vj g" ugo cpve" ercuugu0' Chgt " vj cv." y g" y kmı' f guetłdg" vj g" gzi gt kł gpv' ectłkf " qw" qp" c" ucpf ctf kł gf " NO H' Ctdle" f lewkqct { " cpf " f kweuu" uqo g" qh' vj g" qdv' kpgf " tguwu0' Hkpcmł . " kp" vj g" eqpenwkqp." y g" y kmı' cppqwppeg" uqo g" hwwt" g" y qtmı0' "

2. Related works

Vj ku' ugevkp " ku" f gxqvgf " vq" vj g" tgr tgu' gpvcwkqp" qh' uqo g" tgrv' gf " y qtmı" qh' cxckrdng' ugo cpve" ercuukhlecwkpu" cpf " vj g" ugo cpve" gptlej o gpv' o gvj qf u" qh' ngzlecn' tguqtegu0' "

2.1. Semantic classification

Uxgtcn' ugo cpve" ercuukhlecwkpu" gzi ku' kp" r kgtcwt g0' Y g" ecp" o gpvkp" vj g" xgtdu" ercuukhlecwkqp" *Rłpngt." 3; ; ; = Lcengpf qh" 3; ; 2=Ngxlp." 3; ; 5. " F wdqku(" F wdqku/ Ej ctnıgt." 3; ; 9+0' K' dcugf " qp" dqj " c" eqo o qp" ugo cpve" u" cpf " c" ug' qh' u { pvcve" cngt' pcvkpu" vq" i tqwr gf " ngzlecn' wpku" kpvq" ugo cpve" ercuugu0' Vj ku" v' r g" qh' ercuukhlecwkqp" ku" tgu' tvegf " vq" egt' v' k' ercuu' v' r gu' cpf " tgcw' qp n { " xgtdu0' Uq" pq" eqo r tgj gpukxg" ercuukhlecwkqp" ku' cxckrdng' rko ku' vj g" wghwpguu" qh' vj g" ercuu' hqt" r tcevecn' P NR' v' cunu0' "

O qtgqxt. " y g" ecp" pqvg" vj g" qpvmı lecn' ercuukhlecwkqp" r kng" Y qtf P gv' *O kngt." 3; ; 2+ " vj cv' kpv' p' gf " vq" ercuukł " r j kqu' r j lecn' vj kpi u" cu" vj g { " gzi ku' kp" vj g" y qtr f 0' K' ku' r ct vewct n { " cr r tqr tlcvg" hqt" qdłgev' o qf grkpi . " kpenw' kpi " vj gk" tgrv' kpuj kr u" cpf " r tqr gt v' gu0' Vj gtgłgtg." eqpv' gpv' qh' qpvmı { " f qgu" pqv' kpvtcev' f k' gveł " dw" tcv' gt" y kj " tgrv' kpuj kr u" *k" u { pqp { o { . " cpvq { o . " r ctv' qh" ku' C.ı - 0' Vj ku" ugo cpve" ercuukhlecwkqp" f qgu" pqv' eqpukf gt" vj g" wgdı" eqpv' gv' qh' ngzlecn' wpku." hvtj gt " k' i tqwr u" y qtf " kp" vj g" ercuugu" cu" r tgu' gpv' gf " kp" vj g" tgcı' y qtrf " y kj qw" tghgtłkpi " vq" vj g" r kpi wkułeu' hgcwt gu0' "

Cnuq. " y g" ecp" ekvg" vj g" tghgt' gpv' cn' ercuukhlecwkqp" *I tqwu." 3; 97+ " F lej { . " 4222+ " vj cv' wugf " ugo cpve" hgcwt gu' hng"]- 1' " eqpetgvg_."]- 1' " j wo cp_0' Vj qug" hgcwt gu' ctg" cwcej gf " vq" ngzlecn' wpku" vq" f guetłdg" vj gk" cr r wt' gpv' p' v' vq" vj g"

vj g'tguwvki "o qf gn'y km'dg'f geqtcv'f y kj "vj g'F c'c'Ecvgi qt'lg'u'Tgi kvt { "F ET +4" tgs vkt gf "hqt" 'vj g'o qf grk' cvkqp" qh'vj g'f gcn'rcpi vci g0"

4. Proposed approach

kp'vj ku'ugev'kqp. 'y g'f g'ckl'vj g'r tqr qugf "cr r tqcej "hqt" 'vj g'cwqo cvk" gptlej o gpv'qh'NO H'wcpf ctf k' gf "grgext'qple" f levk'pct'kgu'd { "vj g'ugo cpvle"ercuugu'0'Vj g'hqmqy kpi "hki vt g'3" kmwut'cvgu'vgr u'qh'vj ku'cr r tqcej 0'

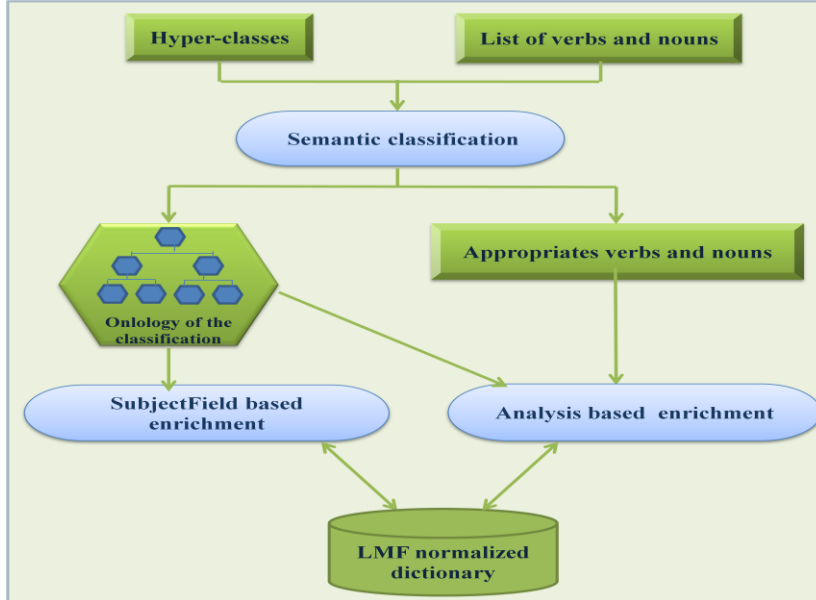


Figure 1: Proposed approach

Vj g'r tqr qugf "cr r tqcej "ku'eqo r qugf "qh'vj tgg"vgr u'c'ugo cpvle"ercuuk'hecvkqp"cpf "vy q"r j'cugu'qh'cwqo cvk" gptlej o gpv'0'Vq'cee'qo r r'kuj "vj g'clo u'qh'vj g'ugo cpvle"ercuuk'hecvkqp. "vj ku'vgr "tgs vkt gu'vj g'j { r gt/ercuugu'qh'vj g' I cuvqp" I tquu'ercuuk'hecvkqp"cpf "c"rku'qh'xgtdu'cpf "pqwpu'qh'vj g'uwf k'gf "rcpi vci g"kp"kp'w0'Vj g'tguwu'qh'vj g' ugo cpvle"ercuuk'hecvkqp"vgr "ctg"vj g'qp'v'q' { "qh'vj g'ercuuk'hecvkqp"cpf "c"rku'qh'cr r tqr t'cv'g'xgtdu'cpf "pqwpu' ej ctcev'gt'k' kpi "vj ku'ercuuk'hecvkqp'0' Y j'gt'gcu." vj g' Uwdl'ge'v'k'gf "dcugf" gptlej o gpv' wugu' vj g' qp'v'q' { "qh'vj g'ercuuk'hecvkqp"vq" gptlej "vj g'NO H'pqto crk' gf "f levk'pct'kgu'd { "kf gp'v'k' kpi "ugo cpvle"ercuugu'0'Vj g'cpcn'f uki"dcugf" gptlej o gpv'tgs vkt gu'cej k'x'kpi "vj g'gptlej o gpv'qh'vj g'NO H'pqto crk' gf "f levk'pct'kgu'd { "dq'vj "vj g'qp'v'q' { "qh'vj g'ercuuk'hecvkqp"cpf "vj g'ku'qh'cr r tqr t'cv'g'u'xgtdu'cpf "pqwpu'k' gf "r t'g'x'k'w'u'0'

4.1. Semantic classification

4.1.1. Basic concept

Qwt'ugo cpvle"ercuuk'hecvkqp"ku'dcugf "qp"vj g'uwf l'gu'qh'vj g'I cuvqp" I tquu"4I tquu."3; ; 6+'ugo cpvle"ercuuk'hecvkqp" *ug'ugev'kqp"4+0'Vj ku'ercuuk'hecvkqp" wugu'vj g'r t'gf'lec'v'cti wo gpv'ut'vewt'g"vq"ercuuk'hecvkqp { "ngz'le'c'n'wpku'0'Vj wa." vj g' uko r ng'ug'v'peg't'gr t'gug'pu'vj g'o k'pko wo "v'p'k'qh'cpcn'f uki'0'k'p'gf g'gf. "vy q'o clqt'ugo cpvle"ercuugu'ej ctcev'gt'k' g'vj ku'ercuuk'hecvkqp"pco gn'<'vj g'ugo cpvle"ercuugu'qh'r t'gf'lec'v'g'u'cpf "vj g'ugo cpvle"ercuugu'qh'cti wo gpv'0'J qy g'x'gt." r t'k'qt"vq"vj g'qdl'ge'v'ercuugu."cpf "dcugf"qp'u'p'v'ev'le"hg'cwt'gu."vj g'ercuuk'hecvkqp"o cl'p'v'k'pu"ercuugu'vj cv't'gi t'qwr "cm' r t'gf'lec'v'g'u'vj cv'vj ct'g'vj g'eqo o qp'u'p'v'ev'le"dg'c'x'k'q'tu'pco gf "J { r gt/ercuugu'0'Vj wa."j { r gt/ercuugu'qh'r t'gf'lec'v'g'u." ur g'ek'k'gf "d { "vj ku'ercuuk'hecvkqp"ct'g'<SCE V'KQP. "GXGP V."UVCVG"cpf "RTGF K'EC V'K'G"J WO CP 0"Y j'k'g"j { r gt/ercuugu'qh'cti wo gpv' ct'g'< \$J WO CP." EQPETGVG." RNCPUVU." CP KOCNU." VIOG." TGPVCN" cpf "CDUVTCEV0"Vj g'ug"j { r gt/ercuugu'ct'g"uwdl'ge'v'vq"uwd"ercuuk'hecvkqp"u'd { "o g'epu'qh'cti wo gpv'r gto w'ev'k'pu" *f'k'w'k'w'k'p'c'n'et'k'g't'c'cr r g'ct'k'pi "kp"q'p'g'qt"o q't'g'r'qu'k'k'p'u'qh'cti wo gpv't'g'rc'v'f"vq"ci'k'x'gp"r t'gf'lec'v'g'u'0'Vj wa."k'i'c" r gto w'ev'k'p'qh'c"pq'w'd { "cp'q'vj g't'eq'p'v'k'w'g'u'vq"c't'w'v'w'g'qh'vj g'o g'ep'k'pi "qh'c"r t'gf'lec'v'g'u'ug'p'ug."vj g'c"p'gy" qdl'ge'v'ercuuk'hecvkqp"ku' t'gs vkt gf "vq" dg" et'g'c'v'f'0'Vj g'ug" qdl'ge'v'ercuugu"cm'qy "j k'j r k'j v'k'pi "vj g'f'k'ht'g'p'v' wugu'qh'c" r qn'ugo q'w'r t'gf'lec'v'g'u"

⁴"y y (u'q'ec'v'q'ti "

4.1.2. Steps of the semantic classification

Y g'r tqr qug'kp'hi wtg'4'vj g'i gpgtci'ugo cpvle"ercuukhlecvkqp'r tqegu0'

"
"
"
"
"
"
"
"
"
"
"
"

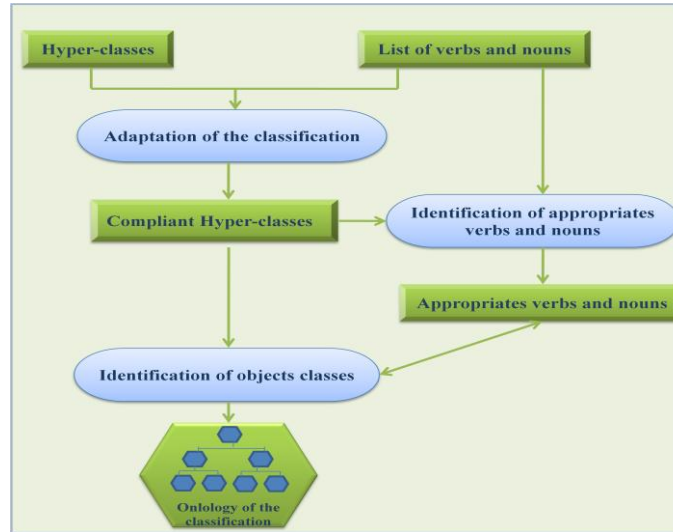


Figure 2: Semantic classification

Vj g'r tqegu'qh'vj g'r tqr qugf "ugo cpvle"ercuukhlecvkqp"ku'tgci' gf "o cpwcm' d { "c"npi wku'K'eqo r qugf "d { "vj tgg" ugr u< "k" Cf cr vvkqp" qh' vj g"ercuukhlecvkqp." "k" K gpvkhlecvkqp" qh' cr r tqr tkvgu" xgtdu" cpf "pqwpu" cpf " "k" K gpvkhlecvkqp"qh'qdlgev'ercuug0'

ko Cf cr vvkqp"qh'vj g"ercuukhlecvkqp<

Hyper-classes"qh'vj g"I cuqp"I tqu'u'wvf lgu"ugg"ugevkqp"600+"cpf "c"list of verbs and nouns"qh'vj g"u'wvf lgf "rcpi wci g'r gthqto "vqi gvj gt"kp"qtf gt"vq"cee'qo r rkuj "vj g"cf cr vvkqp"qh'vj g"cf cr vvkqp"qh'vj g"ercuukhlecvkqp"u'gr 0' Eqpukf gt'kpi "vj cv'vj g"ugo cpvle"ercuukhlecvkqp"ku'r gthqto gf "d { "c"npi wku"vj ku'u'gr "tgs vkt gu"vj g"cdk'k'ku'qh'vj ku' g'zr gt'v'cpf "vj g"u'p'v'le"hg'cwt gu"qh'vj g"u'wvf lgf "rcpi wci g"kp"qtf gt"vq"u'wvf { "vj g'r quukd'k'k'v' "qh'vj g"cf cr vvkqp"qh' vj g"ugo cpvle"ercuukhlecvkqp"qp"vj g"u'wvf lgf "rcpi wci g'0'Qp"vj g"dcuku'qh'u'p'v'le"hg'cwt gu"qh'vj g"u'wvf lgf "rcpi wci g." vj g"gzr gt'v'ecp"kf gpvkh' "pgy "j { r gt/ercuugu"cr r tqr tkvgu"vq"vj g"v'g'v'g' "rcpi wci g."f g'v'g'qt"t'g'p'co g"vj g"gz'k'v'pi " ugo cpvle"j { r gt/ercuugu0'Vj gt'ghq'g'vj g"compliant hyper-classes tgr t'g'v'v'vj g't'g'u'w'v'qh'vj ku'u'gr 0'

kkk K gpvkhlecvkqp"qh'cr r tqr tkvgu" xgtdu"cpf "pqwpu"<

Qp"vj g"dcuku'qh'vj g"pqxgn'rkuv"qh'hyper-classes"kf gpvkh'gf "kp"vj g'r t'g'x'k'w'u'ugr. "t'g'v'g' "vq"vj g"ur g'k'k'le"u'wvf lgf "rcpi wci g." vj g" kf gpvkhlecvkqp" qh' cr r tqr tkvgu" xgtdu" cpf "pqwpu" v'cng" r r'ceg0' Vj ku' u'gr " clo u" vq" f g'v'g'v' vj g" appropriate list of verbs"and nouns"ej ctcev'gt'k' kpi "g'cej "j { r gt/ercuugu"qh'vj g'r tqr qugf "ugo cpvle"ercuukhlecvkqp0'

kkk K gpvkhlecvkqp"qh'qdlgev'ercuugu"<

Vj g"qdlgev'ercuug'eq'p'gr v't'gr t'g'v'v'v'vj g"ej ctcev'gt'k'v'le"qh'vj g'r tqr qugf "ugo cpvle"ercuukhlecvkqp0'Vj wu."vj g"clo "qh' vj ku'u'gr "ku'vj g"identification of object classes hqt"g'cej "ugo cpvle"ercuug0'Vq"cee'qo r rkuj "vj ku'qdlgev'k'g."vj ku'u'gr " tgs vkt gu"vj g"compliant"hyper-classes"qh'vj g"u'wvf lgf "rcpi wci g"cpf "vj g"rkuv"qh'appropriates verbs and nouns" t'geqi p'k'gf "kp"vj g"rcuv'ugr 0'Vj g't'g'u'w'v'qh'vj ku'u'gr "ch'g'v'r t'g'f l'ec'v'gu"ugo cpvle"ercuugu"vq"cu'y g'm'cu'cti wo gpw0' k'p'g'gf. "vj g"gzr gt'v'd'g'p'gh'ku"ht'qo "vj g"u'p'v'le"hg'cwt gu"qh'vj g"u'wvf lgf "rcpi wci g"kp"qtf gt"vq"kf gpvkh' "qdlgev'u'ercuugu t'g'v'k'pi "t'g'v'g'v'g'n"vq"j { r gt/ugo cpvle"ercuugu"qh'r t'g'f l'ec'v'gu"cpf "cti wo gpw0'Cu"j { r gt/ercuugu."vj g" kf gpvkhlecvkqp"qh'vj g"qdlgev'ercuugu"q'w'eqo gu"cu"rkuv"qh'xgtdu"cpf "pqwpu"ej ctcev'gt'k' kpi "g'cej "qdlgev'ercuug0'Vj ku'rkuv' r gthqto u" vq" w'f'v'g' vj g" rkuv' qh' cr r tqr tkvgu" xgtdu" cpf "pqwpu" qh' vj g"ercuukhlecvkqp0' Cp" ontology of the classification"vj cv'u't'gi tqw'u'cm'eqo r r'k'p'v'j { r gt/ercuugu"cpf "qdlgev'ercuugu"t'g'v'g' "vq"vj g"u'wvf lgf "rcpi wci g" t'gr t'g'v'v'vj g't'g'u'w'v'qh'vj ku'u'gr 0'

"

4.2. Enrichment of LMF standardized dictionaries

Chgt "f g x g n r k p i " c " u g o c p v e " e r c u u k h e c v k p . " v j g " g p t l e j o g p v r t q e g u u " q h ' v j g " u c p f c t f k g f " N O H ' f l e v k p c t k g u ' y k j " u g o c p v e " e r c u u g u ' y k m ' v c n g ' r m e g O ' K ' e q o r q u g f " q h ' v y q ' o c k p ' r j c u g u < " k i " v j g " U w d l g e v H g r f " d c u g f " g p t l e j o g p v ' v j c v ' d g p g h k g f " H t q o " v j g " N O H ' f l e v k p c t k g u " u t w e w t g . " r c t l e w r c t n { " H t q o " v j g " w u g u " f q o c k p u " t g r v g f " v q " o g c p l p i u " q h ' r g z l e c n g p v t k g u " k i " v j g " c p c n { u k u " d c u g f " g p t l e j o g p v ' v j c v ' w u g u " h g c w t g u " q h ' v j g " q d v c k p g f " u g o c p v e " e r c u u k h e c v k p o "

4.2.1. Subject Field based enrichment

Vj ku'gptlej o gpv'ku'dcugf "qp"vj g"hgfn "oUwdlgevHgrf o'ceeqt f kpi "vq"vj g"NO H'o qf grO'Cu'uj qy p"lp"hi wt g"5."k' epukuu'qh'vy q'vgr u'f guetldgf "cu'hqmy <"

kuUgctej kpi "ugpugu'y kj "oUwdlgevHgrf o'<"vj g"fo ckpu'qh'wugu"ht"gej "oSenses ""qh'rgzlecn'gpvtkgu"lp"LMF normalized dictionary"ctg"tgr t g u g p v g f " v j t q w i j " c " e r c u u " p c o g f " o U w d l g e v H g r f o ' V j g " c l o " q h ' v j k u " u g r " k u " v j g " g z t c e v k p " H t q o " v j g " f l e v k p c t { , S e n s e s " t g r v g f " v q " t g c v g f " r g z l e c n g p v t { " e q p v c k p i " v j g " o U w d l g e v H g r f o ' H g r f o "

kuUgctej kpi "ugpugu'y kj "oUwdlgevHgrf o'<"vj g"fo ckpu'qh'wugu"ht"gej "oSenses ""qh'rgzlecn'gpvtkgu"lp"LMF normalized dictionary"ctg"tgr t g u g p v g f " v j t q w i j " c " e r c u u " p c o g f " o U w d l g e v H g r f o ' V j g " c l o " q h ' v j k u " u g r " k u " v j g " g z t c e v k p " H t q o " v j g " f l e v k p c t { , S e n s e s " t g r v g f " v q " t g c v g f " r g z l e c n g p v t { " e q p v c k p i " v j g " o U w d l g e v H g r f o ' H g r f o " k o k g p v t k g u " q h ' u g o c p v e " e r c u u g u < " c " r t g t g c v o g p v ' t g c r k g f " q p " " v j g " q d v c k p g f " u g o c p v e " e r c u u k h e c v k p " c p f " v j g " g z l e c n g p v t k g u " o U w d l g e v H g r f o " l p " c p " N O H ' u c f c t f k g f " f l e v k p c t { " e c p " o c f g " c " f k g e v { " e q t t g u r q p f g p e g " " d g y g g p " v j g " j { r g t / u g o c p v e " e r c u u g u " q t " v j g " q d l g e v e r c u u g u ' y k j " v j g " o U w d l g e v H g r f o ' K i ' v j k u " k u " v j g " e c u g . " v j k u " u g r " k f g p v t k g u " v j g " u g o c p v e " e r c u u " H t q o " v j g " o n t o l o g y o f t h e c l a s s i f i c a t i o n " t g r v g f " v q " v j g " h q w p f g f " " S u b j e c t F i e l d " " c p f " w f c v g u " v j g " N O H ' u c f c t f k g f " f l e v k p c t { " d { " v j g " c f f k k p " q h ' v j g " t g v c k p g f " u g o c p v e " e r c u u " v j g " e q t t g u r q p f k p i " S e n s e O "

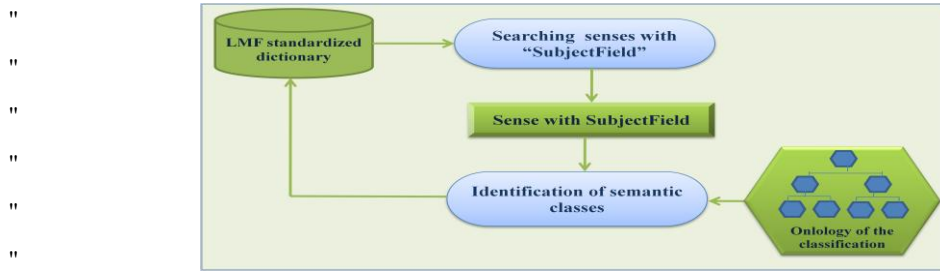


Figure 3: Subject Field based enrichment

4.2.2. The analysis based enrichment

Vj g'cpcn{uku'dcugf "gptlej o gpv'wugu'y g'hcwttgu'qh'vj g'tgcckp g f " u g o c p v e " e r c u u k h e c v k p o " V j g " h q m y k p i " h i w t g " 6 " k n w u t c v g u ' v j g " u g r u " q h ' v j k u " h k p f " q h ' g p t l e j o g p v o "

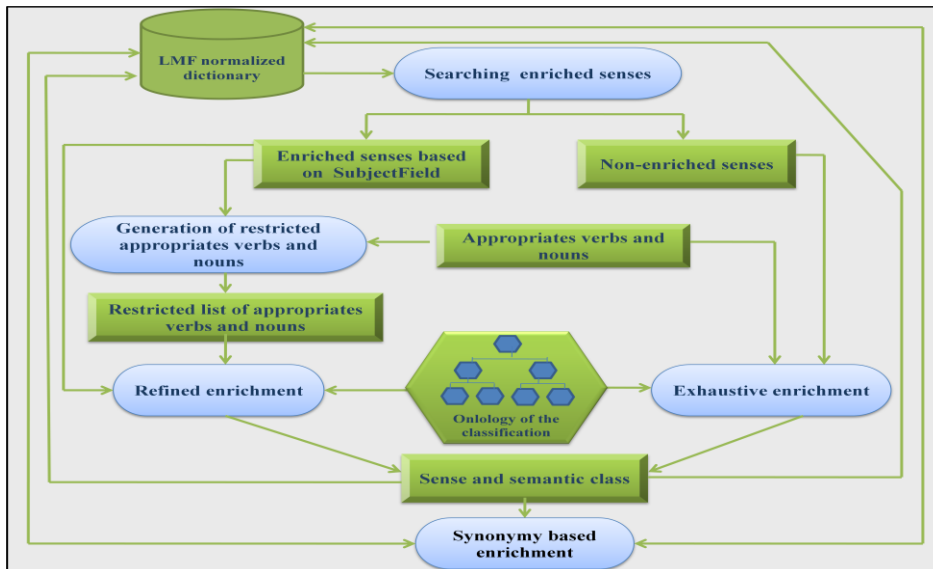


Figure 4: Analysis based enrichment

Vj g'kuu'qh'cr r t q r t l e v g u " x g t d u " c p f " p q w p u " c p f " v j g " q p v m i { " q h ' v j g " e r c u u k h e c v k p " t g r t g u g p v ' v j g " k p r w " q h ' c p c n { u k u " d c u g f " g p t l e j o g p v o " k u " e q o r q u g f " d { " v j g " h q m y k p i " h x g " u g r u < "

j { r gt/ercuu'ku'cnuq'tgckpgr'ht'vj g'Ctdle'rcpi wci g'htqo 'vj g'ercuukhlecwkp'qh'I cuqrp'I tquu'Co qpi 'vj g'qdlgev'ercuu'gdgrpi kpi 'v'vj g'EQPETGVGö'j { r gt'ercuu'y g'pqv'vj g'öErvj guö'ercuu'Kpf ggf. 'vj g'Ctdle'xgd"\$ \$' \$q'y gct\$. "tgr tgupv'vj g'cr r tqr tkvg'xgd"ej ctcevtk kpi "vj ku'qdlgev'ercuu'Vj wu. "ppg"o gcpkpi "qh'vj ku'xgd" f guetkg'cp"\$CEVIQP\$"tgerk gf "d{ "c"htuv'\$J WOCP\$"cti wo gpv'cpf "j ki j rki j vki "cpqv'gt "\$EQPETGVG\$"cti wo gpv'Vj g'gzco r rg'dgrny 'knuw'cvgu'vj tgg'ugpv'pgeu'f g'ckrgf 'vj g'o gcp'qh'vj g'\$ \$'\$q'y gct'\$xgd<

- *3+ "Nedkuc'cn~il.miydu alqub~açaña" Vj g'r wr kl'y gctu'vj g'j cv'
- *4+ "Nedkuc'cn~il.miydu 'cn~ufaHaña" Vj g'r wr kl'y gctu'vj g'cr r rg'
- *5+ "Nedkuc'cm cC 'u alqub~açaña" Vj g'y cvgt'y gctu'vj g'j cv'

Cm'ugpv'pgeu'*3+,*4+cpf '*5+ctg'u{ pvcw'ecmf "eqttge'0'Dw'qpn{ 'vj g'ugpv'pgeu'*3+ku'ugo cpw'ecmf "ceegr vcdrg'0'Kpf ggf. 'kp'ugpv'pgeu'*3+c"öpupil""ecp"öwear"a hat", y j krg'kp'ugpv'pgeu'*4+c""Pupil""ecppqv'y gct'cp""apple"" dgecwug'cp""apple""ku'cp"öCnko gpwö'uq'k'ecp'dg'gcvp'dw'pqv'y gcp'0'Y j gtgcu. 'kp'ugpv'pgeu'*5+vj g""water""ku'cp"öCnko gpv'y cvgtö"cpf "ecppqv'dg'y gcp'0'Vj qug'gzco r rgu'g'zr rkecv'vj g'tgs wktgo gpv'qh'vj g'etgcv'kp'qh'vj g'\$Ervj gu\$'cpf "vj g'\$Cnko gpwö'qdlgeu'ercuu'wpf gt'vj g'\$EQPETGVG\$"j { r gt/ercuu'Vj wu. 'j g'\$Ervj gu\$'qdlgev'ercuu'kpenmf gu'cm'pqwpu'vj cv'ecp'dg'y qtp'd{ "c"\$J WOCP\$'Ctdle'xgtd'uwej "cu<" /xalaça/to undress" " /Air.tady'/to dress", " /labisa/to wear", cpf "pqwpu' rkn<" /kisaA'/cloth", ö hkdCuly gctö, " /haw.jbü/dress" ej ctcevtk g'vj g'öErvj guö qdlgev'ercuu'0

Cti wo gpw' kpu'cpegu' qh' 'vj g' öErvj guö" qdlgev'ercuu' ecp" dg<" /HidaA'/shoes", " /naç.l/sock", " /xuf-û/slipper", " /qlub~açañû/hat", " /sir.waAlû/pant", " /qamiuSû/shirt"0' Vj wu. 'vj g' cr r tqr tkvg'xgtdu'qh'vj g'öErvj guö" qdlgev'ercuu' rkn<" " /xalaça/to undress" " /Air.tady'/to dress"" " /labisa/to wear" ecp'dg'eqttge'v' "k'p'v'w'gu'vj g'cti wo gpw'kpu'cpegu'rkun'dghqtg. Dw'kp'Ctdle'rcpi wci g. " uqo g'xgtdu'ugrgev'htqo "vj g'öErvj guö"cti wo gpw'kpu'cpegu'c"ur gekhle"qpgu'dw'ecppqv'wug'cm'qh'vj go 0'Hqt" gzco r rg.'vj g"* I'Clpöc c'n'I'v'y gct'uj qgu. "ICk'öcf { "k'y gct'uj qgu+xgtdu'ecppqv'r tgegf g'cm'qh'vj g'cti wo gpw'kpu'cpegu'"\$Ervj gu\$"dw'qpn{ ">uj qgu@> @'Vj wu. 'vj g'ugpv'pgeu"* " +*j g'y gct'uj qgu'uj kv'+ku'ugo cpw'ecmf "k'eqttge'v'dgecwug"* I'Clpöc c'n'I'v'y gct'uj qgu+ku'cp'cr r tqr tkvg'xgd"v">uj qgu@> @'ercuu'cpf "*" 1 alqamiuSû/shirt+"f qgu'pqv'tgr tgupv">uj qgu@> @dw'tcv'gt">htkpi wgu@> "@ Vj gtgqtg. "k'ku'pgeu'ct { "v'etgcv'y q'qdlgeu'ercuu'wpf gt'vj g'\$Ervj gu\$"pco gr{ ">uj qgu@> @'cpf ">" htkpi wgu"@> @

Vj g'vcdrg'3'kp'hqmqy kpi "uwo o ctk'gu'vj g'r t'gxlqwu'kf gc<"

Table1: Appropriate verbs for the <shoes> <الحذاء> and <fringues> <ثياب> object class

| | Examples of nouns | verbs | لبس | احتذى | خلع | ارتدى | ارتحل |
|------------------------|---|-------|-------------------|---------------------------|----------------------|-----------------------|-----------------------------|
| | | | labisa to wear | AiH.tady to wear shoes | xalaça to undress | Air.tady' to dress | Ain.taçala to wear shoes |
| <Clothes> object class | <shoes> <الحذاء> خف/xuf-û/ slipper | | | | | | |
| | <shoes> <حذاء> HidaA'/ shoes | | | | | | |
| | <sock> <ناعل> naç.l/ sock | | | | | | |
| | <fringues> <ثياب> قميص/qamiuSû/ shirt | | | | | | |
| | <fringues> <سروال> sir.waAlû/ pant | | | | | | |
| | <fringues> <قبعة> quö~açañû/ hat | | | | | | |

5.2.2. Results of the Arabic semantic classification

Kp'vj ku'ugev'kp'y g'r tgupv'cp'gzco r rg'qh'vj g'ugo cpw'ercuukhlecwkp'qrvqmi { 'ht' 'Ctdle'rcpi wci g'0'Vj g'dgrny " hki wtg" 7" knuw'cvgu' uqo g'tgeqi pk'gf " j { r gt/ercuu'cpf " qdlgev'ercuu'wukpi "vj g' r tqr qugf " r tqegu' qh'vj g' ugo cpw'ercuukhlecwkp'"ugg'ugev'kp'6ö+0'Vj ku'hki wtg'ku'etgcv'f'y kj 'vj g'QY N'qrvqmi { 0

Table3: Examples of available Subject Field in the Arabic standardized dictionary

| Subject Field | | |
|---------------|----------------|------------------|
| In Arabic | Transliterated | In English |
| حَيَوَان | J c{cy cCp" | Cpko crl' |
| حشرة | J c-etc " | Kpugev' |
| نبات | pcdcCv' | Rrnpv' |
| هندسة | j cpff cuc " | I gqo gvt { "" |
| طبّيح | Vcd0z" | Ewkpct { "" |
| جغرافيا | lw 0cCh {cC" | I gqi tcr j { "" |
| موسيقى | o wukf s cC" | O wuke" |
| رياضة | tk{cCF c " | Ur qt v'" |
| طب | Vkd" | O gf lekpg" |
| عسکر | cu0het" | O krkct { "" |

"

Kp"vj g" Ctdle" NO H' pqto crk gf "f lkvqpt {." vj g" \$ \$" \$Cpko crl\$ \$J c{cy cCp\$ cpf " vj g" \$ \$" \$J c-etc \$" \$Kpugev\$ "SubjectFields" ecp" dg" i tqwr gf "kpq" vj g" \$Cpko crl\$ j { r gt/ercuu'K'ku'ko r qtcpv'vq "kpf kecvg'v cv'vj g" \$" \$J c-etc \$" \$Kpugev\$ "SubjectField" eqttgur qpf u" f kgevw " vj g" qdlgev'ercuu' pco gf " oKpugevo" cpf " vj g" \$" \$J c{cy cCp\$ \$Cpko crl\$ "SubjectField" o c { "eqttgur qpf " vj g" qdlgev'ercuu'gu- "oDktf o. "oTqf gpwö. "ötgr vkrguö" cpf " oCs wöle/cpko cniö"cu'uj qy u'lp' hki wtg70'

Vj g'hqmvy kpi 'hki wtg'8' kmwutcvgu'cp'gzco r r ng'qh'vj g'UwdlgevHlgrf "dcugf "gptlej o gpv0'

"

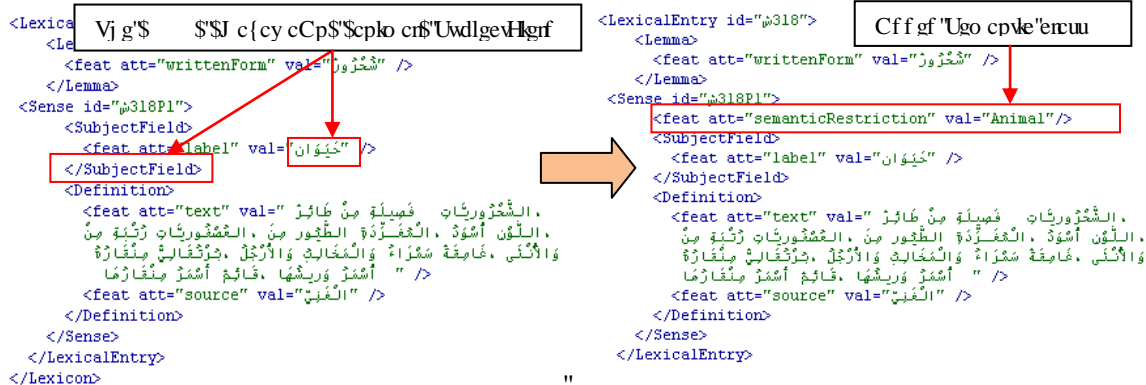


Figure 6: The SubjectField based enrichment applied to a sense of lexical entry

5.3.3. Experimentation of the analysis based enrichment

Vj g"cpnq uku"dcugf "gptlej o gpv'ku" uwdf kxkf gf "kpq" vj q" nkp" qh" gptlej o gpv' Vj g" hkuv' gptlej o gpv' cr r qkpvf "tghkpgf "gptlej o gpv'tgs wktgu'hqt" vj g" r tqi tguu'qh'ku'r tqeguu'vj g"tgutkevfg "rkuv'qh'xgtdu"cpf " pqwpu'kp'qtf gt"vq"tghkpg"vj g" r tko ct { "gptlej o gpv'ectkfg "qw'kp" vj g"UwdlgevHlgrf "dcugf "gptlej o gpv'Qt" vj g"ugeqpf "gptlej o gpv'ku"gzj cwukxg." eqpegtkpi "qpn" pqp/gptlej gf "ugpugu." wugu" vj g" cr r tqr tlcvg" xgtdu"cpf "pqwpu'qh'vj g" ugo cpvle"ercuu'khec'kq"kp"qtf gt"vq"tgcrk'g"vj g" ugo cpvle" gptlej o gpv'qh'vj g" f lkvqpt { 0'

Vj g'cdrg'6' i kxgp'kp'vj g'hqmvy kpi ."eqpvkpu'vj g'tgutkevfg "rkuv'qh'cr r tqr tlcvgu'xgtdu"cpf "pqwpu'tgrvfg " vj g" \$ \$" \$J c{cy cCp\$ \$Cpko crl\$ "cti wo gpw'j { r gt/ercuu'0'

"

"

Table4: Restricted list of appropriate verbs and nouns of “Animal” hyper-class

| Hyper class | Restricted list of appropriate nouns and verbs | Object classes | Restricted list of appropriate nouns and verbs | Sub-Object-classes | |
|-------------|--|----------------|--|------------------------------|----------|
| Animal | " | " | / | / | |
| | " | " | | | Dkf |
| | " | " | | | Tqf gpvu |
| | " | " | Tgr vrgu | | |
| | " | " | " | Hkj | |
| | " | " | Cs wvke/ cpko cnu | حيثان حوت الثدييات | Rluegu |
| | " | " | " | Qvj gtu/cs wvke/ cpko cnu | |

Vj g"cr r rkecvkp"qh'vj g"r tqegu"qh'vj g"tghkpgf "gptlej o gpv'd{ "wukpi "vj g"tgvtkevgf "rkuv"qh'xgtdu"cpf "pqwpu"*cdng7"kp" {gmqy +qp"vj g"ruv"gzvcevgf "htci o gpv'wugf "kp"vj g"UwdlgevHkgrf "dcugf "gptlej o gpv' *hki wtg'8+ecp'i kxg'vj g"gptlej o gpv'r tgvpgv"kp"vj g'hmqy kpi "hki wtg'90"

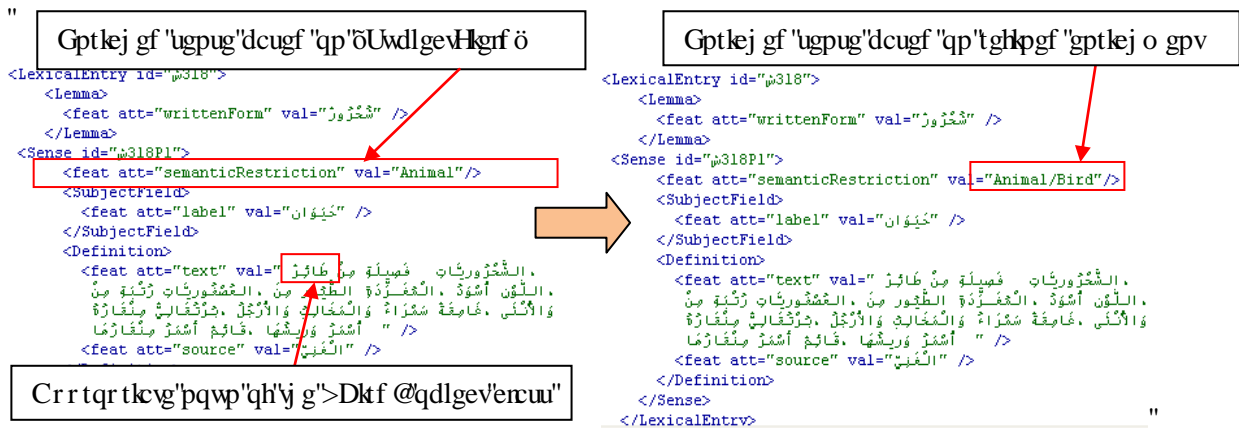


Figure 7: The refined enrichment applied to a sense of lexical entry

Kp"vj g'hmqy kpi .y g"r tgvpgv'cp"gzr gtko gpcvkv"qh'vj g"gzj cwvkvxg"gptlej o gpv'wukpi "vj g"cr r rtr tkcv" xgtdu"cpf "pqwpu"cr r rkef "vq"pqw/gptlej gf "ugpugu"Vj g"cpn{uku"qh"Contexts"cpf "Definitions qh'ugpugu" tgrv"vq" c"ngzkecn'gpx { "kp"vj g"Ctdle"NO H'ucpfc tf k gf "f levkqpc { "d{ "wukpi "vj g"cr r rtr tkcv" xgtdu" cpf "pqwpu"*cdng6+ecp'kf gpxh{ "vj g'tgrgxcpv'ugo cpvke'ercuu'0

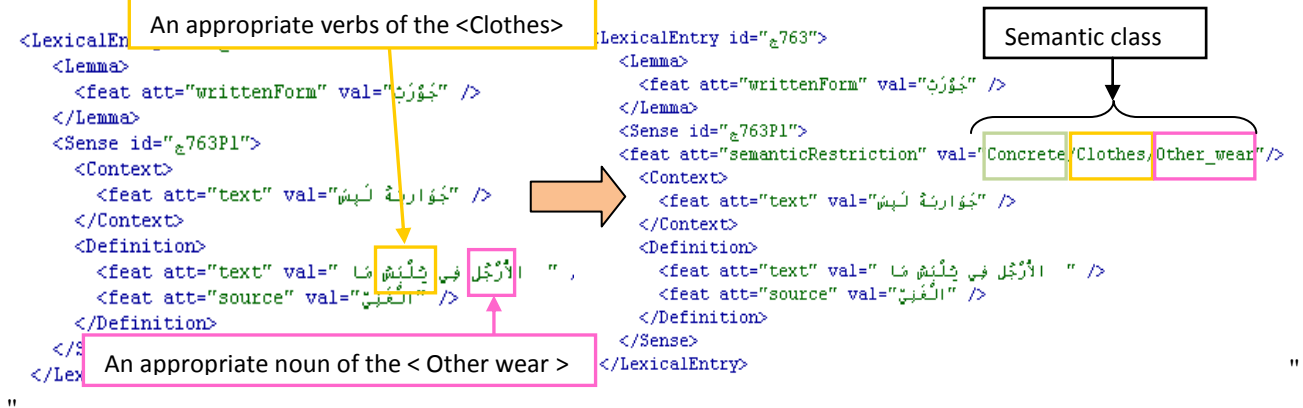


Figure 8: The exhaustive enrichment applied to a sense of lexical entry

5.4. Results

Vq'guv'j g'r gthqto cpeg'qh'j g'ecttkgf "qw'gZR gtlo gpcv'kqp."y g'j cxg'tgcn' gf "c'uc'v'k'ec'n'gxcn'v'k'p0"
 Qw"Ctdle"ucpf ctf k' gf "f'lev'k'pct {"eqpv'k'p"lp"v'cn'56222"ngz'k'ec'n'g'p't'k'g'u"l'p'en'w'f'k'p'i "84379"ugp'ug'u0
 Eqpeg't'p'k'i "j g" Uwdlge'v'k'grf "dcugf "gpt'lej o g'p'v' g'z'r gtlo gpcv'k'p="y g"j cxg" wugf "6"öUwdlge'v'k'grf ö"
 *C'p'k'o cn"K'p'ugev."R'c'p'v'c'p'f "E'w'k'p'c't {"+co q'p'i "j g"3; "cx'k'c'rd'ng"lp"j g"Ctdle" f'lev'k'p'c't {"0'c'p'f "h'q't"j g"
 c'p'c'n'f' u'k'u"dcugf "gpt'lej o g'p'v'y g"j cxg"ej q'leg"q'p'n'f "j g"öEQP ETGVGö"j {r gt"er'cuu"cp'f "ur gek'm'f "j g"
 öEr'q'j guö'q'dlge'v'er'cuu"v'q'g'z'r gtlo g'p'v'j g'r t'q'eg'u'q'h'j k'u'h'k'p'f "q'h'g'p't'lej o g'p'v0"
 Vj g"v'c'd'ng"7" d'g'ny "i k'x'g'u"j g"uc'v'k'ec'n'gxcn'v'k'p" q'h'j g"öUwdlge'v'k'grf ö" c'p'f "j g" c'p'c'n'f' u'k'u" dcugf "
 g'p't'lej o g'p'v0"

Table 5: Evaluation of the enrichment

| | | SubjectField based enrichment | Analysis based enrichment (exhaustive step) |
|--------------------------|----------|-------------------------------|---|
| Number of Subject Field" | Animal | 3; 9" | " |
| | Insect | 3; " | |
| | Plant | 464" | |
| | Culinary | 5; " | |
| Total " | | 497 | |
| Correct assignment | | 676" | ; 2" |
| Incorrect assignment | | 65" | 74" |
| Recall | | ; 3.56" " | 48" " |
| Precision | | ; : ' " | 85" " |

6. Conclusion and perspectives

K'p"j k'u'r cr gt."y g"j cxg"r t'q'r qugf "cp"cr r t'q'cej "h'q't"j g"cw'q'o c'v'e" g'p't'lej o g'p'v'q'h"NO H'uc'p'f c'f k' gf "
 f'lev'k'p'c't'k'u'y k'j "u'g'o c'p'v'e"er'cuu'g'u0Vj k'u'c'r r t'q'cej "k'u'eq'o r qugf "q'h'c"u'g'o c'p'v'e"er'cuu'k'h'ec'v'k'p"dcugf "q'p"
 j g"I cu'q'p"I t'q'u'u'w'f'k'g'u"cp'f "v'y q"v'r gu"q'h"g'p't'lej o g'p'v0Vj g"ht'u'v'g'p't'lej o g'p'v'p'c'o gf "Uwdlge'v'k'grf "
 dcugf "g'p't'lej o g'p'v."v'c'ng'u" c'f x'c'p'c'i gu"ht'q'o "j g"u't'w'ew't'g" q'h"cp"NO H'f'lev'k'p'c't {"y j gt'g" o g'c'p'k'i u"
 eq'p'v'k'p"j g" f'q'o c'k'p" q'h" wug" q'h" c" ngz'k'ec'n' g'p't {"0' Vj g" u'g'eq'p'f "g'p't'lej o g'p'v' ec'ng'f "c'p'c'n'f' u'k'u" dcugf "
 g'p't'lej o g'p'v."w'ug'u"j g"hg'c'w't'g'u"q'h'j g"r t'q'r qugf "u'g'o c'p'v'e"er'cuu'k'h'ec'v'k'p"dcugf "q'p"cr r t'q'r t'k'v'g'u"x'g't'du"
 c'p'f "p'q'w'p'u"ur gek'h'f'k'p'i "g'cej "u'g'o c'p'v'e"er'cuu"cp'f "cr r t'q'gf "v'q"j g"cx'k'c'rd'ng"v'g'z'v'eq'o r q'p'g'p'u"lp"j g"
 f'lev'k'p'c't {"0"

Y g'ecttkgf "qw'gZR gtlo gpcv'k'qp."d {"w'k'p'i "cp"cx'k'c'rd'ng"Ctdle"uc'p'f c'f k' gf "f'lev'k'p'c't {"0'Vj g"q'd'v'k'p'gf "
 t'g'u'w'u" ctg" uc'v'k'h'f'k'p'i " eqpeg't'p'k'i "j g" Uwdlge'v'k'grf "dcugf "g'p't'lej o g'p'v0 Vj g" u'f'p'q'p {"o {"dcugf "
 g'p't'lej o g'p'v'ec'p"t'g'f w'eg"j g"g'p't'lej o g'p'v'g'h'q't'v'c'v'j k't'f u'd'g'ec'w'ug"q'p"cx'g't'c'i g."j g"u'f'p'q'p {"o {"t'g'r'v'k'p"
 eq'p'p'g'ew'j t'g'g'q't' b' q't'g'ug'p'ug'u0"

K'p"j g"hw'w't'g."y g'q'r v'g'f "v'q"cej k'g'x'g"j g"z'r gtlo gpcv'k'qp"q'p"j g"q'v'g't'u'g'o c'p'v'e"er'cuu'g'u'q'h'j g'r t'q'r qugf "
 u'g'o c'p'v'e"er'cuu'k'h'ec'v'k'p"ht"Ctdle"r'p'i w'c'i g'c'p'f "v'q"eq'o r ng'v'j g"t'g'u'v'q'h"Uwdlge'v'k'grf "g'z'k'ng'f "lp"j g"
 Ctdle"NO H' uc'p'f c'f k' gf "f'lev'k'p'c't {"0' K'p" c'f f'k'k'p'."y g" eq'p'k'f'gt" k'o r t'q'x'k'p'i "j g" c'p'c'n'f' u'k'u" dcugf "
 g'p't'lej o g'p'v' d {"c'f f'k'p'i "o q't'g" g'h'k'eg'p'v' u'f'p'c'v'e'k'g'u'g'o c'p'v'e" c'p'c'n'f' u'k'u0' H'k'p'c'm'f'."y g" h'q't'g'ug'g"j g'c'v'j g"
 g'p't'lej o g'p'v'ec'p"q'h'ht"j g"hg'z'k'k'k'v'f "v'q"et'g'v'g'p'gy "q't'k'g'p'v'g'f "x'g't'k'p'u"q'h'j g"u'g'o c'p'v'e"n'p'q'y r'g'f' i g"
 p'g'g'f'gf "h'q't" f'k'h'g't'g'p'v'P RN'cr r r'k'ec'v'k'p'u0"

"
 "
 "

Reference

- Eqj gp."F0*3; ; 8-0Ncy ."uqelcn'r qrlk{."cpf "xlqngpeg<Vj g"ko r cev"qh'tgi kqpcn'ewmwgu0*Journal of personality and Social Psychology*, '92Q 83/; 9: 0"
- "Eqpf co kpgu"C0*4227+0U2 o cpvks wg'gv'eqtr wu."s wmggu't gpeqpt wu'r quukdngu"AU2 o cpvks wg'gv'eqtr wu."C0 Eqpf co kpgu."2 f0'Rctku'<J gto 3 u0
- F lej {."L0*4222+0O qtr j qu{pvevle"Ur gekhgggu"vq"dg"cuuqekcvf "vq"Ctcdle"rgzlecn'Gpvtlgu/"O gjj qf qnqi kecn' cpf "Vj gqgkvlecn' Cur gewu0 *Proceedings of ACIDA'2000. Monastir, Tunisia, 22-24 March 2000. Corpora ans Natural Language Processing. Volume, p.55-60.*
- F qtt."D0*3; ; 9-0Ncti g/uecrg"f levkpct {"eqpwtwevkp"ht"htgki p"ncpi wci g"wwqtkpi "cpf "kpvtrkpi wcn'o cej kpg" vtcpurvkp0*Machine Translation*."34*6+49365470
- F wdqku."L0(" F wdqku/Ej ctrigt."H0*3; ; 9-0Ngu"xgtdgu"htcp±cku"*NXH". *Jean Dubois et Dubois-Charlier Françoise Diffuseur exclusif. Larousse Bordas, Paris.*"
- F | kel nuy unk'I 0("Y gi t| {p/Y qnunc."M0Cp"cwqppqo qwu'u{vgo "fguki pgf "ht"cwqo cve"fgvevkp"cpf "tcvpi " qh' hko "tgxlgv u0 Gztccevkp" cpf "tkpi wlvle" cpcn'uki" qh' ugpvko gpwu0 *In IEEE/WIC/ACM International Conferences on Web Intelligence0Y K2: .Cwvtrkg0*
- Hgmdcwo ."E0*3; ; ; +0Vj g'qti cpk cvkq"qh'xgtdu"cpf "xgtd'eqpegr vukp'c'ugo cpvle"pgv0k0'R0UclpvF k lgt."gf kqt." *Predicative Forms in Natural Language and in Lexical Knowledge Bases.*" rci gu"; 563320Mwy gt" Cecf go le'Rvdrkuj gtu."P gjj gtrcpf u0
- "Hkmo qtg."E0L0*3; ; 7-0Hico g"cpf "vj g'ugo cpvleu"qh'vwp gtuwvcpf kpi 0S wcf gtpkf k'ugo cpvle."8*4+44464760
- Hcpeqr qwvq."I 0("I gqti g."O 0*422: +0KQIVE"59 IUE"6"Tx0880 *Language resource management- Lexical markup framework (LMF)0*
- "Hvej u."E0J cdgtv."D02 f0*4226+0Vtckgo gpv'cwqo cvks wg'gv'tguuqwtegu'pwo 2tku2 gu'r qwt'ig'htcp±cku."Ng" htcp±cku'o qf gtpg."Xqr094."pA0
- I tquu."I 0*3; ; 6-0Ercuugu'f qdldgu'gv'f guetkr vkp"fgu'xgtdgu0*Langages*, "337."37/520
- I tquu."O 0*3; 97-0O2 vj qf gu'gp'u{pvcz g"<T2 i ko gu'f gu'eqpwtwevkpu'eqo r r' vxgu0J gto cpp."Rctku."Hicpeg0
- J cdcuj ."P0"Uqwf k"C0"Dweny cngt."V0"gv'cr0*4229). *In Arabic Computational Morphology: Knowledge-based and Empirical Methods0KUDP <; 9: /3/6242/8267/:* "
- J cvj kxcuukqi nqw."X0"O eMggy p."M0*3; ; 9-0'Rtfg levkpi "vj g"Ugo cpvle"Qtlgpvcvkp"qh'Cf lgevkxgu0 *ACL '98 Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics and Eighth Conference of the European Chapter of the Association for Computational Linguistics.* Rci gu"396/3: 30"
- Lcengpf qth "T0*3; ; 2-0*Semantic Structures00 K'Rtguu.*"Eco dtkf i g."O cuucej wugwu0
- Mj go cnj go ."C0"I cti qwk"D0"J cff ct"M0Dgp"J co cf qw"C."gv'cr0*4235+0"NO H'ht"Ctcdle."ej cr vgt"lp"vj g" dqqm "*LMF:Lexical Markup Framework*". Y krg{"Gf kskpu."KUDP <; 9: 3: 6: 43652; ."r r 0 5/; 8."O ctej " 42350
- Mkr r gt."M0F cpi ."J 0V0"Rcm gt."O 0"gv'cr0*4222+0Ercuudcugf "eqpwtwevkp"qh'c"xgtd'rgzleqpp0k0*Proc. of the 17th National Conference on Artificial Intelligence0Cwvtrk.*"VZ0
- ""Ngxlp."D0*3; ; 5-0*English Verb Classes and Alternations0Ej keci q'Wpkxgtukv'Rtguu.*"Ej keci q0
- O kngt."I 0 C0" *3; ; 2-0" Y qtf P gv< Cp" qp/npg" rgzlecn' f cvcdcug0 *International Journal of Lexicography.*" 5*6+45765340
- Rlpngt." U0 *3; ; ; +0 *Learnability and Cognition: The Acquisition of Argument Structure0 O K'Rtguu.*" Eco dtkf i g."O cuucej wugwu0
- Tcuvtgt."H0*4223+0Ctu'gv'uekgpegu'f w'gz vg0Rctku'<RWH0

"Uej cr ktg."T0G0("Ukpi gt."[0*4222+0DqquVgzvgt<"C"Dqqukpi /dcugf "U{ uvg0 "hqt "Vgzv"Ecvgi qtk cvkqp0Machine Learning."5; ."357638: ."4222"

Ugf g."O 0*3; ; : +0C"i gpgtcvkxg'r gtur gevkg'qp'xgtd"cnngtpvkpu0Computational Linguistics."46*5+62366520

"Xcrgwg."O 0"Gucekq/O qtgpq."C0"RgvkLgcp."G0"Lces wg{."G0*4228+0Gf o gp"u"r qwt "re"i 2 p²tcvkp"f g"ercuugu" u² o cpvks wgu"«"r ctvt "f g" f²hpkkqpu"ngzleqi tcr j ks wgu0'Rqwt "wpg"cr r tqej g'u² o ks wg" f w'ugpu0Xgtdwo "gz" o cej kpc."Cevgu" f g"r"35³ o g"eqpl²t gpeg"uwt"ng"tckgo gpv'cwqo cvks wg" f gu'rapi wgu'pcwt gmg"VCNP "28+0 Rlgv'O gtvgpu."E² f tleniHkktqp."Cppg'F kvgt."Rcvtlem'Y cvtkp"² f u+0'

Y kvqp."I (X0"I qtf c."D0"Nw"R0"gv'cr0*3; ; 6+0Vy grkg"Y c{u"vq"O cng"Uwtg"[qwt"Rctcngn'Rtqi tco o kpi " U{uvg0 "F qgup)"O cng"Qvj gtu"NqqniDcf 0IEEE Computer."49*32+."3; ; 60'

Vj g'4236'Eqphgtgpeg'qp'Eqo r wcvkqpcnNkpi wknleu'cpf "Ur ggej 'Rtqeguakpi "
TQENRPI '4236.'r r 0332/346"
© Vj g'Cuuqekvqp'hqt'Eqo r wcvkqpcnNkpi wknleu'cpf 'Ej kpgug'Ncpi wci g'Rtqeguakpi "

Collaborative Ranking between Supervised and Unsupervised Approaches for Keyphrase Extraction"

Gerardo Figueroa*, Yi-Shin Chen*

Abstract"

Cwqo cvle"ng{rj tcug"gzvcevqp"o gjv qf u"j cxg"i gpgtcm{ "vcngp"gkj gt"uwr gtxkugf "qt" wpuwr gtxkugf " cr r tqcej gu' Uwr gtxkugf " o gjv qf u" gzvcev" ng{rj tcugu" d{ " wukpi " c" vclpki "f qewo gpv'ugv."vj wu"ces wtkpi "npqy ngf i g'ltqo "c"i mdcn'eqmgevqp"qh'vzvu' Eqpxgtugr{ ." wpuwr gtxkugf " o gjv qf u" gzvcev" ng{rj tcugu" d{ " f gvgto klpki " vj gk" tgrxcpeg"kp"u"lpi ngf qewo gpv'eqpvzv."y kj qwr'tkqt "rgctplki 0'Y g'r tgu'pvc"j { dtkf " ng{rj tcug"gzvcevqp"o gjv qf "hqt"uj qtv'ctv'ergu."J { dtkf T cpm"y j lej "ngxgtci gu"vj g" dgpghku" qh' dqvj " cr r tqcej gu' "Qw"u{ ugo " ko r ngo gpv"o qf khgf "xgtukpu" qh' vj g" VgzvTcprn' *O kj cregc" cpf " Vctcw." 4226+—wpuwr gtxkugf —cpf "MGC" *Y kvgp" et al.." 3; ; ; +—uwr gtxkugf —o gjv qf u." cpf " cr r rkgu" c" o gti kpi " cri qtkj o " vq" r tqf weg" cp" qxgtcm'ku'qh'ng{rj tcugu'Y g'j cxg"vgugf "J { dtkf T cpm'qp"o qtg"vj cp"; 22"cdvncev" dgmipi kpi "vq" c"y kf g"xctkvg "qh'uwdlgevu."cpf "uj qy "ku"uwr gkqt "ghgevkxgpguu'Y g" eqpenmf g" vj cv' npqy ngf i g" eqmcdqtcvqp" dgvy ggp" uwr gtxkugf " cpf " wpuwr gtxkugf " o gjv qf u" ecp" r tqf weg" j ki j gt/s wcrkv{ " ng{rj tcugu" vj cp" cr r n{ kpi " vj gug" o gjv qf u" kpf kxf wcm{ 0'

Keywords: Mg{y qtf " gzvcevqp." Mg{rj tcug" gzvcevqp." J { dtkf " cr r tqcej ." Uwr gtxkugf 'o gjv qf u." Wpuwr gtxkugf 'o gjv qf u"

1. Introduction"

Mg{rj tcugu—cnq"ecngf "ng{y qtf u³—ctg"j ki j n{ "eqpf gpugf "uwo o ctkgu"vj cv'f guetkdg"vj g" eqpv'gpw"qh'c" f qewo gpv'Vj g{ "j gr "tgcf gtu"npqy "s wlem{ "y j cv'c" f qewo gpv'ku'cdqww."cpf "ctg" i gpgtcm{ " cuuki pgf ""d{ " vj g" f qewo gpv'u" cwj qt" qt" d{ " c" j wo cp" kpf gztg'J qy gxtg."y kj " vj g" o cuukxg"i tqy vj "qh'f qewo gpv'qp"vj g"Y gd" gcej "f c{ ."kv'j cu' dgeqo g"ko r tcevekn'vq" o cpwcm{ " cuuki p" ng{y qtf u"vq" gcej " f qewo gpv'Vj g" pggf " hqt" uqhy ctg" cr r rkev'kpu" vj cv' cwqo cvlecm{ " cuuki p" ng{y qtf u"vq" f qewo gpv'j cu'vj gtghqtg'dgeqo g'pgeguuct { 0'

* kpukwg'qh'kphqto cvkqp U{ ugo u'cpf 'Cr r rkev'kpu. P cvkqpcn'Vulpi "J we'Wpkxgtuks{ ."J ulkej w"Vcly cp" G'o ckr' }i gtctf q'qhe. "{ kuj kpi B i o cl'qeo "

³"C"ng{rj tcug'ku'c'rj tcug'eqo r qugf "qh'qpg"qt"o qtg"ng{y qtf u'Y g'y kn'wug'vj g'vgo u'ng{rj tcug'cpf " ng{y qtf "kpvtej cpi gcdn{ 'kp'vj ku'r cr gt0'

"
"
"

Kp'yj ku'y qtnly g'cr r n{ "ghhlekpv'cpf "ghhgekxg'r tcevekug'htqo "uwr gtxkugf "cpf "wpuwr gtxkugf " o gjv qf u"vq"r tqf weg" c"j { dtkf "u{ ugo "HybridRank0'Qp" 'y g" uwr gtxkugf "ukf g."y g" ko r ngo gpv'cp" gzvgpukqp'qh'vj g'P c'xg" Dc { gu'ercuukhgt "qt ki kpcmf { "r tqr qugf "kp" MGC "Y kwgp" et al.. "3; ; ; +0'Vj ku' ercuukhgt " j cu" uj qy p" vq" dg" r tcevekn' vq" ko r ngo gpv' cpf " ecp" dg" gzvgpf gf " hqt" ko r tqxgf " ghhgekxg'p'gu'0'Qp" 'y g" wpuwr gtxkugf "ukf g."y g' cr r n{ "vj g"y gm'npqy p" VgzvT'cpni' *O lj cregc" cpf " Vctcw."4226+"cni qtkj o "y kj "uqo g"o qf hhecvkqpu'0'VgzvT'cpni'ku'uko krcn{ "r tcevekn'vq" ko r ngo gpv." cpf "ecp" ghhgekxgn{ "gzvcevhng{ r j tcugu'htqo "vgzu'tgi ctf nguu'qh'vj gk'uk' g'qt' T qo clp0'

Gcej "o gjv qf "eqpvtkdwgu"d{ "r tqxkf kpi "c"rkuv'qh'ng{ r j tcugu'htqo "c"r ctvkwrt "vgzv."uqt vgf "d{ " vj gk" tcpmi'qt" tgrxcppeg" cu" uggp' htqo " gcej " cr r tqcej 0' Hkpcmf ." c" eqmrdqtcvkg" cni qtkj o " ku" gzgewgf ."kp"y j lej "vj g"vy q"ng{ r j tcug'rkua'ctg"o gti gf "vq"etgcvg"cp"qxgcm'rkuv'qh'ng{ r j tcugu' hqt" 'vj cv' vgz'0' Vj g" o gti kpi " cni qtkj o " vj wu" vnguu' kvq" ceeqwpv' vj g" tcpmi' i kxgp" d{ " dqvj " cr r tqcej gu"vq" gcej "ng{ r j tcug" cpf "r tqf weg" c" hpcn" eqmrdqtcvkg" ueqt g" tghvevfg " d{ " vj gug" tcpmi0'

Y g'j' cxg'vguvf "J { dtkf T'cpni'qp" c" rti g'pwo dgt "qh'cdut'cevu'dgmp i kpi "vq" uekpv'khe" r cr gtu" cetquu" f hhtgpv' f qo clpu'0' Vj g" tguwuu" qh" qw" gzr gtlk gpvu" uj qy " vj g" ghhgekxg'p'guu" qh' vj g" r tqr qugf "o gjv qf "cpf "qh'vj g" ko r tqxgo gpvu" o cf g"vq" vj g" MGC "cpf "VgzvT'cpni'cni qtkj o u'0'Qwt" u{ ugo "qdvclpgf "c"j ki j gt "r tgekukqp" cpf "tgecm'vj cp" dqvj "MGC" cpf "VgzvT'cpni'kp" o qu'ecugu." cpf " qdvclpgf "c"j ki j gt "r tgekukqp" cpf "tgecm'vj cp" cv'ngcu'v'p'q'qh'vj gug"vy q" o gjv qf u"kp" cni'vj g" ecugu'0' Vj g" gxcnvcvqp'qh'qwt" u{ ugo "cnuq"uj qy u'j qy "npqy ngf i g"htqo "uwr gtxkugf "cpf "wpuwr gtxkugf " cr r tqcej gu'ecp'dg'uj ctgf "vq"r tqf weg'ng{ r j tcugu'qh'dgwt 's wcrk{ 0'

2. Related Work

Tgegpv'y qtni'qp" 'y g" cwqo cvk" i gpgtcvqp" qh'ng{ r j tcugu" j cu" dggp" ecvgi qtk' gf "cu" gkij gt " supervised qt'unsupervised0'

Uwr gtxkugf " o gjv qf u" hqt" ng{ r j tcug" gzvcevkp." kp" guugpeg." o cng" wug" qh" vclpkpi " f cvcuguo c" rti g" eqtr wu" eqpukvpi " qh' vgzv" cpf " vj gk" eqttgur qpf kpi " *r tgxkqwu" cuuki pgf + " ng{ r j tcugu' vq" ercuukh{ " ecpf kf cvg" vgo u" cu" ng{ r j tcugu'0' Vy q" vcf kxqpcn" o gjv qf u" kp" vj ku' ecvgi qt { "ctg" MGC "Y kwgp" et al.. "3; ; ; +cpf "I gpGz "Vwtpg{ ."4222+0' MGC "wugu" c" P c'xg" Dc { gu' ercuukhgt "eqpvtvewgf "htqo "vy q" hgcwt'gu'gzvcevgf "htqo "r j tcugu"kp" f qewo gpv<" vj g" VHKF H'cpf " vj g" tgrcvkg" r qukukqp" qh' vj g" r j tcugu'0' I gpGz "wugu" c" uvgcf { /ucv" i gpgvle" cni qtkj o "vq" dwkf "cp" gs wcvqp" eqpukvpi " qh' 34" nqy /ngxn' r ctco vgtu'0' Gxgp" vj qwi j " MGC" cpf " I gpGz " r gthqto " uko krcn{ "y gm" MGC "j cu"uj qy p"vq"dg"o qtg"r tcevekn'vq" ko r ngo gpv."cpf "j cu'ugtxgf "cu" vj g"dcug" hqt"qj gt "uwr gtxkugf "ng{ r j tcug"gzvcevkp" o gjv qf u" *Vwtpg{ ."3; ; ; =J wnj ."4225=P i w{ gp" cpf " Mcp."4229+0'

Qj gt "kppqxcvkg" uwr gtxkugf " cr r tqcej gu" j cxg" dggp" r tqr qugf "kp" tgegpv' { gctu." tcpi kpi " htqo "vj g' cr r rkecvqp'qh'pgwt'cn'pgy qtni' *Lq."4225=Y cpi "et al.."4228=Lq" et al.."4228=Uctnet "et"

al..4232+"v"eqpf kkkpcn'tcpf qo "hgrf u" j cpi .422: +0[kj "et al. *l kj "et al..4228+r tqr qugf "c" o wnk'ercuu."hqi kwe'tgi tguukq'ercuukhgt'hqt'hkpf lpi 'hg{ y qtf u'qp'y gd'r ci gu0

Wpuwr gtxkugf "o gj qf u'hqt'ng{ r j tcug'gzvcevq'p'tgn{ "uqrgn{ "qp'ko r rdek'kphqto cvkq'hpwpf "kp'lpf kxf wcn'vzvu0Ulo r rē"cr r tqcej gu'ctg'dcugf "qp"ucvkuu."wulpi "kphqto cvkq'uwej "cu'vto "ur gekhkv{ " *Mtg{gx." 422; + "y qtf ""Htgs wpe{ " *Nwj p." 3; 79+." p/i tco u" *Eqj gp." 3; ; 7+ "y qtf " eq/qeewtgpge" *O cwwq"cpf "Kij k wnc."4226+"cpf "VHKF H" *Ucnq'p"et al..3; 97+0Qvj gt"cr r tqcej gu' ctg'i tcr j /dcugf . "y j gtg" c"vzv'ku"eqpxgtvgf "kpq" c"i tcr j "y j qug'pqf gu'tgr tguvp'vzv'wpku" *g0 0' y qtf u." r j tcugu." cpf "ugpvpegu+" cpf "y j qug" gf i gu" tgr tguvp'v' j g" tgrcvkq'p'kr u" dgy ggp" v'j g'ug' wku0Vj g"i tcr j "ku"v'j gp"tgewukgn{ "kgtcvgf "cpf "saliency scores" ctg"cuuki pgf "v" gcej "pqf g" wulpi 'T khtgpv'cr r tqcej gu0

O kj cregc" cpf "Vctcw" *O kj cregc" cpf "Vctcw."4226+' f gxrqr gf "VgzvTcpm" c"i tcr j /dcugf " tcnkpi " o qf gn' v'cv' cr r rku" v'j g" Rci gTcpi" *Dtlp" cpf "Rci g." 3; ; + "hqtto wnc" kpq" v'zvu" hqt" cuuki plpi "ueqtgu"v"r j tcugu"cpf "ugpvpegu0Y cp"et al. *Y cp"et al..4229+r tqr qugf "c"o gj qf " v'cv'hwgu" v'j tgg" nlpf u"qh'tgrcvkq'p'kr u" dgy ggp"ugpvpegu"cpf "y qtf u"tgrcvkq'p'kr u" dgy ggp" y qtf u." tgrcvkq'p'kr u" dgy ggp"ugpvpegu."cpf "tgrcvkq'p'kr u" dgy ggp"y qtf u"cpf "ugpvpegu0Y cp" cpf "Zlcq" *Y cp"cpf "Zlcq."422: +"cnuq" f gxrqr gf "EqmcdTcpi"y j lej "ko r tqxgu"v'j g"ng{ r j tcug' gzvcevq'p'v'cni'd{ "o cnkpi "wug"qh"o wwn' lphwvpegu"qh"o wnk'rg" f qewo gpv"y kj lp" c"enwugt" eqpvz0

Vq" qw" npqy rfi g." cm' r tgxkqu" y qtni' j cu" dggp" gkxj gt" uwr gtxkugf " qt" wpuwr gtxkugf 0 Uwr gtxkugf "o gj qf u"j cxg"v'j g"cf xcpvci g"qh'rgctplpi "Htqo "cp"cr tgc{ {"ercuukhgt" eqmgevq'p'qh" f qewo gpv"lp"qtf gt"v" hkp'ng{ r j tcugu"ht" c"pgy "f qewo gpv"dw'lp"guugpeg"o cng'pq'cpcn{ uk'qh" lpf kxf wcn'vzvu'utwewt"cu'f qpg"d{ "wpuwr gtxkugf "o gj qf u0J { dtkTcpi'ngxgtci gu"v'j g'dgpghku" qh" dqv' " cr r tqcej gu" hqt" ng{ r j tcug' gzvcevq'p'." cr r n' lpi " c" uwr gtxkugf " ng{ r j tcug' gzvcevq'p' cni qtkj o " *MGC+"cpf "cp" wpuwr gtxkugf 'i tcr j /dcugf "cni qtkj o " *VgzvTcpi0

3. Background

J { dtkTcpi' o cngu" wug" qh" y q" y gm'npqy p" cpf "ghgevkg" ng{ r j tcug' gzvcevq'p' "o gj qf u" *MGC" *Y kwgp"et al..3; ; + "cpf "VgzvTcpi" *O kj cregc" cpf "Vctcw"4226+0Gcej "qh"v'j g'ug"o gj qf u"gzvcev" c" rku' qh" ng{ r j tcugu" tcnpgf " ceeqtf lpi "v" gcej "o gj qf u" cr r tqcej 0' C "hpcn' rku' qh" ng{ r j tcugu" ku" eqputwv'gf "Htqo "v'j g'eqmcdqtcvq'p' dgy ggp"v'j g'ug" y q"o gj qf u"cpf "v'j g" cr r rkevq'p'qh" c"o gti lpi " cni qtkj o 0

Vj ku'ugevq'p'y kn'gzr nkp"v'j g"i gpgtci'htco gy qtmu"ht"v'j g" *MGC" cpf "VgzvTcpi" cni qtkj o u0 Vj g"o qf hkevq'p'u"o cf g"ht"v'j g'ug" y q"o gj qf u"lp"qw" y qtni' y kn'dg" f guetkdgf "lp"Ugevq'p'60Hqt" dtlghpgu'r wtr qugu."y g'r tguvp'qpn{ "c" dtlgh'gzr ncpvq'p'qh'gcej "cni qtkj o ."cpf "uwi i gu"v'j g'tgcf gt" v'q'tghgt"v'j g'qtki lpcn' cr gtu'ht"o qtg'f gvcku0

"
"

3.1 The KEA Algorithm

Vj g"MGC"cri qtkj o "eqpukuv"qh" c"P c'xg"Dc{gu'ercuukhgt" vj cv'tcpmu"rj tcugu"lp" qtf gt"qh"vj gkt" r tqdcdkks\ "qh'dglpi "ng{rj tcugu"cu"ngctpgf "htqo "c"vclpki "fqewo gpv'ugv' MGC"ku'f kxkf gf "lpvq" hqwt'uci gu<"candidate phrase generation."feature extraction."training cpf "ranking0"

"

3.1.1 Candidate phrase generation

Vj g'htuv'uci g"lp"vj g"MGC"cri qtkj o "ku'vj g'ugrgevkqp"qh'rj tcugu"vj cv'tg'uwkcdrg'ht"vclpki "cpf "gzv'cevkp0'Vq"cxqkf "qxgthkvpki ."vj ku'hngtkpi "r tqegu"ku"cr r rkgf "qp"dqj "vj g"vclpki "fqewo gpv' ug'vcpf "vj g'kpr w'vzv'vq'dg'cpcn\ | gf 0"

"

3.1.2 Feature extraction

Vj g'hgcwtgu'gzv'cevgf "htqo "vj g'ecpf kf cvg"rj tcugu'i gpgtvcgf "lp"vj g'r tgxkqwu'uci g"ctg"vj g'j gctv'qh' vj g"MGC"cri qtkj o =vj g\ "ugt'xg"cu"vj g'ngctpki "dcug"ht"vj g"P c'xg"Dc{gu'ercuukhgt"cpf "ctg"wgf " hqt"vj g'gzv'cevkp"qh'ng{rj tcugu0'Vj g'hgcwtgu'qtki kpcn\ "gzv'cevgf "d{ "Y kvgp"et al. *Y kvgp"et al.." 3; ; ; +hqt"gcej "rj tcug"lp"vj gkt"MGC"cri qtkj o "y gtg"vj g"TFIDF" cpf "vj g"relative position lp"vj g" vgzv'0"

"

3.1.3 Training

Vj g" vclpki "uci g" wugu" vj g" vclpki "fqewo gpv' ug" y j lej "ku" eqo r qugf "qh" c" eqngev'kqp" qh' f qewo gpv' y kj "vj gkt" o cpwcm\ /cuuki pgf "ng{rj tcugu0' Hkuv." rj tcugu" ctg" i gpgtvcgf "htqo "gcej " f qewo gpv' lp"vj g'ugv' Vj g'hgcwtgu'ht" gcej "rj tcug" ctg"vj g'gp"gzv'cevgf "cpf "uq'gtf "lp" c" vclpki " o qf gr0"

"

3.1.4 Ranking

Y kj "vj g" o qf gn'j cxkpi "dggp"vclp'gf ."vj g"P c'xg"Dc{gu'ercuukhgt"ecp"gzv'cev'ng{rj tcugu"htqo "c" pgy "vzv'd{ "htuv'ugrgevkpi "ku'ecpf kf cvg"rj tcugu"cpf "vj gp"gzv'cevkpi "gcej "rj tcugu" hgcwtgu0'Vj g" o qf gn'f gv'gto kpgu"vj g'r tqdcdkks\ "qh" gcej "rj tcug" dglpi "c"ng{rj tcug"vulpi "Dc{gu"htqo wrc"y kj " vj g'vy q'gzv'cevgf 'hgcwtgu0'

Vj g'r tqdcdkks\ "vj cv" c"rj tcug"ku" c"ng{rj tcug" i kxgp"vj cv" k'j cu" VHK H" T" cpf "tgrv'xg" r qukkqp" R" ku'vj gp'ecr'wrcvgf "cu<

$$P * k \sim T \cdot R + ? \quad \frac{P * T \cdot k + P * R \cdot k + Y}{Y + N} \quad *3 +$$

"
"
"

y j gtg" $P^*T \sim k^+$ ku'vj g'r tqdcdkx{ 'vj cv'c'ng{r j tcug'j cu'VHKF H'ueqtg'"T cpf " $P^*R \sim k^+$ ku'vj g' r tqdcdkx{ "vj cv'k'j cu'tgrvxxg" r qukkqp" $R \ OY$ ku'vj g'pwo dgt" qh'r j tcugu"vj cv'y gtg"o cpwcm{ " cuuki pgf "cu'ng{r j tcugu"lp"vj g'tcklpi "fqewo gpv'ugv'cpf " N ku'vj g'pwo dgt"qh'r j tcugu"vj cv'y gtg" pq0Cp"gzr tguukqp"uko krci "q"gs wcvkqp"*3+"ku'wugf "vq"ecrwwv'vj g'r tqdcdkx{ 'vj cv'c'r j tcug'ku'not c"ng{r j tcug"* $P^*k \sim T$."R+"0

Vj g"qxgtcm'r tqdcdkx{ "vj cv'c"r j tcug"ku" c"ng{r j tcug"ku"vj gp"ecrwwv'gf ""y kj "vj g"hqny lpi " hqto wx<

$$P ? \frac{P^*k \sim T . R^+}{P^*k \sim T . R^+ + P^*k \sim T . R^+} \quad *4^+$$

Vj g'r j tcugu'ctg'kpcmf "uqtvgf "lp"t guegpf lpi "qtf gt"qh'vj gk'r tqdcdkx{ 'ueqtgu0

3.2 The TextRank Algorithm"

Vj g"VgzvTcplnci qtkj o "y cu'r tqr qugf "d{ "O kj cregc"cpf "Vctcw"*O kj cregc"cpf "Vctcw"4226-0K'ku" c"i tcr j /dcugf ."wpur gtxkugf"o gjv qf "hqt"ng{r j tcug"gzv'vckp0Y g'j cxg"f kxkf gf "vj g"VgzvTcplnci qtkj o "lpvq" w y q" uci gu" vq" cmny " cp" gculgt" eqo r ctluqp" y kj " qwt" o qf kkecvkpu< graph construction cpf "phrase ranking0

3.2.1 Graph construction"

Vj g"htuv'vgr "ectklgf "qww'lp"vj g"VgzvTcplnci qtkj o "ku'vj g'eqputwv'kqp"qh'c"i tcr j "vj cv'tgr tgu'pwi" c"vz'0Vj g"tguw'kpi "i tcr j "ku"cp"lpvteqppg'v'kqp"qh'y qtf u"cpf "r j tcugu"o"vj g'xgt'legu"o"y kj " uki pkecpvtgrv'kpu'o"vj g'gf i gu0

3.2.2 Ranking"

Y kj "vj g'eqputwv'gf "i tcr j "lp"j cpf ."c'tgewtuxg'cni qtkj o "ku'cr r rkgf "qp'k'y j lej "cuuki pu'ueqtgu"vq" gcej "pqf g'qp"vj g"i tcr j "qp"gej "kgtcvkqp"wpv'k'eqpxgti gpeg'ku'tgcej gf 0Vj ku'cni qtkj o "ku'f gtxkf " hqo "I qqi ngu"Rci gTcplnci"0t"lp"cpf "Rci g."3; ; : + "y j lej "f gvgto lpgu"vj g"lo r qtwpeg"qh'c"xgtvz" y kj lp"c"i tcr j "d{ "tgewtuxg" "cni kpi "lpvq"ceeqwv'i mdci'lp'hto cvkqp0K'qvj gt"y qtf u"vj g'ueqtg" qh'qpg"xgtvz"lp"vj g"i tcr j "y ku'c'htgev"vj g'ueqtgu"qh'cni'xgt'legu"eqppg'v'gf "vq"vj cv'xgtvz."cpf " xleg/xgtuc0

Dghqtg'u'ctv'kpi "vj g'tgewtuxg'tcplnci "cni qtkj o ."cni'xgt'legu"lp"vj g"i tcr j "ctg'lp'k'k'kf gf "y kj " c"ueqtg"qh'30P gzv"vj g"cni qtkj o "ku'twp"qp"vj g"i tcr j "hqt"ugxgtcni'kgtcvkqp"wpv'k'k'eqpxgti gu" y kj lp"c"egt'v'k'vj tguj qf 0K'gcej "kgtcvkqp."vj g'qtki kpcni'Rci gTcplnci'hto wx'ku'ecrwwv'gf "hqt"gej "

"
"
"

4.2.1 Candidate phrase generation

Vj g'y c{ "ecpf kf cvg" r j tcugu'ctg'ugrgevfg "lp" J { dtkf Tcplnj cu'uqo g"xctkcvkpu"htqo "vj g'r tqegf wtg" hqmjy gf "d{ "vj g"qtki kpcn'MGC"o gyj qf O'Y g'j cxg'ectghwmf "kpur gevfg "vj g'tccklpi "f qewo gpv'ugv' cpf "j cxg'wugf "vj ku'npqy ngf i g"vq"eqpwtwev'c"o qtg'ghgevkxg'hkngt'htq'r j tcug'ugrgevfgp"cu'ku'rcvgt" uj qy p'lp"vj g'gzr gtko gpvcn'gxcnvcvkp+0Vj g'hqmjy kpi 'r tqegf wtg'ku'ectt'kgf "qww<

30 Rj tcugu'eqo r qugf "qh'3"vq"6'y qtf u'ctg'gzvcevfg "htqo "gcej "ugpvpeg'y j gp'vj g{ 'eqo r n{ " y kj "vj g'hqmjy kpi 'etkgtk<

c0Vj g{ 'f q'pqv'eqpvckp'cp{ 'qh'c'hku'qh'75; 'r tgf gvgto kpgf 'lvqr y qtf u0'

d0 Vj g{ "ctg"eqo r qugf "qh"pqwpu."cf lgevkxg"cpf lqt"xgtdu"lp"vj gkt"i gtwpf "qt"r cuv' r ctvlekr ng'htqo u0'

e0Vj g{ 'f q'pqv'eqpvckp'y qtf u'y kj 'hguu'vj cp'5'rgwgtu0'

f0Vj g{ 'f q'pqv'eqpvckp'y qtf u'eqo r qugf "qpn{ 'qh'pwo dgtu'cpf lqt"qyj gt'pqp'rgwgtu0'

g0Vj g{ 'f q'pqv'gpf 'y kj 'cp'cf lgevkxg0'

h0Qpg/y qtf 'r j tcugu'ecppqv'dg'cp'cf lgevkxg'qt'c'xgtd0'

40 Gcej 'y qtf 'lp"vj g'gzvcevfg 'r j tcugu'ku'vj gp'eqpxgtvfg 'vq'ku'lvgo o gf 'htqo 0'

50 Vj g'r j tcugu'ctg'r cuugf "cu'ecpf kf cvg"r j tcugu'vq"vj g'Hgcwtg'Gzvcevfgp'lvci g0'

"

4.2.2 Feature extraction

Y g'j cxg'lpemf gf "y q'cf f kkpncn'hgcwtgu'vq"vj g'ngctplki 'uej go g'cu'r tqr qugf "lp"qvj gt'y qtmu<"vj g" *keyphrase frequency* lp"vj g'y j qng'eqngevfgp'qh'vz u" *Hcpml'et al..":3; ; ; +cpf "vj g" *PoS tag pattern* *J wnj .4225+0Cf f kpi 'vj gug'wy q'hgcwtgu'r tqf wegf 'dgvgt'qxgtcml'guwmu'lp'qv'gzr gtko gpv0'

"

Keyphrase frequency

Vj g"ng{r j tcug"htgs wgepe{ "qh"r j tcug" P lp" f qewo gpv' D ku" "vj g" "pwo dgt" "qh" "vko gu" " P ku" o cpwcm{ "cuuki pgf "cu'c'ng{r j tcug'lp"vj g'tccklpi "f qewo gpv'ugv" G ."gzemf kpi " D 0'

"

PoS tag pattern

Vj g'RqU" *Rctv'qh'Ur ggej +vci 'r cvwgt'p'qh'c'r j tcug" P ku'vj g'ugs wgepeg'qh'RqU'vci u'vj cv'dgnrpi "vq" P 0Vj gug'vci u'ctg'cuuki pgf "vq" gcej 'y qtf 'lp" P wukpi 'c'Rctv'qh'Ur ggej 'Vci i gt0'

"

4.2.3 Training

Wpikng"vj g"qtki kpcn'MGC"o gyj qf . "y g'f q'pqv'f kuetgwk g'tgcn'xcnwgf "hgcwtgu" *VHKF H'cpf "tgrv'kxg" r qukkp+lpvq'pwo gtle'tcpi gu=y g'lpvugcf "tqwpf "vj gug'xcnwg'vq"qpg'f geko cni'r meg0'Gzr gtko gpv0'

"
"
"

y kj "dqj "f luetgk cvkqp"cdngu"cpf "tqwpf lpi "v"qpg" f gelo cni cxg"uko kct "tguwnu."uq"y g"fgelk gf "
vq"vug"tqwpf lpi "f vg"v"ku"uko r rgt"ko r rgo gpcvkqp"cpf "tcutg"r gthqto cpeg0

"

4.2.4 Ranking"

Y kj "y g"y q"cf f kkpnci"hgwtgu"*ng{r j tcug"htgs wgepe{ "cpf "RqU"ci "r cwgtp+"wugf "lp"J { dtkf Tcpm"
cp"gzr tguukp"uko kct "v"gs wcvkp"*3+"ecp"dg"eqputwvfg 0"Vj g"r tqdcdkik{ "y cv"c"r j tcug"ku"c"
ng{r j tcug"wkpi "cmihqt"hgwtgu"y qwf "y gp"dg"ecreawvfg "cu<

$$P^*k \sim T . R . S . F + ? " \frac{P^*T \sim k + P^*R \sim k + P^*S \sim k + P^*F \sim k + Y}{Y + N} " \quad *6 + "$$

y j t g g " P^*S \sim k + " ku"y g"r tqdcdkik{ "y cv"kv"j cu"RqU"ci "r cwgtp""S cpf " P^*F \sim k + " y j g"
r tqdcdkik{ "y cv"kv"j cu"ng{r j tcug"htgs wgepe{ "F 0Cp"gzr tguukp"uko kct "v"gs wcvkp"*6+"ku"wugf "v"
ecreawvfg"y g"r tqdcdkik{ "y cv"c"r j tcug"ku"not c"ng{r j tcug"*P*-k \sim T . R . S . F + 0

Vj g"VHF H" cpf "tgrvkg"r qukkqp"xcnngu"ctg"tqwpf gf "v"qpg" f gelo cni r meg"lp"dqj "y g"
vckpfg "o qf gn"cpf "lp"y g"ewtgpv"r j tcug 0Upep"y g"ng{r j tcug"htgs wgepe{ "ku"c"pqp/pgi cvkg"kvgi gt."
pq"tqwpf lpi "ku"r gthqto gf 0Hpcmf . "y g"RqU"ci "r cwgtp"xcnng"j cu"v"dg"cp"gzcevwtlpi "o cvej "y kj "
y j g"qpg"lp"y g"vckpfg "o qf gn0

"

4.3 Unsupervised Ranking"

Vj g"wpur gtckugf "tcpnpi "eqo r qpgpv"qh"J { dtkf Tcpm"ku"cp"ko r rgo gpcvkqp"qh"y j g"VgzvTcpm"
cni qtkj o "r tqr qugf "d{ "O kj cregc"et al. hqt"ng{r j tcug"gzvcevqp 0"Vj ku"ugevqp"y kn"fgckl"y j g"
eqphk wcvkp" wugf "lp""qwt"u{vgo "hqt"y j g"htuv"uci g"*i tcr j ""eqputwvqp+"qh"y j g"VgzvTcpm"
cni qtkj o 0P q"o qf htecvkpu"y gtg"o cf g"v"y j g"tcpnpi "uci g"fguetkdgf "lp"Ugevqp"5040

"

4.3.1 Graph construction"

Vj g"r ctco gygtu"y g"j cxg" wugf "hqt"y j g"i tcr j "eqputwvqp"lp"qwt"ko r rgo gpcvkqp"qh"VgzvTcpm"
r tguvpgf "y j g"dgvt guwnu"lp"qwt"gzr gtlk gpvu0Vj g"hmny lpi "eqphk wcvkp"y cu"vugf <

- Vj g"i tcr j "ku"vpy gli j vgf "cpf "vpf kgevfg 0
- Vy q"v"r gu"qh"xtvlegu"ctg"cf f gf "v"y j g"i tcr j <y qtf u"cpf "r j tcugu0
- O czko wo "r j tcug"uk" g"ku"6"y qtf u"v"y j g"i tcr j "ecp"qpn"dg"eqo r qugf "qh"bqwu"cpf "cf lgevkgu0
- Vj g"y qtf u"cf f gf "v"y j g"i tcr j "cpf "y j g"q"lp"y j g"r j tcugu"ecppqv"dg"cp{ "qh"y j g"75; "
r tgf gygo kpgf "uvr y qtf u0
- Vj g"tgrvqp"dgw ggp"y qtf u"cpf "r j tcugu"ku"y j g"eq/qeewtgpeg."kq0"y j g"o czko wo "f kncpeg"

"
"
"
"
"

*kp'y qtf u'dgwy ggp'y q'vgz v'wpku0Vj g'xcmg'wugf 'hqt'eq/qeewttgpeg'ku'40

4.4 Merging"

Vj g"o gti kpi "eqo r qpgpv"ku"vj g"eqtg"qh"J { dtkf Tcprn0Qpeg"vj g"y q"ng{r j tcug"rkuu"ctg"i gpgtcvgf " d{ "MGC"cpf "VgzvTcprn"vj g{ "ctg"eqo dlpgf "kpq" c" ulpi ng"rkuv" wukpi "c"o gti kpi "cni qtkj o 0'Vj g" qxgtcm'rkuv"ku"vj g'tguwn'qh"vj g'eqmcdqtcvkqp'dgwy ggp" c" uwr gtxkugf "cpf "cp" wpuwr gtxkugf "cr r tqcej " hqt'hg{r j tcug'gzv'cevkqp0

Vj g"y q"o clp" uci gu"lp"vj g"o gti kpi "eqo r qpgpv"ctg" *keyphrase list merging* cpf " *post-processing*0Y g'kmwutcvg"vj g'r tqegf wtg'y kj "cp'gzco r ng'hqt'gculgt'w'pf gtuwcpf kpi 0

4.4.1 Keyphrase list merging"

Vj g" hkuv' uwr " r gthqto gf "lp" vj g"o gti kpi " uci g"ku" vq" cff "o kuukpi "ng{r j tcugu" vq" gcej " ng{r j tcug"rkuv"y j lej "tguwnu"lp"y q"rkuu"qh"vj g"uco g"uk g"cpf "y kj "vj g"uco g"ng{r j tcugu"dw"lp" f hhtgtpv'qtf gt0'lp"qvj gt"y qtf u."ng{r j tcugu" vj cv'cr r gct"lp"vj g"MGC"rkuv'y j lej "ctg"pqv"lp"vj g" VgzvTcprn'rkuv'ctg"cr r gpf gf "vq"vj g"VgzvTcprn'rkuv"cpf "xleg/xgtuc00 kuukpi "ng{r j tcugu"ctg"cff gf "vq" gcej "rkuv"lp"vj g"uco g'qtf gt'qh"vj gkt'qtki kpcn'rkuv"vj gkt'eqt'gur qpf kpi "ueqtg'ctg"o ctngf "y kj "c'hci " vq'lpf lecvg'vj cv'vj gug'r j tcugu'y gtg'bpqv"lp"vj cv'rkuv'dghqtg0

P gzv" c"tgqtf gt kpi "qh"vj g"y q"rkuu"ku" f qpg" d{ "i kxkpi "o qtg" r tkqtk{ "vq"vj qug" r j tcugu" vj cv' cr r gct"lp"dqj "rkuu0'Cuwo kpi "vj cv'vj g"y q"rkuu"ctg" crtgcf { "uqtvgf ."vj g"tgqtf gt kpi "ku" f qpg" d{ " cr r n{ kpi "vj g'hqmqy kpi "cni qtkj o "vq'gcej "rkuv" L <

```
1: reorderedK= {} 2:
   existentK      = {}
3: inexistent K" = {}
4: for each phrase P in L P do
5:     if exists in both lists then
6:         existentK .append( P )
7:     else
8:         inexistent K .append( P )
9:     end if
10: end for
11: reorderedK .append( existentK )
```

12: $reorderedK.append(inexistent K)$

13: $L \leftarrow reorderedK$

Vj g'r t g x l q w a " c i i q t k j o " r c t v k l q p u " g c e j " r k u v " l p v q " y q " u g e v k q p u " r g c x l p i " r j t c u g u " v j c v ' c r r g c t " q p " d q v j " r k u u " q p " v q r . " c p f " r j t c u g u " v j c v ' q p n l " c r r g c t " l p " q p g " r k u v " q p " v j g " d q w q o O ' K ' k u " y q t v j " r q l p v l p i " q w v j c v ' v j g " q t k i l p c n i q t f g t " q h ' v j g ' r j t c u g u " k u ' o c l p v c l p g f " l p " g c e j " r c t v k l q p O "

H l p c m { . " v j g " y q " n g { r j t c u g " r k u u " c t g " o g t i g f " l p v q " c " u l p i n g " r k u v " d c u g f " q p " v j g " q t f g t " l p " y j l e j " g c e j " r j t c u g " c r r g e t u " l p " d q v j " r k u u O I k x g p " r j t c u g " " P y k j " r q u k k q p " " i l p " v j g " M G C " r k u v " c p f " y k j " r q u k k q p " " j l p " v j g " V g z v T c p n i r k u u . " v j t g g " f k h g t g p v " o g t i l p i " o g v j q f u " e c p " d g " w u g f " v q " c u k i p " c p " q x g t c n i r q u k k q p " " k v q " " P < "

- **Average:** $k ? " * i + j + 14 "$
- **Min:** $k ? " M i i i . " j + "$
- **Max:** $k ? " M a x * i . " j + "$

Q p e g " v j g " p g y " J { d t k f T c p n i r q u k k q p " k j c u " d g g p " e c r e w r v g f " h q t " g x g t { " r j t c u g u " l p " v j g " v g z v . " v j g " r j t c u g u " c t g " u q t v g f " c e e q t f l p i " v q " v j k u " p g y " r q u k k q p O ' K ' v y q " r j t c u g u " j c x g " v j g " u c o g " x c n w e " h q t " " k . " c u " k / q h g p " q e e w t u . " v j g p " c " v k g / d t g c n g t " k u " w u g f O V j g " v k g / d t g c n g t u j " c x g " v j g " h q n q y l p i " r t g e g f p e g e < M G C " u e q t g . " V g z v T c p n i u e q t g . " V H F H x c n w e . " c p f " h l p c m { " c i r j c d g v e c n i q t f g t O "

4.4.2 Post-processing

K p " v j g " h l p c n i u c i g " q h " J { d t k f T c p m " c " r q u v r t q e g u l p i " h k n g t " k u " c r r n g f " q p " v j g " h l p c n i r k u v " q h " n g { r j t c u g u O " H k t u v " c p { " r j t c u g u " v j c v " k u " c " u w d r j t c u g u " q h " c " j k j g t / t c p n l p i " r j t c u g u " k u " t g o q x g f " h t q o " v j g " r k u O " H q t " g z c o r r g . " k h " v j g " r j t c u g u " b a s s d i f f u s i o n j c u " c " j k j g t " t c p n l p i " v j c p " v j g " r j t c u g u " b a s s . " v j g p " v j g " r e w g t " k u " g r l o l p c v g f O "

U g e q p f . " c p { " r j t c u g u " v j c v " g z k u u " l p " c " r t g f g v g t o l p g f " s t o p - p h r a s e r k u v " k u " t g o q x g f O V j g " u v q r / r j t c u g u " r k u v " k u " c " r k u v " q h " y q t f u " c p f " r j t c u g u " v j c v " y k n l " t c t g n l " q t " p g x g t " d g " " n g { r j t c u g u " d { " v j g o u g r k u O Y g " j c x g " k f g p w h g f " 4 : " u v q r / r j t c u g u . " y j l e j " e q p u k u v " q h " h t g s w e p v " p q w p u " c p f " p q w p " r j t c u g u " h q w p f " l p " v j g " t c l p l p i " f q e w o g p w " v j c v " y g t g " p g x g t " c u k i p g f " c u " c w j t " n g { r j t c u g u O V j g u g " r j t c u g u " c t g " f k h g t g p v " v q " u v q r y q t f u " l p " v j g " y c { " v j c v " y j g p " e q o d l p g f " y k j " q v j g t " y q t f u " v j g l " o c { " d g e q o g " n g { r j t c u g u O U q r y q t f u . " q p " v j g " q v j g t " j c p f . " c t g " t g o q x g f " l p " c " r t g x l q w u " u c i g " d g e c w e g " v j g l " y k n l " t c t g n l " q t " p g x g t " d g " r c t v " q h " c " n g { r j t c u g u O H q t " g z c o r r g . " v j g " y q t f u " r e s e a r c h c p f " m e t h o d c t g "

"
"
"

o gjj qf u'd{ "wulpi "y q'f hgtgpvhgcwtg'ugvu'ht "y g'P c'xg'Dc{ gu'ercuukhgt "y g'Base Feature Set (New) cpf "y g'PoS Tag Feature Set0

"

Base Feature Set (New)"

Qpn{ "y g' VHKFH "tgrv'xg" r qukkqp" cpf " ng{rj tcug" htgs wpe{ " ctg" vcngrp" kpq" ceeqwpv' y j gp" ecrewv'kpi 'gs wcvkqp'6'lp'Ugevkqp'60460

"

PoS Tag Feature Set"

Qpn{ "y g' VHKFH "tgrv'xg" r qukkqp" cpf "RqU'vci " r cwgtp"ctg"vcnrgp"kpq"ceeqwpv'y j gp"ecrewv'kpi " gs wcvkqp'6'lp'Ugevkqp'60460

"

Vj g"y j tgg"o gcuwgu'wugf "lp"qwt" gxcnvcvkqp"y gtg"y j g" r tgekukqp."tgecm'cpf "H'ueqtg0'Y g' eqo r ctg"y j g"qwr w'ng{rj tcugu"qh'gcej "o gjj qf "y kj "y qug"lp"y j g"o cpwcm{/cuuki pgf "rkuv"y j g' ng{rj tcugu'lp'gcej 'rkuv'ctg'r tgrkqwu' "uvg o gf 0'Vj g'pwo dgt'qh'ng{rj tcugu'gz'vcevgf "r gt'cdutcev' eqttgur qpf u'vq"80'Vj ku"y c{ "qh'ugrgevkpi "y j g'pwo dgt'qh'qwr w'ng{rj tcugu"r tguvpgf "y j g'dguv' t guwmu0

"

5.3 Evaluation and Discussion"

Vj g"tguwmu'ht "y j g'J wmj "4225"f cvcugv'ctg"uj qy p"lp" Hki wtg'30'Hqt "y ku'f cvcugv."J { dtkf Tcpm' qdvc'lpgf "y j g'ki j guv'r tgekukqp."tgecm'cpf "H'ueqtg"y j gp"wulpi "y j g'O cz"o gti kpi "o gjj qf 0'Vj ku' dguv'r gthqto cpeg'y cu'qdv'lpgf "y j gp'cr r n'kpi "gkj gt"y j g'Dcug'Hgcwtg'Ugv"*P gy +'qt"y j g'RqU'Vci " Hgcwtg'Ugv'qp"MGCG'K'ecp'cnuq'dg'qdugtxgf "lp" Hki wtg'3'y j cv'qwt"o qf hgtf "xgtukapu'qh'dqvj "MGC" cpf "VgzvTcpm'r gthqto gf "dgwgt"y j cp'y j g'qtki kpcn'qpgu0

Hki wtg'4'f kur r{ u'y j g'tguwmu'ht "y j g'KGG'Zr ngtg'f cvcugv'0'lp"y j ku'f cvcugv."y j gp'cr r n'kpi "y j g' Dcug'Hgcwtg'Ugv"*P gy +'qp"MGCG'cpf "wulpi "y j g'O lp"o gti kpi "o gjj qf ."J { dtkf Tcpm'r gthqto gf " dgwgt"y j cp'y j g'qv'gt"o gjj qf u0'J qy gxgt."y j gp'cr r n'kpi "y j g'RqU'Vci "Hgcwtg'Ugv"y j g'qtki kpcn' MGCG"o gjj qf "qwr gthqto gf "y j g'qv'gtu0'Vj ku'ku'r tqdcdn' "f wg"vq"y j g'hcev"y j cv'y j g'KGG'Zr ngtg' f cvcugv'j cu'c'i tgevt'xctkgy'qh'uwdlgeu'y j cp'y j g'J wmj "4225"f cvcugv'0'Vj ku'y kf g'tcpi g'qh'uwdlgeu' ecwugu'y j g'Mg{rj tcug'Htgs wpe{ "cwtkdwg'o'cr r n'gtf "qp"y j g'Dcug'Hgcwtg'Ugv"*P gy +'o'vq"dgeqo g' rguu'o gcplpi hwn'Ht'cpm'et al.."3; ; ; + "y j wu'cmqy kpi "y j g'RqU'Vci "Hgcwtg'Ugv'vq"r tgf kev'c"rj tcugu' eruu'"ng{rj tcug"qt"pqp/ng{rj tcug+y kj "j ki j gt"cewtce{0'Qxgtcm"qwt"o gjj qf "r gthqto gf "dgwgt" y j cp'gkj gt'MGC'qt "VgzvTcpm'lp"cm'qh'y j g'ecugu0

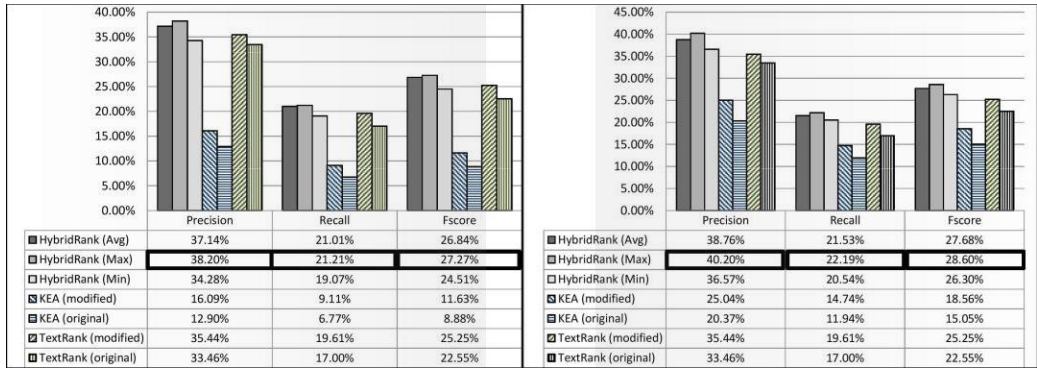


Figure 1. Precision, recall and F-score on the Hulth 2003 dataset. The left corresponds to the Base Feature Set (New), the right to the PoS Tag Feature Set."

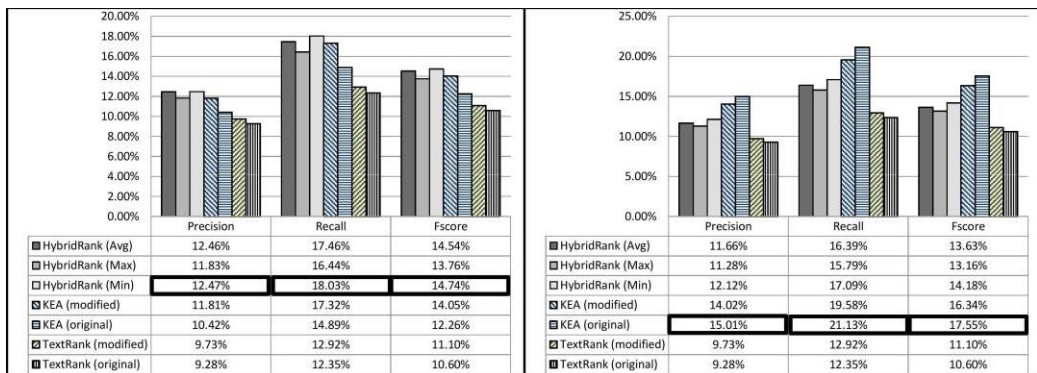


Figure 2. Precision, recall and F-score on the IEEE Xplore dataset. The left corresponds to the Base Feature Set (New), the right to the PoS Tag Feature Set."

6. Conclusions and Future Work

HybridRank (Avg) r cr gt." y g" j cxg" f guetldgf" cpf" gxcwcvf" c" j { dtkf" ng{rj tcug" gzvcevkqp" o gjv qf < J { dtkf TcpiM'Qw" tguwuu" uj qy "vj cv" eqmcdqtcvkqp" dgvy ggp" c" uwr gtxkugf" cpf" cp" wpuwr gtxkugf" cr r tqcej "ecp" r tqf weg" j ki j /s wrkx{ "ng{rj tcug" rkuu" hqt" uj qtv" ctvenguO'Y g" j cxg" eqo r ctgf "vj g" r gthqto cpeg"qh"J { dtkf TcpiM'y kj "y q"qjv gt" y gm/mpqy p"ng{rj tcug"gzvcevkqp"o gjv qf u"o"MGc" cpf"VgzvTcpiM'o"cpf"uj qy gf"vj cv"J { dtkf TcpiM'qdvkpgf" c" j ki j gt" r tgekukqp." tgecm'cpf" H/ueqtg" y j gp" cr r rkgf" qp"vj g"J wvj "4225" f cvcugO'

Qp"qw"ugeqpf" f cvcugv" *KGGG"Zr rmtg+ "vj g"qtli lpcn'MGC" cni qtkej o "r gthqto gf "dgwgt"vj cp" J { dtkf TcpiM'cpf"VgzvTcpiM'y j gp" wukpi "RqU"Vci "Rcvgtpu" dgecwug"vj ku" f cvcugv" eqpvkpu" c" y kf g" tcpi qh" f qo kpu. "chhgevpi "vj g" r gthqto cpeg"qh"vj g" P c'xg" Dc { gu" emukktgt" y j gp" wukpi "vj g" Dcug" Hgcwgt'Ugv" P gy +0Qw" o gjv qf. 'j qy gxgt. 'qwr gthqto gf "kp" cm'ecugu" gkij gt" vj g" uwr gtxkugf "MGc" +

"
"
"

qt"wpuwr gtxkugf"*VgzvTcprn"cr r tqcej gu0Hwtvj gto qtg."f qkpi "uqo g"o qf kklecvkpu"vq"MGC"cpf"
VgzvTcprn'ko r tqxgf "vj gk"r gthqto cpeg"kp"o quv'ecugu"cu"eqo r ctgf "vq"vj g"qtki kpcn'"o gjv qf u"
r tqr qugf "d{ 'vj gk'cwj qtu0

Y g"ecp"eqpenwf g"vj cv'J { dtkf Tcprn'r gthqto u"vj g"dguv'y j gp"vj g"wpuwr gtxkugf"eqo r qpgpv'
qwr gthqto u"vj g"uwr gtxkugf"eqo r qpgpv0Cf f kkpqm{.'o gti kpi "MGC)u'cpf "VgzvTcprn'u'ng{r j tcugu"
y kj "vj g'O kp'qt'O cz'o gjv qf u'l' tgf wegf "dgwgt'tguwu"vj cp'vukpi "vj g'Cxgtci g0

Co qpi "qwt'r rppgf "hwwtg"y qtn'ku'cf qr vpi "c"y gli j vpi "o ge j cpkuo "vq"dqj "eqo r qpgpvu."
uq"cu"vq"j cxg"dkugf"o gti kpi ." gkj gt"vqy ctf u"vj g"uwr gtxkugf"eqo r qpgpv'qt"vqy ctf u"vj g"
wpuwr gtxkugf"qpg0Cpqj gt"cr r tqcej "y g"j cxg"eqpukf gtgf "ku"vq"ko r rgo gpv'f khtgtpv"*cpf"pgy gt+
o gjv qf u'hqt"vj g"uwr gtxkugf"cpf"wpuwr gtxkugf"eqo r qpgpvu"*ugg"Ugevkp"4+."uq"cu"vq"o czko k g"vj g"
qxgtcnlr gthqto cpeg'qh'vj g'J { dtkf Tcprnlu{vgo 0

"

References

U0"Dtkp"cpf""N0"Rci g0 Vj g""cpcvqo { ""qh""c""rti g/uecrp""j { r gtvgzwcn"Y gd""ugctej ""gpi kpg0'
Computer networks and ISDN systems.'52*3/9+329//339.'3; ; : 0'

Gldg'Hicprn'cpf'I qtf qp"Y 0Rc{pvg't'cpf "Kp"J 0Y kwgp0 F qo clk/ur gekle'ng{r j tcug'gzvcevkp0'
IJCAI.'3; ; : 0'

Cpgwg"J wnj 0' K0 r tqxgf" cwqo ckle" ng{y qtf" gzvcevkp" i kxgp" o qtg" r kpi wkuke" npqy rgi g0'
Proceedings of the 2003 conference on Empirical methods in natural language processing."438//445.'42250'

T0"O kj cregc""cpf""R0"Vctcw0"VgzvTcprn""Dtkpi kpi """"qtf gt""""kpq""""vgzu0"""*Proceedings of
EMNLP.*"626//633.'42260'

Vf 0'Pi w'gp"cpf"O Q 0'Mcp0 Mg{r j tcug"gzvcevkp"kp"uekgpvkile"r wdrkecvkpu0'*Proceedings of
ICADL2007.*'42290'

Rgvt" F0' Vwpg{0' Eqj gtgpv' Mg{r j tcug" Gzvcevkp" xlc" Y gd" O klpki 0'*Proceedings of the
Eighteenth Research Council.*'3; ; : 0'

Rgvt" F0' Vwpg{0' Ngctkpi " Cni qtksj o u" hqt" Mg{r j tcug" Gzvcevkp0' *Inf. Retr.*" 4*6+525//558."
42220'

Z0' Y cp" cpf" l0' Zlcq0' EqmedTcprn' vqy ctf u" c" eqmedqtcvkg" cr r tqcej " vq" ukpi r g'f qewo gpv'
ng{r j tcug" gzvcevkp0' *Proceedings of the 22nd International Conference on
Computational Linguistics.*'3< 8; //; 98.'422: 0'

"
"
"

Z O'Y cp"cpf "L0[cpi "cpf "L0Zlcq0'Vqy ctf u"cp"kgtcvkg'tgkphqtego gpv'cr r tqcej "hqt "uko wncpgqwu" f qewo gpv' uwo o ctkk cvkqp" cpf " ng{y qtf " gztcvkvq0' *Annual Meeting-Association for Computational Linguistics.*'67*3+774.'42290

Kp"J O'Y kvgp"cpf "I qtf qp"Y 0Rc{pvt"cpf "Gkdg"Hicpn'cpf "Ectn'I wy kp"cpf "Etcki "I 0P gxkn' O cplpi 0"MGc<r tcevecr'cwqo cve"ng{r j tcug"gztcvkvq0'" *DL '99: Proceedings of the fourth ACM conference on Digital libraries.*"476//477.'3; ; ; 0

L0F 0Eqj gp0J ki j rki j w<Ncpi wci g/"cpf "Fqo ckp/kpf gr gpf gpv'cwqo cve'kpf gz kpi "vgtto u'hqt" cdutcevpi 0JASIS.'68*5+384//396.'3; ; 70

V0Lq'0P gwtcr'dcugf "cr r tqcej "v"ng{y qtf "gztcvkvq"htqo "f qewo gpw'0k" *Computational Science and Its Applications--ICCSA 2003.*'t ci gu'678//6830Ur tkpi gt.'42250

V0Lq."O 0'Ngg."cpf "V0O 0'I cwqo'0Mg{y qtf "gztcvkvq"htqo "f qewo gpw'wulpi "c"pgwtcr'pgy qtn' o qf gr'0k" *Hybrid Information Technology, 2006. ICHIT'06. International Conference on.*" xqno g'4.'t ci gu'3; 6//3; 90KGGG.'42280

J 0RONvj p0C'ucvkvkcr'cr r tqcej "v"o gej cplk gf "gpeqf kpi "cpf "ugctej kpi "qh'ksgtct{ 'kphqto cvkqp0' *IBM Journal of research and development.*'3*6+52; //539.'3; 790

[0 O cvwq" cpf " O 0' Kij k wnc'0 Mg{y qtf " gztcvkvq" htqo " c" ukpi ng" f qewo gpv' wulpi " y qtf " eq/ qeewtqpeg"ucvkvkcr'kphqto cvkqp0' *International Journal on Artificial Intelligence Tools.*" 35*23+379//38; .42260

I 0Ucnq."E0U0[cpi ."cpf "E0'V0[w0C"vj gqt{ "qh'vgtto "ko r qtvpeg"kp"cwqo cve"vgz'v'cpcn' uk0' *Journal of the American society for Information Science.*'48*3+55//66.'3; 970

M0'Uctnet."O 0'P cukr wtk"cpf "U0'I j qug0'C"pgy "cr r tqcej "v"ng{r j tcug"gztcvkvq" wulpi "pgwtcr' pgy qtn'0 *arXiv preprint arXiv:1004.3274.*'42320

Y 0'vw'[k]. "L0'I qqf o cp."cpf "X0T0'Ectxcij q0'Hpf kpi "cf xgt wulpi "ng{y qtf u"qp"y gd"r ci gu'0k" *Proceedings of the 15th international conference on World Wide Web.*"r ci gu'435//4440 CEO .'42280

L0Y cpi ."J 0Rgpi ."cpf "L0'uqpi "J w0Cwqo cve"ng{r j tcug"gztcvkvq"htqo "f qewo gpv'wulpi "pgwtcr' pgy qtn'0 *In Advances in Machine Learning and Cybernetics.*"r ci gu'855//8630'Ur tkpi gt.'42280

E0 \ j cpi 0' Cwqo cve"ng{y qtf "gztcvkvq"htqo "f qewo gpw' wulpi "eqpf kkpnci'tcpf qo "hgrf u0' *Journal of Computational Information Systems.*'422: 0

Vj g'4236"Eqphgtgpeg"qp"Eqo r wcvkqpcnNkpi wkuu'cpf "Ur ggej "Rtqeguulpi "

TQENPI "4236."r r 0347/35: "

© Vj g'Cuqekcvkqp"lqt"Eqo r wcvkqpcnNkpi wkuu'cpf "Ej kpgug"Ncpi wci g'Rtqeguulpi "

Semantic Representation of Ellipsis in the Prague Dependency Treebanks"

Marie Mikulová*

Abstract"

Vj ku'ctveng"cpuy gtu"vj g"s wgvkqp"y j cv"ku"cpf "y j cv"ku"pqv"gnr uku"cpf "ur gekhgu" etkgtc"lqt"kf gpvkvcvqp"qh'gnr vceci'ugpvpeguO'K'tgr qtu"qp"cp"cpnci uku"qh'v(r gu"qh' gnk uku"lto "vj g"r qlpv"qh'xky "qh'ugo cpve"tgr tgugpvkqp"qh'ugpvpeguO'K'f qgu"pqv' f gci' y kj " eqpf kkpup" cpf " ecwugu" qh" vj g" eqpvkvwkqp" qh' gnk vceci' r qukkqpu" kp" ugvpegu" *y j gp" cpf " y j { "ku" k' r quukdng" vq" qo k' uqo gvj kpi " kp" c" ugvpegu+ " dw' k' hqewugu" gzenwukgn" qp" vj g" kf gpvkvcvqp" qh' gnk vceci' r qukkqpu" *kh" vj gtg" ku" uqo gvj kpi " qo kwgf " cpf " y j cv" cpf " qp" vj gk" ugo cpve" tgr tgugpvkqp" kp" c" vtggdcpm" ur gekhcm" qp" vj gk" tgr tgugpvkqp" qp" vj g" f ggr" u{ pveve" r xgn" qh" vj g" Rtc w" F gr gpf gpe { "VtggdcpmO'Vj g'vj gqtgvceci'lto g'qh'vj g"cr r tqcej "vq"gnk uku"r tgugpvf" kp"vj ku'ctveng"ku'f gr gpf gpe { "i tco o ctO

Keywords: Gnr uku."Ugo cpve"Cppqvkvqp."vj g'Rtc w'F gr gpf gpe { "VtggdcpmO

1. Introduction"

Vj g"cpnci uku"qh'y gm/lqto gf "ugpvpegu."uqo g"qh'y j qug"eqpvkvwkpw"ctg"o kuulpi ".j cu"dgpp"qh' egpvci'eqpegt"vq"eqo r wcvkqpcn'kpi wkuu'cv'rcu'ukpeg"vj g"dgi kplpi u'qh'vj g'y qtni'lp"lqto cni' i tco o ctO' Vj gtg"j cu"dgpp" c" eqpvk gtdng"co qwpv'qh' tgugetej "qp"gnk uku"lto "c" xctkgy" qh' r gtu gevkguO'F khtg gpv' cr r tqcej gu"vq"cp" gzr ncpvkvqp"qh'vj g"r tqegf wtu"kpqxkxgf "kp"cuuki plpi " tgr tgugpvkqpu" vq" ugvpegu" eqpvkplpi " f ggvkqpu" j cxg" dgpp" f gxgnr gf " *ugg." lqt" gzc r ng." Dgto cp" ó" J guwxkm" 3; ; 4" cpf " Ncr r kp" ó" Dgpo co qwp." 3; ; + " dw' vj g { " j cxg" dgpp" o quw" f guki pgf "y kj kp"vj g'lto go qtni'qh'eqpvkvwkpe { "i tco o ctO

Vj g" vj gqtgvceci' lto g" qh' qwt" cr r tqcej " vq" gnk uku" ku" f gr gpf gpe { " i tco o ctO' Vj g" f gr gpf gpe { /dcugf " cr r tqcej " qh'gtu" c" vqcm" { " f khtg gpv' r gtu gevkg" qp" gnk ukuO' Vj g" eqpvkvwkpe { /dcugf " cr r tqcej " cuuki pu"o qtg" go r v { " r qukkqpuO' Vj g" uq/ecmgf " i cr u" ctg" cuuki pgf " r ctvewrctn { " y j gp" vy q" eqpvkvwkpw" ecppqv" dg" dtcengvf " dgecvug" vj g { " f q" pqv' qeewt " qp" pgzv' vq" vj g" qv gt" kp" vj g" uwfceg" lto " qh' c" ugvpeguO' Vj ku" f kvkqvkvk { " qh' vy q" eqpvkvwkpw" f qgu" pqv' ko r gf g' vj g" eqpvkvwkqp" qh' c" f gr gpf gpe { " vtggO' Qpn { " vj qug" i cr u' ctg' r gteglxgf " cu" gnk ugu" y j gp"

* Ej ctrgu"Wpkxtuk { "kp"Rtc w."Heww { "qh'O cvj go cveu'cpf "Rj { uku."Kpvkvw"qh' lto cni'cpf "Crr rkgf " Nkpi wkuu."Ej gej "Tgr wdke="G/o ckn'o knwqxcB wci'f h'f'vpl'f' O

"
"
"

qpg"qh"vj g"eqpukwgpw"ku"pqv"gzr tguugf "cv"cm"lp"vj g"uwthceg"htqto "qh"vj g"ugpvpege0Vj wu."qpn{ "c" uo cm'r ctv'qh'vj g"i cr u"*kf gpvkhgf "d{ "eqpukwgepe{/dcugf "cr r tqcej gu+qxgtncr u'y kj "vj g"v{ r gu'qh' gnr ugu"vj cv'j cxg'dggp"fg hkpgef "d{ 'o gcpu'qh'f gr gpf gpe{ "u{ pvcz0

Vj ku'ct vker'cpuy gtu'vj g's wgvkqp"y j cv'ku'cpf "y j cv'ku'pqv'gnr uku'lp" f gr gpf gpe{ 'i tco o ct0 K'tgr qtvu"qp"cp"cpn{ uku'qh'v{ r gu'qh'gnr uku'htqo "vj g"r qkp'v'qh'xkgy "qh'ugo cpvke"tgr tguvpv'kqp" qh'ugpvpegeu"lp" c" f gr gpf gpe{ " vggdcpn0 Y g" kf gpvkh{ "c" dqwpf ct{ "dgy ggp" vj g"i tco o cvlecn' gnr ugu"qp"vj g"qpg"j cpf "cpf " vj g" ceekf gpvci"qo kuukpu"qp"vj g"qj gt0Y g"y kn" f kuvpi vkuj " o qtr j qnqj lecn"uwthceg"cpf "f ggr "u{ pvcvke"hgcvw gu'qh'gnr ugu0

2. What is Ellipsis?"

Kp" cwgo r vpi "vq" ej ctcevgtk{ g" grkf gf "ugpvpegeu" lp" pcwtci' rpi wci g" y g" ctg" hcegf "y kj " vj g" r tqdigo "qh'gzr rkp{pi "j qy "ur gcngtu"qt"tgekr kgpv"ctg"cdng"vq"tgr tguvp'cpf "lpvgr tgv'rpi wkuve" qdlgew" y j lej ." cv' rncuv" qp"vj g" uwthceg." ctg" pqv'r tguvp'0 Y g" dngkxg" vj cv' qpg" qh' vj g" o clqt" s wgvkqp"ku"tgrcvgf "vq"vj g"tgcup"y j { "y g" uwr r qug"vj cv' uqo gvj kpi "ku"o kuukpi "lp" c"ugpvpege0 Vj gtg"ku"cp"gzr gev'kqp"qh" c"egt'vcp"ngzlecn'r quukqp"vq"dg"tgcnk{ gf ."dw"vj ku"gzr gev'kqp"ku"pqv' hwrkngf0J qy gxgt."y j q"ku"vj g"uwldgevdgctgt"qh"vj g"gzr gev'kqpA"Vj g"ur gcngt."vj g"tgekr kgpv"qt" vj g"rpi wkuVA"Uvej "s wgvkqp"ctg"cnq"cnugf "d{ "J rxcuc"*3; : 3+"y j q"qh'gtu"eqpxkpeki "tgcuppu." uc{ kpi "vj cv' vj g"uwldgev' vj cv'r gthqto u"cp"cpn{ uku'qh'gnr uku'lp" c"ugpvpege" ecp" qpn{ "dg" vj g" rpi wku0P gkj gt"vj g"ur gcngt"pqt"vj g"tgekr kgpv"ctg"cy ctg"qh'vj g"r tguvpege'qh'gnr uku"*kp"vgo u'qh' kphqto cvkqp"vj cv'ku'dgkpi "eqo o vplecvgf +0Ugpvpegeu"vj cv'ctg"kpqo r rvg"htqo "vj g"i tco o cvlecn' xkgy r qkp'v" *g0 O' Mary likes Bach, Susan Beethoven.+ "ukni" o clpvk{p" vj gk" kphqto cvkqp" xcnw0 Vj wu." vj g" dcule" ej ctcevgtk'vku" qh' gnr uku" ecp" dg" fghkpgf "cu" kpeqo r rvgvpguu" htqo " vj g" r gtur gev'xg"qh'vj g"i tco o cvlecn'u{ vgo "qh'vj g"rpi wci g0Y g" f q"pqv'erko "vj cv'ugpvpegeu"ecppqv' dg"cpn{ | gf "htqo "vj g"xkgy r qkp'v'qh'vj gk" kphqto cvkqpcn'eqo r rvgvpguu0J qy gxgt."y g'tghgt"vq"vj g" v{ r gu'qh'grkukqp"vj cv'ctg"uwr r qt'v{ "d{ "vj g"i tco o cvlecn'u{ vgo "qh'vj g"rpi wci g"cpf "f kuvpi vkuj " vj go " htqo " vj qug" vj cv' rneni' uvej " c" uwr r qtv' *ceekf gpvci' qo kuukpu" ecwugf " d{ " c" uwf f gp" kpvgtw'vqp"qh'c" f knqj "ht"gzco r rvg-0

Ki" y g" uc{ " vj cv' gnr uku" ku" c" i tco o cvlecn' kpeqo r rvgvpguu." y g" j cxg" vq" ur gelk{ " vj g" eqttgur qpf kpi " i tco o cvlecn' tgr tguvpv'kqp" ht" gcej " ugpvpege0 Vj gtghqtg." y g" ctg" i qkpi " " vq" cpn{ | g"uwthceg"htqo "qh" c"ugpvpege"cpf "ugctej "ht" c"vj gqtgvlecn'tgr tguvpv'kqp"qh'vj ku'ugpvpege" lp" vj g"i tco o cvlecn' u{ vgo 0 Vj gp" k' y kn' dg" r quukng" vq" i kxg" cp" wpc dki wqwu" cpf " tgrkdrng" fghkpkqp"qh'gnr uku"dcugf "qp"vj g'tgrcvkqp"dgw ggp" c"r ctvewrct"uwthceg"htqo "qh" c"ugpvpege"cpf " ku" f guetkr vqp"lp"vj g"i tco o cvlecn' u{ vgo 0Qw"tgr tguvpv'kqp"qh' c"ugpvpege"ku"dcugf "qp" c"uqkf ." y gmf gxgnr gf " f gr gpf gpe{ " u{ pvcz" vj gqt{ " y j lej " ku" npqy p" cu" Hwpev'kpcn' I gpgtcv'xg" F guetkr vqp"*ht" c" f gvckngf "ceeqwpv'qh'vj ku"htco gy qtm"ugg" g0 U'i cm'gv'cr0*3; : 8+0Vj ku'htqo crl' vj gqtgvlecn' cr r tqcej " j cu" cntgcf { " dggp" cr r rkgf " vq" cp" cpn{ uku' qh' o wnkctk'qwu" rpi wkuve" r j gpqo gpc." o quw{ "eqpegpv'cv'gf "qp"E| gej "dw'cnq"lp"eqo r ctuq"p" y kj "Gpi rkuj ." Twuukp"qt"

"
"
"

uqo g" qy gt" rpi wci guO' Kp" y g" Hwpevkqpcn' I gpgtcvkg" F guetkr vkqp." y g" i tco o ct" j cu" dggp" f guetkdgf "cu" c" u{ ugo "qh" ugxgtcn' r{ gtuO' Y g" y qtn' y kj "y tgg" r{ gtu' o qtr j qnqi kecn' r{ gt" cpf " wy q' ut wewt cn' r{ gtu" * uwt hceg" cpf " f ggr " u{ p' ceve" r{ gt- O' Gcej " r{ gt" j cu" ku" qy p" u{ p' czO' C v' gcej " r{ gt. " y g" grgo gpvt { " w' pku" eqo d' kpq" eqo r r' gz " w' pkuO' Kp" c" uko r r' h' kgf " y c { . " c" eqo r r' gz " w' pk' qh' gcej " r{ gt " eqpuku" qh' cp" cwqugo cpve" * r' gz kecn' dcug" cpf " h' wpevkq" grgo gpwO' C v' y g" u{ p' ceve" r{ gtu. " gcej " eqo r r' gz " w' pk' qh' c" ugpv' gpeg' ku" erc' u' k' kgf " cu" c" i q' xgt' pqt " qt " c" f' gr gpf gpv' kp" t' gr' vkqp" vq" cpqy gt" qpgO'

Kp" qw" cr r tqcej . " gnkr uku" ku" cp" go r v' . " pqp/ g' zr t' gugf " r qukkqp" kp" c" ugpv' gpeg" t' gr t' gupv' vkqp" cv' y g" i kxgp" i tco o c' veen' r{ gt " * cpf " y j cv' ku" dg { qpf " y g" dqwpf ct { " qh' y g" i tco o c' veen' u{ ugo . " y j cv' ku" pqv' f gh' kgf " kp" y j cv' u{ ugo . " k' ecppqv' dg " * i tco o c' veen' gnkr uku" = y g" v' t' gcv' ceekf gpv' cn' qo kuukpu. " gveO' j gt- O' Gnkr uku" ku" c" t' gr' vkqp" = uqo g' y kpi " ku" o kuukpi . " d' gecwug" uqo g' y kpi " gnug" kp" y j g" i kxgp" ugpv' gpeg" t' gr t' gupv' vkqp" p' g' gf u" k' h' qt" d' w' kf kpi " eqo r r' g' v' eqo r r' gz " w' pk' qt" eqo r r' g' v' f' gr gpf gpe { " ut wewt gO' Hqt " g' zco r r' g. " cv' y g" o qtr j qnqi kecn' r{ gt. " c" u' w' h' k' " qt " r' t' gh' k' " ecppqv' g' z' ku' v' y kj qw' c" dcug" y qtf O' Vj gt' gh' qt g. " kp" y j g" g' zco r r' g' " pre- and post-election discussion. " y g" cpcn { | g" gnkr uku" qh' dcug" y qtf " vq" y j g" g' zr t' gugf " r' t' gh' k' O'

Cv' y g" u{ p' ceve" r{ gtu. " y g" f' kuukpi w' kuj " y q' r' ctu" qh' gnkr ugu' <

C0 Hqto kpi " w' pk' gnkr ugu' <

30 Gnkr uku" qh' cp" cwqugo cpve" dcug" vq" c" h' wpevkq" grgo gpv'

40 Gnkr uku" qh' c" h' wpevkq" grgo gpv' vq" cp" cwqugo cpve" dcug"

D0 F gr gpf gpe { " ut wewt g" gnkr ugu' <

30 Gnkr uku" qh' c" i q' xgt' p' kpi " w' pk' vq" c" f' gr gpf gpv' * ut wewt cn' gnkr ugu' +

40 Gnkr uku" qh' * cp" q' d' r' k' cvqt { + f' gr gpf gpv' vq" c" i q' xgt' p' kpi " w' pk' * x' crgpe { " gnkr ugu' +

Y g" f' gh' kg" h' qto kpi " w' pk' gnkr ugu" cv' dqy " u{ p' ceve" r{ gtuO' C v' y g" uwt hceg" u{ p' ceve" r{ gt. " c" eqo r r' gz " w' pk' eqpuku" qh' cp" cwqugo cpve" xgt d" y kj " ku" cwz' k' r' ct { " xgt d" = r t' gr qukkqp" ct g" eqppg' evgf " y kj " p' q' w' puO' Uq. " y g" cpcn { | g" gnkr uku" qh' cp" cwqugo cpve" xgt d" vq" cp" cwz' k' r' ct { " xgt d" * v { r g" C O + " cpf " gnkr uku" qh' cp" cwz' k' r' ct { " xgt d" vq" cp" cwqugo cpve" xgt d" * v { r g" C O + " = gnkr uku" qh' c" p' q' w' vq" c" r t' gr qukkqp" * v { r g" C O + " cpf " gnkr uku" qh' c" r t' gr qukkqp" vq" c" p' q' w' * v { r g" C O + " O' Cu" q' p' g' f' ggr " eqo r r' gz " w' pk: " y g" v' t' gcv' k' k' qo u" cpf " eqppg' evqpu" qh' cp" cwqugo cpve" xgt d" y kj " c" o qf' cn' t' r' j' cug" xgt dO'

Vj g" gnkr uku" qh' c" i q' xgt' p' kpi " w' pk' vq" c" f' gr gpf gpv' * v { r g" D O + " ku" y j g" o c' kp" v { r g" qh' gnkr ugu" cv' y j g" uwt hceg" u{ p' ceve" r{ gtO' Vj g" gnkr uku" qh' cp" * q' d' r' k' cvqt { + f' gr gpf gpv' vq" c" i q' xgt' p' kpi " w' pk' * v { r g" D O + " ku" y j g" o c' kp" v { r g" qh' gnkr ugu" cv' y j g" f' ggr " u{ p' ceve" r{ gtO' J gt g. " y j g" gnkr uku" ku" c" o cwgt" qh' x' crgpe { " * y g" y qtn' y kj " y j g" y g' qt { " qh' x' crgpe { " y j cv' y cu" v' t' gcv' f' r' ct' v' ew' r' ct n' d { " Rcp' g' x' q' x' a " * ugg" Rcp' g' x' q' x' a . " 3; : 2 = ugg' c' nuq' W' g- q' x' a . " 4233. " cdq' w' x' crgpe { ' kp" y j g' R' t' ci w' g' F' gr gpf gpe { " V' t' gg' d' c' p' m' O'

"
"
"
"
"

3. Types of Ellipses"

Gnrk ugu"ctg"emcuukhgf "*cu"o qtrj qnqi kcn"uwthceg"cpf "f ggr "u{pvcvle+"ceeqt f kpi "vq"vj g"rc{ gt"qh" vj g"i tco o cvlecriu{ ugo "vj cv"j cu"dggp"cuuki pgf "vq"vj g"go r v{ "r qukskqp0"Vcdng"3."4"cpf "5"lpenwf g" c" dtlgh"uwo o ct{ "qh" vj g" f ghkpgf " v{r gu" qh'gnrk uku'y kj "uqo g" uko r ng" gzco r ngu" vj cv" i kxg" c" uwthelgpv'kmwvckvqp"qh'vj gug'v{r gu"*vj g"gnrk gf "vgz v'ku'lp"us wct g'dtcengw+0"

Table 1. Morphological ellipses."

| Word ellipses" | |
|--|---|
| Gnrk uku'qh'gpf "qh'y qtf "vq"r tghkz" | <i>pre-]grgevkqp_"and post-election discussion"</i> |
| Gnrk uku'qh'dgi lppkpi "qh'y qtf "vq"uwthkz" | <i>in terms of the context(]eqpvzvs)"</i> |
| Gnrk uku'qh'r ctv'qh'eqo r quksk'y qtf" | <i>two-]ugcvgt_"or four-seater car"</i> |

Table 2. Surface syntactic ellipses."

| Forming unit ellipses" | |
|---|---|
| Gnrk uku'qh'cwzkrict { 'xgtd'vq'cwqugo cpvle'xgtd" | <i>Peter will go to Prague and]y km_"visit his mother there."</i> |
| Gnrk uku'qh'cwqugo cpvle'xgtd'vq'cwzkrict { 'xgtd" | <i>(Will you go?) Yes, I will]i q_."</i> |
| Gnrk uku'qh'pqwp"vq"r tgr qukskqp" | <i>in front of]eco gtc_"and behind the camera"</i> |
| Gnrk uku'qh'r tgr qukskqp"vq"pqwp" | <i>in Prague and]kp_"Pilsen"</i> |
| Gnrk uku'qh'xgtd'vq'uwdqtf lpcvki "eqplvpevkqp" | <i>We know, when]uj g'eco g_"and why she came."</i> |
| Gnrk uku'qh'uwdqtf lpcvki "eqplvpevkqp" lp"fr gr gpf gpv'erwug" | <i>He said, that they would come and]vj cv' vj g{ 'y qwf_"stay the night."</i> |
| Gnrk uku'qh'ugeqpf "r ctcvleve"grgo gpv" | <i>Go away otherwise...]A_"</i> |
| Structural ellipses" | |
| Gnrk uku'qh'i qxgtplkpi "xgtd" | <i>Mary likes Bach, Susan]rnngu_"Beethoven."</i> |
| Gnrk uku'qh'i qxgtplkpi "pqwp" | <i>Central]Gwtqr g_"and Eastern Europe"</i> |

Table 3. Deep syntactic ellipses."

| | |
|---|--|
| Forming unit ellipses" | |
| Gnr uku'qh'b qf cn'xgtd'vq'cwqugo cpv'e'xgtd" | <i>Peter wants to relax and j y cpv'u'vq_'listen to music."</i> |
| Gnr uku'qh'c'cwqugo cpv'e'xgtd'vq'b qf cn'xgtd" | <i>(Stay the night.) I could not jvc{ _ "</i> |
| Gnr uku'qh'r ctv'qh'kf kqo " | <i>They buried j y j' g'j' cvej gy_'and then dug up the hatchet again."</i> |
| Valency ellipses" | |
| Vgz wcn'gnr uku" | <i>(Did the shop assistant pack the book?) Yes, he did j y j cv'v' g'dqgm_ "</i> |
| I gpgtcr'cti wo gpv" | <i>Jane sells at Bata j y j cv_'jv'y j qo _ "</i> |
| Eqptqri' | <i>The company planned j y j q/v'j g" eqo r cp{ _"to increase production."</i> |
| Tgekr tqekv' " | <i>John a Mary met j y j q/gcej "qv'j gt_ "</i> |

4. The Prague Dependency Treebanks"

Vj g" Rtc i wg" F gr gpf gpe{ " Vtggdcpm' *RF V+ " ku" vj g" Htuw' eqo r rgz" rki vku'ecm' " o qv'xcv'f " vtggdcpm'ecr wtkpi "cnuq" vj g" f ggr "u{pv'ev'e" utwewtg"qh'ugpv'pegu0'K'ku'dcugf "qp" c" f gr gpf gpe{ " u{pv'ev'e" vj ggt{." qp" vj g" Hwpe'v'kpcn' I gpgt'cv'xg" F guetk'v'kq'0' Vj g" vtggdcpm' eqpuku" qh' eqpv'kpv'v'u'E| gej "vgz'u'o quv'v' "qh'v'j g"lqwt'p'c'v'v'e" ut'ng"cp'cn' | gf "cv'v'j tgg"rc {gtu'qh'cppq'cv'kq'p< o qtr j qnq' k'cn' *o /rc {gt+ " utw'c'eg" u{pv'ev'e" *cp'cn'v'ec'n" c/rc {gt+ " cpf " f ggr " u{pv'ev'e" *v'ge'v'qi tco o c'v'ec'n" v'rc {gt-0' k'p" cf f k'k'q" vq" vj g'ug" vj tgg" cppq'cv'kq'p" rc {gtu' vj gt'g" ku" cnuq" qpg" pqp/cppq'cv'kq'p" rc {gt." tgr t'gug'v'kpi " vj g" o'tcy /vgz'v'0' C'v' vj ku" rc {gt." ecn'gf "y qtf "rc {gt" *y /rc {gt+ " vj g"vgz'v'ku'ugi o g'p'v'f "k'p'v' f q'ewo g'p'u'cpf "r c'tci t'c'r j u'cpf "k'p'f k'k'f w'cn'v'q'ng'p'u'c't'g't'ge'q'i p'k' g'f "cpf " cuu'q'ek'v'f " y k'j " w'p'k'w'g" k'f g'p'v'k'g'tu'0' C'v' vj g" o /rc {gt" gcej " ug'p'v'peg" ku" " r t'q'x'k'f g'f " " y k'j " o qtr j qnq' k'cn' ec'v'gi q't'kgu" *r'go o c." v'ci -0' C'v' vj g" c/rc {gt." c" f gr gpf gpe{ " vtgg" ecr w't'gu" utw'c'eg" u{pv'ev'e" t'gr'v'k'p'u'w'ej "cu'U'w'd'g'ev." Q'd'g'ev." c'p'f "C'f x'g't'd'k'c'r'0' Vj g" j k' j g'u'v'rc {gt" eq'p'v'k'p'u'c'm' vj g" k'p'h'q'to c'v'k'p" vj c'v'ku" g'p'eq'f g'f "k'p" vj g" utwewtg"qh'v'j g'ug'p'v'peg" c'p'f "ku'ngz'k'c'n'k'go u'0'U'q." vj ku"rc {gt" ecr w't'gu" vj g" f ggr ." ugo c'p'v'eq/u{pv'ev'e" utwewtg." vj g" h'w'p'v'k'p'u" qh' ku" r c't'u." vj g" o'f g'gr o' i tco o c'v'ec'n' k'p'h'q'to c'v'k'p." eq't'g'h't'g'peg" c'p'f " v'q'r'k'e/h'q'ewu" c't'v'w'v'v'k'p" k'p'cn'f'k'p'i " vj g" f ggr " y qtf " qtf g't'0' H'k'i w'g" "3" "f'k'ur rc {u" vj g" t'gr'v'k'p'u" d'g'y g'g'p" vj g" p'g'k'i j d'q't'k'p'i "rc {gtu"cu"cppq'cv'f "cpf "

"
"
"
"

tgr tɔgɔpɔvɔf "kɔ" jɔ g" f c v c 0' V j w u " h q t " g z c o r r g . " j ɔ g " E | g e j " u g p v p e g " B y l b y ŝ e l d o l e s a . " r k s g t c m l " k ɔ " G p i r k u j < - J g / y c u ' y q w r f " y g p v ' q h q t g u 0 ' e q p v k p u ' r c u v ' e q p f k k q p c n ' q h ' j ɔ g " x g t d " j t * b y l b y ŝ e l - j g / y c u ' y q w r f " y g p v 0 ' c p f " c ' v l r q " * d o l e s a - j q h q t g u 0 ' "

"
"

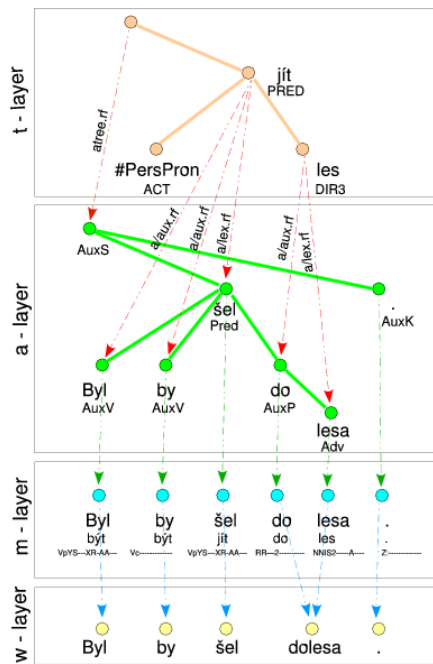


Figure 1. Linking the layers in PDT."

V j g " v c n l ' p w o d g t " q h ' u g p v p e g u " c p p q v c v f " c v ' c m l ' j ɔ g " j ɔ t g g ' r c { g t u " k u " 6 ; . 6 6 4 . " c o q w p v k p i " v q " : 5 5 . 5 7 9 " * q e e w t g p e g u " q h " p q f g u 0 ' V j g " R F V " x g t u k q p " 3 0 " * y k j " j ɔ g " c p p q v c v k p " q h ' j ɔ g " h k t u v " y q " r c { g t u = J c l k " g v ' c r 0 " 4 2 2 3 + " k u " c x c k r d r g " h t q o " j ɔ g " N k p i w k u l e " F c v c " E q p u t w o . " c u " k u " j ɔ g " x g t u k q p " 4 0 " * y k j " j ɔ g " c p p q v c v k p " q h ' j ɔ g " j k f . " f g g r " u { p v c v k e " r c { g t = J c l k " g v ' c r 0 " 4 2 2 8 + 0 ' V j g " r e v u v ' x g t u k q p " R F V " 5 0 " * D g l g m l g v ' c r 0 " 4 2 3 5 + " y k j " u q o g " c f f k k q p u " * g z w c n l ' e q t g h g t g p e g . " f k u e q w t u g " t g r v k p u . " i g p t g " u r g e k h e c v k p . " o w n k y q t f " g z r t g u l k q p u + " k u " c x c k r d r g " h t q o " j ɔ g " N k F C V I E N C T R " t g r q u k q t { 0 C " u k o k r c n l " d c u g f " c p p q v c v k p " j c u ' d g g p " w u g f " h q t " q j g t " R t c i w g " t g g d c p m 0 ' V j g " R t c i w g " E | g e j / G p i r k u j " F g r g p f g p e { " V t g g d c p m l " * J c l k " g v ' c r 0 " 4 2 3 3 + " e q p v k p u " r c t c m g n l " R F V / r k n g " c p p q v c v k p u " q h " G p i r k u j " v z w u " * Y c m l ' U t g g v ' l q w t p c n l ' r c t v ' q h ' R g p p " V t g g d c p m l " c p f " q h ' j ɔ g k t " r t q h g u l k p c n l ' t c p u r v k p " v q " E | g e j 0 ' V j g " R t c i w g " F g r g p f g p e { " V t g g d c p m l " q h ' U r q n g p " E | g e j " * k v ' k u " r r c p p g f " v q " d g " t g n g c u g f " c v ' j ɔ g " g p f " q h ' 4 2 3 6 + " e q p v k p u " u r q p v c p g q w a " f k c n j i w g " u r g g e j . " t c p u e t k d g f . " t g e q p u t w e v g f " c p f " h w t j g t " c p p q v c v f " k ɔ " j ɔ g " R F V " u v l g 0 "

"
"

5. Capturing Ellipsis Techniques"

Kp"vj g"Rtci wg" F gr gpf gpe{ "Vtggdcpmu." gnrk uku"ku"tgcvgf "cv"vj g"j ki j guv"vrc{gt0' Cm' v r gu"qh' gnrk ugu"i kxgp"kp"Ugevkqp"5"*gzemf kpi "uqo g"o kpat"gzegr vkpu+ctg"ecr wtgf "cv"vj ku"rc{gt0'Vj g" r tlpkr ngu"qh"vj g"dwrk/w "qh"vj g"my gt"rc{gtu"kp"RF V" *pco gn{ "vj cv"vj g"pwo dgt"qh"pqf gu"cv" vj gug"rc{gtu"ku"kf gpvkecn'vq"vj g"pwo dgt"qh"vngpu"kp"vj g"ugpvpeg+f q"pqv'cmqy "ecr wtkpi "gnrk uku" d{ "cp"cf f kxqp"qh" c"pqf g" *y j lej "ku"vj g"o quv"kpwkxg"y c{ "qh"ecr wtkpi "gnrk uku"kp"eqtr qtc+0' Vj gtghqtg."cv"vj g"o /rc{gt"cpf "c/rc{gt."gnrk uku"ku"kp f lcvgf "qpn{ "d{ "c"ur gekcn'cwtkdwg"cv"uwej "cp" gztgtuugf "pqf g"vj cv'ecppqv'dg"cppqvcvgf "ceeqt f kpi "vq"uwcn'twrgu"dgecvug"qh'cp"gnrk uku"kp"vj g" ugpvpeg0'c'v'j g"vrc{gt."gnrk ugu"ctg"tgr tguvpvf "d{ "vj g"nkpki "vj g"pqf gu"qh'vj g"j ki j guv"vrc{gt" y kj "vj g"pqf gu"qh'vj g"my gt" c/rc{gt" *Ugevkqp"70+cpf "d{ "cf f gf "pqf gu"cv"vj g"vrc{gt" *Ugevkqp" 70+0'

"

5.1 Links of the t-layer to the a-layer"

Vj g"tgrvkqp"dgw ggp"vj g"vrc{gt"cpf "c/rc{gt"ku"tcvj gt"eqo r rgz0'Vj g"nkpki "dgw ggp"vj g"rc{gtu" ku"pqv'qpn{ "c"vgej plecn's wgnkqp."dw"kw"ecttkgu" c"r kgeg"qh'nkpi wknk "kphqto cvkp0'K'ecr wtgu"vj g" vtcpkkqp"htqo "vj g" *nkpi wknk+ "o gcpkpi "qh" c" ugpvpeg"vq"ku"htqo . "vtcpkkqp"htqo "vj g" f ggr " u{pvcke"utwewt"vq"vj g"utwheg"gzr tguakp0'Vj g"dcule"r tlpkr ng"qh'nkpki "ku"cu"hmjy u<vj gtg" ku" c" nkpki"htqo "vj g" v'pqf g" *pqf g" cv"vj g" vrc{gt+"vq" gcej "c/pqf g" *pqf g" cv"vj g" c/rc{gt+"vj cv' kphwgegu"vj g"xcnwg"qh'uqo g"qh'ku"cwtkdwg0'Dcugf "qp"vj g"v r g"qh'cwtkdwg"vj cv'ku"lphwgegf " d{ "vj g" c/pqf g"y q'v r gu"qh'nkpki"vq"vj g" c/rc{gt"ctg" f khtg gpvkevgf <

/ lex ? "nkpm"vq"vj g" c/pqf g"htqo "y j lej "vj g" v'pqf g"i qv'ku"rgzkecn'o gcpkpi " *qt"ku" dki i guv" r ctv0'Vj g" c/pqf g" lphwgegu"vj g"xcnwg"qh'vj g" lemma cwtkdwg0'Wuwcm{ "k'ku"cp" c/pqf g" vj cv'tgr tguvpw'cp"cwqugo cpvke'y qtf " *pqwp."cf lgevkxg."xgtd."cf xgtd+0'

/ aux ? "nkpm"vq"tgo clpki "c/pqf gu"vj cv'lphwgegu"xcnwg"qh'cwtkdwgu"qh'vj g"i kxgp"v'pqf g0' V{r kecm{ ."vj g" c/pqf gu"chge'v'xcnwg"qh'u{pvcke"cpf "o qtr j qnqi kecn'cwtkdwgu"uwej "cu" functor *f ggr "u{pvcke"hwpevkp+."subfunctor *f gvckgf "emukhkecvkqp"qh'hwpevtu+"cpf " gram *vj g" utwewt" g"qh' cwtkdwgu" vj g" uq/ecmgf "i tco o cvgo gu." y j lej "ecr wtg" f ggr " i tco o cvkecn'eqttgrvgu"qh' vj g" o qtr j qnqi kecn'ecvgi qtkgu+0' Vj gug"ctg" wuwcm{ " c/pqf gu" tgr tguvpki "r tgr qukkqpu."eqplwpevkqpu."cwzkrkt{ "xgtdu."uwr r qt vki "gzr tguakp0'

Hqt"gzco r rg."lp" Hki wtg"3."vj gtg"ctg"vj tgg"nkpm"htqo "vj g"v'pqf g"tgr tguvpki "vj g"xgtd"j i qai q0'vq" vj g" c/pqf g" tgr tguvpki "y qtf " sel -y gpv0'vq" vj g" c/pqf g" tgr tguvpki "y qtf " by -y qwf 0'cpf "vq"vj g" c/pqf g"tgr tguvpki "vj g"y qtf "byl -j g/y cu0'Vj g"nkpm"vq"vj g"htu'v'o gpvkpgf " c/pqf g" *sel -y gpv0'ku"lex v r g=vj g"qj gt"v y q'nkpm'ctg"aux v r gu0'

"
"
"

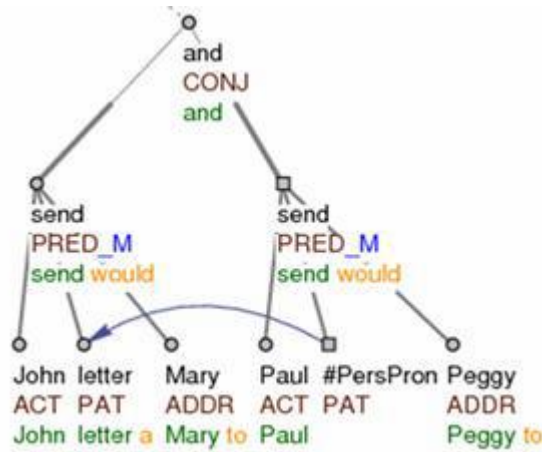


Figure 2. PDT-annotation: *John would send a letter to Mary and Paul to Peggy.*"

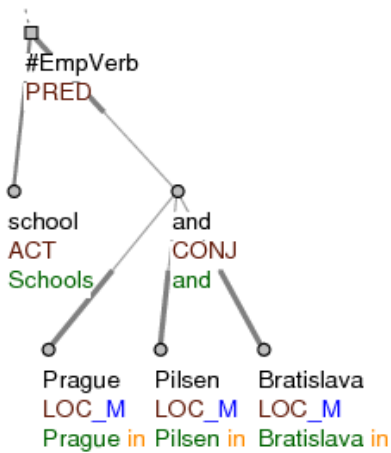


Figure 3. PDT-annotation: *Schools in Prague, Pilsen and Bratislava.*"

Vj gtg"o c {"dg"o qtg"rkpm"v"qpg"c/pqf g"ó"ltqo "xctkqu"v/pqf gu"K"j cr r gpu"lp"ecug"qh"cp" gnr uku"lHqt"gzco r ng."lp"vj g"ugeqpf"erwug"qh"vj g"ugpvpeg"John would send a letter to Mary, Paul to Peggy.."c"v/pqf g"y kni"dg"cf f gf"lqt"vj g"grkf gf"i qxgtplpi "xgtd"cpf "vj gtg"y kni"dg"c"rgz/rkpm"ltqo "vj ku"cf f gf"v/pqf g"v"vj g"c/pqf g"tgr tguvpvpi "vj g"y qtf"send cpf"cwz/rkpm"v"vj g" c/pqf g"tgr tguvpvpi "vj g"y qtf"would"Vj g"kf gp"ecnrkpm"y kni"cnq"dg"wgf"lp"ecug"qh"vj g"v/pqf g"vj cv"tgr tguvpvpi"i qxgtplpi "xgtd"lp"vj g"htuv"erwug"Vj wu"vj gtg"ctg"y q"fqwdng"rkpm"v"c/pqf gu"tgr tguvpvpi "vj g"y qtf u"send cpf"would"ltqo "vj g"vrc{gt"cpf"vj ku"vj g"y c{"y g"ecr wtg"vj g"vgzwcni"gnr uku"qh"c"i qxgtplpi "wplv"*cnpi"y kj"vj g"cf f kkp"qh"vj g"v/pqf g"l"K"vj ku"ecug."vj g"o clp"lpf lecvq"qh"gnr uku"lp"vj g"ugpvpeg"ku"qh"eqwtug"vj g"cf f gf"v/pqf g"J qy gxgt."vj g"vgzwcni"gnr uku"ku"gnr uku"ecr wtg"gf"d{"vj g"fqwdng"rkpm"v"vj g"rc{gt"*d{"vj g"rgz/rkpm"v"vj g"c/pqf g"ltqo "

"
"
"

vj g"cf f gf "vppf g+0Ugg"Hi wtg"4=vj g"rnpnpi "c/pqf gu"ctg"y tkwgp"lp"i tggp"*rgz/rkpm"cpf "qtcpi g"
*cwz/rkpm"eqmqt"d{ "gcej "vppf g0"

Uqo g"v{r gu"qh"gnkr ugu"ctg"qpn{ "ecr wtgf "xlc"rkpm"vq"vj g" c/m{ gt0Hqt"gzco r ng."vj g"hev"
vj cv'r tgr qukkqpu"lp"vj g"itci o gpv'Schools in Prague, Pilsen and Bratislava0ctg"pqv'tgr gcvgf ."ku"
tgr tgugpvf "qpn{ "d{ "cp"cwz/rkpm"vq"vj g"gzr tguugf "r tgr qukkqp"itqo "gcej "qh"vj g"vppf gu"vj cv"
tgr tgugpv"vj g"pco gu"qh"ekkgu"*vj gtg"ku" c"rkpm"vq"vj g"r tgr qukkqp" in itqo "vj tgg"vppf gu="ugg"
Hi wtg"5+0"

"

5.2 Addition of Nodes"

Kp"qtf gt"vq"ecr wtg"vj g"y j qrg"o gcplpi "qh" c"ugpvpeg. "k'ku"uqo gwko gu"pgeguuct { "vq"cf f "vppf gu"
vj cv'meni"vj gk" f kgev"eqwvgr ctv"cv"uwthceg"itqo "qh"ugpvpeg"*cv"vj g" c/m{ gt+0Y g" f kkpki wkuj "
wy q"dcule"v{r gu"qh"cf f gf "vppf gu"

/ copy ?" c" vppf g" vj cv" j cu" vj g" uco g" xcnwgu" qh" egtvklp" cwtkdwgu" *rgo o c." uqo g"
i tco o cvgo gu"cpf "xcrpe{ "itco g+cu"cp"gzr tguugf "vppf g0Y g"ecni"vj ku"cf f gf "vppf g" c"
oeqr {o" qh" vj cv" vppf g0 Vj gtg" ctg" f qwdrg" rgz/rkpm" vq" vj g" c/pqf g" vj cv" tgr tgugpv"
gzr tguugf "eqr kgf "y qtf 0"

/ substitute ?" cp" ctvklcn"vppf g"vq"y j lej "qpg"qh" vj g"hmjy lpi "rgo o c" uwdukwgu" j cu"
dggp"cu"ki pgf <
c0#EmpNoun=#EmpVerb"
d0#Cor. "#Gen. "#Oblfm. "#PersPron. "#Rcp"

Xctkqwu"v{r gu"qh"gnkr uku"ctg" ecr wtgf "d{ "xctkqwu"v{r gu"qh"cf f gf "vppf gu0 Vj g"vgzwni"
gnkr uku"qh"i qxgtplpi "xgtd"qt"pqwp"ku"ecr wtgf "wulpi "oeqr {o"vppf gu0Vj g"u{ ugo "gnkr uku"*y kj "
pq"cpvgeg gpv"qh"i qxgtplpi "xgtd"qt"pqwp"ku"ecr wtgf "wulpi "ouwdukwgo"vppf g"y kj "c"rgo o c"
uwdukwgo"o gpvklpgf "lp"vj g"rkuv" c+"cdqxg0 Vj g"xctkqwu"v{r gu"qh" xcrpe{ "gnkr ugu"ctg"ecr wtgf "
wulpi "vppf gu"y kj "c"rgo o c"uwdukwgo"o gpvklpgf "lp"vj g"rkuv" d+0Kp"ecug"qh"vgzwni"gnkr uku."vj gtg"
ku" c"eqtghgtpeg"*cpcr j qtc+"cttqy "itqo "vj g"cf f gf "vppf g"vq"vj g"vppf g"tgr tgugpvpi "gzr tguugf "
wpk"*cpcr j qt+0"

Kp"vj g"tgg."vppf gu"tgr tgugpvpi "gzr tguugf "y qtf u"ctg" f tcy p"cu"ektergu="cf f gf "vppf gu"ctg"
f tcy p"cu" uswctgu0 Kp"Hi wtg" 4." vj gtg" ku" cp"gzco r ng" qh"ecr wtlpi " c" vgzwni"gnkr uku" qh" vj g"
i qxgtplpi "xgtd0Kp"vj g"ugeqpf "erwug."vj gtg"ku"cp"cf f gf "vppf g"tgr tgugpvpi "vj g"i qxgtplpi "xgtd"
cpf "vj ku"vppf g"ku" c"eqr { "qh"vj g"vppf g"tgr tgugpvpi "vj g"gzr tguugf "xgtd"lp"vj g"htuv"erwug0Vj g"
u{ ugo "gnkr uku"qh"vj g"i qxgtplpi "xgtd"ku"ecr wtgf "lp"Hi wtg"50Vj gtg"ku"cp"cf f gf "vppf g"y kj "
#EmpVerb"rgo o c"uwdukwgo"qp"vj g"i qxgtplpi "xgtd"r qukkqp0"

"
"
"

kp'Hki wtg'4."vj gtg'ku'cnuq'cp"gzco r rg"qh'vgz wcn'xcrgpe{ "gnkr uku0Vj g'grkf gf "Rcvkpv'letter
*kp"vj g"ugeqpf "ercwug+"ku'ecr wtgf "wukpi "cp"cf f gf "vppqf g'y kj "#PersPron rgo o c0Vj gtg'ku"c"
eqtghgtgpeg'cttqy "Itqo "vj g'cf f gf "vppqf g"vq"vj g'gzr tguugf "Rcvkpv'kp"vj g'htuv'ercwug0C"uko krcr"
gzco r rg"qh'vgz wcn'xcrgpe{ "gnkr uku'ku'cnuq'kp"Hki wtg'60

6. Summary of Representation of Ellipsis"

Vj ku'ugevqp"dtlpi u'cp'qxgtxky "qh'vj g'y c{u'vj g'fghkpgf "gnkr ugu'*ugg'vj g'uwo o ct{ 'lp"Ugevqp"5+"
ctg'tgr tguugvgf "cv'vj g'vrc{ gt0Utwewtcn'gnkr ugu'*v{r g'D0'lp"Ugevqp"4+'ctg'cny c{u'ecr wtgf "d{ "
cp"cf f kkpq"qh'c"oeqr { "o"qt"c"ouwdunxwgo"vppqf g"*ugg"Vcdng"6="kh'vj g'gnkr uku'ku'ecr wtgf "d{ "c"
vppqf g'y kj "c"rgo o c'uwdunxwgo."qpn{ "vj g'rgo o c'ku"o gpv'kpgf "kp"vj g'vcdng+0Xcrgpe{ "gnkr ugu"
*v{r g'D0'lp"Ugevqp"4+'ctg'cny c{u'ecr wtgf "d{ "cp"cf f kkpq"qh'c"vppqf g'y kj "c"rgo o c'uwdunxwgo=
vj g'qxgtxky "qh'xcrgpe{ "v{r gu'qh'gnkr ugu'cpf "vj gkt'rgo o cu'uwdunxwgo'ku'r tguugvgf "kp"Vcdng'70

"
"

Table 4. Representation of structural ellipses in PDT."

| | Textual ellipsis" | System ellipsis" |
|-------------------------------|------------------------------|---------------------------|
| Gnkr uku'qh'i qxgtplpi "xgtd" | eqr { "vppqf g" rgz/rkpm' | #EmpVerb" pq'rgz/rkpm' |
| Gnkr uku'qh'i qxgtplpi "pqwp" | eqr { "vppqf g" rgz/rkpm' | #EmpNoun" pq'rgz/rkpm' |

Table 5. Representation of valency ellipses in PDT."

| | |
|---------------------------|-----------------------------------|
| Qdrki cvqt { 'cti wo gpv' | #PersPron" eqtghgtgpeg'cttqy " |
| Eqpvtqmgf "cti wo gpv' | #Cor" eqtghgtgpeg'cttqy " |
| Tgek' tqekv{ " | #Rcp" eqtghgtgpeg'cttqy " |
| I gpgtcn'cti wo gpv' | #Gen" pq'eqtghgtgpeg'cttqy " |

Cv' vj g" vrc{gt." y g" cmq" ecr wtg" dqj " v{r gu" hqto lpi " wplk' gnrk ugu' uwt hceg" u{pvele" #hqto lpi " wplk' gnrk ugu"lp" Vcdng"4+"cpf "f ggr "u{pvele" #hqto lpi " wplk' gnrk ugu"lp" Vcdng"5+0'Vj g" o gj qf "qh" ecr wt lpi " ku" uko krc" hqt" dqj " v{r gu' Vj g" gnrk uku"qh" cp" cwqugo cpve" dcug" vq" c" hwpvklp" grgo gpv* v{r g" C08" lp" Ugevkp"4+"ku" ecr wt gf " ceeqtf lpi " vq" vj g" uco g" twgu" vj cv' cr r n{ " vq" vj g" gnrk uku"qh" c" i qxgtplpi " xgt d" qt" pqwp. " k0y kj " vj g" j gr " qh" vj g" oeqr { o" vppqf g" #kp" ecug" qh" vgz wcn' gnrk uku+ " qt" ouwdunkwgö" vppqf g" y kj " vj g" #EmpVerb qt" #EmpNoun rgo o c" #kp" ecug" qh" u{vgo " gnrk uku+0'Vj gt g" ctg" cwz/ rkpni" vq" c/ pqf gu" tgr tguvp lpi " gzi tguuf " hwpvklp" grgo gpw0' Hqt" c" erget" qxgt xky . " ugg" Vcdng"80' Cp" gzco r rg" qh" ecr wt lpi " cp" gnrk uku"qh" cp" cwqugo cpve" dcug" vq" c" hwpvklp" grgo gpv' ku" lp" Hki wt g" 60' Vj gt g" ku" cp" gnrk uku"qh" vj g" xgt d" obtain vq" vj g" o qf cn' xgt d" want lp" vj g" ugeqpf " erwug0

Gnrk ugu"qh" c" hwpvklp" grgo gpv' vq" cp" cwqugo cpve" dcug" #v{r g" C04" lp" Ugevkp"4+" ctg" cny c { u" ecr wt gf " qpn{ " d { " o gcpu" qh' rkpni" vq" vj g" c/ r { gt0' Ugg" vj g" gzco r rg" Schools in Prague, Pilsen and Bratislava lp" Hki wt g" 5" y j gt g" gnrk uku"qh" c" r tgr quklp vq" pqpwu" ku" ecr wt gf 0

Table 6. Ellipsis of autosemantic base to function element in PDT."

| Textual ellipsis" | System ellipsis" |
|-------------------|---------------------|
| eqr { 'vppqf g" | #EmpVerb !#EmpNoun" |
| rgz/ rkpni' | pq' rgz/ rkpni' |
| cwz/ rkpni" | cwz/ rkpni" |

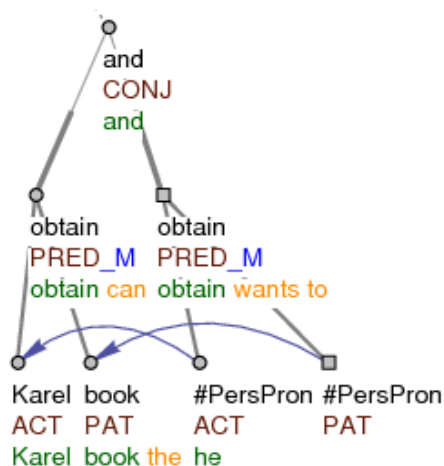


Figure 4. PDT-annotation: Karel can obtain the book and he wants to."

Vj g"o clqtka{ "qh"vj g" fghkpgf "v{r gu"qh"gnkr ugu"*ugg"vj gkt"uwo o ct{ "kp"Ugevkqp"5+"ku"ecr wtgf "kp"vj g"RF V0J qy gxgt. "vj g"o qtr j qnqi lecn'y qtf "gnkr uku"*ugg"Vcdrg"3+"uvm't go clpu"wpcppqvvgf 0 Y g"uwi i guv"vj cv'ur gekn'o qtr j qnqi lecn'vci u'uj qwf "dg"ko r rgo gpvgf "vq"t guqkkg"vj g"ecugu'y j gp" qpn{ "c"r ctv'qh" c"y qtf "r tghkz. "uwhkz. "gpf kpi . "r ctv'qh" c"eqo r qukg" y qtf "+"ku"ugr ctvvgf "cu"cp" kpf kxf wcn'rgzlecn'wpk'dcugf "qp"vj g"r tpekr ng"qh"vqngpk cvkqp"oitqo "c"ur ceg"vq" c"ur cegö0P gzv." y g"uwi i guv'ecr wtgpi "vj g"y qtf "gnkr ugu"kp" c"uko krt"y c{ "vq"vj g"tgr tgugpvkqp"qh"gnkr uku"qh"cp" cwqugo cpvle"dcug"vq" c"hwpevkqp"grgo gpv"*ugg"Vcdrg"8+0"Vj g"rgo o c"qh"vj g"v'pqf g"tgr tgugpvkpi " vj g"grkf gf "y qtf "uj qwf "i gv" c"pqp/grkf gf . "hwn'rgzlecn'hqto "cpf "vj g"rgzlecn'cpf "o qtr j qnqi lecn' cwtkdwgu"qh"vj ku"v'pqf g"uj qwf "eqttgur qpf "vq" c"pqp/grkf gf "tgeqpuw wevgf "hqto 0Vj gtg"uj qwf "dg" c"rgz/rkpn"vq" c/pqf g"tgr tgugpvkpi "c"pqp/grkf gf"y qtf "kp"ecug"qh"gz wcn'gnkr uku+"qt"vj gtg"uj qwf " dg"pq"rkpn'cv'cmi"kp"ecug"qh"u{vgo "gnkr uku+0"Vj g" c/pqf g"tgr tgugpvkpi "vj g"gzr tguvgf "r ctv'qh" c" y qtf "r tghkz. "uwhkz. "gpf kpi . "r ctv'qh" c"eqo r qukg" y qtf "+"uj qwf "dg"ecr wtgf "cu"cp"cwz/rkpnkpi " c/pqf g0

7. Conclusion

Kp"vj g"ctveng."y g" fgrko kcvg"vj g"vj gqtgvlecn'dcuku"qh"gnkr uku."dtkpi "c"r tqf wevkg"emcukhlecvkqp" cpf "f guetkdg"ku"ugo cpvle'tgr tgugpvkqp"kp"vj g"Rtci wg" F gr gpf gpe{ "Vtggdcpm0

Vj g"Rtci wg" F gr gpf gpe{ "Vtggdcpm"j cu'citgcf { "dggp"vugf "kp"o cp{ 'r tqlgew."dqvj "kp"rkpi wkuveu" tgugctej "cpf "hq" "vj g" f gxgnr o gpv"qh"PNR"ngctpkpi "cni qtkej o u"cpf "uqhy ctg"vqqu0"Vj g" cpcn{ vlcni'cpf "rcvt"vj g"vgevqi tco o vlcni'rc{gt"j cu'dggp"vugf "cu" c"dcuku"ht"eqpxgtukqp"vq"vj g" y gm'npqy p"EqP NN"hqto cv'kp"vj g"4229"cpf "422; "uj ctgf "vcumi"*j clk "gv'cni"422; +0"kp"vj gug" uj ctgf "vcumi."qxt"52"r ctugtu"cpf "UTN"vqqu"j cxg"dggp"etgcvgf ."cu" f guetkdgf "kp"vj g"r tqeggf kpi u" qh"EqP NN"4229"cpf "EqP NN"422; ""Uj ctgf ""Vcumi0"Rctcmgn"E| gej /Gpi rkuj ""F gr gpf gpe{ " Vtggdcpni"j cu'dggp"cmq"vugf "hq"r gctpkpi "r ctugtu"cpf "ugo cvle"tqrg"rdgrtu."y j lej "j cxg"kp"wtp" dggp"vugf "hq"o cej kpg'tcpurvkqp"r tqlgew" *Dqlct"gv'cni"422; =Vco ej {pc"gv'cni"4236=O ctg gm" Rqr gn"ficdqntum . "4232+0"dcugf "qp"vj g" f cvc"cpf "gzr gtlgpeg"y kj "vj g"RF V."c"pgy "o wnktpi wcn' vggdcpniJ co rGF V"y kj "52"vggdcpni"kp" f khtgtpv'rcpi wci gu"j cxg"dggp"eqpxgtvgf "vq" c"wpkhqto " hqto cv'cpf "o cf g"cxkkrdrg"*Tquc"gv'cni"4236+0

"

Acknowledgments

Vj g"tgugctej "tgr qtvgf "kp"vj g"ctveng."y cu'uw r qtvgf "d{ "vj g"O lkurt { "qh"Gf wecvkqp."[qwj "cpf " Ur qt v'qh"E| gej "Tgr wdike"y kj kp"vj g"r tqlgew"NR F CV/Ertkp"NO 42322350"Vj g"r ctvlekr cvkqp"qp" TQENR I "4236"eqphgtgpeg"y cu'cmq"uw r qtvgf "d{ "vj g"v'qf c"Cwq"eqo r cp{0

"
"
"
Reference

Dgl gm"G0'gv'crl" *4235-0'Prague Dependency Treebank 3.00'F cvc luqhy ctg."Kpukswg"qh"Hqto cri' cpl" Cr r rkgf" Nlpi wkuu." Ej ctrgu" Wpkxgtukf " lp" Rtcu wg." Rtcu wgo' WTN< j wr <lhvcr0 hfwpk0 l r f v0l"

Dgto cp." U0'6" J guwkm" C0' *gf u0" *3; ; 4-0' Proceedings of the Stuttgart Ellipsis Workshop' Ctdgur cr lgtg" f gu" Uqpf gthqtuej wpi udgtglej u" 562." Dgtlej v' P q0' 4; ." Wpkxgtukf " qh' Uwwi ctv."Uwwi ctv'

Dqlct" Q0" Ecrkuqp/Dwtej " Ej 0" J clk " L0" Mqgj p" R0' *422; +< Ur gekri' Kuwg" qp" Qr gp" Uqwtg" Ocej lpg" Vtcurvqpp" Vqqu."The Prague Bulletin of Mathematical Linguistics."Ej ctrgu" Wpkxgtukf " lp"Rtcu wg."Rtcu wg."P q0'; 3."KUDP "; 9: /: 2/; 26397/3/; ."KUP "2254/87: 70

D3/4 o qx^."C0"J clk ."L0"J clk qx^."G0"J rcf n1."D0" *4225-0'Vj g"Rtcu wg" f gr gpf gpe{ "vtggdcpn^c" 5/rxgn' cppqcvqpp" uegpctk0' k0"Treebanks: Building and using parsed corpora, gf 0' C0' Cdgm? ."Mwy gt."F qtf tgej v."325/34: 0'

J clk "L" Ektco kc" O 0" Lqj cpuuqp" T0" Mcy cj ctc" F 0" O ctv" O 0' C0" O cts wgl " N0" O g{ gtu" C0" P kctg" L0" Rcf > "U0" r^ a pgn' L0" Utc ^ n' R0" Uwf gcpw' O 0' Z wg" P kcy gp. \ j cpi " [k^ *422; +< Vj g" EqP NN'422; " Uj ctgf " Vcum^ U{ pcvle" cpl " Ugo cpvle" F gr gpf gpekgu" lp" O wnr rg" Ncpi wci gu' k0' < Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL): Shared Task."Cuqekvqpp" hqt" Eqo r wvqpcn' Nlpi wkuu." Dqwf gt."EQ."WUC."KUDP "; 9: /3/; 54654/4; /; ."3/3: 0'

J clk ."L0" gv" crl" *4228-0' The Prague Dependency Treebank 2.0. EF/TQO 0' Ecw0' P q0' NFE4228V23." Nlpi wkuu" F cvc" Equqt vwo ." Rj krf gr j kc0'

J clk ." L0" J clk qx^ ." G0" Rcpqxq^ ." L0' gv' crl" *4234-0' Cppqwpekpi " Rtcu wg" El gej /Gpi rkuj " F gr gpf gpe{ " Vtggdcpn' 400' k0" Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC 2012)." Gwtqr gcp" Ncpi wci g" Tguqwtegu" Cuqekvqpp."Kncpdwn"KUDP "; 9: /4/; 73962: /9/9."5375/53820'

J clk ." L0' gv' crl" *4233-0' Prague Czech-English Dependency Treebank 2.00' F cvc luqhy ctg." Kpukswg"qh"Hqto cri' cpl" Cr r rkgf" Nlpi wkuu."Ej ctrgu" Wpkxgtukf " lp" Rtcu wg." Rtcu wgo'

J clk ."L0'gv'crl" *4223-0'Prague Dependency Treebank 1.0 (Final Production Label)0'EF/TQO 0' Ecw0' P q0' NFE4223V32." Nlpi wkuu" F cvc" Equqt vwo ." Rj krf gr j kc." KUDP " 3/7: 785/ 434/20'

J rxcu." \ 0' *3; ; 3-0' Qp" U{ pcvlecri' Gnr uku' k0" Satzsemantische Komponenten und Relationem im Text."gf u0' F cpge-."H0' Xlgy gi gt." F 0"—L " UCX." Rtcu wg."33; /34: 0'

Ncr r kp." U0'6" Dgpo co qwp." G0' *gf u0" *3; ; ; 4-0' Fragments. Studies in Ellipsis and Gapping." Qzhqtf " Wpkxgtukf " Rtguu." P gy " [qtm" KUDP "2/3; /734524/80'

O ctg gni' F 0" Rqr gni' O 0' fiedqntum " \ 0' *4232-0' O czko wo " Gpvtqr { " Vtcurvqpp" O qf gni' lp" F gr gpf gpe{ /Dcugf " O V" Hco gy qtn0' k0' < Proceedings of the Joint Fifth Workshop on Statistical Machine Translation and MetricsMATR." Cuqekvqpp" hqt" Eqo r wvqpcn' Nlpi wkuu." Wf r ucrc." Uy gf gp." KUDP "; 9: /3/; 54654/93/; ."423/4230'

"
"
"

- O kmwqx^a. " O' gv' cř' *4228+< Annotation on the tectogrammatical level in the Prague Dependency Treebank. Annotation manual' Vgej plecri' tgr qt v' pq' 4228 52." ěřvkwg" qh' Hqto cř'cpf' Cr r rġf' Nġpi wkuġeu.'Ej ctrgu'Wpkxgtuk{ "ġ"Rtci wg.'Rtci wg' O ctewu." O' gv' cř' *3; ; 7+0' Treebank 20' EF/TQO' Ecv' P q' NFE; 7V9." Nġpi wkuġeu" F cv' Eqpuqt vkw. 'Rj křf gr j k' O' Tquc" T0" O c-gm' L0" O ctg gm' F 0" Rqr gr' O 0" \ go cp" F 0" fcdqntvum "\ 0' *4236+0' J co ġF V" 40< Vj kv{ " F gr gpf gpe{ " Vtggdcpnu" Ucphtf k ġf 0' ġ< Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC 2014)." Gwtqr gcp" Ncpi wci g" Tguqtegu' Cuuġek vġp. "Tg{ nġx ġm' Ěgrcpf. "KUDP "; 9: /4/; 73962: /: /6." 4556/45630' U cm" R0' gv' cř' *3; ; 8+0' The Meanings of the Sentence in Its Semantic and Pragmatic Aspects." F qtf tgej v. " F 0' T ġf gr' Rwdġuj ġpi 'Eqo r cp{ "/ "Rtci c." Cef go k. "Rtci wg' Vco ej { pc" C0" Rqr gr' O 0" Tquc" T0" Dqġct" Q0' EWP Kġp" Y O V36" *4236+0' Ej ko gtc" Uġm' Cy cku" Dgnġtqr j qp' ġ< Proceedings of the Ninth Workshop on Statistical Machine Translation." Cuuġek vġp" ġt" Eqo r wcvġpcř' Nġpi wkuġeu. "Dcġko qtg. "O F. "WUC. "KUDP "; 9: /3/; 63865/ 39/6." 3; 7/4220' Wg-qx^a "\ 0' *4233+< Valenġnġ slovnġk Prařskġho zavislostnġho korpusu (PDT-Vallex). ěřvkwg" qh' " Hqto cř' " cpf "" Cr r rġf "" Nġpi wkuġeu. "" Ej ctrgu "" Wpkxgtuk{ "" "" ġ "" "" Rtci wg. "" "" Rtci wg' " KUDP "; 9: /: 2/; 26793/3/; 0'

Vj g'4236'Eqphgtgpeg'qp'Eqo r wcvkqpcn'Nkpi wku'leu'cpf 'Ur ggej 'Rtqegu'kpi "
 TQENRPI '4236.'r r 035;/374"
 © Vj g'Cuuqek'v'k'p'ht'Eqo r wcvkqpcn'Nkpi wku'leu'cpf 'Ej kpgug'Ncpi wci g'Rtqegu'kpi "

Sketching the Dependency Relations of Words in Chinese"

Meng-Hsien Shih* and Shu-Kai Hsieh*

Abstract"

Y g'r tqr qugu'c'ncpi wci g'tguqwtg'd{ 'cwqo cvkcm{ 'ungvej kpi "i tco o cvkcn'tgrv'kpu" qh'y qtf u'dcugf "qp" f gr gpf gpe { "r ctugu'ht qo "wpvci i gf "vgz u0'Vj g'cf xcpvci g'qh'y qtf " ungvj "dcugf "qp" r ctugf "eqtr qtc" ku." eqo r ctgf "vq" Ungvej "Gpi kpg" *Mkri cttkh "T { ej n{ ." Uo t| ."("Vvi y gm"4226+ "vq" r tqxkf g" o qtg" f gvcku'cdqww'v'j g' f khgtgp'wuci g'qh'gcej " y qtf "uwej "cu" xctkqu" v' r gu'qh' o qf khkcv'kpp." y j kej "ku" cnuq "ko r qtvcv' k'p" ncpi wci g' r gf ci qi { 0' Cmj qwi j "uqo g'ncpi wci g'tguqwtg" qh'q'v' gt "ncpi wci gu" j cxg" cwgo r vgf "vq" ungvj "y qtf u'dcugf "qp" r ctugf "f cvc." k'p" Ej kpgug" y g' j cxg" pqv' uggp" c" tguqwtg" hqt " f gr gpf gpe { " ungvj "qh' y qtf u' k'p" ewxqo k' gf "vgz u0' Vj gtghqtg." y g' r tqr qug" uvej "c" tguqwtg" cpf "gxcncv" y kj "Ej kpgug" Ungvej "Gpi kpg" *J wcp "gv' cr0"4227+ "k'p" vgo u'qh' eqttgur qpf kpi 'v'j gucvw' hmpcv'kpp0

Keywords: F gr gpf gpe { 'i tco o ct. 'I tco o cvkcn'tgrv'kpp. 'P NR'v'qqu'it guqwtg u0

1. Introduction"

U{ pvc i o cvk "tgrv'kpcn' kphqto cvkpp" j cu' dggp" v' j g' hqewu' qh' v' j g' k'p'v'gthceg" uwf lgu' qh' u{ pvcz " cpf " ugo cpv'leu0' Y kj " v' j g' tcr kf " f gxgnr o gpv' qh' eqtr qtc." xctkqu" eqtr wu" s wgt { ." r tqh'k'kpi " cpf " xkuwck' cvkpp" v'qqu" j cxg" go gti gf "s wlem{ "qxgt" v' j g' r cuv" { gctu0' Co qpi "v' j gug" v'qqu. "Y qtf "Ungvej " Gpi kpg" *Mkri cttkh" gv' cr0"4226+ "J wcp " gv' cr0"4227+ "j cu" r tqxkf gf "cp" gh'gcv'kpg" cr r tqcej "vq" s wcpv'kcv'kgn{ " uwo o ct k' g" i tco o cvkcn' cpf " eqm'qecv'kpp" dgj cxkqt "3'0' Vj g' r tqxkf gf " hmpcv'kpu" k'p'ncv' g' Eqpeqtf cpegt. "Y qtf "Ungvej . "Ungvej "F kh" "Vj gucvw." cpf "qv' j g' y gd" eqtr wu' etcy k'pi " cpf 'r tqegu'kpi "v'qqu0

Rt gxkqu" rkgctw'gtu" j cxg" t'gxg'crgf " v' j cv' eqtr wu" k'pi wku'leu" j cu' d'gpg'k'hsf " i tgcw{ " ht qo " Ej kpgug" Ungvej "Gpi kpg" *J qpi "(" J wcp " . "4228+0' Cmj qwi j " r tqr t'kgw{ . " Y qtf " Ungvej " Gpi kpg" u{ ugo "ku" r qr w'ct" co qpi "eqtr wu" k'pi wku'leu" cpf "ncpi wci g' v'gcej gtu' d'gecv'wug" qh' ku' hmpcv'kpu" k'p" ncpi wci g' cpcn{ uku0' Jy gxgt. "v' j g' eqp'ut wcv'kpp" qh' Ungvej "Gpi kpg" ku' v'ko g' eqpuwo k'pi "f w'g" v' j g' "

* I tcf wcv' k'p'v'k'w'g' qh' Nkpi wku'leu. 'P cvkqpcn' Vcly cp' W'p'k'x'gt'uk{ . 'Vck' gk' Vcly cp' "

G' o ckn' } uko qp'0' kcp. "v' j w'ck' B i o ckn'0' qo "

³ j w' <1y y y 0' ungvj gpi kpg'0' q'0' w' "

"
"
"

o cpwcm{ " gf kqf " ungvj " i tco o ct0 J gtg" y g" r tqr qug" cp" cngtpcvkxg" cr r tqcej " vq" ungvj " vj g" i tco o ct' r tqhkg'qh'y qtf u'cwqo cvkcm{ 'htqo "c'vzv'eqtr wu0

Vj g" r cr gt" ku" qti cpk{ gf " cu" hqmjy u< Ugevkqp" 4" tgvkgy u" vj g" ewtgpv' f guki p"qh' tgrvqf " rpi wci g" tguqwtegu0 Ugevkqp" 5" f guetkdgu" vj g" r tqr qugf " o gvj qf "qh' ungvj kpi " y qtf u" kp" c" r ctugf " eqtr wu0 Ugevkqp" 6" r tguvpu" vj g" tguwmu" htqo " vj g" r tqr qugf " cr r tqcej " cpf " gxcnvcvkp0' k{ vj g" hpcn' ugevkqp. " y g' j cxg" c" dtkgh' eqpenwukqp" hqt " vj ku' r cr gt0

2. Review"

Y qtf "Ungvej "Gpi kpg" *Y UG+ " r tqxf gu" c" ugv" qh' eqtr wu" s wgt { "vqnu" vj cv' clo u" vq" j gr " wgtu' tgvkgn' rpi wkuke " r cwgtpu" kp" rpi wci g" wug0' Co qpi " vj gug" vqnu. " y qtf " ungvj " hwpevkp" i ckpu" vj g" o quv' r qr wrtk{ " cpf " j cu" y kf gr{ " cr r rkgf " kp" vj g" uwf lgu" qh' eqtr wu" rpi wkukeu" cpf " rpi wci g" r gf ci qi { " *Mki cttkh' 4229:0

I kxgp" vj g" r tgr tqeguqf " eqtr wu" f cvc. " vj g" cxckrdng" Y UG" u{ vgo " kp" o quv' rpi wci gu" o cngu" wug" qh' tgi wrt " gztguukpu" vq" gztcev' i tco o cvkcn' kphqto cvkqp" htqo " c" RQU" vci i gf " eqtr wu0 Vj g" uq/ecnf " sketch grammars. " o quv{ " o cpwcm{ " etchgf " d{ " rpi wkuu. " f guetkdg" vj g" tgrvkqp" dgvy ggp" c" vti gv' y qtf " cpf " ku" f gr gpf gpv. " eqputckpof " qp" vj g" uwtqwpf kpi " eqpvz0' k{ ku" f guki p" qh' i tco o ct" gpi kpggtkpi . " vj g" ungvj " i tco o ctu' ctg' wugf " hqt' hpkq/ucv' vj g" cmjy " r ctukpi " vq" gztcev' vj g" f hhtgpv' i tco o cvkcn' tgrvkpu" 4" 0' V{ r lecn' tgrvkpu" kp" Gpi rkj " Y UG" kpenf g<]QDLGE VaQH_ "]CF LaO QF KHGT_ "]P QWP aO QF KHGT_ "]O QF KHGU_ "]CP F IQT_ "]RRa kP VQ_ " gve0

k{ vgtu u" qh' eqtr wu" rpi wkukeu. " vj g" " sketch hqt" c" " y qtf " r tguvpu" c" " ecpf kf cvg" " ugv" qh' ku" collocates qti cpk{ gf " d{ " vj gk" i tco o cvkcn' tgrvkpu" vj g{ " ucpf " kp" vq" vj g" vti gv' y qtf 0' Vj gug" eqmtecvu' ctg' uqtvgf " ceeqtf kpi " vq" egtvcp" ucvkuke" o gcuwg" qh' eq/qeewtgpqeg. " cu" kmwutcvgf " kp" vj g" ecug' qh" 打j öj kō<

"
"
"
"
"
"
"
"
"
"
"
"
"
"

⁴ j wr <ly y y Ungvej gpi kpg0eq0mif qewo gpcvkp ly knkUnGU gr IETgcvEqtr wu"
⁵ j wr <ly qtf ungvj 0kpi 0kplec0gf w0y "

打 sinica freq = 2695

| PP 在 59 5.4 | Object 1834 3.4 | SentObject_of 122 3.3 | Modifier 776 2.6 | Subject 1173 2.2 |
|-------------|-----------------|-----------------------|------------------|------------------|
| 臉 8 19.07 | 電話 92 32.92 | 敢 13 21.23 | 去 49 19.95 | 武松 13 28.39 |
| 身 6 14.45 | 折 30 32.3 | 開始 12 16.41 | 要 55 18.61 | 棍子 5 16.55 |
| | 籃球 42 32.29 | 喜歡 10 15.84 | 愈 9 17.32 | 球 9 14.15 |
| | 高爾夫球 23 31.69 | 怕 6 14.05 | 就 41 16.82 | 我 99 13.91 |
| | 零工 14 29.41 | 繼續 7 13.95 | 先 18 16.78 | 你 49 13.54 |
| | 仗 15 26.76 | 知道 8 11.88 | 再 26 15.86 | 他 66 11.62 |
| | 招呼 16 24.74 | | 不會 15 15.2 | 爸爸 9 11.06 |
| | 折扣 18 24.03 | | 該 9 14.87 | 雨 5 9.77 |
| | 交道 7 23.57 | | 一起 10 14.11 | 電話 8 9.63 |
| | 呵欠 8 23.5 | | 不能 14 13.88 | 人 47 8.49 |
| | 勝仗 6 21.95 | | 會 31 13.7 | 人家 5 8.11 |
| | 太極拳 9 21.74 | | 連 7 13.35 | 老師 10 7.46 |
| | 冷顫 6 21.33 | | 一直 10 13.35 | 妳 7 7.22 |
| | 寒顫 6 21.33 | | 一 17 13.17 | 他們 16 7.1 |
| | 寒噤 5 20.96 | | 亂 5 13.11 | 她 22 7.05 |
| | 高爾夫 11 20.16 | | 不要 10 12.84 | 門 5 6.98 |
| | 嗝 6 19.68 | | 不 42 12.58 | 誰 5 5.86 |
| | 預防針 6 19.53 | | 各 6 11.37 | 媽媽 5 5.26 |
| | 敗仗 6 19.53 | | 能 20 11.23 | 同學 5 5.21 |
| | 盹 5 19.44 | | 可以 18 10.93 | 學生 10 4.71 |
| | 虎 14 19.35 | | 還 15 10.39 | 我們 15 4.67 |
| | 場 33 19.22 | | 別 5 9.76 | 孩子 6 4.55 |
| | 羽毛球 6 19.11 | | 又 11 9.6 | 自己 10 3.44 |
| | 排球 9 19.11 | | 只 11 9.18 | 時候 5 3.28 |
| | 強心針 5 18.92 | | 都 15 8.98 | |

Figure 1. Word sketch of 打“hit”

Vj g"eqtg"eqo r ppgpv"kp"Y UG"u{ ugo "ku"vj g" sketch grammar." y j kej "f ghkpgu"vj g" rkpget" r cwgtpu"y kj "tgi wrct"gzr tguakp"lqt"vj g"u{ ugo "q"cwqo cvkcm"kf gpvkh{ "r quikdrg"tgrcvkpu"vq"vj g" vcti gv'y qtf 0Hqt"kpucpeg."ppg"qh"vj g"ungvej "i tco o ct"twrgu"f ghkpgf "kp"vj g"j wi g"Ej kpgug"eqtr wu" *j VgpVgp33."y kj "408"dlarkp"vqngpu+"r tqxkf gf "d{ "Y UG"ctg"eqpegtpgf "y kj "o qf khlecvkp0Vj cv" ku."y g"ecp"kf gpvkh{ "vj g"ecugu"qh'o qf khlecvkp"tgrcvkpu"y j gtg"vj g"vcti gv'y qtf "kpf kcvgf "d{ "vj g" r tghk"ö3-ö+"ecp"dg"cp{ "pqpw"lqmny gf "d{ "ppq/pqwpu0Cpf "vj g"eqmqecvgu"lq0"vj cv'y qtf u"y g" y cpv"q"ecr wtg"o ctngf "y kj "vj g"r tghk"ö4-ö+ku"vcngp"vq"dg"cp{ "xgtd"lqmny gf "d{ "e'y qtf "的"

*DUAL"

=A_Modifier/Modifies"

2:]vci ?\$XQ \$_]y qtf ?\$f\$_]vci ?\$P Q \$_]2.1:]vci ?\$P Q \$_]vci #?\$P Q \$]"

"
"
"
"
"
"

Vj g"ungvej "i tco o ct"ecp"dg"o qtg"eqo rñecvqf "y kj "vj g"i tcpwætkv "qh"RQUV Vj g"hqmqy lpi "
i tco o ct"uj qy u"vj g"ercuukñecvqf"tgrvqpp"fgxgnr gf "d{"J wcpv "gv'cn0*4227+"cpf "ko r ngo gpvqf "kp"
vj g"Ej kpgug"Y qtf Ungvej "u{ungo ⁶."kq0"vj g"vci gv'y qtf "ecp"dg"c"pqwp"r tgegfg gf "d{"c"o gcuwtg"
y qtf "vci i gf "d{"P h<

"
"

? O gcuwtg"

"

4-\$P h0 \$""*\$C\$XJ 33\$XJ 35\$XJ 43\$XQ \$""\$F G\$+""\$vci ?\$P jcddef _Q \$""(""\$vci #?\$P ef \$"
3-\$vci ?\$P jcddef j h_0 \$" (" vci #?\$P de0 \$" (" vci #?\$P ef 0 \$" (" y qtf #?\$" 者 \$" (" y qtf #?\$" 們 \$"
jvci #?\$P jcddef j gh_0 \$vci ?\$P de0 \$vci ?\$P ef 0 \$"

"
"

J qy gxt."vj g"y tkkpi "qh"i tco o ct"ku"vko g/eqpuwo lpi ."twppkpi "tkun'qh"-ny "tgecmæ"uq"y g'wtp"vq"
gzt nkv"vj g"fg gr gpf gpe{"r ctugt"ht"gpolej lpi "vj g"tgrvqpcn'kphqto cvkqp0'Wpikng"r j tcug/utwewtg"
i tco o ct."f gr gpf gpe{"i tco o ct"eqpegv"vcgu"qp"vj g"typed dependency dgw ggp"y qtf u."tcvj gt"
vj cp"eqpukwgpv'kphqto cvkqp0'k'ku"j ki j n{"cf xcpvci gqwu"vq"qwt"uwwf {"ht"kv'ku"tkpi wkuwecm{/tlej "/"
ecr wt lpi "pqv"qpn{"u{pvcvke" kphqto cvkqp"uwej "cu"nsubj *pqo kpcn' uwdlgev" dw' cnq" cdutcev"
ugo cpvke"qpgu"uwej "cu"loc *ñecrk gt+"/"cpf "ecp"dg"htvj gt"cr r ñgf "vq"qvj gt"u{pvcvke/ugo cpvke"
kpgthceg'cumu"Ej cpi ."Vugpi ."Lwtchun{.('O cpplkpi .'422; +0'

Vj g"Ucphqtf "ñzlecnk gf "r tqdcdkñvke"r ctugt" *Ngx{" (" O cpplkpi ."4225+"y qtnu"qw"vj g"
i tco o cvlecn' utwewtg" qh' ugpvpegu" y kj " c" hcvqtf "r tqf wev" o qf gr' ghñekgpvñ" eqo dlpi "
r tghgtpegu"qh'REHI "r j tcug" utwewtg"cpf "ñzlecn'f gr gpf gpe{"gzt gt u0' k' cf f kvkqp" vq"r j tcug"
utwewtg""tgg."vj g"r ctugt""cnq""r tqxkf gu"Ucphqtf ""F gr gpf gpeku""UF +7""vj cv"ctg""npqy p""cu"
i tco o cvlecn'tgrvqpu"dgw ggp"y qtf u"kp"c"ugpvpeg0'Vcnq"vj g"hqmqy lpi "Ej kpgug"ugpvpeg"ht"
gzco r r<我很喜歡兩則惜福與惜緣的故事。 Vj g"head 喜歡 j cu"c"dependent qh'
我 cu'ku'pqo kpcn' uwdlgev"cpf "cpqvj gt "f gr gpf gpv'qh"故事 cu'f k gev'qdlgev"Hi 04+0'

"
"
"
"

⁶""j wr <ly qtf ungvj 0kpi 0kplecQf w0y "

⁷""j wr <lpn 0ncphqtf Qf wluqhy ctg lncphqtf /f gr gpf gpeku0j vo n'

"
"
"
"
"

| | |
|--|---|
| <pre> (ROOT (IP (NP (PN 我)) (VP (ADVP (AD 很)) (VP (VV 喜歡) (NP (DNP (NP (NP (NR 兩)) (NP (NN 則) (NN 惜福) (NN 與) (NN 惜緣))) (DEG 的)) (NP (NN 故事)))))) (PU 。))) </pre> | <pre> nsubj(喜歡-3, 我-1) advmod(喜歡-3, 很-2) root(ROOT-0, 喜歡-3) nn(惜緣-8, 兩-4) nn(惜緣-8, 則-5) nn(惜緣-8, 惜福-6) nn(惜緣-8, 與-7) assmod(故事-10, 惜緣-8) assm(惜緣-8, 的-9) dobj(喜歡-3, 故事-10) </pre> |
|--|---|

"

Figure 2. Dependencies in a Chinese sentence with PCFG: 我很喜歡兩則惜福與惜緣的故事。

Vj g"UF "j cu'dggp'y kf grŋ "wugf "lp"PNR/tgrcvf "hgrf u'uwej "cu'ugpŋo gp'cpcnŋ uku"*O ggpc"(" Rtcđj cncŋt."4229+."vzwcni' gpvcŋo gpv' *Cpf tqwuqr qwqu"(" O cncncukqku."4232+0' Vj g' Ej kpgug" xgtukqp"qh"UF "Ej cpi "gv'cŋo"422; +!u'cncq"cxckŋdrg"qp"vj g"Ucphqtf "F gr gpf gpeku'r ci g⁸0'Vj g"UF " ecp'gxp'f kŋkpi wŋj '67'ŋ' r gf 'f gr gpf gpeku'co qpi 'Ej kpgug'y qtf u'cu'lj qy p'lp"Vcdrg'30'

"
"
"
"
"
"
"
"
"
"
"
"

⁸""j wr <lpn ūcphqtf ūf wluqhy ctg lŋcphqtf /f gr gpf gpeku'ŋij vo n/Ej kpgug"

Table 1. Chinese dependency relations (Chang et al., 2009)"

| abbreviation | short description | Chinese example | typed dependency | counts | percentage |
|--------------|---|--------------------|------------------|--------|------------|
| nn | noun compound modifier | 服务中心 | nn(中心, 服务) | 13278 | 15.48% |
| punct | punctuation | 海关统计表明, | punct(表明, ,) | 10896 | 12.71% |
| nsubj | nominal subject | 梅花盛开 | nsubj(盛开, 梅花) | 5893 | 6.87% |
| conj | conjunct (links two conjuncts) | 设备和原材料 | conj(原材料, 设备) | 5438 | 6.34% |
| dobj | direct object | 浦东颁布了七十一件文件 | dobj(颁布, 文件) | 5221 | 6.09% |
| advmod | adverbial modifier | 部门先送上文件 | advmod(送上, 先) | 4231 | 4.93% |
| prep | prepositional modifier | 在实践中逐步完善 | prep(完善, 在) | 3138 | 3.66% |
| nummod | number modifier | 七十一件文件 | nummod(件, 七十一) | 2885 | 3.36% |
| amod | adjectival modifier | 跨世纪工程 | amod(工程, 跨世纪) | 2691 | 3.14% |
| pobj | prepositional object | 根据有关规定 | pobj(根据, 规定) | 2417 | 2.82% |
| rcmod | relative clause modifier | 不曾遇到过的情况 | rcmod(情况, 遇到) | 2348 | 2.74% |
| cpm | complementizer | 开发浦东的经济活动 | cpm(开发, 的) | 2013 | 2.35% |
| assm | associative marker | 企业的商品 | assm(企业, 的) | 1969 | 2.30% |
| assmod | associative modifier | 企业的商品 | assmod(商品, 企业) | 1941 | 2.26% |
| cc | coordinating conjunction | 设备和原材料 | cc(原材料, 和) | 1763 | 2.06% |
| clf | classifier modifier | 七十一件文件 | clf(文件, 件) | 1558 | 1.82% |
| ccomp | clausal complement | 银行决定先取得信用评级 | ccomp(决定, 取得) | 1113 | 1.30% |
| det | determiner | 这些经济活动 | det(活动, 这些) | 1113 | 1.30% |
| lobj | localizer object | 近年来 | lobj(来, 近年) | 1010 | 1.18% |
| range | dative object that is a quantifier phrase | 成交药品一亿多元 | range(成交, 元) | 891 | 1.04% |
| asp | aspect marker | 发挥了作用 | asp(发挥, 了) | 857 | 1.00% |
| tmod | temporal modifier | 以前不曾遇到过 | tmod(遇到, 以前) | 679 | 0.79% |
| plmod | localizer modifier of a preposition | 在这片热土上 | plmod(在, 上) | 630 | 0.73% |
| attr | attributive | 贸易额为二百亿美元 | attr(为, 美元) | 534 | 0.62% |
| mmod | modal verb modifier | 利益能得到保障 | mmod(得到, 能) | 497 | 0.58% |
| loc | localizer | 占九成以上 | loc(占, 以上) | 428 | 0.50% |
| top | topic | 建筑是主要活动 | top(是, 建筑) | 380 | 0.44% |
| pccomp | clausal complement of a preposition | 据有关部门介绍 | pccomp(据, 介绍) | 374 | 0.44% |
| etc | etc modifier | 科技、文教等领域 | etc(文教, 等) | 295 | 0.34% |
| lccomp | clausal complement of a localizer | 中国对外开放中兴起的明星 | lccomp(中, 开放) | 207 | 0.24% |
| ordmod | ordinal number modifier | 第七个机构 | ordmod(个, 第七) | 199 | 0.23% |
| xsubj | controlling subject | 银行决定先取得信用评级 | xsubj(取得, 银行) | 192 | 0.22% |
| neg | negative modifier | 以前不曾遇到过 | neg(遇到, 不) | 186 | 0.22% |
| rcomp | resultative complement | 研究成功 | rcomp(研究, 成功) | 176 | 0.21% |
| comod | coordinated verb compound modifier | 颁布实行 | comod(颁布, 实行) | 150 | 0.17% |
| vmod | verb modifier | 其在支持外商企业方面的作用 | vmod(方面, 支持) | 133 | 0.16% |
| prtmod | particles such as 所, 以, 来, 而 | 在产业化所取得的成就 | prtmod(取得, 所) | 124 | 0.14% |
| ba | "ba" construction | 把注意力转向市场 | ba(转向, 把) | 95 | 0.11% |
| dvpm | manner DE(地) modifier | 有效地防止流失 | dvpm(有效, 地) | 73 | 0.09% |
| dvpmod | a "XP+DEV(地)" phrase that modifies VP | 有效地防止流失 | dvpmod(防止, 有效) | 69 | 0.08% |
| prnmod | parenthetical modifier | 八五期间 (1990 - 1995) | prnmod(期间, 1995) | 67 | 0.08% |
| cop | copular | 原是自给自足的经济 | cop(自给自足, 是) | 59 | 0.07% |
| pass | passive marker | 被认定为高技术产业 | pass(认定, 被) | 53 | 0.06% |
| nsubjpass | nominal passive subject | 镍被称为现代工业的维生素 | nsubjpass(称作, 镍) | 14 | 0.02% |

Qp'yj g'qj gt'j cpf.'o quv'ugo cp'le'tguwtegu'rkng'Rtqr Dcpm'Rcm gt.'I kf gc.'("Mpi udwt {. " 4227+'cpf'Hco gP gv'Dcngt.'Hkm qtg.'("Nqy g.'3; ; : +'gkj gt'r tqxf g'eqctug/i tckpgf'lpqto cvkqp" qt'y kj "rko kqf"eqxgtci g'k'j' y ku'r cr gt.'y g'r tqr qug" c"ngzlecn'tguwteg"vqni'vq'f guetkdg'o qtg" f gwckrgf"lpqto cvkqp"ht'cmi'y qtf u'k'p" c"vgz v'eqtr wu'Y g"ej qqg"Ukplec"Eqtr wu"*Ej gp."J wpi ." Ej cpi.'("J uw'3; ; 8+'cu'qwt'vgz w'cpf'gxcn'v'j g'tguw'u'y kj'Ej kpgug'Ungvej "Gpi kpg'k'p'vgt o u' qh'eqttgur qpf kpi 'j gucv w'u'hwpevkp0

3. Method

Ɔp"y ku'ecug'uwf {."wpci i gf "gzvu"qh"789.924"ugpvgpegu"htqo "Ukplec"Eqtr wu"509"y gtg"r ctugf " y kj "f gr gpf gpe { "tgrvƆqpu"d { "y g"Ucphqtf "Retugt" *Ej cpi "gv'cn0"422; +0'Y g"qdvclpƆf "796.774" f gr gpf gpe { "tgrvƆqpu"qh"45"v { r gu"dgvy ggp'66.479'y qtf u0'

Vq"ungvej "c"y qtf ."y g"o cnƆ"wg"qh'y j g" f gr gpf gpe { "wr ngu"htqo "y j g"r ctugf "eqtr wu"uƆg"y j g" tki j v'r cpgi'qh"Hi 0'4+"vq"gzvcev'y j g"tgrvƆqpu"qh"gej "y qtf "y kj "ku'f gr gpf gpvu."cpf "qdvclp"y j g" unvej "qh'y qtf u'wej "cu"打 0j k0'uj qy p'dgruy <

Table 2. Dependency sketch of 打 “hit”

(Matches with Chinese Sketch Engine are marked in red bold face)

| prep | dobj | advmod/mmod | nsubj | asp | conj |
|------|------|-------------|-------|-----|------|
| 在 | 電話 | 去 | 武松 | 了 | 重建 |
| 到 | 折 | 要 | 棍子 | 著 | 是 |
| 自 | 籃球 | 就 | 球 | | 鬧 |
| | 高爾夫球 | 先 | 我 | | |
| | 硬仗 | 不會 | 你 | | |
| | 招呼 | 該 | 他 | | |
| | 折扣 | 一起 | 爸爸 | | |
| | 哈欠 | 會 | 雨 | | |
| | 太極拳 | 連續 | 人 | | |
| | 麻藥針 | 一 | 老師 | | |
| | 盹兒 | 能 | 他們 | | |
| | 虎 | 可以 | 她 | | |
| | 羽毛球 | 還要 | 學生 | | |
| | 排球 | 都 | 自己 | | |
| | 蛇 | 雖然 | 湖人 | | |
| | 起來 | 仍然 | 來 | | |
| | 秋千 | 而 | 政 | | |

9"y y y UkplecƆf w0y UkplecEqtr wu"

"
"
"
"
"

Upep" y g" Ucphtf "Rctugt" ukni' uwhgtu" htqo " r ctukpi " f khlwnt "kp" Ej kpgug" *Ngx { " (" O cpkpi ."
4225+ " y g" i tco o cvecln' tgn' vkpu" cwqo cvecln' " tgs wktgf . " y qwi j " ko r tguukxg . " o c { " eqpvckp "
j gvtqi gpgwu" gttqtu" qtki kpcvpi " htqo " o kuci i kpi " gttqtu" : " . " u { pveve " co dki wklgu" cpf " qvj gt "
f gr gpf gpe { " r ctukpi " kuwgu . " uq " y g " j cxg " qdugt xgf " uqo g " o kpat " ungewj " gttqtu " kp " yj g " tguwn0 "
J qy gxgt . " kwl' j ctf " vq " gxcnvcg " yj g " tguwn' kp " cp " cwqo cve " y c { " cu' eqpxgpvkpcnk' gf " kp " yj g " hgrf " qh' "
P NR0Vj g' o clp' tgcuppu' ctg <

"
"

13_0Ewtgpvt . " yj gtg " ku' pq " i qnf / uvcpf ctf " *kp " Ej kpgug " 0K' ku' r ctvevntn' " j ctf " vq " o gcuwg' tgecm' hqt "
yj g' ugv' qh' = eqttgev' cpuy gt0ku' pqv' cxckndrg0 "

14_0Cp " qxgtcm' gxcnvcvqp " qh' yj g' ungewj " r gthqto cpeg " y km' j cxg " vq " tgn' " qp " yj g " cuuguuo gpv' qh' gcej " "
o qf wrg " *y qtf " ugi o gpvcvqp . " RQU " vci i kpi . " ungewj " i tco o ct " cpf lqt " f gr gpf gpe { " r ctukpi . " gve0 "
ugr ctcvgn' 0C " eqo r ctcvkg " vcdrg' ku' uj qy p' kp " Vcdrg' 50 "

"
"

Table 3. Comparison of Different Word Sketch Systems"

| Y qtf " Ungvej " | y qtf " ugi o gpvcvqp " | r qu' vci i kpi kci ugv' " | ungvej " i tco o ct " | f gr gpf gpe { " r ctugt " |
|------------------|--|---|------------------------------|----------------------------|
| EY UG0kplec " | EMR " | EMRICUDE " | j cpf / etchgf " twgu " | , " |
| j VgpVgp03 " | Ucphtf 'Ej kpgug' " Y qtf 'Ugi o gpvt " | Ucphtf 'Nqi / rkpget " Rctv' Qh' Ur ggej 'Vci i gt' " " Ej kpgug' Rgpp' Vtggdcpni' " uvcpf ctf " | j cpf / etchgf " *4' twgu+ " | , " |
| Rtqr qugf " | Ucphtf 'Ej kpgug' " Y qtf 'Ugi o gpvt " | , " | , " | Ucphtf " f gr gpf gpeku " |

: "kp" yj ku' uwf { . " ukpep " Ucphtf " Rctugt " vngu " o cpwn' " vngpk' gf " kpr w' htqo " Uplec " Eqtr wu . " yj g " ugi o gpvcvqp " gttqt " o c { " dg' ngu " yj cp " yj cv' htqo " cp " cwqo cve " ugi o gpvt " cpf " dg' qo kwgf " j gtg0 "

"
"
"

Ķ" cf f klqp." Itqo " vj g" r gtur gevkg" qh' rpi vci g" tguqwegu" eqputwekqp" cu" y gm' cu" cr r rkgf" rnzleqi tcr j {." cu"vj g" u{ uvgo " clo u"vq" kf gpvkh{ " j ki j n{ " ucrlgpv'ecpf kf cvg" r cwtgtpu." vj g"pqku{ " f cvc" uj qwf "pqv'eqpukswg" c" ugtkqu" r tqdrgo " hqt" vj g"vcun0'Vj g" r qukkqp" ku"cuq" y gm'ctvkwrcvfg" cpf" r tqr qugf " Ķ" *Mri cttkh" Mqx^a ." Mrgm" Uf cpqkx ." (" Vldgtku" 4232+." y j gtg" c" xctkcpv' qh' gxcnvcvqp't ctf ki o "wugt lf gxgnr gt/qtkepvfg 'r ctf ki o +ku'tgs wkt gf 0'

F hgtgpv'ltqo "Co dcvk'Tgff {."cpf "Mri cttkh"*4234+'cpf 'Tgff {.'Mtr chku.'O eEctj {.'cpf " O cpcpfj ct" *4233+' y j gtg" cp" gzvgpcn' gxcnvcvqp" vcunl uvej " cu"topic coherence qt" semantic composition y gtg""cf qr vfg." y g" gxcnvcvfg " vj g""r tqr qugf " o gvj qf "y kj " vj g""vcunl' qh' cwqo cvk" eqputwekqp"qh'vj gucwtu."hqt"qw"o clk"eqpegtp"ku"vj g"eqputwekqp"qh'rpi vci g"tguqweg"tcvj gt" vj cp'PNR'u{ uvgo 'r gthqto cpeg0'

Vj g" vj gucwtu" Ķ" Y UG" ku" ecngf "distributional thesaurus." cpf ""ecp" dg" dwkx" hqt" cp{ " rpi vci g"kh'vj g"y qtf "ungvej gu'f cvc"qh'vj g"rpi vci g"ku'cxckrdrg0' Vj g"vj gucwtu'ku'eqputwevfg" d{ " eqo r wkpi "vj g"uko kctkv{ "dgy ggp"y qtf u'dcugf "wr qp"vj g"qxgtncr r Ķpi "tcv"qh'vj gkt"y qtf "ungvej gu0' Qw"o gvj qf "Ķpvvfgf . "hmqy u"vj g" distributional semantic model *F Ķpw"Rj co ."("Detqpk"4235= Vwtpg{ "("Rcpvgn"4232+' cpf" cpej qtu"qp"y q"o cpwcm{ "eqputwevfg" tguqwegu" qh' vj g" Chinese Wordnet: "" cpf 'Ej kĶp" *Ej cq{ 'Ej wpi .'4235+³²⁰

4. Evaluation"

Vj g" f gr gpf gpe{ " f cvc"qh' Hxg"ugrvevfg" u{ pqp{ o "ugv" *經常." 原因." 按照." 相當." cpf "快樂+"ltqo " Ej Ķpgug"Y qtf pgv'y gtg"eqpxgtvfg "Ķpv"o wnk'f Ķo gpukqpcn' *v"cxqkf "ur ctugpguu."qpn{ " f gr gpf gpv" uj ctgf "d{ "dqj" u{ pqp{ o u'y gtg"Ķpvnf gf +Ķ"qtf gt"vq"ecrwcvg" f kvtkdwkqpcn'uko kctkv{ "dgy ggp" u{ pqp{ o u0' Hxg" u{ pqp{ o " ugvu" ltqo " Ej Ķpgug" Y qtf pgv' y gtg" gzco Ķpgf 0' Hqt" gzco r rĶ." vj g" f gr gpf gpe{ " f cvc"qh' 高興 cpf "快樂 ctg"eqpxgtvfg "cu" hmqy u" *Ķ kt gi ctf Ķpi " vj g" f gr gpf gpe{ " v{r g+<

| | | | | | |
|----|----|----|----|----|----|
| | 不 | 也 | 了 | 他 | 可以 |
| 高興 | 9" | 3" | 7" | 7" | 3" |
| 快樂 | 3" | 6" | 3" | 5" | 4" |

"
"

¹ j wr <ltqr g0Ķpi wkvleupw0gf v0y ley p4"

^{32j} wr <leqf g0 qqi rĶ0qo lr ly /u{ pqp{ o u/ej kĶp"

"
"
"

Vj gp"y g"cf qr v"ppg"qh' yj g"eqo o qp"o gcuwtgu" hqt" uko kctk\{ "lp" f kurtkdwkqpcn' o qf gnu." cosine similarity." vq"ecrewrvcg" yj g" uko kctk\{ "dgy ggp" y q" y qtf u" *g0 0" 高興 cpf "快樂" Vj g" o gcplpi " qh'c" y qtf "ku" f gyto kpgf "d\{ "ku" eqmjecvqp." cpf "tgr tgugpvf "cu" c" xgevqt "qh'ku" eq/qeewttgpeg" y kj " qy gt" y qtf u" lp" o wnr rg" f ko gpukqu0'kp" yj ku" o qf gn" yj g" uko kctk\{ "dgy ggp" y q" y qtf "xgevqtu." w1" cpf "w2." ecp" dg" ecrewrvcg" f "d\{ "yj gkt" equlpg" xcwv<

"
$$CosSimilarity(w_1, w_2) = \frac{w_1 \cdot w_2}{|w_1| |w_2|}$$
 *3+

Vq" kmwutcvg." eqpukf gt" qpn\{ "yj g" hktuv" y q" f ko gpukqu" qh' 高興 cpf "快樂." yj g" equlpg" uko kctk\{ " dgy ggp" yj g" y q" y qtf u" y qwf "dg" *9.3+ *3.6+ 1\sqrt{(9^4- 3^4)+1\sqrt{(3^4- 6^4)+? 2099." cpf "yj g" ecrewrvcg" ecp" dg" gz vpf gf "vq" gxgp" o qt g" f ko gpukqu0'ki" y q" y qtf u" j cxg" uko kct" eqmjecvqp" y kj "qy gt" y qtf u." yj g" xcwv" qh' equlpg" uko kctk\{ "y kn' cr r tqcej " yj g" wr r gt" dqwpf "qh' 30" cpf "eqwf "dg" eqpukf gt gf "c" r ckt" qh' u\{ pqp\{ o u0'

Hpcm\{ "vq" qdvclp" c" u\{ pqp\{ o " rkuv" yj g" f gr gpf gpw" qh' c" vti gv" y qtf "ctg" tcpngf "d\{ "yj gkt" uko kctk\{ "y kj " yj g" vti gv" y qtf ." tgi ctf ngu" qh' yj gkt" f gr gpf gpe\{ " tgrvclpu0' Vj g" tguwmu" hqt" yj g" ugrgevgf "hxg" u\{ pqp\{ o " ugw' lp" Ej kpgug" Y qtf pgv' cpf "Ej kkp" ³³ " ctg' uq y p' lp" Vcdrng' 60'

"

Table 4. Comparison of the results with Sketch Engine"

| | 快樂 | 經常 | 原因 | 按照 | 相當 |
|--------------------------|-------------------|-------------------|------------------------|---------------------|---------------------|
| Eklp" | 高興,愉快,樂,... | 時常.常常." 時時.常0' | 因.故.緣故.緣 由.í " | 依照.比照.遵照.í " | 頂.相當於.í " |
| Ej kpgug" Y qtf pgv' | 樂.愉快.愉 | 時常.勤.常 常.常.恆 | 關係.肇始.因."故. 導因.緣 | 按.依照.依據." 根據.í " | 很.相當於.具體 |
| Rtqr qugf " O gj qf " | 有趣.愉快." 美好.í " | 時常 | 關係 | 按.依照 | 具體.í " |
| Ungvej " Gpi kpg" | 愉快.美 好.í " | p k" | 因素.背景.條件." 環境.理由.0' | p k" | 莫大.重大.重 要.直接.í " |

³³Cnj qwi j "Y qtf P gv'ku" o qtg" wugf "lp" pcwtcr'ncpi wci g' r' tqeguulpi . "Ej kkp" ku" eqpukf gt gf "c" o qtg" cr r tqr tkvg' tguqwtg' f guki pgf "hqt" yj gucwtwu0J gtg" y g' r' tgugpv' yj g" eqo r ctkuqp" y kj "dqj 0'

''
''
''

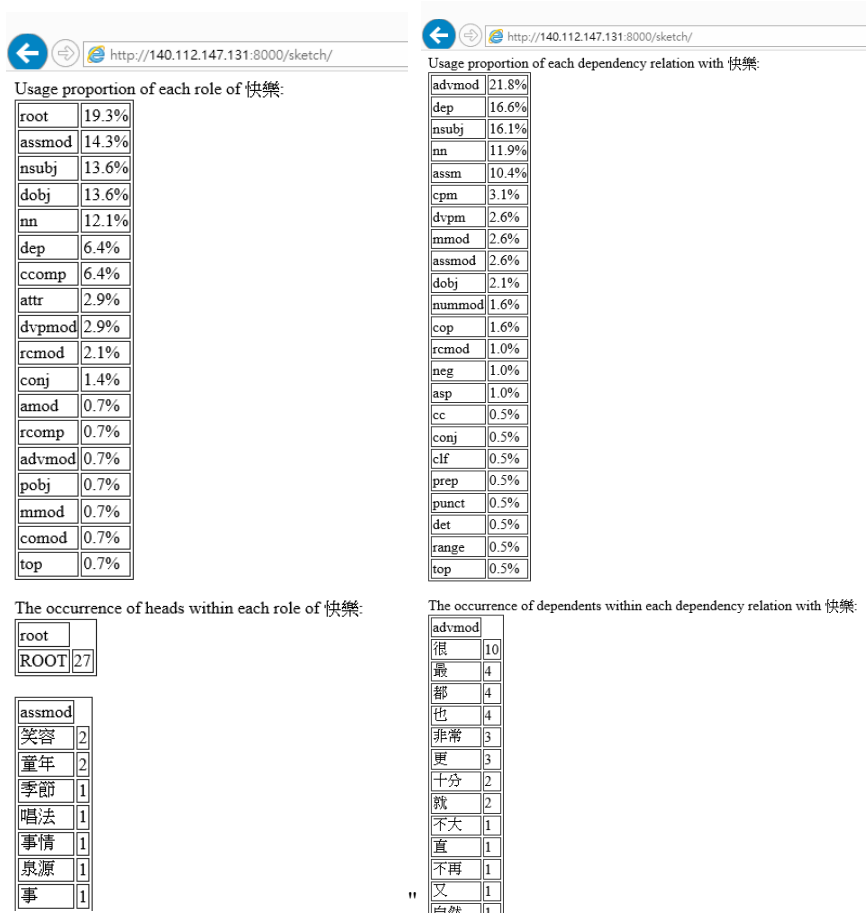


Figure 3. Snapshot of the sketch function''

5. Conclusion''

Y qtf "ungvej "ku" c" eqtr wu/ dcugf "cwqo c\le" uwo o ct { "qh" c" y qtf \u" i tco o c\lecn/ cpf " eqmjec\kqpcn' dgi cxlqt0' Dcugf " qp" \j g" j cpf / etchnf " h\pk\g/ uc\vg" ungvej " i tco o ct " qxgt" c" RQU/ wci i gf " eqtr wu. " y qtf " ungvej " u{ ugo " ecp" kf gp\wh{ "\j g" eqmjec\vu" y kj " i tco o c\lecn' tgr\kqpu" vq" \j g" vti gv" y qtf 0' J qy gxgt. " \j g" i tco o ct " gpi lpggtkpi " ku" \ko g/ eqpuwo kpi " cpf " tgs \k\gu" g\zr gtvu. " k\p" \j ku" r cr gt. " y g" r tqr qug" cp" cngt\pc\kxg" d{ " r\gxgtci kpi " cp" g\zk\kpi " f gr gp\ gpe { " r ctugt0' Vj g" t\guwu" y gtg" gxc\m\c\vgf " dcugf " qp" \j g" eqo r ct\k\qp" qh' f k\ut\kdw\kqpcn' \j guc\w\w\y kj ' \ki p\k\ec\peg0'

''

Vj ku' r cr gt " ugtxgu" cu" \j g" h\k\uv" cwgo r v' vq" etgc\vg" cp" qr gp/ uq\w\egf " y qtf " ungvej / r\kng" eqtr wu' r tqh\k\kpi " u{ ugo " h\q" " E\j lpgug" " kpi w\k\k\eu" cpf " " Vgcej kpi " E\j lpgug" cu" " Ugeqpf " " N\cpi wci g0' Vj g" ''

"
"
"

r tqr qugf "o gj qf "ku"r kr grkpgf "cpf "ecp"dg"cr r rkgf "v" wugt/etgcvgf "eqtr wu"Vj g"gzvcevgf "tgrvkvq"
vkr ngu">y 3."T."y 4@ecp"dg" wugf "v" gptlej "qwt" qp/i qkpi "Ej kpgug" DK NGZ "f cvdcug0 Hwwtg"
y qtnu"lpenwf g"gzr nrtkpi "qvj gt "f gr gpf gpe{ "r ctulpi "cni qtkj o ."lpeqtr qtcvki "cf xcpegf "ucvkvku"
vq"ulpi ng"qaw"ucrkpgv"eqnqecvkvpu "cpf "cp"qr gp"gxcnkvkvq"r rlvhto "hqt"vj g"lmtvj gt "ko r tqxgo gpv"
qh'vj g'tguvteg'ctg'lp'r tqi tguo'

"
"

Acknowledgements"

Vj g"cwj qtu"y qwf "rkg"v"vj cpn'tgxky gtu"qh" TQENKI "4236" hqt"vj gk"luki j vhwieqo o gpu"
vj cvj gr "ko r tqxg"vj g'o cpwuetkr v0'

"
"

References"

Co dck"D0T0" Tgf f {."U0"("Mki cttkh"C0*4234-0)Word Sketches for Turkish. Rcr gt"r tguvpgf "cv"
vj g"NTGE0'
Cpf tqwuqr qwqu."K0"("O cncnkvkv."R0*4232-0C"uvxg{ "qh'r cter j tculpi "cpf "vgz wcn'gpvckn gpv"
o gj qf u0J. Artif. Int. Res., 38*3+:357/3: 90'
Dcngt."E0HD" Hkm qtg."E0L0"("Nqy g."L0D0'*3; ; : -0)The Berkeley FrameNet Project. Rcr gt"
r tguvpgf "cv" vj g"Rtqeggf kpi u" qh' vj g" 58vj " Cppwcn' O ggkpi " qh' vj g" Cuuqekvkvq" hqt"
Eqo r wcvkpcn' Nkpi wkvku" cpf " 39vj " kvgtpcvkvpcn' Eqphgtgpeg" qp" Eqo r wcvkpcn'
Nkpi wkvku"/Xqno g'3.'O qpv gcn'S wgdge.'Ecpcf c0'
Ej cpi ."R0E0"Vugpi ."J 0"lvtchun{.'F0"("O cplkpi ."E0F 0*422; -0)Discriminative reordering with
Chinese grammatical relations features0Rcr gt"r tguvpgf "cv"vj g"Rtqeggf kpi u"qh'vj g"Vj kf "
Y qtnuj qr "qp"U{pvz'cpf 'Utvewtg'lp'Ucvkvkcn'Vtcpuvcvkvq.'Dqwf gt.'Eqnqtcf q0'
Ej cq."H0"C0"("Ej vpi .""U0H0"*4235-0"C" F ghpkkvq/dcugf "Uj ctgf/eqpegr v' Gzvtcekvq"y kj kp"
I tqw u" qh' Ej kpgug" U{pqp{ o u<C"Uwf {" Wkkl kpi " vj g""Gzvgpf gf "Ej kpgug" U{pqp{ o "
Hqtguo0IJLCLP, 18*4-0'
Ej gp."M0L0" J wpi ." E0T0" Ej cpi ." N0R0" (" J uw" J 0N0' *3; ; 8-0) Sinica corpus: Design
methodology for balanced corpora. Rcr gt"r tguvpgf "cv""vj g""Vj g""33vj " Rcekhe""Cukc"
Eqphgtgpeg"qp'Ncpi wei g.'kvhhto cvkvq'cpf 'Eqo r wcvkvq""RCENIE/33-0'
F kpw" I 0'Rj co ."P0"("Detqpk" O 0' *4235-0) DISSECT-DIStributional SEMantics Composition
Toolkit. Rcr gt" r tguvpgf "cv" vj g" 73uv' Cppwcn' O ggkpi " qh' vj g" Cuuqekvkvq" hqt"
Eqo r wcvkpcn'Nkpi wkvku0'

"
"
"

J qpi .L0HD("J wpi ."E0T0*4228+0Using chinese gigaword corpus and chinese word sketch in linguistic research. Rcr gt" r t g u g p v g f " c v ' j g " V j g " 42 j " R c e k h e " C u l c " E q p h g t g p e g " q p " N c p i w c i g . k p h t o c v k p p ' c p f ' E q o r w c v k p p " * R C E N E / 42 + 0

J wpi ."E0T0" M k i c t t k h " C 0 " Y w [' 0 ' E j k w " E 0 O 0 " U b k j . " U 0 ' T { e j n . " R 0 ' 0 0 0 ' E j g p . " M 0 L 0 * 4227 + 0 Chinese Sketch Engine and the extraction of grammatical collocations. Rcr gt " r t g u g p v g f " c v ' j g " R t q e g g f l p i u ' q h ' j g " H q w t j " U K J C P " Y q t m u j q r " q p ' E j k p g u g " N c p i w c i g " R t q e g u l p i 0 " M k i c t t k h " C 0 * 4229 + 0 W u l p i " e q t r q t c " l p " r p i w c i g " r g c t p l p i < j g " U n g v e j " G p i k p g 0 " O p t i m i z i n g t h e r o l e o f l a n g u a g e i n T e c h n o l o g y - E n h a n c e d L e a r n i n g , 22 . " 430

M k i c t t k h " C 0 " M q x a . " X 0 " M t g m " U 0 " U f c p q x k . " R 0 " (" V l d g t k w u . " E 0 * 4232 + 0 A quantitative evaluation of word sketches. Rcr gt " r t g u g p v g f " c v ' j g " R t q e g g f l p i u " q h ' j g " 36 j " G W T C N G Z " k p v g t p c v k p c n E q p i t g u u . " N g g w y c t f g p . " V j g " P g y g t r c p f u 0

M k i c t t k h " C 0 " T { e j n . " R 0 " U b t | . " R 0 " (" " V w i y g m " F 0 * 4226 + 0 K r k 26 / 2 : " j g " u n g v e j " g p i k p g 0 Information Technology, 105 . " 3380

N g x { . " T 0 " (" O c p p l p i . " E 0 * 4225 + 0 Is it harder to parse Chinese, or the Chinese Treebank? Rcr gt " r t g u g p v g f " c v ' j g " R t q e g g f l p i u " q h ' j g " 63 u v " C p p w c n " O g g v l p i " q p " C u u q e k c v k p p " h q t " E q o r w c v k p c n N p i w k u l e u / " X q n w o g 3 . " U c r r q t q . " L c r c p 0

O g g p c . " C 0 " (" R t c d j c n c t . " V 0 X 0 * 4229 + 0 Sentence level sentiment analysis in the presence of conjuncts using linguistic analysis0 Rcr gt " r t g u g p v g f " c v ' j g " R t q e g g f l p i u " q h ' j g " 4 ; j " G w t q r g c p ' e q p h g t g p e g " q p " K T " t g u g c t e j . " T q o g . " K c n 0

R c w g t . " O 0 " I k f g c . " F 0 " (" M p i u d w t { . " R 0 * 4227 + 0 V j g " R t q r q u k k p " D c p n x " C p " C p p q v c g f " E q t r w u " q h ' U g o c p v e " T q r g u 0 Computational Linguistics, 31 * 3 + . " 93 / 3280 f q k 320838402 ; 3423275852486 "

T g f f { . " U 0 " M r r c h k u . " R 0 " O e E c t j { . " F 0 " (" O c p c p f j c t . " U 0 * 4233 + 0 Dynamic and Static Prototype Vectors for Semantic Composition. Rcr gt " r t g u g p v g f " c v ' j g " K E P N R 0

V w t p g { . " R 0 F 0 " (" R c p v n " R 0 * 4232 + 0 H t q o " H t g s w g p e { " v j " o g c p l p i < x g e v q t " u r c e g " o q f g m " q h ' u g o c p v e u 0 J . A r t i f . I n t . R e s . , 37 * 3 + " 363 / 3 : : 0

Vj g"4236"Eqphgtgpeg"qp"Eqo r wcvkqpcr/Nkpi wku'eu"cpf "Ur ggej "Rtqeguulpi ""
TQENRPI "4236."rr0375/384" ""
©"Vj g"Cuuqekvqp"lqt"Eqo r wcvkqpcr/Nkpi wku'eu"cpf "Ej kpgug"Ncpi wci g"Rtqeguulpi "

使用中文字筆畫構形資料庫校正字形相似之別字

Using Chinese Orthography Database to Correct Chinese Misspelling Words With Graphemic Similarity

張道行¹ 陳學志² 鄭健良¹

摘要

中文別字自動偵測與校正是個相當重要的工具，許多分析別字類型的研究指出，「字音混淆」、「字形混淆」與「字義混淆」是別字產生的主要原因，因此近年來許多別字自動校正的研究也都採取分別針對字音、字形、字義造成的混淆進行探討。但對於字形相似別字的校正正確率仍不夠好，主要原因之一是因為在中文字字形結構的資訊不夠完整。本文的目的是利用中文部件組字與形構資料庫的筆畫結構資料提出一個演算法，計算兩個中文字筆畫結構序列的相似程度，並用於字形相似類別字的偵測與校正。實驗結果顯示筆畫結構用在偵測與校正字形相似別字的效能較原先以部件的方法來得有效。"

關鍵字：別字偵測；別字校正；字形相似；筆畫；中文部件組字與形構資料庫；LCS演算法

1. 緒論

中文別字自動偵測與校正是個相當重要的工具，在資訊應用中也具有相當的價值。例如許多文書處理軟體都有提供別字檢查與校正建議。但由於中文別字的偵測與校正較困難，這些軟體在中文別字校正效能上仍有不足。因此除了英文的別字校正研究[3]_9_之外，一直以來也有許多研究[6]_8_]：[33]_35_]37_]38_]提出不同的方法試圖解決中文別字偵測與校正的問題。"

早期別字自動偵測方法大多是透過建立混淆字集方式搭配語言模型運作。例如Ej cpi j4_將欲處理語句中的每個字視為別字，並以該字混淆字集中的字逐一替換原字，組合出多種可能的句子。再利用 dk i tco_ 語言模型計算所有句子的分數，選出最可能的句子。

³" F gr ctvo gpv'qh'Eqo r wgt "Uelgpeg"cpf "lphqto cvkq"Gpi kpggtkpi . "P cvkqpcr/Mcqj ukxpi "Wpkxgtuks{ "qh' Crr r'ngf "Uelgpegu."Mcqj ukxpi . "Vcly cp"

⁴" F gr ctvo gpv'qh'Gf wecvkqpcr/Ru{ej qm i { "cpf "Eqwpu'gkpi . "P cvkqpcr/Vcly cp "P qto cr/Wpkxgtuks{ . "Vclr gk" Vcly cp"

Go ckn'ej cpi vj B i o (mucu'gf w0y "

許多研究也提出了類似的方法，但改用不同的語言模型，如 *tki tco* [39]、混合規則式與機率式的模型 [34]。此類方法相當容易實作，也可同時適用於字音相似與字形相似兩類別字的偵測與校正。但這種方法的主要缺點為適當的混淆字集建立不易，如果廣泛蒐集所有可能的混淆字，很容易造成誤報率提高。但若只蒐集常見別字，則會遺漏不在混淆字集中的字。"

許多分析別字類型的研究 [31; 32] 指出，「字音混淆」、「字形混淆」與「字義混淆」是別字產生主要原因，因此近年來許多別字自動校正的研究也都採取分別針對字音、字形與字義造成的混淆進行分析、再加以整合的策略。這些研究在字音造成的別字部分都有相當好的校正正確率，但對於字形與字義造成的別字的校正正確率仍不夠好。對於字形造成的別字校正效能有限，主要原因之一是因為在中文文字字形結構的資訊不夠完整，導致建立某些中文字字形混淆字集時容易產生相似度誤差使得字集不夠精確。"

陳學志等人 [44] 建立的「中文部件組字與形構資料庫」是一個解決這個問題的可能線索。該資料庫於 4232 年起提供每個中文字的部件組字資料，又於 4235 年進一步擴充，將每個部件進一步拆解為筆畫結構的組合，得到每個字以筆畫為基本單位的構形資料，我們稱為筆畫結構序列。本文的目的就是利用這個中文字筆畫結構資料庫處理字形相似度問題，並用於字形相似類別字的偵測與校正。但是筆畫結構資料是描繪基本筆畫單元間的空間關係，要如何使用於計算兩字間的字形相似度並不是個容易的問題。本文的主要貢獻在於提供一個演算法可計算兩個中文字筆畫結構序列的相似程度。"

本文其餘內容將組織如下。第 4 節說明字形相似別字的相關研究以及所遭遇的問題。第 5 節介紹「中文部件組字與形構資料庫」中筆畫結構資訊的表達方式。第 6 節說明本文所提的字形相似度計算方法。第 7 節說明如何將字形相似度計算結果整合在別字預測工具中。第 8 節展示應用此方法於字形錯別字的實驗結果。最後提出未來可能的研究方向。"

0

2. 相關研究

早期研究多使用混淆字集偵測與校正別字，然而這種方法未必能將真正字形相似字完整涵括。因此有些研究針對字形相似性設計自動建立混淆字集的方法。最早的相關研究是使用與字形相關的倉頡輸入法偵測字形相似別字。如 *Nlp* 等人 [1] 利用倉頡碼建立混淆字集，再採用類似 *Ej cpi* [4] 的方法偵測與校正。實驗結果顯示校正的成功率有明顯提升且有較少的誤報率。*Nkw* 等人 [32] 也是利用倉頡碼對形構相近字提出計算相似度的方法。 [32] 利用倉頡輸入法將中文字編碼，透過倉頡碼 48 個基本單位的組合比較兩中文字的相似性。然而

倉頡輸入法為了維持每一個漢字輸入不超過五個碼的效率，因此某些部件較多或較複雜的漢字倉頡碼會被簡化，使得字與字之間的相似證據難以找出。於是]32_建構了一套「倉頡詳碼」完整地還原字的構形，將每個字以倉頡詳碼表示，可提供較精確的相似度訊息。"

基本上使用倉頡碼建立混淆字集的概念是認為漢字可拆解成數個基本單元、或稱為部件**tcf lecn+*，比對部件的使用可判斷字形相似性。而有一些研究]44_]45_提出了更完整的中文字部件構形資料庫，例如]45_提出的「漢字構型資料庫」，以及"]44_建構的「中文部件組字與形構資料庫」。以後者為例，其蒐集了 82; 9 個常用字，從內拆解出 65; 個基礎中文部件，並歸納出 33 種字形空間的結構關係。透過以上指標與部件跟結構關係的組合，可以更加瞭解字形的構成，進而觀察字與字之間的相似關係。與倉頡碼相比，使用中文字部件資料庫比較字形相似度是更為合理且精確的作法。"

先前 E]j cpi 等人]7_的研究也利用「中文部件組字資料庫」中的部件偵測與校正字形別字。該研究也發現使用部件偵測及校正字形相似字的確有效，但該研究也指出有些部件由於未能進一步拆解導致某些字的相似度計算有相當大的誤差。例如西本身是一個部件，因此當他和西計算相似度時因為沒有相同部件導致相似度非常低。這個問題突顯了以部件量測字形相似性的侷限所在，但也指出改良的方向，也就是應該以更基本的字形組成單元測量相似性。"

"

3. 中文字筆畫結構表示法

在「中文部件組字與形構資料庫」中，所有漢字可歸納出 63 個基礎筆畫與 33 個字形結構關係。筆畫中包含基本的橫}一_i、豎}丨_i、撇}丿_i等。其中部分筆畫無法透過系統 *wpleqf g* 編碼呈現，便會將其透過組合的方式表達，例如「國」的第二劃為橫豎，但橫豎無法以一般 *wpleqf g* 編碼呈現，因此以筆畫}口 4_i*代表「口」字書寫時的第二劃，也就是橫豎+表示。有些更複雜的筆畫必須透過結構關係描述，例如「大」的第三劃為撇捺，同樣無法以 *wpleqf g* 呈現，因此會以筆畫}尺 1_i表示，其中筆畫內置於「尺」右側的「1_i」代表右下方的筆畫，即「尺」字右下方的那一劃，也就是撇捺。"

字形結構部分則包含 32 種組合及其它結構。32 種組合包含：垂直組合*例如「二」字結構是由上下垂直組合而成+、水平組合*例如「林」字是由左右水平方式組合+、封閉組合*例如「國」字是由四面包圍方式組成+等。這些組合以符號表示，例如垂直組合為「/」、水平組合為「~」、封閉組合為「2」等。大多數字都可以用這 32 種組合結構表示，但有些字的結構不屬於任何一種組合，故這些字歸類為無字形組合之結構，例如「一」、「乙」字

都為單獨存在的形體結構，不屬於任何一種組合。”

在資料庫中每一個漢字都是由筆畫、筆畫間連接關係、與字形結構三項所組合，稱為「基礎筆畫組合」。例如「二」字，其基礎筆畫組合為「/ * } - u_i . } - i_j +」，從中可看出垂直組合 * + 是主結構，這個組合包含了兩個筆畫 } - u_i 與 } - i_j，其中 } - u_i 代表書寫時較短的 } - i_j。「二」字無筆畫連接關係，對於有筆畫連接關係的字也會在基礎筆畫組合以符號組合「 - } c < d B e_i 」表示連接關係。此符號組合一般都會緊鄰在一筆畫之後，例如「 } | i - } c < d B e_i 」，表示會以筆畫 } | i 為基準，本文稱為基本筆畫。符號組合中的 c 代表整個基礎筆畫組合中由左而右數來第 c 個筆畫和基準筆畫 } | i 有連接關係，這第 c 個筆畫本文稱為連接筆畫。符號組合中的 d 代表基準筆畫在連接筆畫的 d 位置相接，「 B e_i 」則代表連接筆畫在基準筆畫的 e 位置相接。此資料庫將每一個基礎筆畫依照比例平均劃分為 32 等分，用來表示前述相接位置。符號組合中的「 - 」表示筆畫間的連接關係為相互交錯；若筆畫間為相接但非交錯的關係，則以「 ϕ 」符號表示。筆畫間連接關係只有以上兩種。”

圖 3 是一個基礎筆畫組合的範例。範例中的「十」字其基礎筆畫組合為「 } | i - . } | i - * 3 < B 7 + 」。 「十」的筆畫連接關係表示為 } - i_j 在其 7 的位置 * 也就是正中間 + 與 } | i 在其 7 的位置相互交錯。然而「十」的結構為單獨存在，故在筆畫組合裡面不會看到結構符號、只會有筆畫。但是兩個筆畫有連接關係，因此僅以「 } | i - * 3 < B 7 + 」表示。”

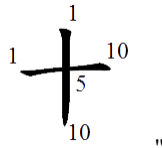


圖 3 基礎筆畫組合範例字「十」

4. 中文字筆畫結構相似度

由於前述的筆畫結構表示法不容易直接用來計算兩中文字筆畫結構的相似度，因此必須先將整體筆畫結構轉換為筆畫結構配對。轉換方法是先將每個漢字的基礎筆畫組合中以有連接關係的一組筆畫配對作為相似度計算的基本單位。例如「大」字，基礎筆畫組合為「 [{ - } , { 月 1 } + (1 : 5 @ 3) , { [尺 /] } ~ (1 : 5 @ 0) ~ (2 : 3 @ 0)] 」，將其透過兩兩筆畫連接關係配對成 ({ - } , { 月 1 }) 、 ({ - } , { [尺 /] }) 及 ({ 月 1 } , { [尺 /] }) 三個筆畫配對，並以漢字書寫時的筆畫順序將配對排序，形成一個配對序列。

接著便可使用最長公共子序列 (Nqpi guv" Eqo o qp" Uwdugs wgegu". "NEU) 演算法 [9]_36_ 計算兩個配對序列的相似度。NEU 演算法可以經由比對得到兩序列中最長相同的子序列，

因此廣泛被使用於計算兩字串或序列之間的相似度。其計算方法是將序列中的所有子序列依照順序一個個與另一序列之子序列做比對，若相同則比較下一個子序列，不同則記錄比較至目前最長的相同子序列，接著再從頭比對另一序列的下一個子序列，重複上述動作直到所有子序列比較完成即可取得最長相似的子序列組合。例如比較兩字串「Nqqm」與「Dqqm」，會先以「Nqqm」的 N 與「Dqqm」的 D 比較，不同則比較「Dqqm」下一個子字元 q，重複以 N 比對完「Dqqm」所有字元後，再以「Nqqm」的第二字元 q 重新比對「Dqqm」。藉由不斷地循環比對便能得到最長的子字串為「qqm」，長度為 5。”

由於 NEU 演算法具備「擷取序列必須順序相同」的特性，符合漢字是依固定筆畫順序書寫的規則，因此我們採用 NEU 演算法作為兩字間筆畫配對序列相似計算的方法。以「大」和「太」兩字之相似度比較為例，其配對序列分別為({一},{月 1})、({一},{[尺/])、({月 1},{[尺/])與({一},{月 1})、({一},{[尺/])、({月 1},{[尺/])、({、})，兩字間除了配對({、})不一樣之外，其餘三個配對均相同，且其配對經由 NEU 演算法計算後，最長連續相似配對的配對數為 5。

上述的配對方法雖然依照順序的特性能有效計算大多數字之間的相似度，但有兩種特殊情況必須考慮。第一種情形是部件相同但部件位置不同的鏡寫字會造成相似度不足的問題。例如「部」與「陪」兩字，其筆畫配對完全相同，但出現的順序不同，若如果利用原本配對序列經由 NEU 演算法計算並無法得到合理的相似度。因為有鏡寫障礙的學生很容易寫出這類的別字，故我們修正前述配對序列產生的方式以改善 NEU 演算法的問題。”

我們將欲比較的兩序列其中之一複製後連結在原先序列的後面，形成一個長度為原先兩倍的新序列，稱為複製倍增序列，再將複製倍增序列與另一序列做比較。以「cdef gh」與「f ghcde」兩序列為例，由於兩序列僅「cde」與「f gh」順序顛倒，若使用原先 NEU 演算法計算，其最長相似長度為 5。因此我們將序列「cdef gh」倍增為「cdef ghcdef gh」後，再與另一序列「f ghcde」比較，可得最長共同序列為「f ghcde」，長度為 8。但由於這種做法會使一個字與自己的最長共同序列的相似度會和它與鏡寫字間的相同，因此當兩字為鏡寫字時，必須在使用最長共同序列長度作為相似度評估時略做調整，才能有較合理的結果。

第二種情形則是配對序列省略了字形結構資訊導致某些字的相似性計算不合理。例如「昌」與「田」兩字，其筆畫配對序列完全相同，若以 NEU 方法比較兩字，兩序列完全相同，但兩字並不是相同字，是屬於主字形結構不同的相異字。因此相似度公式必須考慮比較兩字之間的主結構關係，根據研究經驗給予不同的字形主結構的結構相似係數。在 NEU 計算相似性時須另外再乘上結構相似係數，作為兩字最後的字形相似分數。例如「昌」

與「田」字，其主字形結構分別屬於「垂直組合*+」與「水平組合*+」，這兩個結構關係差異較大，較容易分辨，因此兩結構的相似係數為 206。。

另外，本文使用筆畫評估字形相似性的原始動機之一，是因為有些相似的部件無法被有效判斷相似性。雖然用筆畫能將部件拆解成更小的比較單位，但筆畫與筆畫之間仍有相似的問題，例如筆畫}一_i和筆畫}刁_{4_i}、}L_{4_i}及}┐_i相當相似。為解決這個問題，我們使用J44所提供的筆畫相似度資訊，J44將所有筆畫分類為43個筆畫相似集。因此本文在計算筆畫配對序列相似度時，對於兩個筆畫配對，若對應的筆畫在同一個筆畫相似集中，即使筆畫不同仍視為相同筆畫配對。

綜合以上討論，對於兩個中文字 C_1 與 C_2 ，其筆畫配對序列分別為 S_1 與 S_2 ，則兩字的字形相似度公式如下：

$$sim(C_1, C_2) = sms(C_1, C_2) \times \frac{2 \times lcs(rep(S_1), S_2) - mirror}{num(S_1) + num(S_2)} \quad *3+''$$

其中 rep 函數可產生輸入序列的複製倍增序列； lcs 函數計算兩序列的最長共同子序列的長度； num 函數計算該序列的配對數量； sms 函數計算兩字主字形結構的相似性；參數 $mirror$ 在 lcs 函數值分別等於兩個 pwo 函數值且兩字相異時，其值為 3，否則為 2。

5. 別字校正模型

本文使用的別字訓練模型是以 Ej cpi 等人J7所提的方法為基礎。J7利用訓練資料中別字所產生的候選詞的四個參數*字音相似度、字形相似度、字頻機率比值、詞性機率比值+建立一個線性迴歸公式作為預測模型，用來判別候選詞是否可以對原句的可疑字組進行替換校正。本文與J7的差別就在於字形相似度的計算方法不同。為了建立預測模型，本文準備 872 筆含有別字的資料，在計算有別字的可疑字組與有正確字的候選詞間的四個參數後，成為一個 r quklxg 樣本集合；另外準備 872 筆句子包含無別字句子與含有別字句子，在計算可疑字組與無正確字的候選詞間的四個參數後成為一個 pgi cvkxg 樣本。使用這兩個樣本集合共 3522 句可以建立線性迴歸公式。這個線性迴歸公式如下：

$$y = \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 \quad *4+''$$

其中 X_i 為候選詞與可疑字組間的四個參數， β_i 為四個參數的迴歸係數，將候選詞對應

可疑字組的四個參數及各項迴歸係數代入上述公式可得到 { 值。當候選詞代入公式所得之 { 值越大時，其作為應替換的正確詞彙的可能性越高。相對地，若 { 值越小其作為可替換詞的機會也相對減少，因此可以利用 { 值作為候選詞應替換原字組可能性的依據。

本文從訓練的 3522 句 rqukkxg 及 pgi ckkxg 樣本集合裡候選詞與可疑字組間四個參數代入公式所得之 { 值找出最佳分類正確率，並以此 { 值作為預測模型的門檻值。若候選詞是應替換原字組的可能性分數大於門檻值即會被判定為可替換之正確詞彙而進行校正，反之若一個候選詞是應替換原字組的可能性分數小於門檻值，該候選詞則會被判定為無法替換詞彙而被捨棄。

"

6. 實驗

本文採用 UK J CP /9" Dcng/qth'4235 2Ej lpgug'Ur gnrkpi 'Ej genj37_所提供的資料集作為實驗的測試資料。資料集內包含 UwdVcun3 與 UwdVcun4 各 3222 份文本。其中 UwdVcun3 只提供含有別字的文本與別字位置的資料，因此本文僅採用 UwdVcun4 資料作為量測中文別字偵測與校正效能使用。UwdVcun4 的 3222 份文本是來自於蒐集學生寫作常見的錯誤類型，每一份文本內至少含有一個以上的別字及數句未包含別字的句子，全部共包含 3487 個別字。本文保留含有別字的句子作為測試資料，合計共 343; 筆句子。由於我們是測試字形相似度的效能，因此只採用字形相似以及形音相似實驗的別字。另外，實驗使用的別字偵測工具只能擷取 Y GECp]5 斷詞系統將詞彙中別字分成單字詞者，因此我們只使用符合上述兩項條件的 673 個別字作為測試資料。"

本文別字校正實驗將以第 7 節所提方法設計兩個預測工具，一個是使用 Ej cpi 等人 J7_所提出的部件相似度計算方法，另一個則是使用本文所提筆畫結構相似度計算方法。兩個工具分別針對 673 筆測試資料進行測試。另外由於測試的別字有些為字音與字形皆相似的別字，為了評估字形相似計算之結果差異，在系統提出可替換候選字時會淘汰僅字音相似的正字詞彙，只比較字形混淆與形音同時混淆的別字。結果於表 3 所示。"

"

表 1、使用部件相似與筆畫相似方法之校正結果

| | 校正字數" | T gecni' | Rt gekukqp" | Hcng'Crcto " |
|-------|-------|----------|-------------|--------------|
| 部件相似法 | 250 | 77065' | 8504; ' | 20 5' |
| 筆畫相似法 | 417 | ; 4068' | 9; 07: ' | 2082' |

" 由表 3 的結果可知，在只比較形別字與形音別字的情況下，使用本文所提之筆畫相似方法能有效提升校正成功率。從本文所提方法的三項評估指標均遠優於使用組字部件的方法，可以得知筆畫計算字形相似的方法確實能改善部件相似性不足的問題。 "

"

7. 結論

本文提出了一個使用中文字筆畫結構資料的字形相似評估方法，並將其與現有別字偵測方法整合。從實驗結果可以證實本文所提方法用於字形相似別字的偵測與校正效能較原先以部件相似的方法來的好，其原因可以推論是本文所提方法讓字形相似度計算變得更加精確所導致。但受限於實驗材料沒有明顯區分字形別字以及形音均相似別字，因此無法證實本文所提方法在次分類別字上的效果，後續研究可以進一步加以證實。 "

" 另外，字形相似度計算的正確性評估可以考慮先建立真人相似度評估的標準，用以檢驗各種方法所得字形相似度與人類認知結果的一致性。我們目前正在進行這項工作。若是這項工作得以確認本文所提方法在計算字形相似度的正確性，本文所提方法可廣泛應用於漢字教學以及提供漢字習得理論研究相當有效的工具。 "

"

誌謝

本文作者感謝科技部計畫編號 PUE"324/4733/U/373/224 及 PUE"325/4; 33/K225/523 的支持，同時也感謝教育部及國立台灣師範大學「邁向頂尖大學計畫」的支持。

參考文獻

- 13_ Dtguucp."U" *4226+0'O qtr j qm j le" pqp/y qtf " gttqt" f gvevqp0'Rtqeggf lpi u'qh" " KGGG"37j " Kpvtpevqpcn'Y qtmuj qr "qp'F cvdcug'cpf 'Gzr gtv'U{ ugo u'Cr r rkecvqpu053/57"
- 14_ Ej cpi ."E0'J 0' *3; ; 7+0' C" P gy " Crr tqcej " hqt" Cwqo cvle" Ej kpgug" Ur gnkpi " Eqttgevqp0' Rtqeggf lpi u'qh'P cwtrcn'Ncpi wci g"Rtqegu lpi "Rcekle"Tk o "U{o r qukw o æ 7."Ugqvn"Mtgc0' 49: /4: 50'
- 15_ Ej cpi ."V0J 0'Uwpi ." [0'V0'("Ngg." [0'V0'*4234+0'C'Ej kpgug"y qtf "ugi o gpvqvp"cpf "RQU" wci i lpi "u{ ugo "hqt"tgcfc dklv{ "tgugctej 0'Rcr gt"r tguvpgf "cv'64pf "Cpwwcn'O ggwpi "qh'vj g" Uqelgv{ 'hqt'Ego r wgtu'lp'Ru{ ej qm j {.'O kppgcr qrku.'O P ."WUC0'
- 16_ Ej cpi ."V0'J 0'Uw"U" [0'("Ej gp."J 0'E0' *4234+0'Cwqo cvle" Eqttgevqp" hqt" I tcr j go le" Ej kpgug" O kur gmfg " Y qtf u0' Rtqeggf lpi u' qh" vj g" 46vj " Eqplhtgpeg" qp" Eqo r wcvqpcn' Nlpi wku'cu'cpf "Ur ggej "Rtqegu lpi ."Vc{ wcp."Vcly cp0347/3620'
- 17_ Ej cpi ."V0J 0'Ej gp."J 0'E0'Vugpi ." [0J 0'("\ j gpi ."L0'N0'*4235+0'Cwqo cvle" F gvevqp"cpf " Eqttgevqp" hqt" Ej kpgug" O kur gmfg " Y qtf u" Wulpi " Rj qpqm j kecn' cpf " Qtvj qi tcr j le" Uko krctklgu0' Rtqeggf lpi u' qh' vj g" Ugxgvj " UK J CP " Y qtmuj qr " qp" Ej kpgug" Ncpi wci g" Rtqegu lpi "UK J CP /9+ 'P ci q{c.'Lcr cp0; 9/3230'

- 18_ Ej kw"J 0'Y0"Y w"l0'E0"("Ej cpi ."l0'U0*4235-0'Ej kpgug"Ur gmkpi "Ej gengt"Dcugf "qp" Ucvkwnecn" O cej kpg" Vtcpuvcvqp0' Rtqeggf kpi u" qh' yj g" Ugxgpvj " UK J CP" Y qtmuj qr " qp" Ej kpgug"Ncpi wci g"Rtqeguulpi "UK J CP/9+P ci q{c."lcr cp06;/750
- 19_ Hkfg o cp."E0"("Uk grk"t0' *3; ; 4-0'Vqrgtcvki "ur gmkpi "gttqtu"fwtkpi "r cvkppv"xcrkf cvkqp0' Lqwtpcn'qh'Eqo r wgtu'cpf "Dkqo gf kecn'Tgugctej ."47*7+*6: 8672; 0
- 1]_ J wcpj ."E0'O 0'Y w"O 0'E0"("Ej cpi ."E0'E0*4229+0'Gttqt "F gvgevkqp"cpf "Eqttgevkqp"Dcugf " qp"Ej kpgug"Rj qpgo ke"Crr j cdgv'kp"Ej kpgug"Vgzv0'Rtqeggf kpi u'qh'yj g"Hqwtvj "Eqphgtgpeg"qp" O qf grkpi "F gekukpu'ht" Ct vklekcn'kpgvni gpeg*O F CKX +0'685/6980
- 1]_ Nkp."l 0L0'J wcpj ."H'N0"("l w"O 0'U0*4224+0C"Ej kpgug"Ur gmkpi "Gttqt"Eqttgevkqp"U{ ugo 0' Rtqeggf kpi u" qh' yj g" Ugxgpvj " Eqphgtgpeg" qp" Ct vklekcn' kpgvni gpeg" cpf " Crr rkecvkpu0' 429/4340
- 132_ Nkw"E0N0"Nck'O 0J 0'Vkgp."M0Y 0'Ej wcpj ."l 0J 0'Y w"U0J 0"("Ngg'E0[0*4233+0'Xkuwcmf " cpf " Rj qpqmj kecmf " Uko krt" Ej ctcevgtu" kp" kpeqttgev" Ej kpgug" Y qtf u<" Cpcn'ugu." Kgpvklecvkqp." cpf " Crr rkecvkpu0' Lqwtpcn' qh' CEO " Vtcpuvcvqp" qp" Culcp" Ncpi wci g" kphqto cvkqp'Rtqeguulpi ."32*4+*32-3/5; 0
- 133_ Nkw"Z0'Ej gpi ."H'Nwq."l 0'F vj ."M0"("O cwuwo qvq."l 0*4235+0C'J { dtkf "Ej kpgug"Ur gmkpi " Eqttgevkqp"Wulpi "Ncpi wci g"O qf gn'cpf "Ucvkwnecn"O cej kpg"Vtcpuvcvqp"y kj "Tgtcprkpi 0' Rtqeggf kpi u" qh' yj g" Ugxgpvj " UK J CP" Y qtmuj qr " qp" Ej kpgug" Ncpi wci g" Rtqeguulpi " *UK J CP/9+P ci q{c."lcr cp076/7: 0
- 134_ Tgp."H'Uj k'J 0"("\ j qw"S 0*4223+0C'j { dtkf "crr tqcej "vq'cwqo cvk"Ej kpgug"vgz v'ej genkpi " cpf "gttqt"eqttgevkqp0'Rtqeggf kpi u'qh'4223"KGGG"kvgtpcvqpcn'Eqphgtgpeg"qp"U{ ugo u." O cp."cpf "E { dgtpgvku038; 5/38; : 0
- 135_ Xctqn" E0' *4233+0' Rcvwtp" cpf " Rj qpgle" Dcugf " Utggv" P co g" O kur gmkpi " Eqttgevkqp0' Rtqeggf kpi u'qh'4233"KGGG"Eqphgtgpeg"qp"K/P I . "Ncu'Xgi cu."P X0775/77: 0
- 136_ Xcej j ctclcpk" X0' *4234+0' Ghgevxg" Ncdgn' O cej kpi "ht" Cwqo cvk" Gxcnvcvqp" qh' Wug" o" Ecug" F kci tco u0' Rtqeggf kpi u" qh' 4234" KGGG" Hqwtvj " kvgtpcvqpcn' Eqphgtgpeg" qp" Vgej pqmji { "ht" Gf vecvqp."J { f gtcdf 0394/3970
- 137_ Y w" U0' J 0' Nkp." E0' N0" (" Ngg." N0' J 0' *4235+0' Ej kpgug" Ur gmkpi " Ej geni' Gxcnvcvqp0' Rtqeggf kpi u" qh' yj g" Ugxgpvj " UK J CP" Y qtmuj qr " qp" Ej kpgug" Ncpi wci g" Rtqeguulpi " *UK J CP/9+P ci q{c."lcr cp0576640
- 138_ [gj ."l0H0"Nk"U0H0"Y w"O 0T0'Ej gp."Y 0[0"("Uw"O 0'E0*4235+0'Ej kpgug"Y qtf "Ur gmkpi " Eqttgevkqp" Dcugf " qp" P/i tco "Tcprngf " kvxgtvqf " kvf gz" Nku0' Rtqeggf kpi u'qh' yj g" Ugxgpvj " UK J CP" Y qtmuj qr " qp" Ej kpgug" Ncpi wci g" Rtqeguulpi " *UK J CP/9+P ci q{c."lcr cp0 65/6: 0
- 139_ \ j cpi ."N0'J wcpj ."E0'\ j qw"O 0"("Rcp."J 0*4222+0C'wqo cvk" f gvgevkpi keqttgevkpi "gttqtu" kp"Ej kpgug"vgz v'd{ "cp" crr tqzko cvg"y qtf /o cej kpi "cni qtkj o 0'kv<"Rtqeggf kpi u'qh'yj g"5: yj " Cppwcn'O ggkpi "qh'yj g" Cuuqekcvkqp'ht'Eqo r wcvkpcn'Nkpi vkuveu."46: 64760
- 13]_ 丘慶鈴 *民; 4+。避免小學生寫錯別字之教學策略" *碩士論文+。國立新竹師範學院，新竹市。"
- 13]_ 陳思彧" *民; 9+國民小學高年級別字矯誤教學研究" *碩士論文+ 臺北市立教育大學，臺北市。"
- 142_ 林佳儂" *民; ; +。國中錯別字教學策略研究" *碩士論文+。國立彰化師範大學，彰化縣。"
- 143_ 張嘉惠、林書彥、李淑瑩、蔡孟峰、李淑萍、廖湘美、孫致文和黃鏢" *民; ; +。以最佳化及機率分佈標記形聲字聲符之研究。中文計算語言學期刊，37*4+，92/：6。"
- 144_ 陳學志、張璣勻、邱郁秀、宋曜廷、張國恩 (民 322)。中文部件組字與形構資料庫之建立及其在識字教學的應用。教育心理學報。48;/4; 2。"

145_ 莊德明、謝清俊（民；6）。漢字構形資料庫的建置與應用。漢字與全球化國際學術研討會，台北。”

學術論文簡介的自動文步分析與寫作提示

黃冠誠 吳鑑城 許湘翎 顏孜曦 張俊盛

Abstract

近年來，英文逐漸變成全世界學術研究最主要的溝通的媒介。而學術英文寫作，也成為非常重要的研究與教學的領域。學者也很重視，如何透過電腦的輔助，幫助一般性的語言學習，甚或特定性的學術論文寫作。學術寫作包含許多的文節類型，包括研究論文、計畫申請書、回顧與評論文節等 (Swales, 1990)。其中，研究論文占有最重要的角色。

在學術論文中，「簡介」是絕大部分論文都有的第一個節。現今，幾乎沒有學術論文，沒有「摘要」與「簡介」，而直接詳細地描述研究的目的、方法、結果。而且，對寫者和讀者而言，「簡介」在學術論文中都扮演非常重要的角色。一篇好的簡介，要能為整篇論文定調，抓住讀者的興趣，提供論文的扼要資訊。換言之，「簡介」肩負重大責任——吸引讀者注意，讀完全文。

因此，有一些研究開始分析論文簡介如何達成其溝通的任務。Graetz (1985) 發現論文簡介似乎有共同的「問題—解法」修辭結構，依序包括問題 (problem)、方法 (solution)、評估 (evaluation)、結論 (conclusion) 等部分。Swales (1990) 分析大量的論文簡介，分析歸納出一套修辭的動機與模式：「創造研究空間」(Create A Research Space, CARS)。Swales 認為論文爭取研究得到讀者的認同，有如環境中生物爭取生存空間。為此，大部分作者依循三個修辭的步驟——也就是文步 (moves)——來說服讀者。

| ECTU ³ 文步 | 子文步與資訊內容 |
|-----------------------|--|
| 文步 K 界定範圍 | 30 聲明研究領域的重要性，及 I 或 40 聲明研究課題的廣泛性與普及性，及 I 或 50 回顧與評論前人研究 |
| 文步 KK 建立利基 | 3C0 提出與前人不同的聲明，或 3D0 指出前人研究的缺口 (i cr)，或 3E0 提出本論文的研究議題 (tgugctej 's wguvkqp)，或 3F0 說明本研究所根據的典範與傳統 |
| 文步 KKK 佔據利基 | 3C0 概述本論文的目的，或 3D0 概述本論文的方法 40 宣布本論文的主要結果與發現 50 指出本論文的結構 |

圖 30³ Uy crgu³; ; 2⁴ 提出的 ECTU³ 模式的文步與資訊內容

如圖 1 所示，這三個文步包括了「界定研究範圍」、「建立利基」、「佔據利基」。在每一個文步下，又需要描述若干必要或選項的內容。另外，美國國家醫學圖書館，也主張醫學論文作者，應提供分段有標題（labeled sections）的結構化摘要（structured abstract）。

目前已經有許多學術寫作教材，透過文步分析來教導英文非母語的學生，如何寫作學術論文（如 Swales and Feak, 2004; Glasman-Deal, 2010）。也有研究者開發軟體系統（例如，Marking Mate: writingtools.xjtlu.edu.cn:8080/mm/markingmate.html），分析學生的作文，自動產生批改的建議與評分。但是很少有系統能夠在學生寫作中，依照文步的推進，適時地提供寫作提示與輔助。直覺上，如果我們能將大量的論文簡介加以處理，自動化分析其中每句的文步，繼而分析特定文步句子的常見片語或句型，我們將可以在寫作的過程，有效地協助學生。

然而，過去所提出的自動文步分析方法，都需費時費工標註大量論文。有鑑於此，我們提出新方法，以降低人工標註的工作量，運用於訓練統計式分類器，來預測論文簡介的句子的文步，並藉以開發一個線上輔助寫作系統 *WriteAhead*。在 *WriteAhead* 的開發過程，我們採用了比 CARS 更簡單的文步分類，如圖 2 所示。用了這一套分類方式，系統容易自動分類文步，而使用者比較容易掌握使用於寫作過程。

我們期望此一自動文步分析工具，以及 *WriteAhead* 系統，有助於提升英文非母語者（non-native speakers, NNS），寫作學術論文的能力。在本論文中，我們提出了一套監督式機器學習的方法，能夠自動地學習如何將語料庫內的簡介句子，大略地分類為幾個文步。有了分類的句子之後，我們就可以統計各文步的連續詞頻率（ngrams）。在 *WriteAhead* 系統，就參考使用者選擇的文步，以及游標之前的幾個字，提示接續片語。

WriteAhead 能夠提供與排列這些提示，是因為 *WriteAhead* 透過大量的論文原始資料以及少量的人工標示，學習如何辨識 OWN 文步的句子，以及這些句子內的常見片語及其頻率。我們將在第三節詳述 *WriteAhead* 所運用的文步分類器的訓練過程。

致謝詞

本研究承蒙科技部補助研究經費，計畫標號 NSC 100-2511-S-007 -005 -MY3

以二維共振峰分布建立語者音色模型及其在語者驗證上之應用

Using 2D Formant Distribution to Build Speaker Models and Its Application in Speaker Verification

呂嘉毅¹，蕭志濱²，李明慶²，蒲長恩³，吳家隆^{2,*}

¹ 國立台北大學資訊工程學系

² 法務部調查局鑑識科學處，³ 法務部調查局通訊監察處

* 通訊作者

摘要

語音是重要的生物特徵之一，也是鑑識科學上的重要工具。在鑑識實務上常遭遇到的一個挑戰，就是通訊線路及錄音裝置的多元性。不同的裝置與線路特性會對語音證物的頻譜產生相當的影響，從而也會影響到鑑識的正確性。共振峰是語音中重要的要素，並且較不易受到通道及裝置之頻率響應的影響。在本論文中我們提出一個從分析較長時間語料，所得之二維共振峰的分布，來建立起一個語者之音色模型的方法。這個方法對於相同語詞及相異語詞方式的語者驗證工作均適用。在實驗的部分，我們報告了對約七十人規模的語料分別進行數位錄音及電話錄音的語者驗證測試。

關鍵詞: 語者驗證，線性預測方法(LPC)，共振峰，語者音色模型

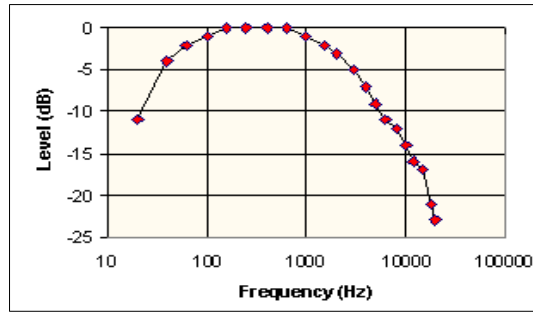
1、目的與背景

語音是人類彼此間溝通最方便也最首要的方式。語音不但是用於傳播信息，也是一項重要的生物特徵 (biometrics)，可以用來做身份識別之用。對於利用電腦來分析語音這方面的研究，大致可分為兩個領域：一是語詞識別 (speech recognition)，一是語者識別 (speaker recognition) [1-4]。若是要分辨某一個語音樣本是否來自某一個特定的語者，則又稱為語者驗證 (speaker verification 或 speaker authentication) 語者驗證又可細分為限定語詞 (text dependent) 與非限定語詞 (text independent) 兩種方式[5,6]。在限定語詞的方式中，用來比對的兩段語音樣本，其語音之內容須為相同或相似。而在非限定語詞的方式下，其語句之內容可為不同。後者之處理難度較高，但在取樣上較不受限，其應用也較為廣泛。本研究之內容是屬於語者驗證性質，同時包括了限定語詞與非限定語詞的方式。

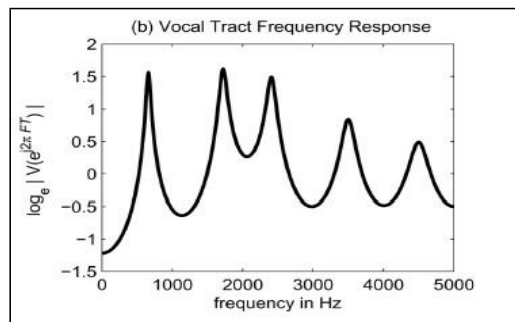
語音分析最基本的技術就是頻譜分析。由於每個人的口腔構造及發音習慣均有所不同，所以發出聲音的共鳴結構就會有所不同，語音中較強的共鳴成分就會形成頻譜中的波峰，稱為共振峰(formants)。不同的語者，因生理構造或口音上的差異，即使是發出同一個字音，其頻譜的形狀也會有所差異。所以藉著分析頻譜，我們可以分辨語者，也可以辨識字音中之韻母。

在語音分析的工作上一個經常遇到的問題就是裝置或是通訊線路(channel)所帶來的影響。如果裝置或線路的頻率響應為已知，我們尚可藉著演算還原出原音訊之頻譜，但是對於鑑識工作方面的實案而言，裝置或線路的特性通常為未知。一般而言，裝置及線路的頻率響應呈現平滑的變化，通常是中頻部分高起而高頻與低頻部分低下的情形。如下圖一所示。相對而言，發聲道的頻率響應，特別是共振良好時，其

變化較為複雜，從低頻到高频會出現若干個高峰點，如圖二所示。



圖一、裝置或線路之頻率響應示意圖



圖二、發聲道頻率響應例圖

當我們把較緩慢變化的裝置或通道的頻率響應，與較快速變化的發聲道頻率響應相乘(或相加)時，其結果就是前者會影響到後者的整體起伏變化，但是較不會影響到個別高峰點的出現及位置。也就是說，雖然語音之頻譜會受到裝置的影響，但是其中共振峰的位置相對穩定。本研究所提出之語者音色模型將以共振峰的位置為主。我們將從語音中擷取出共振峰，然後以共振峰的分布來建立一語者的音色模型。

近年來有越來越多的研究指出觀察長時間共振峰分析(LTF, Long term formant analysis)在語者驗證上的重要性。所謂長時間共振峰分析就是累計整段語料中各共振峰出現的位置，通常是前四個共振峰。因為是包含眾多不同的字音在內，所以各共振峰的位置並非固定，一段時間累積下來，就會得到一個分布曲線。因為是來自於相同語者，所以語者的因素自然包含在其中。又因為是來自於許多不同的字音的混合，所以字音的因素就會被淡化掉。因為語者驗證的重點是在語者的音色特徵而非語詞內容，所以長時間共振峰分析會是一個可利用的工具。

英國的 Nolan 與 Grigoras 在 2005 年的一篇論文中報告，紀錄語音中前四共振峰的長時間分布，在語者鑑識實案上十分有效[7]。在後續的研究中他們進一步報告各共振峰的長時間分布多呈現出不對稱(skewed)的情形，並且其分布最高點(mode)之位置在鑑識上的重要性超過其平均位置[8]。歐洲學者 Becker、Jessen、及 Grigoras 在 2008 年提出將長時間共振峰分析所得之參數值套用到高斯混合模型(Gaussian mixture model)來進行語者識別[9]。他們假定各共振峰的長時間分布為高斯分布，並自各段語料估計出高斯分布的平均值與標準差值，以進行 likelihood 計算。他們對 68 位男性語者的語料，以前三個共振峰的位置及頻寬為參數(共六個)，達到了 EER 為 0.03 的驗證成績。

德國學者 Moos 對 71 位男性語者的行動電話錄音語料進行 LTF 分析[10]，他發現 F2 與 F3 合用時有優良的語者鑑別效果。他同時發現，F3 較 F2 有著更好的穩定性，也就是對同一個語者其變異性較低。在

文章中也指出，LTF 也具有一些其他良好的特性，例如不易受到說話速度快慢及音調高低等因素的影響。中國學者 Xu 與 Kong 在 2012 年的一篇文章中報告他們以 LTF 分析進行跨語言的語者驗證[11]。他們以前四個共振峰的分布之 peak, kurtosis, 與 skewness 作為特徵值，發現能夠成功的以三種不同語言(中、英、韓)的語料進行跨語言的語者驗證。歐洲學者 Jessen 與 Becker 在 2010 也曾報告他們對德語、俄語、及阿爾巴尼亞語所進行的實驗，也有著相似的結論[12]。

前述學者所提出的長時間共振峰分析，多係對個別共振峰的分布一一的來進行，也就是屬於一個維度的分析。本論文提出的方法是將前幾個共振峰做成對的分析，也就是求得二維的共振峰分布來進行分析。又因為前二共振峰的分布與幾個主要單音韻母有很明確的對應，我們進一步將 F1-F2 平面分割為若干區域，並分別分析落在這些區域中的音框以建立更細緻的語者音色模型。在下一節中我們將詳細介紹本研究所提出建立音色模型的方法。在第三部分中我們會將這個音色模型應用到語者驗證的實驗上。

2、研究方法

本論文所提出的方法大致可分為以下幾個步驟。首先我們找出一段語音中具有共振的部分，也就是其中的有聲字音(voiced sounds)部分。其次我們以線性預測方法，逐一分析這些有聲字音的音框，找出其中的共振峰。再根據所找出的共振峰的分布建立起該位語者的音色模型。最後，我們藉比對兩組共振峰分布的相似度，來比對兩段語料之音色相似度。這些步驟分別敘述如下。

2.1、找出語料中之有聲字音部分

因為本方法是要找出語料中的共振峰分布，以建立起一語者之音色模型，所以首先我們就要找出語料中具有明顯共振的部分，即是語料中的有聲字音。在本研究中，我們先將語料切割為 20ms 大小的音框，相鄰的音框有 10ms (即為 50%)的重疊。我們對每一個音框計算出一個音量大小值，以及求取其 autocorrelation function (ACF)曲線，並找出其在合理之週期範圍所能達到的最高值。如果一個音框具有足夠的音量大小以及夠大的 ACF 峰值，我們就接受此一音框為一個有聲字音的音框。在語料量足夠的情形下，上述兩項門檻值可以做較嚴格之設定，以確保所找出的音框均有不錯的共振品質。

2.2、以線性預測法(LPC)找出音框之共振峰

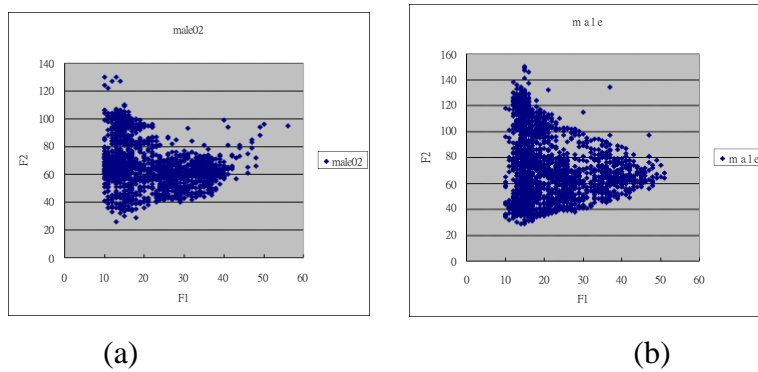
本研究聲音樣本的取樣率(sampling rate)為 11,025 Hz，依據文獻的建議，階數 p 的設定大約是在 14 至 16 左右。在這個取樣率之下，音訊中的頻寬大約保留到 5K Hz，其中共振峰的數目因字音而異，大約有五到六個。本研究在建立語者音色模型時，將會利用到前四個共振峰，即 F1 到 F4。在推導這 p 個 LPC 係數方面我們所用的方法是常見的 Levinson- Durbin 演算法。這個演算法首先自一音框求出 p 個 autocorrelation function 的值，然後再藉由一個遞迴式的演算法解出模型的 p 個係數值來。在對一個音框求出一組係數值之後，我們會再將音框之音訊值帶入模型，並計算出預測值與實際值之誤差。倘若誤差值過大，則表示其所找出的共振峰並不準確，或是該音框的共振仍為不佳，或是該音框受到了較大的雜訊。當有此情形發生時，我們就會略過這些音框不用。一般而言在此一階段會被淘汰的音框大約占有聲字音的 5% 以內。

2.3、建立語者之音色模型

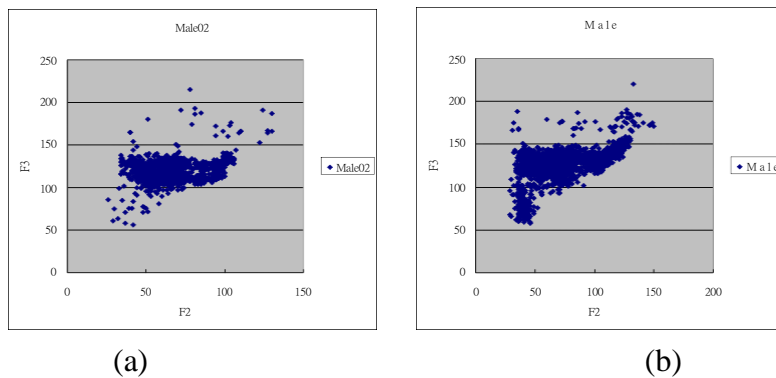
自一語者之語料找出有聲字音之音框，並推算出其中的共振峰之後，我們就可開始建立該語者之音色模型。如果是使用於相同語詞式的語者驗證工作，無論語料多寡，我們都可以建立起音色模型，只是在語料量少時，所建立起的音色模型也是與語詞內容相關。而當使用於相異語詞式的語者驗證時，我們就需要有較多的語料才能建立起一個較為完整的音色模型。

在我們所建立的語者模型中，第一部分就是共振峰分布的情形。在上一個小節中我們提到，我們自每一個音框找出其前四道的共振峰(F1- F4)。在此我們將自語料產生出三個二維(2D)的共振峰分布圖，分別是 F1 對 F2，F2 對 F3，以及 F3 對 F4。每一個有效的音框將會對應到這些圖中的一個點，語料的時間越長，則分布圖中的點也會越多。因為每個語者的音色有所不同，即使是在發出同一個音，其共振峰的位置也會有所差異，反映在這些二維的分布圖上，就是這些點的集中位置會有所不同。

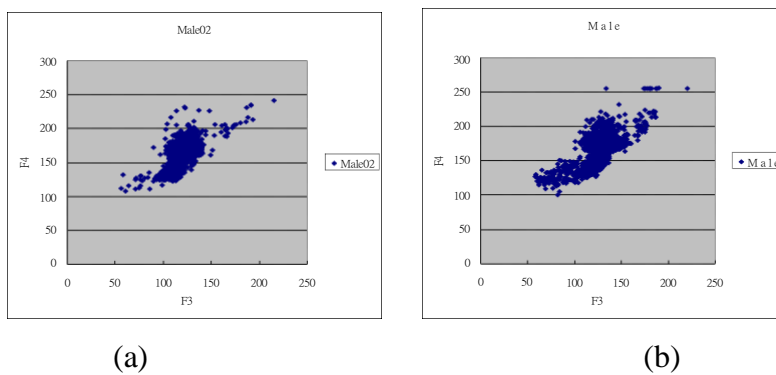
在圖二中我們顯示了兩位不同男性語者在約 60 秒的相同語詞語料所呈現出的 F1-F2 共振峰分布圖。從圖中我們可以很清楚的看出兩位語者的 F1 對 F2 在分布上的差異。因為是依據相同的語詞，所以這裡所反映出的差異主要是來自於二人在音色上的不同。圖四及圖五分別顯示出此二位語者之 F2-F3 與 F3-F4 的分布圖。



圖三、兩位男性語者之 F1-F2 分布圖

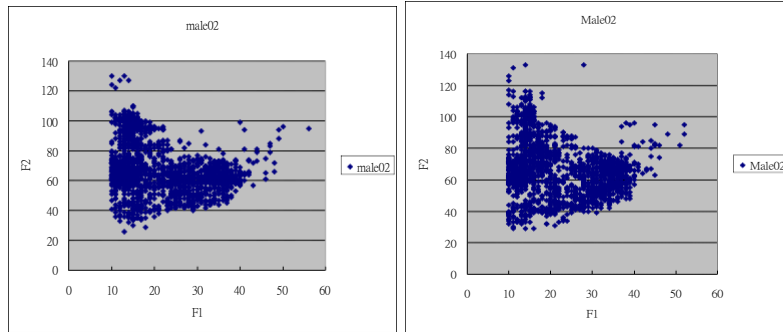


圖四、兩位男性語者之 F2-F3 分布圖

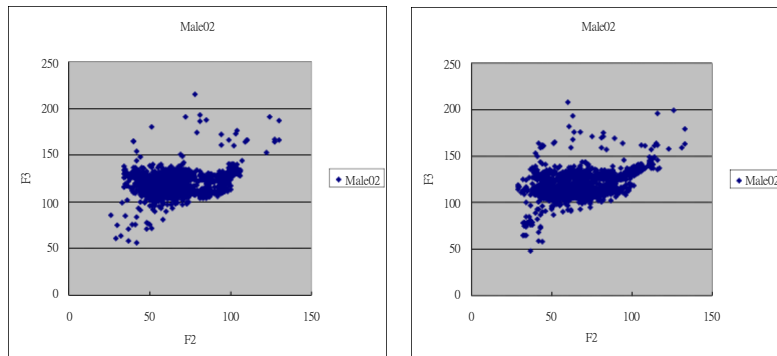


圖五、兩位男性語者之 F3-F4 分布圖

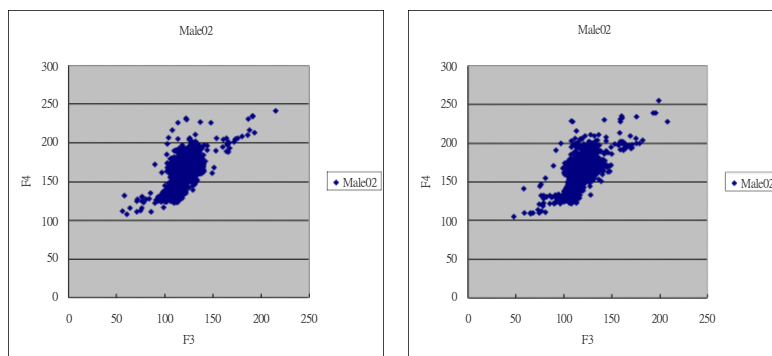
在上二圖中我們同樣可以觀察到，此二位語者在共振峰分布上的差異。接下來我們要顯示的是，相同一位語者對於不同的語詞其所呈現之共振峰分布情形。在這裡我們將以 **Male02** 這位男性語者在兩段各約 60 秒，但語詞內容為不同的語料所得之共振峰分布圖來相比較。在下面的三個圖中我們分別比較此位語者在這兩段語料在 F1-F2，F2-F3，與 F3-F4 分布上的差異。



(a) (b)
圖六、同一位男性語者不同語詞語料之 F1-F2 分布圖



(a) (b)
圖七、同一位男性語者不同語詞語料之 F2-F3 分布圖



(a) (b)
圖八、同一位男性語者不同語詞語料之 F3-F4 分布圖

自以上三對共振峰分布圖我們可以清楚看出它們之間的相似性。即使在語詞內容為不同的情形下，因著語料長度夠長，亦即字音重疊程度提高，相同語者之共振峰的累進分布也就趨近於相似。基於以上對不

同語者相同語詞，及相同語者不同語詞之共振峰分布所做的觀察，我們認為以共振峰之分布來建立語者音色模型，將會是一個有效區分語者的工具。

除了以共振峰來建立語者音色模型外，我們也會為語者建立若干個頻譜模型。其原因為頻譜中能量變化上下起伏，但共振峰只標示出其中高峰點的位置，而忽略了頻譜中下凹部分變化的資訊。這些下凹的部分包含了實際發聲道系統中的零點(zeros)，但是 LPC 方法假定發聲道模型中只有極點，所以其所包含的資訊並不完整。

我們為每個語者建立起十個平均頻譜作為其音色之頻譜模型。在此我們先將整個 F1-F2 平面分割為十個區域，其中之九個區域大致對應到一般(F1,F2)會落入的區域，而第十個區域就是對應(F1,F2)不可能出現的區域，對應到此一區域即表示在共振峰的推算上發生了錯誤。對於每一個有聲字音的音框，我們除了求取其前四個共振峰之外，我們也求出其對數功率頻譜(log power spectrum)。我們依其 F1 與 F2 的值將其頻譜累計至前述之十個區域的平均頻譜之中。

本研究的目標是要發展出，一種能夠同時適用於相同語詞式、與相異語詞式語者驗證的方法。共振峰的分布圖仍多少會受到語詞內容的影響，特別是當語料的量較少時。但是這九個區域中的平均頻譜，則較不會受到語詞內容不同的影響。其原因為，當語詞內容為不同時，只會影響到落到各區域中音框的數目，但是對於取自各區中音框之平均頻譜的影響則較輕。也就是說，由這些平均頻譜所組成的模型，能夠更全面性的表現出一位語者的音色特徵。

下表顯示出對男聲我們分割 F1-F2 平面的方式。因為男女聲在共振峰分布上有著相當顯著的差異，所以在區域的分割上也有所不同。在我們的分割方式中，有部分的區域剛好與幾個主要的韻母有所對應，我們也一併標示於表中。

表一、將男性語者之 F1-F2 平面分為 10 個區域

| 區域編號 | F1 範圍 (Hz) | F2 範圍 (Hz) | 所對應之單韻母音 |
|------|------------|-------------|----------|
| 1 | 258 – 387 | 580 – 1010 | ㄨ |
| 2 | 258 – 387 | 1010 – 1913 | |
| 3 | 258 – 387 | 1913 – 2687 | 一 |
| 4 | 387 – 688 | 580 – 1075 | ㄛ |
| 5 | 387 – 688 | 1075 – 1720 | ㄜ |
| 6 | 387 – 688 | 1720 – 2472 | ㄝ |
| 7 | 688- 1075 | 580 – 1075 | |
| 8 | 688- 1075 | 1075 – 1483 | ㄩ |
| 9 | 688- 1075 | 1483 – 2257 | |
| 0 | 其他頻率 | 其他頻率 | 屬錯誤情形 |

再接下來我們就對落在每一個區域內的音框進行分析，並擷取出一些描述語者聲音特色的參數。計有以下各項：

A、各區域內音框之 FFT 平均頻譜

我們利用 LPC 頻譜來找出 F1 及 F2 共振峰，藉以將音框分群，其原因為 LPC 頻譜的長處就是在於找出主要之共振峰。但是聲音中除了共振峰外尚其他的特徵及變化，FFT 頻譜就有較全面的紀錄。在此我們先求出各音框的 FFT 頻譜，再將這些頻譜作平均，以得到一個平均 FFT 頻譜。因為音框已經大致依照韻母經過分類，所以落在同一區域中的音框，其發音大致相近，頻譜也會相似。當我們將這些頻譜加以平均時，其個別之差異將會淡化，而共有之特徵將會得到增強。未來在比對兩個語者模型時，我們將對應區域中之平均頻譜加以比對，以計算其相似度。

B、各區域內音框之 LPC 平均頻譜

LPC 頻譜係由線性預測之係數所推出，其特色在於較能顯出共振峰之位置。在比對語者之音色時，共振峰仍為最主要之資訊，因為共振峰位置是由發聲器官及發音習慣所決定。我們在此將一個區域內所有音框之 LPC 頻譜加以平均，以得到一個平均頻譜。這個頻譜在高峰的部分與 FFT 平均頻譜十分相近，但是在波谷的部分則有較大的差異。我們可藉比對 LPC 頻譜特別來檢視兩位語者在共振峰位置上的相似程度。

C、各區域內音框之共振峰的累加曲線

我們做完 LPC 分析之後，對於每個音框我們得到了一組的共振峰。目前我們是找出前五至六個，就是 F1 到 F6。不過在有效的頻寬範圍中我們可能只會看到四個共振峰，此時第五個共振峰就經常在有效頻寬之外，在比對時可忽略不看。在這裡我們將一個區域內，所有音框的所有共振峰，全部投到同一個頻率軸線上。因為區域本身就是按 F1 與 F2 來劃分的，自然這些音框在 F1 與 F2 的頻率範圍會有相當高的一致性。但是我們發現這些來自同一語者的聲音，在 F3, F4, 甚至於 F5 也會有著相當好的一致性。這一項特徵參數也是用來反映出一位語者在發出不同字音時，其共振峰分布的情形。但是與前一項不同之處是，在前一項中，我們是對 LPC 頻譜之強度值做平均，所以一個音強的音框，會較一個音弱的音框有著更大的影響。但是在這裡我們對每一個音框的共振峰都是以同一強度值紀錄，所以在意義上略有不同。

2.4、藉比對音色模型決定二段語料之音色相似度

在前面的部分我們說明了如何自一段語料建立起其語者的音色模型，一個模型中包含了三個二維的共振峰分布圖，以及將 F1-F2 平面分為十個區域後，在每個區域所累積出的 FFT 平均頻譜、LPC 平均頻譜、以及共振峰分布累加曲線。在比，我們將比對兩個音色模型的內容，以估計兩段語料之語者在音色上的相似程度。

我們進行比對兩個音色模型的基本方法為計算其相關係數值。在比對二維共振峰分布時，我們分別就三對的二維共振峰分布(F1-F2, F2-F3, F3-F4)兩兩計算出其間之相關係數值。在比對兩個音色模型中對應之 FFT 頻譜、LPC 頻譜、及一維之共振峰累加曲線時，我們則兩兩計算其間之一維相關係數。在完成以上之計算後，我們將得到六個相關係數值。因為這些相關係數均具有不同的特性，之後我們可以依語料的特性顯選擇使用，或是將這些相關係數做加權平均，以得到一個綜合相似指標值。

3、實驗與結果

在以相異語詞進行語者驗證之實驗部分，我們分別針對了男聲及女聲，並以數位錄音及電話錄音兩種方式進行實驗。每一種的錄音方式又可分為比對同次錄音中之不同語句，以及比對不同次錄音中之不同語句兩種方式。在實驗中，我們使用了 72 人的語音樣本，其中有男生 38 人及女生 34 人，均為 18 歲以上之成年人。採樣分兩次進行，時間上的間隔為兩個月。實驗中所用到的國語語句每組的句數均是六十句，每句有六至十個字不等。

每次的錄音因語者說話速度快慢不同，大約有三分鐘的長度。我們再將每份語料分為前後兩段，每段的長度大約在 90 秒左右，其中包含一句與一句之間停頓的時間。倘若扣除掉語句間停頓所花的時間，每段錄音的長度約為 60 秒左右。因為前段與後段有著不同的語詞內容，所以我們可以用同次錄音中的前段與後段進行相異語詞之語者驗證。因為是取自於同一次的錄音，無論是錄音裝置、錄音環境、或是語者的生理狀況都會極為相似，所以我們預期比對的結果(即驗證正確率)將會較好。

我們也將利用不同次錄音中的前後段落進行交叉的比對，例如將第一次錄音中的前半段，與第二次錄音

中的後半段進行比對；或是將第一次錄音的後半段與第二次錄音的前半段進行比對。這種的比對方式不單是語詞不相同，就連錄音的裝置、錄音的環境、通訊的線路、以及語者的生理狀況都有可能不同。我們預期驗證的正確率也將會下降。

在驗證所用的參數部分，如在前面一節中所述，我們有以下幾個，分別給予編號P1到P9：

- P1 九個特徵區域中音框之 FFT 平均頻譜
- P2 九個特徵區域中音框之 LPC 平均頻譜
- P3 九個特徵區域中音框之共振峰分布曲線
- P4 全域之 F1-F2 分布
- P5 全域之 F2-F3 分布
- P6 全域之 F3-F4 分布
- P7 P1-P3 之綜合
- P8 P4-P6 之綜合
- P9 P1-P6 之綜合

其中 P1 到 P6 是個別的參數曲線或是分布圖，P7 是把前三個特徵曲線(P1-P3)加以綜合的結果。而 P8 是將 P4-P6 這三個分布圖加以綜合的結果，而 P9 則是再進一步把 P7 和 P8 加以綜合。接下來我們就依序表列不同條件下所得之驗證正確率。

A、以同一次數位錄音中之不同段落進行語者驗證

相對於電話錄音，數位錄音有著較大的頻寬，保存語者音色的能力較佳。又因為比對用的語料取自同一次的錄音，在錄音裝置、錄音環境、通訊線路、以及語者生理狀態等多方面均為最相近，所以驗證的準確率最高。

表二、以同次數位錄音中之相異段落進行語者驗證所得之驗證等錯誤率 EER (%)

| 參數 | 男聲 | 女聲 |
|----|-----|-----|
| P1 | 0.3 | 0.0 |
| P2 | 0.3 | 0.0 |
| P3 | 1.3 | 0.1 |
| P4 | 0.3 | 0.7 |
| P5 | 5.4 | 4.3 |
| P6 | 5.6 | 4.5 |
| P7 | 0.2 | 0.1 |
| P8 | 1.8 | 0.1 |
| P9 | 1.1 | 0.0 |

從上表中可以看到，無論是男聲或是女聲，我們都達到了相當高的正確率。這就反映出這些特徵值確實能夠掌握到一個語者的音色特徵。仔細的比較，我們可以發現 P1-P3 的表現略為優於 P4-P6，但是 P4 仍然是有著相當高的驗證正確率(EER 1%以下)。女聲的正確率略優於男聲，這也是因為女聲所使用的頻域較男聲為寬，其音色可有較大的差異度。

B、以不同次之數位錄音中之不同段落進行語者驗證

如前所述，兩次錄音之間間隔了兩個月的時間，在裝置及人員方面均有所變化，所得到的語者驗證正確率也就有所下降。

表三、以不同次數位錄音中之不同段落進行語者驗證所得之驗證等錯誤率 EER (%)

| 參數 | 男聲 | 女聲 |
|----|------|------|
| P1 | 15.1 | 14.3 |
| P2 | 6.2 | 9.2 |
| P3 | 7.0 | 11.2 |
| P4 | 19.4 | 21.3 |
| P5 | 24.2 | 16.0 |
| P6 | 18.8 | 21.0 |
| P7 | 8.4 | 10.5 |
| P8 | 12.4 | 11.5 |
| P9 | 6.9 | 10.6 |

從上表中可以看到，所有特徵值的驗證 EER 值都有明顯的上升。但其中部分的參數，尤其是 P2 與 P3，表現相對穩定。因此之故，P7 與 P9 也有著較佳的表現。在這個部份我們看到對男聲的正確率略為優於女聲，一個可能的原因就是男聲相對受到隨時間變化因素的影響較小。

C、以同次之電話錄音中之不同段落進行語者驗證

電話錄音之頻寬為 3.5k Hz 左右，相當程度低於數位錄音之約 5.5k Hz 的頻寬。這個減少的頻寬對於語者的音色會產生一定程度的影響，這樣的影響也些微的反映在驗證的 EER 之上。

表四、以同次電話錄音中之相異段落進行語者驗證所得之驗證等錯誤率 EER (%)

| 參數 | 男聲 | 女聲 |
|----|-----|-----|
| P1 | 0.1 | 0.2 |
| P2 | 0.3 | 0.4 |
| P3 | 0.3 | 0.1 |
| P4 | 1.9 | 2.1 |
| P5 | 1.8 | 1.4 |
| P6 | 3.0 | 1.5 |
| P7 | 0.2 | 0.2 |
| P8 | 0.2 | 0.4 |
| P9 | 0.1 | 0.2 |

比較上表與表二，我們看到在男聲的部分差異不大，但是在女聲的部分正確率有略為下降。這個下降的情形主要是受到頻寬被壓縮的緣故，對女聲音色的影響較比對男聲明顯。與前兩表(表二及表三)相似的是，P1-P3 的表現仍是優於 P4-P6，但是當我們將 P4-P6 綜合起來用(即為 P8)，仍然是有著不錯的正确率。

D、以不同次之電話錄音中之不同段落進行語者驗證

在四種組合之中，此一組合之情形最接近鑑識工作的實際情況。此處所利用到的電話線路十分多元，包括固網、手機、長途等等。語者的發話環境也各為不同，年齡的範圍也較廣。所以這個部分的語料的品質較為接近實案中的情形。

表五、以不同次電話錄音中之不同段落進行語者驗證所得之驗證等錯誤率 EER (%)

| 參數 | 男聲 | 女聲 |
|----|------|------|
| P1 | 11.5 | 8.1 |
| P2 | 12.5 | 9.0 |
| P3 | 12.5 | 8.4 |
| P4 | 7.7 | 11.4 |
| P5 | 16.7 | 8.9 |
| P6 | 24.6 | 15.6 |

| | | |
|----|------|-----|
| P7 | 12.2 | 8.7 |
| P8 | 15.5 | 6.3 |
| P9 | 14.4 | 6.2 |

將上表與表四比較，我們可以發現驗證之 EER 有所上升。從 P9 可以看出來，女聲部分大約上升了六個百分點，但是男聲則上升約十個百分點。值得注意的是女聲部分 P8 的表現超越 P7，這表示出當錄音之線路及裝置多元時，共振峰分布之特徵比頻譜特徵有著較佳的表現。此外男聲部分 P4 的表現也相對較佳。P6 在電話錄音的部分表現較差的原因是，共振峰 F4 常常是高於截止頻。所以找到的 F4 經常並非真正的 F4。不同次(相隔兩個月)的電話錄音，又是在非實驗室的環境之下錄製，能夠達到 90% 以上的驗證正確率，顯示本方法具有實用之潛力。

4、結論

在本論文中我們提出了一種可以以相異語詞語料進行語者驗證的方法。因為用來比對之語料的語詞可能為不同，我們無法逐句地來進行比對。我們所提出的方法，乃是自語料分析長時間共振峰的分布以建立起語者音色模型，然後再就兩個音色模型加以比對。這樣的一個音色模型，我們認為至少具有以下兩方面的優點。第一是頻譜本身的形狀相當容易受到通訊線路或錄音裝置的影響，從而影響到比對的結果，但是共振峰的位置卻是相對較為不會受到裝置或線路的影響。因為在鑑識實務上，線路及錄音裝置相當多樣化，也難以取得其頻率特性資料。使用共振峰特徵將有助於提升鑑識的穩定性。第二方面的優點是來自於我們將 F1-F2 平面分區之後，再就落在各區之中的音框，分別求取其平均頻譜和共振峰分布曲線，以建立語者音色模型。這些區域大致與不同的單韻母音對應，在語詞不同的情形下，各韻母出現的次數或為不同，但所建立起的語者音色模型仍屬完備。

5、參考文獻

1. R.D. Peacocke and D.H. Graf, "An introduction to Speech and Speaker Recognition," *IEEE Computer Magazine*, pp. 26-33, August 1990.
2. J.P. Campbell, "Speaker Recognition: A Tutorial," *Proceedings of the IEEE*, Vol. 85, pp. 1437-1462, September 1997.
3. Sadaoki Furui, *Digital Speech: Processing, Synthesis, and Recognition*, 2nd. Edition, Marcel Dekker, New York, New York, 2001
4. Thomas F. Quatieri, *Discrete-Time Speech Signal Processing principles and Practice*, Prentice Hall, 2002.
5. T. Dutta, "Text Dependent Speaker Identification based on Spectrograms," *Proceedings of Image and vision Computing New Zealand 2007*, pp. 238-243, December 2007.
6. F. Bimbot, J.-F. Bonastre, C. Fredouille, G. Gravier, I. Magrin-Chagnolleau, S. Meignier, T. Merlin, J. Ortega-Garcia, D. Petrovska-Delacretaz, and D. A. Reynolds, "A Tutorial on Text-Independent Speaker Verification," *EURASIP Journal on Applied Signal Processing*, Vol. 4, pp. 430-451, 2004.
7. F. Nolan and C. Grigoras, "A case for formant analysis in forensic speaker identification," *International Journal of Speech Language and the Law*, Vol. 12, No. 2, pp. 143-173, 2005.
8. K. McDougall, P. Harrison, F. Nolan, and C. Kirchhubel, "Voice Similarity and Long Trem Formant Analysis," University of Cambridge report.
9. T. Becker, M. Jessen, and C. Grigoras, "Forensic Speaker Verification Using Formant Features and Gaussian Mixture Model," *Proceedings of Interspeech 2008 Special Session: Forensic Speaker Recognition – Traditional and Automatic Approaches*, Brisbane, Queensland, Australia, September, 2008.
10. A. Moos, "Long-Term Formant Distribution (LTF) based on German spontaneous and read speech," *Proceedings of IAFPA 2008*, Swiss Federal Institute of Technology, Lausanne, 2008.
11. Y. Xi, "Vocal tract characteristics on long-term formant distribution," *Proceedings of the 2012 International Conference on Computer Science and Network Technology (2012 ICCSNT)*, pp. 207-211, Dec. 29-31, 2012.
12. M. Jessen and T. Becker, "Long-term formant distribution as a forensic-phonetic feature," *Journal of Acoustical Society of America*, Vol. 128, No. 4, p. 2378, 2010.

Unsupervised Approach for Automatic Keyword Extraction from Arabic Documents

Arafat Awajan *

Abstract

In this paper, we present an unsupervised two-phase approach to extract keywords from Arabic documents that combines statistical analysis and linguistic information. The first phase detects all the N-grams that may be considered keywords. In the second phase, the N-grams are analyzed using a morphological analyzer to replace the words of the N-grams with their base forms that are the roots for the derived words and the stems for the non-derivative words. The N-grams that have the same base forms are regrouped and their counts accumulated. The ones that appear more frequently are then selected as keywords. An experiment is conducted to evaluate the proposed approach by comparing the extracted keywords with those manually selected. The results show that the proposed approach achieved an average precision of 0.51.

Keywords: Keyword extraction, Keyphrase extraction, Arabic Language, N-gram.

1. Introduction

Keyword extraction is the process of identifying a short list of words or noun phrases that capture the most important ideas or topics covered in a document. Keyword extraction has been used in a variety of natural language processing applications, such as information retrieval systems, digital library searching, web content management, document clustering, and text summarization (Rose et al. 2010). Although keywords are very useful for a large spectrum of applications, only a limited number of documents with keywords are available on-line. Therefore, appropriate tools that can automatically extract keywords from text are increasingly needed with the continually growing amount of electronic textual content available online.

In this paper, an unsupervised two-phase approach for keyword extraction from Arabic

* Princess Sumaya University for Technology – Department of Computer Science, Amman – Jordan
E-mail: awajan@psut.edu.jo

documents is described. The proposed method combines the document's statistics and the linguistic features of the Arabic language to automatically extract keywords from a single document in a domain-independent way. In the first phase, all the N-grams are extracted and those considered as potential candidate keywords are retained. In the second phase, the candidate keywords are analyzed linguistically by a morphological analyzer that replaces each term with its base form, which are the roots of the derived words and the stems of the non-derivative words. The candidate keywords are then grouped in such a way that the keywords extracted from similar roots and stems are put together and their counts accumulated.

This paper is organized as follows. In section 2, we present related works and the main approaches to keyword extraction. Section 3 highlights the main Arabic language features used in our technique. A detailed description of the proposed technique and its two phases provided in Section 4 and Section 5. Section 6 consists of the experimental results and the main findings of the evaluation of the proposed method.

2. Related Work

Existing automatic keyword extraction methods can be divided into two main approaches: supervised and unsupervised (Pudota et al. 2010; Hasan and Ng 2010). In the supervised approach, the keyword extractor is trained to determine whether a given word or phrase is a keyword or not. An annotated set of documents with predefined keywords is always used in the learning phase. All the terms and noun phrases in the text are considered as potential keywords, but only those that match with keywords assigned to the annotated data are selected. The main disadvantages of this approach are its dependency on the learning model, the documents used as the training set, and the documents' domains. Furthermore, training data and learning processes are usually time-consuming (Turney 2000; Turney and Pantel 2010; Frank et al. 1999; Hulth 2003; Hulth 2004).

The unsupervised approach for keyphrase extraction avoids the need for annotated documents. It uses language modeling and statistical analysis to select the potential keywords. A candidate keyword is often selected based on features such as its frequency in the document, the position of its first occurrence in a document, and its linguistic attributes, such as its stem and part-of-speech (POS) tag (Matsuo and Ishizuka 2004; Mihalcea and Tarau 2004; Liu et al. 2009). The unsupervised methods are in general domain-independent and less expensive since they do not require building an annotated corpus.

Keyword extraction algorithms from both approaches have been successfully developed and implemented for documents in the European languages (Rose et al. 2010; Liu et al. 2009; Matsuo et al. 2004). However, despite the fact that Arabic is one of the major international languages making up about 4% of the Internet content, not many studies about extracting

Arabic keywords have been performed. El-Shishtawy and Al-Sammak (2009) presented a supervised method that uses linguistic knowledge and machine learning techniques to extract Arabic keywords. The system uses an annotated Arabic data set of 30 documents from a specific domain, compiled by the authors as a training data set. The keywords from the documents' data set used to evaluate their system were assigned manually.

An unsupervised keyphrase extraction system (KP-Miner) was proposed by El-Beltagy and Rafea (2008). This system was basically developed for the English language and then adapted to work with the Arabic language. Statistical analysis of the texts was conducted in order to determine the most weighted terms. Two main conditions are considered; the first states that a phrase has to have appeared at least n times in the document from which the keywords are to be extracted, and the second condition is related to the position where a candidate keyphrase first appears within an input document. The linguistic analyses performed on the texts are limited to stop word removal and word stemming.

The hypothesis defended in this work is that using the linguistic features of the Arabic language — mainly its rich and complex morphological structure — may present an attractive paradigm to improve the extraction of keywords. The proposed approach is designed to work on a single document without any prior knowledge about its content or domain. Typically, a generic unsupervised keyphrase extractor features two steps; the first is to extract as many candidate words as possible, and the second is to apply the linguistic knowledge of the text language to tune the final list of extracted keywords.

3. The Features of Arabic Language

Arabic is a Semitic language with rich morphology that is a combination of non-concatenative morphology and concatenative morphology. Regarding the concatenative aspect, an Arabic word is composed of a stem, affixes, and clitics. The affixes are concatenative morphemes that mark the tense, gender, and/or number of the word (Al-Sughaiyer and Al-Kharashi 2004). A clitic is a symbol consisting of one to three letters that can be attached to the beginning or the end of a word. It represents another part of speech, such as a preposition, a conjunction, the definite article, or an object pronoun (Habash 2010; Awajan 2007; Diab et al. 2007). In terms of their formation, most of the stems obey non-concatenative rules and are generated according to the root-and-pattern scheme. In general, an Arabic word may be decomposed in its components according to the structure shown in figure 1. For example, the word “والعبون”, or “and the players” in English, consists of the clitics “و” and “ال”, the stem “العب”, and the postfix “ون”. Its stem is generated from the root “لعب”, according to the pattern “فاعل”. Figure 2 shows the steps for a word formation.

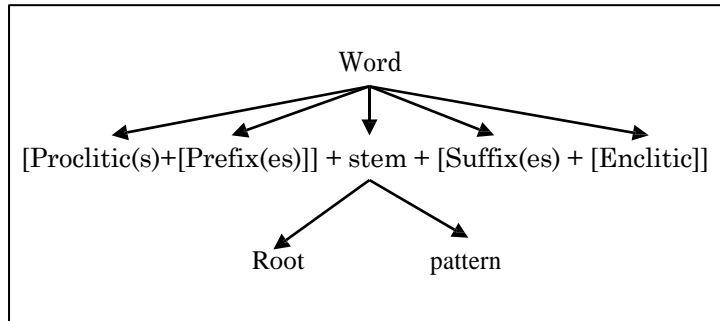


Figure 1. Arabic derivative word structure

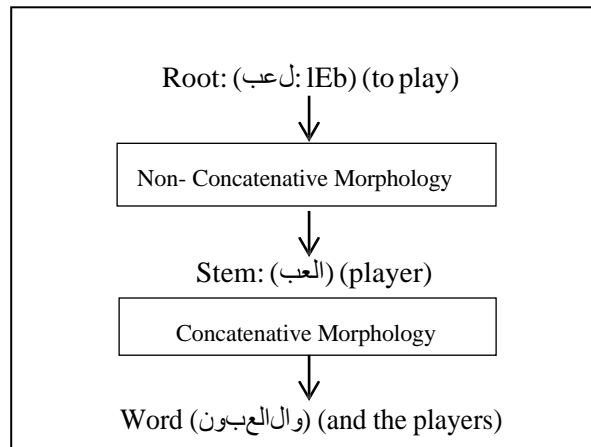


Figure 2. Arabic word formation (Example)

Arabic words are classified into two categories: derivative words and non-derivative words. The stems of derivative words are generated from the roots according to standard patterns or templates. These standard patterns represent the major spelling rules governing Arabic words. Based on the above, a derivative Arabic word can be represented by its root along with its morphological pattern, and its roots carry its basic conceptual meaning.

Non-derivative words include two sub- categories: fixed words and foreign words. Fixed words are a set of words that do not obey the derivation rules. These words are generally stop words, such as pronouns, prepositions, conjunctions, question words, and the like. The foreign words are nouns borrowed from foreign languages.

The combinatory nature of the Arabic language morphology creates an important obstacle for different natural language processing applications, including keyword extraction. This property, generally known as “data sparseness”, results in a large number of words generated from the same root but with different stems (Benajiba et al. 2009). Consequently, the grouping of words according to their surface or stems cannot give keywords that

accurately reflect the content of the document.

In order to tackle this problem, we need to conduct a deeper morphological analysis to extract the roots and to consider their properties in order to group related words and increase the weight of those representing the main ideas covered by the text. The linguistic analysis we are proposing will be applied at two different levels of the keyword extraction. The input text is preprocessed to assign each word with its POS in order to detect all the possible N-grams. The detected N-grams are then post-processed to extract the roots, and to group the N-grams generated from the same roots, and to accumulate their weights.

4. N-Gram Extraction

4.1 Part-of-Speech Tagging

This phase consists of several operations: sentence delimiting, tokenization, and POS tagging. The input text is processed to delimit sentences, following the assumption that no keyphrase parts are located separately in two or more different sentences (Pudota et al. 2010). Punctuation marks, such as commas, semicolons, and dots, are used to divide the input documents into sentences.

Tokenization aims at turning a text into a list of individual words or tokens (Manning et al. 2009). As the clitics attached to a word always refer to other entities, such as pronouns, prepositions, conjunctions, and the definite article, a tokenizer is applied to separate all the clitics except the definite article from the word. The tokenizer is repeatedly applied until the word stops changing.

We then assign a POS tag to each token using the Stanford Arabic parser (Green and Manning. 2010). The assigned POS tags are later used to select the possible N-grams, remove the verbs, and remove meaningless terms, such as the stop words.

4.2 N-gram Extraction and Filtering

A keyword is typically a combination of nouns and/or adjectives. Furthermore, the number of terms that are allowed in a keyword is often limited to three words. Thus, each sentence is processed to extract all the possible N-grams that constitute a sequence of adjacent words with a maximum length of three words. All the N-grams that contain verbs, stop words, or clitics are removed. Only the N-grams that have their members labeled with one of the POS tags marking nouns or adjectives are retained. In addition, the unigrams that are not labeled as nouns are removed from the N-gram list. Figure 3 shows the detected unigrams, bi-grams, and trigrams from a sentence.

| | |
|-----------------------------------|--|
| Input Sentence in Arabic: | هي هاشميا ن يقاالرد كئل مءلا الءى زءارة كءى المرءنء ىءلا مءقا |
| Input Sentence in English: | The American president visited the Hashemite Kingdom of Jordan. |
| Tokenization: | :ءاءة هاشميا ن يقاالرد كئل مءلا الءى ءارة ءب كءى مءلا ءسء ىءلا |
| Unigrams: | كئل مءلا - ءءمء هءا - ن يقاالرد - كئل مءلا - ءارة - كءى مءلا - |
| Bi-grams: | كءى المرءنء ىءلا مءلا - هي هاشميا ن يقاالرد - ن يقاالرد |
| Tri-grams: | هى هاشميا ن يقاالرد كئل مءلا |

Figure 3. N-Grams Extraction

5. Keywords Selection

5.1 N-gram Normalization

Normalizing N-grams is the process of reducing the words of an N-gram into their base forms. This process will allow the clustering of N-grams carrying the same information, hence reducing the sparseness of the text's potential keywords. To achieve this objective, a word morphological analyzer is developed based on the Alkhalil Morpho-Syntactic System (Boudlal et al. 2010). It is applied individually to the words on the list of N-grams. The morphological structures produced by the analyzer are used to determine the category of words, derivative or non-derivative. The derivative words are represented by their root along with their morphological pattern, and the non-derivative words are represented by their stem, permitting different N-grams that have common base forms to reinforce each other in scoring and to reduce the number of redundant terms and concepts. Each N-gram is associated with its list of base forms called the normalized N-grams (NNG) at the end of this step.

5.2 N-gram Clustering and Weighting

All the N-grams generated from the same base forms are grouped together, their counts accumulated, and represented by their NNG. A vector representation of the text is produced where each detected NNG and its frequency are listed. In this work, we define the frequency of a normalized N-gram NG_i noted $Freq(NG_i)$ as the sum of all the N-grams having the same base forms of NG_i .

Each normalized N-gram should be assigned a weight that represents its relevance to be selected as a keyword. The keyword frequency and the keyword degree are generally considered for scoring potential keywords (Rose et al. 2010, Mihalcea and Tarau 2004). The weight of a normalized N-gram NG_i is given by the following formula:

$$\text{Weight}(\text{NG}_i) = \text{Freq}(\text{NG}_i) / \sum_{j=1}^m (\text{Freq}(\text{NG}_j))$$

where m is the number of Normalized N-grams.

As the unigrams are generally more frequent than the bi-grams and bi-grams are more frequent than tri-grams, we need to correct the weight of N-grams by introducing a new measure called score. The N-gram score takes into account the relevance of individual components forming the N-gram. The score of a unigram is equal to its weight since a unigram has one component. The score of other N-grams (bigrams, trigrams, ...) is given by the following formula:

$$\text{Score}(\text{NG}_i) = \text{Weight}(\text{NG}_i) + \sum_{j=1}^N (\text{Weight}(\text{T}_j)),$$

where the $\text{T}_1, \text{T}_2, \dots, \text{T}_N$ represent the N roots/stems of the normalized N-gram NG_i .

The degree of an N-gram is calculated as the sum of its Weight and the Weights of all the higher structures containing this N-gram. Thus, the degree favors terms occurring frequently in longer candidate keywords, and the score favors the frequent terms regardless of their co-occurrence with other terms.

5.3 Keywords Selection

The list of N-grams is reordered according to their scores since the highest scores determine the potential candidate keywords. The number of extracted keywords is set by the user. The selection of keywords is done according to the following rules.

- If two N-grams have the same score, the longer one will be selected.
- If two candidate keywords have the same number of components and the same score, we select the higher degree.
- If an N-gram is selected, all the possible combinations of its components will be removed from the list of N-grams to guaranty that an extracted keyword will not be included in another one.

The list of keywords is then built by replacing each selected normalized N-gram by the most frequent of its surface N-gram in the original text. Therefore, the list of keywords that will be associated with the document will have more readable form.

6. Experiments and Evaluation

In order to evaluate the performance of the proposed system, an experiment was carried out to test it by comparing the extracted keywords against the manually assigned ones. A collection of 70 journal articles and article abstract selected from six journals and covering different domains was used. The dataset is divided into three groups according to their size [table 1]. The average number of words per article is 3406. Each one of these articles was assigned a list of keywords. The number of keywords varies from 2 to 14, with an average of 5.14 keywords per document. The number of extracted keywords is set to the same number of keywords assigned manually to the documents, so the number of false positive detections and false negative detections will be equal, and the three measures P, R, and F will be identical.

Table 1 shows the main results of the conducted experiment. An average precision of 0.51 was achieved. Since the primary analysis of the dataset showed that only about 73% of the human-generated keywords appear in the document texts, this result can be considered as a good result. The results have shown also that better results are achieved with larger documents.

Table 1: Results

| Dataset | Number of Documents | Average of words per article | Precision |
|---------|---------------------|------------------------------|-----------|
| 1 | 22 | 6523 | 0.56 |
| 2 | 28 | 3238 | 0.54 |
| 3 | 20 | 212 | 0.41 |
| All | 70 | 3406 | 0.51 |

7. Conclusion

This paper proposed an unsupervised two-stage approach for keyword extraction from Arabic texts that avoids the necessity of annotated data. The conducted experiments showed that the proposed method can extract keywords from single documents in a domain-independent way. The linguistic analysis of the texts and the grouping of N-grams according to their linguistic features improve the quality of extracted keywords. An average precision of 0.51 was achieved in despite the fact that that only about 73% of the human-assigned keywords appear in the document texts.

Reference

Al-Sughaier, I., Al-Kharashi, I. (2004). Arabic morphological analysis techniques: A comprehensive survey. *Journal of The American Society for Information Science and Technology (JASIST)*, 55(3), 189-213.

- Awajan, A. (2007). Arabic Text Preprocessing for the Natural Language Processing Applications. *Arab Gulf Journal of Scientific Research*, 25(4), 179-189.
- Benajiba, Y., Diab, M., Rosso, P. (2009). Arabic Named Entity Recognition: A Feature-Driven Study. *IEEE Transactions on Audio, Speech, and Language Processing*, 17(5), 926-934.
- Boudlal, A., Lakhouaja, A., Mazroui, A., Meziane, A., Ould Abdallahi, M., Shoul, M. (2010). Alkhalil Morpho Sys: A Morphosyntactic analysis system for Arabic texts, *International Arab Conference on Information Technology (ACIT)*. Riyadh, Saudi Arabia.
- Diab, M., Hacioglu, K., JURAFSKY, D. (2007). Automatic Processing of Modern Standard Arabic Text. *Chapter in Arabic Computational Morphology*. Springer Ed. 159-179.
- El-Beltagy S., & Rafea A. (2008). KP-Miner: A keyphrase extraction system for English and Arabic documents, *Information Systems*. 34(1), 132-144.
- El-Shishtawy, T., & Al-Sammak, A. (2009). Arabic Keyphrase Extraction using Linguistic knowledge and Machine Learning Techniques, In *Proceedings of the Second International Conference on Arabic Language Resources and Tools*, The MEDAR Consortium, Cairo, Egypt.
- Frank, E., Paynter, G.W., Witten, I.H., Gutwin, C., Nevill-Manning, C.G. (1999). Domain-Specific Keyphrase Extraction. *Proceedings of the International Joint Conference on Artificial Intelligence*, Stockholm, Sweden, 668–673.
- Green, S., and Manning, C. D. (2010). Better Arabic Parsing: Baselines, Evaluations, and Analysis. In *COLING- Beijing*. 394–402.
- Habash, N. (2010). *Introduction to Arabic Natural Language Processing*. Morgan & Claypool Publishers, USA.
- Hasan, K.S., NG, V. (2010). Conundrums in unsupervised Keyphrase Extraction: Making Sense of the State-of-the-Art. *COLING 2010*, 365-373.
- Hulth, A. (2003). Improved automatic keyword extraction given more linguistic knowledge. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, Sapporo, Japan,
- Hulth, A. (2004). *Combining Machine Learning and Natural Language Processing for Automatic Keyword Extraction*. Doctoral dissertation. Department of Computer and Systems Sciences, Stockholm University.
- Liu, Z., Li, P., Zheng, Y., Sun, M. (2009). Clustering to Find Exemplar Terms for Keyphrase Extraction. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, Singapore. 257–266.
- Manning, C. D., Raghavan, P., Schütze, H. (2009). *Introduction to Information Retrieval*. Cambridge University Press.
- Matsuo, Y., Ishizuka, M. (2004). Keyword Extraction from a Single Document using Word Co-occurrence Statistical Information. *International Journal on Artificial Intelligence Tools*, 13(1), 157-169

- Mihalcea, R., Tarau, P. (2004). TextRank: Brining order into texts. In Proceedings of *EMNLP 2004*, Association for Computational Linguistics, Barcelona, Spain. 404-411.
- Pudota, N., Dattolo, A., Baruzzo, A., Tasso, C. (2010). A New Domain Independent Keyphrase Extraction System. *Digital Libraries: Communications in Computer and Information Science*, 91, 67-78.
- Rose, S., Engel, D., Cramer, N., Cowley, W. (2010). Automatic keyword extraction from individual documents. *Text Mining: Applications and Theory* edited by Michael W. Berry and Jacob Kogan, John Wiley & Sons, Ltd. 3-20
- Turney, P. D. (2000). Learning Algorithm for Keyphrase Extraction. *Information Retrieval*, 2(4), 303-336.
- Turney, P. D., & Pantel, P. (2010). From Frequency to Meaning: Vector Space Models of Semantics. *Journal of Artificial Intelligence Research*, 37, 141-188.

Testing Distributional Hypothesis in Patent Translation

Hsin-Hung Lin* and Yves Lepage*

Abstract

This paper presents a wordlist-based lexical richness approach to testing distributional hypothesis for genre analysis in translation studies. In recent years, there has been continuing interest in patent translation. However, there are only a few lay their interests on comparison between native and non-native writing. The proposed approach to terms distribution of technical words contained in United States Patent and Trademark Office (USPTO) and Japan Patent Office (JPO) in terms of lexical variation, lexical density and lexical sophistication, in brief, highlights distributional similarity of technical genre, and in particular, distributional difference of academic and general genres.

Keywords: Patent Translation, Native Characterization, Corpus, Co-Occurrence.

1. Introduction

As globalization has resulted in rapid greater economic growth, the challenges of interdisciplinary interaction in pursuit of precise patent writing have incredibly increased.

In Lin and Hsieh (2010a), English patent documents were statistically extracted and computationally examined from LexisNexis Academic, a database for legal professionals. They compiled a reference corpus of independent claim texts and lay the focus on their collocation features. Mutual information is attainable with the help of selectional collocation features underlining specific clausal types represented in natural language processing of patent specification.

While their work appears to fill a niche in the ESP (English for Specific Purposes) field (and particularly in the English for Occupational Legal Purposes), Lin and Hsieh (2010b) further compiled a modern patent language technical term list with statistical-retrieval methodologies as a mandatory

* Graduate School of Information, Production, and Systems, Waseda University, Japan.

E-mail: nobuhiro602@toki.waseda.jp; yves.lepage@waseda.jp

approach. The research content and statistical investigations assist patent attorneys expand the vocabulary size for the advancement of patent writing at an international level.

Lin and Hsieh (2011) proposed a mixed-method approach to detecting scholarly discourse in patent technical documents. The Patent Technical Word Corpus (hereafter PTWC), containing 16 million word tokens, was compiled to elucidate the underpinning principles in identifying discourse elements, text-structure components, and the location of references. Whereas most existing IPR (intellectual property rights) databases accessible for information retrieval, the creation of PTWC, based on corpus-statistics and text-processing technology, refines more decisive characteristics of terminological knowledge as potential contribution for evaluation of technical documents.

To characterize technical genre in translation studies, we use lexical richness based on technical wordlist to test distributional hypothesis.

2. Technical Terms Distribution

We firstly conduct a quantitative survey based on USPTO Glossary to rank the distribution of technical terms used in United States Patent and Trademark Office (USPTO) and Japan Patent Office (JPO) within the time period from year 2010 to 2013. Table 1 below presents the statistical results.

‘Comprising’, a term of art used in claim language which means that the named elements are essential in describing the invention, ranked the first in USPTO. According to USPTO Glossary, it is a transitional phrase that is synonymous with "including," "containing" or "characterized by;" is inclusive or open-ended and does not exclude additional, unrecited elements or method steps. On the contrary, ‘consisting of’, a transitional phrase that is closed and excludes any element, step, or ingredient not specified in the claim, ranked the 6th.

To characterize transitional phrases of technical genre in translation studies, we retrieved co-occurring information of ‘comprising’ and ‘consisting of’ to compare it with academic and general genres.

Table 1. Distribution of patent technical words in USPTO

| Rank | Term | Frequency | Rank | Term | Frequency |
|------|-------------------------|-----------|------|-------------------|-----------|
| 1 | comprising | 3785213 | 11 | specification | 854667 |
| 2 | scope | 2459656 | 12 | continuation | 738785 |
| 3 | patent | 1603882 | 13 | dependent claim | 625886 |
| 4 | Group | 1306808 | 14 | composed of | 617353 |
| 5 | element | 1245265 | 15 | independent claim | 587926 |
| 6 | consisting of | 1165427 | 16 | representative | 518762 |
| 7 | drawing | 1015261 | 17 | benefit claim | 437599 |
| 8 | disclosure | 919881 | 18 | person | 383784 |
| 9 | application (patent) | 884470 | 19 | priority claim | 381352 |
| 10 | patent application | 884470 | 20 | interference | 341173 |

We give the survey of terms used in JPO in Table 2. It is noted that “comprising” ranked the first in distribution of USPTO and JPO, whereas “consisting of” ranked the 6th.

Table 2. Distribution of patent technical words in JPO

| Rank | Term | Frequency | Rank | Term | Frequency |
|------|-------------------------|-----------|------|-------------------|-----------|
| 1 | comprising | 629750 | 11 | applicant | 60293 |
| 2 | composed of | 371852 | 12 | drawing | 53469 |
| 3 | element | 272496 | 13 | person | 48893 |
| 4 | POWER | 272088 | 14 | IDS | 24946 |
| 5 | Group | 176103 | 15 | Control No. | 22905 |
| 6 | consisting of | 136992 | 16 | interference | 22445 |
| 7 | PAIR | 122746 | 17 | RE | 19777 |
| 8 | representative | 72519 | 18 | specification | 18102 |
| 9 | Request (PCT) | 70606 | 19 | classification | 16977 |
| 10 | application (patent) | 62027 | 20 | independent claim | 15513 |

3. Methodology

3.1 The Distributional Hypothesis

Sahlgren (2008:33) maintains that distributional approaches to meaning acquisition utilize distributional properties of linguistic entities as the building blocks of semantics. This hypothesis is often stated in terms like words which are similar in meaning occur in similar contexts (Rubenstein & Goodenough, 1965). In other words, words that occur in the same contexts tend to have similar meanings (Pantel, 2005).

3.2 Corpus Preparation

Transitional phrases in patent application were used to specify whether the claim is limited to only the elements listed, or whether the claim may cover items or processes that have additional elements. The most common transitional phrase used is the open-ended phrase "comprising". However, many claims use closed-ended language such as "consisting of".

In this regards, we retrieve co-occurring information containing "comprising" and "consisting of" from LexisNexis Academic for corpus preparation. Table 3 shows the structure for the corpus creation.

Table 3. Genre-based co-occurrence corpus of transitional phrases

| Genres | Native Writing | Non-Native Writing |
|------------------------|-------------------------|--------------------|
| Technical (Patent) | USPTO | JPO |
| Academic (Law Journal) | Canadian Legal Journals | HK Law Journal |
| General (Newspapers) | US Newspapers | Non-US Newspapers |

3.3 Lexical Richness

Lexical richness is a concept about one's lexical uses, which can be measured by lexical density, sophistication, and variation (Kao and Wang, 2014:54).

Kojima and Yamashita (2014:23) suggest that lexical richness measures primarily assess learners' vocabulary use. Lexical variation, the proportion between different words (types) and the total number (tokens) of words used in the text, is known as the type-token ratio (TTR).

Lexical density is defined as the percentage of lexical words in the text, for

example, nouns, verbs, adjectives, and adverbs (Laufer and Nation, 1995:309). Since only content words carry semantic meanings, a greater lexical density indicates more semantic information conveyed in a text.

Read (2000: 200) distinguishes dimensions of lexical richness, and one of these is lexical sophistication, which he defines as ‘the use of technical terms and jargon as well as the kind of uncommon words that allow writers to express their meanings in a precise and sophisticated manner’. The proportion of words used at different frequency levels, in terms of K1, K2, AWL (Academic Word List), and off-list words, in the text. K1 and K2 words are the most commonly used first 1000 and 1001 to 2000 words, respectively, in English. Words beyond these K1, K2, and AWL are placed into the off-list level, where proper nouns, rare words, special terms, acronyms, abbreviations, incompletions, and even misspellings may be found.

4. Results and Discussion

In terms of lexical density, non-natives employed more semantic information than the natives, among all genres. In terms of lexical variation, non-natives employed more lexical diversity than the natives in technical and academic genres.

Academic Genre, in particular, HK Law Journal, containing most semantic information (83%), among the all, whilst general genre, Non-Us Newspapers, containing least lexical diversity, as we excluded technical genre for analysis.

4.1 Technical Genre

In technical genre, in particular, JPO (Patent Abstract of Japan), containing least advanced words (15.03%) in the texts, among all.

Table 4. Lexical sophistication of “comprising” in technical genre

| Word Level (%) | USPTO | JPO |
|----------------|-------|-------|
| K1 Words | 50.35 | 50.37 |
| K2 Words | 2.61 | 3.61 |
| AWL Words | 23.74 | 21.78 |
| Off-List Words | 23.31 | 24.24 |

Less vocabulary knowledge in K2, AWL, and Off-list words were employed in “consisting of”, compared with that of “comprising”. The natives employed more academic words in Table 4, more off-list words in Table 5.

Table 5. Lexical sophistication of “consisting of” in technical genre

| Word Level (%) | USPTO | JPO |
|----------------|-------|-------|
| K1 Words | 62.37 | 64.89 |
| K2 Words | 0.29 | 0.34 |
| AWL Words | 19.43 | 19.74 |
| Off-List Words | 17.92 | 15.03 |

4.2 Academic Genre

In academic genre, HK Law Journal, containing most advanced words (33.28%) in the texts, among all.

Table 6. Lexical sophistication of “comprising” in academic genre

| Word Level (%) | Canadian Legal Journal | HK Law Journal |
|----------------|------------------------|----------------|
| K1 Words | 59.02 | 51.02 |
| K2 Words | 5.88 | 3.07 |
| AWL Words | 13.47 | 12.63 |
| Off-List Words | 21.63 | 33.28 |

As shown in Table 6 and Table 7, non-natives employed more off-list words in academic legal genre, whereas the natives employed more K1 and K2 words in academic legal genre.

Table 7. Lexical sophistication of “consisting of” in academic genre

| Word Level (%) | Canadian Legal Journal | HK Law Journal |
|----------------|------------------------|----------------|
| K1 Words | 59.74 | 50.74 |
| K2 Words | 6.12 | 2.97 |
| AWL Words | 13.07 | 14.24 |
| Off-List Words | 21.07 | 32.05 |

4.3 General Genre

As can be seen in Table 8 and Table 9, the non-natives employed more K2, AWL, and Off-list words but less K1 words in general genre.

Table 8. Lexical sophistication of “comprising” in general genre

| Word Level (%) | US Newspapers | Non-US Newspapers |
|----------------|---------------|-------------------|
| K1 Words | 63.34 | 51.97 |
| K2 Words | 3.69 | 5.47 |
| AWL Words | 11.66 | 16.73 |
| Off-List Words | 21.30 | 25.84 |

Table 9. Lexical sophistication of “consisting of” in general genre

| Word Level (%) | US Newspapers | Non-US Newspapers |
|----------------|---------------|-------------------|
| K1 Words | 63.37 | 53.48 |
| K2 Words | 4.03 | 7.71 |
| AWL Words | 12.61 | 16.42 |
| Off-List Words | 19.99 | 22.09 |

In short, K1 words were employed more by the natives in academic and general genres, whilst less used in technical genres.

5. Conclusion and Future Work

There is a correlation between distributional similarity and meaning similarity, which allows us to utilize the former in order to estimate the latter (Sahlgren, 2008:33). In terms of distribution statistics, the technical genre reveals more distributional and meaning similarity.

In summary, lexical richness is a valid and reliable measure to characterize genres. For future research, we seek to investigate the origin differences between syntagmatic and paradigmatic relations to further refine the preliminaries of the present study.

References

- Kao, S. M., & Wang, W. C. (2014). Lexical and organizational features in novice and experienced ELF presentations. *Journal of English as a Lingua Franca*, 3(1), 49-79.
- Kojima, M., & Yamashita, J. (2014). Reliability of lexical richness measures based on word lists in short second language productions. *System*, 42, 23-33.
- Laufer, B., & Nation, P. (1995). Vocabulary size and use: Lexical density in FL written production. *Applied Linguistics*, 16, 307-322.
- Lin, H. H., & Hsieh, C. Y. (2010a). Collocation Features of Independent Claim in US Patent Documents: Information Retrieval from LexisNexis. *ROCLING XXII: Conference on Computational Linguistics and Speech Processing*, pp. 296-310. Taiwan: Academia Sinica.
- Lin, H. H., & Hsieh, C. Y. (2010b). The Specialized Vocabulary of Modern Patent Language: Semantic Association in Patent Lexis. *PACLIC 24: Pacific Asia Conference on Language, Information, and Computation*, pp. 417-424. Japan: Waseda University Press.
- Lin, H. H., & Hsieh, C. Y. (2011). Characteristics of Independent Claim: A Corpus-Linguistic Approach to Contemporary English Patents. *International Journal of Computational Linguistics and Chinese Language Processing*, 16(3-4), 77-106.
- Pantel, P. (2005). Inducing ontological co-occurrence vectors. In *Proceedings of the 43rd Conference of the Association for Computational Linguistics*, pp. 125–132). Morristown, NJ, USA: Association for Computational Linguistics.
- Read, J. (2000). *Assessing vocabulary*. Cambridge: Cambridge University Press.
- Rubenstein, H., & Goodenough, J. (1965). Contextual correlates of synonymy. *Communications of the ACM*, 8 (10), 627-633.
- Sahlgren, M. (2008). The distributional hypothesis. *Italian Journal of Linguistics*, 20(1), 33-53.

Spectrum Analysis of Cry Sounds in Preterm and Full-Term Infants

Li-mei Chen¹, Yu-Hsuan Yang¹, Chyi-Her Lin², Yuh-Jyh Lin², Yung-Chieh Lin²

Abstract

Long-time average spectrum (LTAS) was used to analyze the cry phonations of 26 infants under four months old; 16 of them are full-term and the other 10 infants are preterm. The results of first spectral peak, mean spectral energy, spectral tilt, high frequency energy were used to compare the cry phonatory between full-term and preterm infants. In addition, cry duration and percent phonation is also compared. According to previous studies, full-term and preterm infants' crying behavior show significant differences because immature neurological development of preterm infants. Major findings in this study are: (1) There was no significant difference in unedited cry phonation across groups; (2) There was no significant difference in percent phonation across groups; (3) There was no significant difference in first spectral peak across groups, and no significant difference within groups could be found. However, full-term infants have higher first spectral peak than that of preterm infants; (4) There was no significant difference in mean spectral energy across groups, yet there was a significant main effect for partition; (5) There was no significant difference in spectral tilt across groups. Post hoc comparisons identified higher spectral tilt in P2 than in P3 in the full-term infants; (6) There was no significant difference in high frequency energy across groups. Significant differences were observed across partition, and in both groups, P1 had higher HFE than P3. The differences in the measures of crying behavior between full-term and preterm infants can help to estimate health condition of infants who are under 4 months old.

Keywords: Long-time average spectrum, Infant cry, Preterm infants

1. Introduction

Previous studies show that preterm infants are prone to immaturity of neurological development which leads to their sensitiveness toward pain stimulation, and the greater pain they suffer would reflect on crying behavior. If a set of distinctive measures can be identified, it might be possible to differentiate infant cries in the spectrum of normative behavior and cries due to organic pathology. The measures

¹ Department of Foreign Languages, National Cheng Kung University

² Department of Pediatrics, Medical College, National Cheng Kung University

would be helpful for doctors and caregivers to identify if the unknown cries are caused by just infant colic or other more complicated factors. Long-time average spectrum (LTAS) of crying behavior were analyzed in two groups of newborn infants in this study. LTAS was used to analyze the infant's non-partitioned crying episode (NP), as well as the 3 equal-length partitions (P1, P2, P3). First spectral peak (FSP), mean spectral energy (MSE), spectral tilt (ST), and high frequency energy (HFE) were measured.

Colic strikes infants who are under four months old, and it makes the infants cry in the evening on a daily bases or at the moment of waking up (Lester, Boukydis, Gracia-Coll, & Hole, 1990). The cause of this pain is still unknown (Zeskind & Barr, 1997). Colic occurs when infants are around one month old and it often disappears without a reason when infants are older than three months (Clifford, 2002). It is a universal and commonly-seen phenomenon which is the cause for excessive cry behavior. Though previous studies have suggested that higher fundamental frequency and a larger percentage of dysphonation in crying behavior can be found in the pain cries of infants who suffer from colic, no standard acoustic features in cry vocalization of infants with colic was established (Zeskind & Barr, 1997). Long-time average spectrum might provide an option to see if there are any significant variations in the cries of infants with colic and those who are healthy from first spectral peak, mean spectral energy, spectral tilt, and high frequency energy.

Though infants are not able to talk, they can express their feelings and emotions through crying, facial expression, and body movements. Diseases are able to be discovered by some characteristics in crying behavior (Radhika, Chandralingam, Anjaneyulu, & Satyanarayana, 2012). For example, different pain stimuli would lead to different fundamental frequencies in infant cry vocalization (Radhika et al., 2012). If more specific characteristics are found in certain diseases, it would be more effective in prescribing and curing. Sometimes parents can differentiate why their babies cry by their various crying behavior (Soltis, 2004). As for the way of eliciting cries, Johnston, Stevens, Craig, & Grunau (1993) had proposed two different ways: the heel-stick procedure and injection. In this current study, injection was used as the only standard method of eliciting cries to avoid any nuances that might caused by the different types of pain stimuli. However, even though there are some scientific ways of detecting the pain intensity infants endure, the experience of pain is quite subjective and is not merely related to physiological but also psychological factors (Qiu, 2006). Moreover, since infants use crying to arouse caregivers' attention, it can be expected that infants' crying behavior differs with and without their caregivers around them (Green, Gustafson, Irwin, Kalinowski, & Wood, 1995). Usually, the responses from caregivers bring cry behaviors to a halt (Green et al., 1995). Thus, crying is regarded not only an independent behavior but also plays an important role in social interactions between infants and their caretakers (Green et al., 1995). Furthermore, crying is a way of drawing other people's attention to help infants get rid of the uncomfortable situation or meet their needs (LaGasse, Neal, & Lester, 2005).

Because of the immature development of nervous systems caused by premature birth, preterm infants' crying behavior is believed to reveal different characteristics from that of full-term infants whose nervous systems are comparatively well-developed (Goberman & Robb, 1999). Premature infants are reported to have higher *f₀* in their cry phonation, and it might be due to the immature, shorter vocal folds (Johnston et al., 1993). Or as Zeskind (1983) stated that high-risk infants are not able to perfectly control their crying behavior and that they tend to react more intensely towards pain stimuli than low-risk infants. Infants react differently to the same stimulus pain whether they are healthy or born at risk. However, while some studies have shown that preterm infants are more sensitive to pain stimuli, others found that some premature infants have less intense reactions towards pain than normal infants (Qiu, 2006).

The main objective of this current study is to find out how the crying behavior between full-term and preterm infants differs from each other. The findings might help in detecting infants' health conditions. Moreover, if the differences of the crying behavior can be systematically characterized, the measurements can be further applied to identify features in neonate cries due to infant colic.

2. Method

2.1 Participants

Previous studies indicated that gender did not lead to significant differences in first spectral peak, mean spectral energy, spectral tilt, and high frequency energy (Goberman & Robb, 1999; Goberman, Johnson, Cannizzaro, & Robb, 2008). Therefore, gender was not controlled in this study. There were 26 infant participants; 16 were full-term infants and the other 10 were preterm infants. The infants were all under four months old for both full-term infants and preterm infants according to their gestational ages. All of the infants in this study were considered to have normal hearing according to interview with parents.

2.2 Data Collection

For data collection, TASCAM wave recorder and RODE uni-directional microphone were used while recording the cry phonation of both preterm and full-term infants. The microphone was held near the infants' mouth. All infants were in the supine position while receiving the injection because acoustic properties (e.g. fundamental frequency) might be influenced by postures of infants (Lin & Green, 2007). Data were collected in the hospital. The cry phonation of both groups of infants was recorded during and after they receive the injection. The pain stimulus was the same in both groups of infants.

2.3 Acoustic Analysis

Based on Goberman and Robb (1999), a crying episode of infants was defined as the duration of the continuous crying activity, beginning with the first audible cry sound after the pain stimulus, and an episode was completed as soon as the infants stopped crying. The non-voiced parts of a crying episode

were first edited out in the cry vocalization, making a “non-partitioned crying episode” (Goberman & Robb, 1999). In this current study, all the inspiratory cry was eliminated from pain stimulus, only the phonatory parts were analyzed. Then, a non-partitioned episode was divided into three partitions with the same length of durations (P1, P2, P3). P1, P2, P3 are regarded as the early, middle, and late sections of the crying episode, respectively, corresponding to the attack, cruise, and subdual phases of a crying episode as suggested by Truby and Lind (1965). First spectral peak, mean spectral energy, spectral tilt, and high frequency energy were measured. First spectral peak was identified as the first amplitude peak across the LTAS display. Mean spectral energy was measured with the mean amplitude value from 0 to 8000 Hz. Spectral tilt was the ratio of energy between 0-1000 Hz, and 1000-5000 Hz. High frequency energy was the sum of amplitudes from 5000 Hz to 8000 Hz.

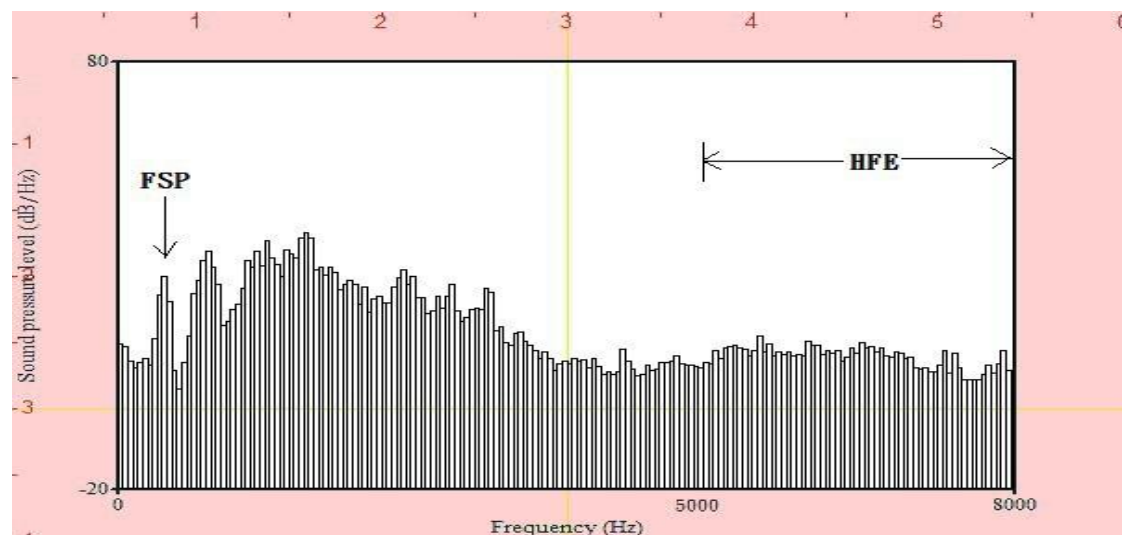


Figure 1. Typical LTAS display showing the location of the first spectral peak (FSP) and high frequency energy (HFE) between 5000Hz and 8000Hz.

3. Results & Discussion

3.1 Unedited Cry Duration

The average duration of crying episodes for the 16 full-term infants was 43.26s (SD=31.27), and for the 10 preterm infants was 36.21s (SD=30.93). Although full-term infants have longer average duration of crying episodes, there are no significant differences between the two groups, $t(24) = 0.48$, two-tailed, $p > .05$. The result is the same as that of Goberman and Robb (1999).

3.2 Percent Phonation

The percentage of cry phonation in a long-term non-partitioned, unedited crying episode was calculated. The average percent phonation across the crying episodes of the 16 full-term infants and the 10 preterm infants was 67% (SD=17.04) and 67% (SD=13.98) respectively. That is, 67% of the unedited crying

episode contained cry phonation. Like what was found in Goberman and Robb (1999), there is no significant differences across groups in the percentage of cry phonation, $t(24) = 0.39$, two-tailed, $p > .05$.

3.3 First Spectral Peak (FSP)

The non-partitioned and partitioned first spectral peak values of the 16 full-term and the 10 preterm infants are listed in Table 1. A two-way analysis of variance (ANOVA) was performed to check if there were significant differences in FSP values between the two groups, and whether there was significant variation between the three equal-length cry durations (P1, P2, P3) in each group. The results indicated no significant term by partition interaction ($p > .05$), no significant main effect for term status ($p > .05$), and no significant main effect for partition ($p > .05$). Despite the fact that there were no significant differences in statistical tests, we found that full-term infants have higher non-partitioned FSP (182.07Hz) than that of preterm infants (130.44Hz). From overall observation, full-term infants demonstrated higher FSP in non-partitioned and the three partitioned episodes. Full-term and preterm infants displayed different trends of FSP in P1, P2, and P3.

While the infants were receiving injections, the sharp pain stimulated them and all the infants burst out crying. According to the previous studies (Johnston et al., 1993), preterm infants are expected to have higher FSP because preterm infants were thought to be more sensitive and would react more intensely to pain. Intensive crying behavior causes the increase of the subglottal pressure and the stiffness of the vocal folds. However, in this current study, the mean FSP of the full-term infants turned out to be higher than that of the preterm infants, in both the non-partitioned episode and the three equal-length episodes. Moreover, full-term infants' crying episode involved more distinct phases with decrease of FSP in P3. In preterm infants, FSP kept increasing from P1 to P3.

Table 1. First spectral peak from the non-partitioned episodes (NP) and three partitioned crying episodes with equal length (P1, P2, P3) from the full-term and preterm groups of infants

| Group | FSP (Hz) | | | | |
|-----------|----------|--------|--------|--------|--------|
| | | NP | P1 | P2 | P3 |
| Full-term | M | 182.07 | 135.88 | 184.79 | 149.46 |
| | SD | 139.06 | 113.24 | 142.45 | 119.31 |
| Preterm | M | 130.44 | 104.35 | 117.40 | 139.14 |
| | SD | 71.74 | 52.06 | 67.36 | 82.11 |

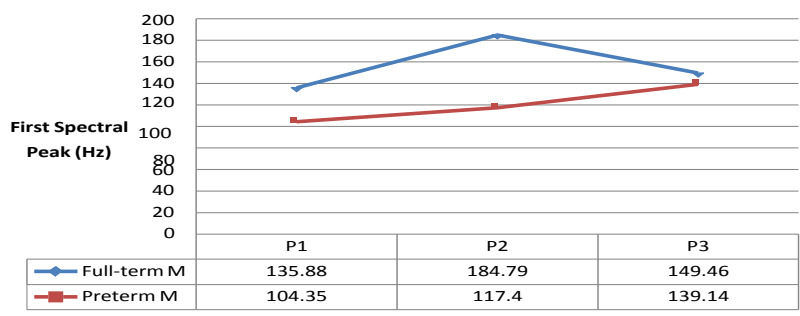


Figure 2. First spectral peak of full-term and preterm infants over time (P1, P2, and P3 are three equal-length partitioned crying episodes.)

3.4 Mean Spectral Energy (MSE)

The mean spectral energy of non-partitioned and partitioned episodes of the 16 full-term and the 10 preterm infants are compared and listed in Table 2. A two-way analysis of variance (ANOVA) was performed to investigate if there was significant variation between the two groups, and whether there was significant variation between the three equal-length cry durations (P1, P2, P3) in each group. No significant term by partition interaction ($p > .05$) was found. There was a significant main effect for partition ($F = 6.47$, $p < .05$), yet there was no significant main effect for term, $p > .05$. One-way ANOVA tests were then performed in each group to check the changes in MSE (in P1, P2, P3). In full-term infants, P2 was significantly higher than P3. In preterm infants, P1 showed significantly higher energy than P2 and P3. The changes are illustrated in Figure 3.

Premature infants, compared to full-term infants, are reported to have higher f_0 in their cry phonation, and it might be due to tension of the larynx (Johnston et al., 1993). Greater pain stimulus makes laryngeal muscles tighten. In this current study, although no significant differences could be identified, preterm infants showed higher mean of MSE in non-partitioned episode and the three equidurational crying episodes. This shows that during the cry duration, the preterm infants' laryngeal muscles are tighter and they have a more severe reaction toward pain stimulus. The tighter laryngeal muscles suggest a more intense cry behavior. Moreover, a decrease in MSE could be observed in both full-term and preterm infants. This might suggest that the laryngeal muscles of both groups of infants loosen by phase, especially in preterm infants. There is sharper decrease of MSE from P1 to P3 in preterm infants.

Table 2. Mean spectral energy from the non-partitioned episodes (NP) and three partitioned crying episodes with equal length (P1, P2, P3) from the full-term and preterm groups of infants

| | | MSE (dB) | | | |
|-----------|----|----------|--------|--------|--------|
| Group | | NP | P1 | P2 | P3 |
| Full-term | M | 19.368 | 19.982 | 19.507 | 14.323 |
| | SD | 9.627 | 9.523 | 11.158 | 11.266 |
| Preterm | M | 22.801 | 25.201 | 18.695 | 15.628 |
| | SD | 5.785 | 6.409 | 7.963 | 6.153 |

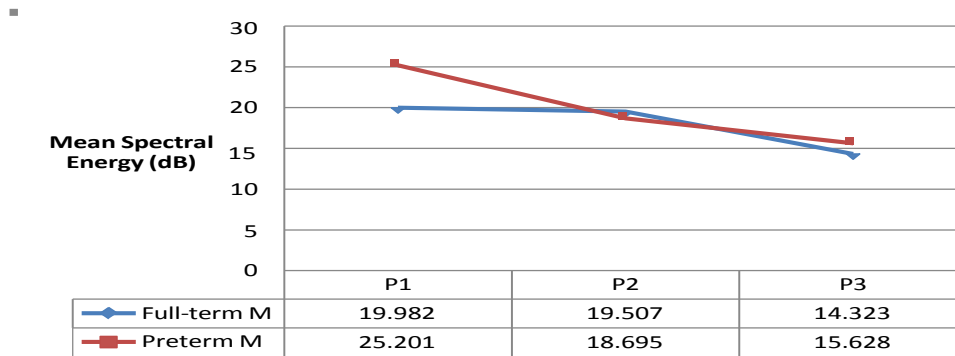


Figure 3. Mean spectral energy of full-term and preterm infants over time (P1, P2, and P3 are three equal-length partitioned crying episodes.)

3.5 Spectral Tilt (ST)

The spectral tilt of non-partitioned and partitioned crying episodes of the two groups are compared and listed in Table 3. In order to identify if there was significant differences of ST between the two groups and whether there was significant variation between the three equal-length cry durations (P1, P2, P3) in each group, a two-way analysis of variance (ANOVA) was performed. There were no significant term by partition interaction ($p > .05$), no significant main effect for partition, and no significant main effect for term. To investigate changes in ST across partitions within each group, separate one-way ANOVA tests were performed for full-term and preterm infant groups. In full-term infants, post hoc comparisons identified a significantly higher ST for P2 than for P3 ($p < .05$), but no significant differences in ST across partitions for the preterm infants ($p > .05$). Overall, the full-term infants showed higher ST values at the onset of crying vocalization which decreased over time; whereas the preterm infants had lower ST values at the onset, which increased over time. The changes are demonstrated in Figure 4. For full-term infants, higher P2 than P3 in ST was observed. The ST of full-term infants does not increase over time as mentioned in Goberman and Robb (1999); on the contrary, the ST of full-term infants

decreases over time. However, the increase of ST can be found in preterm infants.

The variation of ST of full-term and preterm infants is shown in table 3. Preterm infants have a slightly higher ST than full-term infants. ST refers to how quickly the amplitude of the harmonics decreases. A higher ST value was related to hypoadduction of the vocal folds (Mendoza, Munoz, & Naranjo, 1996). In this current study, hyperadduction is observed in the decrease ST of full-term infants, whereas, hypoadduction is observed in the increase ST of preterm infants.

Table 3. Spectral tilt from the non-partitioned episodes (NP) and three partitioned crying episodes with equal length (P1, P2, P3) from the full-term and preterm groups of infants

| | | ST | | | |
|-----------|----|-------|-------|-------|-------|
| Group | | NP | P1 | P2 | P3 |
| Full-term | M | 1.381 | 2.242 | 1.423 | 1.118 |
| | SD | 0.307 | 3.207 | 0.387 | 0.300 |
| Preterm | M | 1.839 | 1.935 | 2.811 | 3.218 |
| | SD | 0.685 | 0.659 | 2.795 | 5.326 |

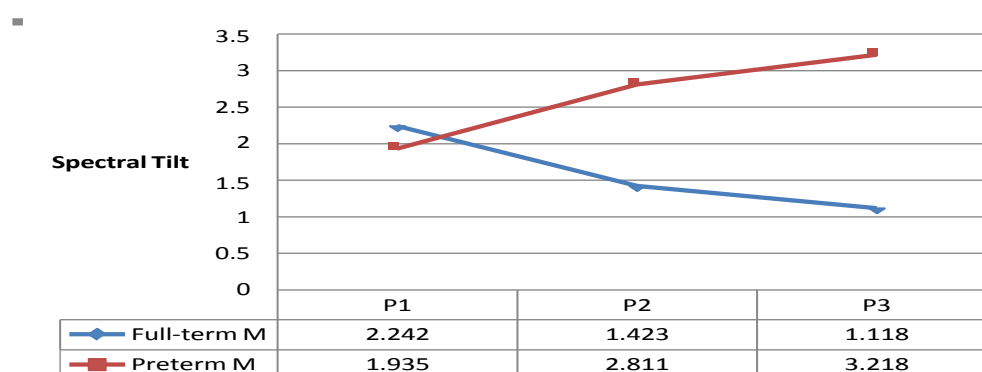


Figure 4. Spectral tilt of full-term and preterm infants over time (P1, P2, and P3 are three equal-length partitioned crying episodes.)

3.6 High Frequency Energy (HFE)

The high frequency energy of non-partitioned and partitioned crying episodes of the two groups are compared and listed in Table 4. In order to identify if there was significant variation between the two groups, and whether there was significant variation between the three equal-length cry durations (P1, P2, P3) in each group, a two-way analysis of variance (ANOVA) was performed. No significant term by partition interaction ($p > .05$) was found. There was no main effect for term ($p > .05$), stating that there were no significant differences in HFE across the two groups. There was significant main effect for partitions ($F = 8.29, p < .05$). Furthermore, one-way ANOVA tests were performed to check changes in HFE across partition within each group. Significant differences in HFE were found across

partition for both full-term infants ($F = 3.91, p < .05$) and for preterm infants ($F = 4.57, p < .05$). There was a significantly higher P1 in HFE than P3 in both infant groups ($p < .05$). In both groups, HFE decreases over time. The HFE of full-term infants does not change drastically over time; however, in preterm infants, the HFE shows a steep descent. The changes are demonstrated in Figure 5. HFE in both full-term infants and preterm infants decreases overtime. However, HFE of preterm infants has a wider range, crossing from 1225 to 1807, while the range HFE of full-term infants is about 300.

Table 4. High frequency energy from the non-partitioned episodes (NP) and three partitioned crying episodes with equal length (P1, P2, P3) from the full-term and preterm groups of infants

| | | HFE (dB) | | | |
|-----------|----|----------|------|------|------|
| Group | | NP | P1 | P2 | P3 |
| Full-term | M | 1672 | 1703 | 1543 | 1272 |
| | SD | 582 | 552 | 720 | 590 |
| Preterm | M | 1737 | 1807 | 1511 | 1227 |
| | SD | 469 | 514 | 546 | 509 |

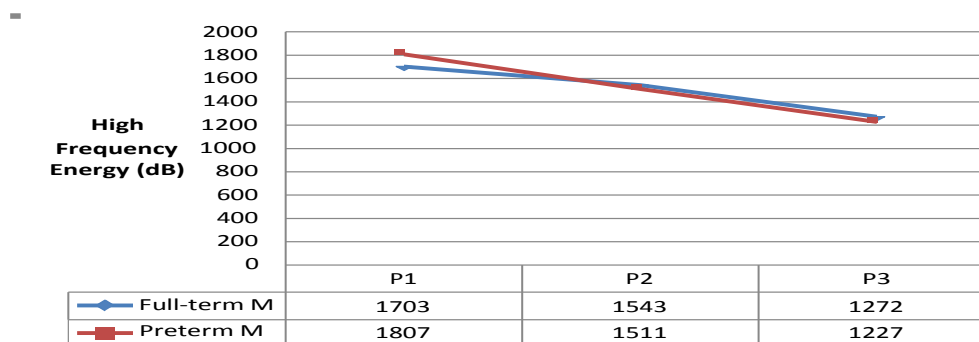


Figure 5. High frequency energy of full-term and preterm infants over time (P1, P2, and P3 are three equal-length partitioned crying episodes.)

Some of the results in this current study did not match the findings in previous studies. The differences could be due to a few discerning variables. First, although the uni-directional microphone was used in this study, the environmental noises could not be completely controlled because the nurses were required to explain the procedure to the caregivers. Moreover, there was unavoidable overlapping from noises of other infants' cry sound. Once the infant's cry vocalization was overlapped with adults' voice or other infants' cry sound, the partition could no longer be used for further analysis. Second, all the infants receiving injections had their caregivers around them. This caused inevitable interaction between adults and infants, bringing unexpected nuances to the results. Both full-term and preterm

infants might use more strength in crying, hoping their caretakers would alleviate their pain. Third, some caretakers tended to soothe the infants as soon as they started crying, which would significantly change the natural cry episode since the soothing and consolation from the caretakers might influence their cry phonation. The infants might feel safe and stopped crying. This might cause incomplete early, middle, and late sections in a cry episode, as Goberman and Robb (1999) mentioned. In further studies, the interaction of infants and their caregivers is probably one of the variables that should be strictly controlled. Moreover, video recording should be implemented in order to identify whether the infants stopped crying spontaneously or their attention was drawn by other things. Lacking complete three sections of cry episode might be the main reason of the discrepancy in the findings of this current study and previous studies.

4. Summary

Cry phonations of neonates from 16 full-term infants and 10 preterm infants were analyzed with long-time average spectrum (LTAS). Major findings were: (1) Full-term infants had higher first spectral peak than that of preterm infants; (2) Mean spectral energy in both groups of infants decreased over time; (3) Spectral tilt in full-term infants decreased, but increased in preterm infants over time; (4) High frequency energy in both full-term and preterm infants decreased over time. In further study, the environmental noise (e.g., from nurses, parents, and other infants around) should be controlled in order to acquire sufficient data to identify more systematic distinctions in the patterns of cry phonation between full-term and preterm infants.

Acknowledgements

This research was supported by a grant from National Science Council (NSC102-2410-H-006 060) and the Ministry of Education, Taiwan, R.O.C. The Aim for the Top University Project to the National Cheng Kung University. We thank the infant participants and their families for their time and cooperation.

References

- Clifford, T. (2002). Infant colic: A prospective, community-based Examination (Unpublished doctoral dissertation). The University of Western Ontario, Canada.
- Goberman, A. M., Johnson, S., Cannizzaro, M. S., & Robb, M. P. (2008). The effect of positioning on infant cries: Implications for sudden infant death syndrome. *International Journal of Pediatric Otorhinolaryngology*, *72*, 153-165.
- Goberman, A. M. & Robb, M. P. (1999). Acoustic examination of preterm and full-term infant cries: The long-time average spectrum. *Journal of Speech, Language, and Hearing Research*, *42*, 850-861

- Green, J. A., Gustafson, G. E., Irwin, J. R., Kalinowski, L. L., & Wood, R. M. (1995). Infant crying: Acoustics, perception, and communication. *Early Development and Parenting*, 4(4), 161-175.
- Johnston, C., Stevens, B., Craig, K., & Grunau, R. (1993). Developmental changes in pain expression in premature, full-term, two- and four-month-old infants. *Pain*, 52, 201-208.
- LaGasse L. L., Neal A. R., & Lester B. M. (2005). Assessment of infant cry: Acoustic cry analysis and parental perception. *Mental Retardation and Development Disabilites*, 11, 83-93.
- Lester, B., Boukydis, C., Gracia-Coll, C. & Hole, W. (1990). Colic for developmentalists. *Infant Mental Health Journal*, 11(4), 321-333.
- Lin, H. C., & Green, J. A. (2007). Effects of posture on newborn crying. *Infancy*, 11(2), 175-189.
- Mendoza, E., Munoz, J., & Naranjo, N. (1996). The longtime average spectrum as a measure of voice stability. *Folia Phoniatica*, 48, 57-64.
- Qiu, J. (2006). Does it hurt? *Nature*, 444, 143-145.
- Radhika, R. L., Chandralingam, S., Anjaneyulu, T. & Satyanarayana, K. (2012). A suggestive diagnostic technique for early identification of acyanotic heart disorders from infant's cry. *International Journal of Electrical and Electronics*, 1(3), 32-38.
- Soltis, J. (2004). The signal functions of early infant crying. *Behavioral and Brain Sciences*, 27, 443-458.
- Truby, H., & Lind, J. (1965). Cry motions of the newborn infant. In J. Lind (Ed.), *Acta Paediatrica Scandinavica: Newborn Infant Cry* (Suppl.163), 7-58.
- Zeskind, P. (1983). Production and spectral analysis of neonatal crying and its relation to other biobehavioral systems in the infant at-risk. In T. Field & A. Sostek (Eds.), *Infants born at-risk: Physiological and perceptual processes*. New York: Grune & Stratton.
- Zeskind, P. & Barr, R. (1997). Acoustic characteristics of naturally occurring cries of infants with "colic". *Child Language Development*, 68, 394-403.

Web-Based Recording and Visualization Framework for Moving Trajectories

Po-An Yang¹ LiJung Chi¹

Seth Chen² Jonathan Tsai²

Kun-Ta Chuang¹

Yung-Chung Ku²

¹Dept. of Computer Science and Information

²Innovative DigiTech-Enabled Applications &

Engineering

Services Institute, IDEAS

National Cheng Kung University

Institute for Information Industry

Tainan, Taiwan, R.O.C.

Taipei, Taiwan, R.O.C.

ktchuang@mail.ncku.edu.tw

{seth, jonathan, piperku}@iii.org.tw

Abstract

隨著智慧型手機的普及，個人軌跡資料的收集越來越方便，有關軌跡的應用也越來越多，像是 Foursquare [1]、Human [2]等 app 已陸續成為最火紅的手機應用。不過這些應用通常都只有很簡單的呈現打卡地點或軌跡，讓資料以靜態的出現，並沒有進一步的將資料視覺化，即使有，通常也要等到官方自行整合後釋出資料。我們提出的 Web-based 的系統架構，不僅可以即時動態呈現軌跡，更能跨平台，不受限於裝置環境，使用者可以透過瀏覽器來記錄自己的軌跡。軌跡儲存後，我們也提供一個即時呈現大家軌跡的移動動態呈現機制，不僅可以看見自己的軌跡，也可以看見他人的軌跡。另外有關於個人軌跡紀錄，通常使用者會有個人隱私的擔憂，在我們的架構中，會給每個使用者 UUID，透過 UUID 模糊了使用者與個人軌跡之間的連結，只有使用者的 UUID 及軌跡會留在 Database 中，系統不會知道使用者的 identity 資訊。最後我們也實作一些資料視覺化的技術，讓軌跡得以動態方式呈現。

Keywords: Trajectory, GPS, Data Visualization.

1. Introduction

有關於地圖軌跡的應用，已經是現在人不可缺少的一環，但多數的系統都偏向做導航或打卡等，鮮少有系統像世界迷霧 [3]一樣，讓使用者以地圖為日誌，記錄自己的軌跡。而即使有這些系統，通常這樣的應用程式都無法跨平台，包含 iOS 及不同的 android 等作業系統，均須特別撰寫不同的應用程式。在本文中，我們提出並實作了一個 Web-based 的軌跡擷取及呈現的系統架構，讓使用者可以簡單的透過瀏覽器或手機 app 記錄自己的

軌跡，並且以資料視覺化的技術將軌跡呈現。

我們實作的平台是 Github Pages [4]，透過 Github Pages，在此設計中，我們完全使用現在熱門的雲端服務，完全不用架設自己的 server，只要將 code 放到指定的 branch，使用者就能很方便的使用我們的服務，並且讓未來需要這個系統架構的開發者或研究員，可以很方便的 fork 我們的作品，增加他們需要的功能。

另外在記錄軌跡的部份，有些裝置沒有 GPS 定位的晶片，因此我們的系統支援 HTML5 的技術透過 wifi 或行動網路來定位；如果裝置有 GPS，我們的系統會用 GPS 來定位較精準的位置。有關軌跡資料的儲存，我們系統用 Firebase [5]當作資料庫。Firebase 是目前矽谷很熱門的雲端資料庫系統，Firebase 讓我們可以專注在前端的開發，透過 Firebase 簡便的 API，我們的系統可以很容易的存取使用者的軌跡資料。在軌跡呈現的部份，地圖是用 Open Street Map [6]，OSM 是一個 Crowdsourcing 平台，由於圖資較其他開放式地圖完整，所以我們系統選用 Open Street Map 當底層地圖；軌跡是用 Leaflet [7]，這是一套 Open-source JavaScript Library，除了有很方便的 API 可以在地圖上畫路徑，更可以讓開發者在他們的架構下，增加自己需要的套件，以疊圖層的方式，將資料呈現在地圖上，我們會有兩層的 layer，一層是靜態的路線，一層是動態的路線移動。

以下我們將分別介紹系統的整體架構以及實作的細節，最後也會探討一些未來的加值應用。

2. Framework Overview

在 Framework Overview 中會介紹我們系統架構的流程，如下圖(一)。我們全部的架構都是以 Web 為基礎，如此可以很容易的跨平台。

首先使用支援 HTML5 的瀏覽器或是手機 app 來取得使用者的位置，由於現今瀏覽器大多支援 HTML5，所以可以方便我們跨平台做定位。接下透過我們系統架構中的 GPS-Logger¹記錄使用者的軌跡，當使用者按下 Save 後，GPS-Logger 會把軌跡傳至 Firebase。將 Firebase 獨立出來，可以很清楚的切割我們的架構，透過 Firebase 連結前面的 GPS-Logger 收集軌跡資料，及後面的 Trajectory Visualization 軌跡呈現，Firebase 讓我們不需要任何伺服器編碼，就能將整套系統即時呈現。最後使用者進入 Trajectory Visualization，便能馬上看到自己動態的軌跡。

在 Part3 中，我們會介紹 GPS-Logger 如何運作，如何記錄使用者的軌跡。在 Part4 中，我

¹ GPS-Logger : <http://frankyang0529.github.io/gps-log/>

們會介紹資料儲存的格式為何，以及如何存入跟讀取軌跡。在 part5 中，我們會介紹 GPS-Path-View² 如何以新穎的方式，將軌跡資料呈現。在 Part6 中，我們提出一些未來展望，會介紹幾個有潛力的研究方向，除了前端的 Data Visualization 還有後端的 Data Mining，讓 Trajectory 的應用更豐富。最後在 Part7 中，我們會對我們的系統下結論，並且提出我們對未來 LBS 應用趨勢的看法。



圖(一) Our Framework

3. Implementation of GPS-Logger

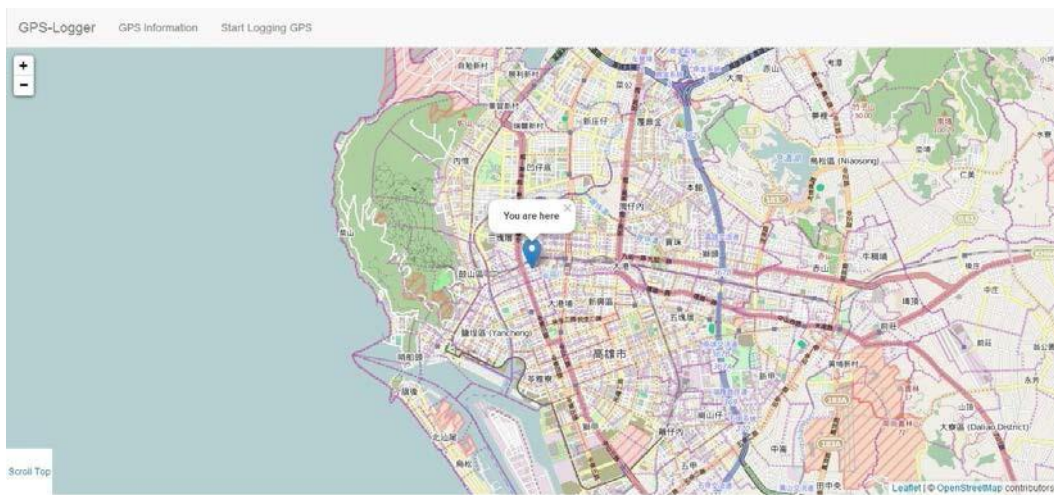
在 GPS-Logger 中，我們定位是使用 HTML5 Geolocation API [8]，這是由 World Wide Web Consortium(W3C) [9]所訂定的標準，由於這套 API 的支援度很廣泛，所以可以滿足我們跨平台的需求。定位前會先詢問使用者是否願意讓我們的系統存取位置資料，如果使用者同意，Geolocation API 會透過 IP address、Wi-Fi、Bluetooth MAC address、radio-frequency identification(RFID)、Global Position System(GPS)或 GSM/CDMA cell ID 去查詢使用者的位置，並且回傳精準度最高的位置資訊。目前多數瀏覽器都有支援這套 API，像是 Internet Explorer 9、Firefox 3.5、Chrome、Safari 5 及 Opera 10.6 以後的版本都有支援這套標準。Geolocation API 有很多細節可以操作，可以讓使用上更為靈活，像是有個參數是 enableHighAccuracy，讓開發者可以選擇是否取得較精準的位置，如果選擇 true，定位上會比較久，但位置資訊較為精準。

使用者進入我們系統架構的當下，系統就會透過 random 函數製造一個長度為 36 的 Universally Unique Identifier(UUID)，UUID 是通訊唯一識別碼，藉此來隱藏使用者的真實資訊，達到匿名的效果，當使用者重新開啟 GPS-Logger 時，都會重新分配 UUID，讓使用

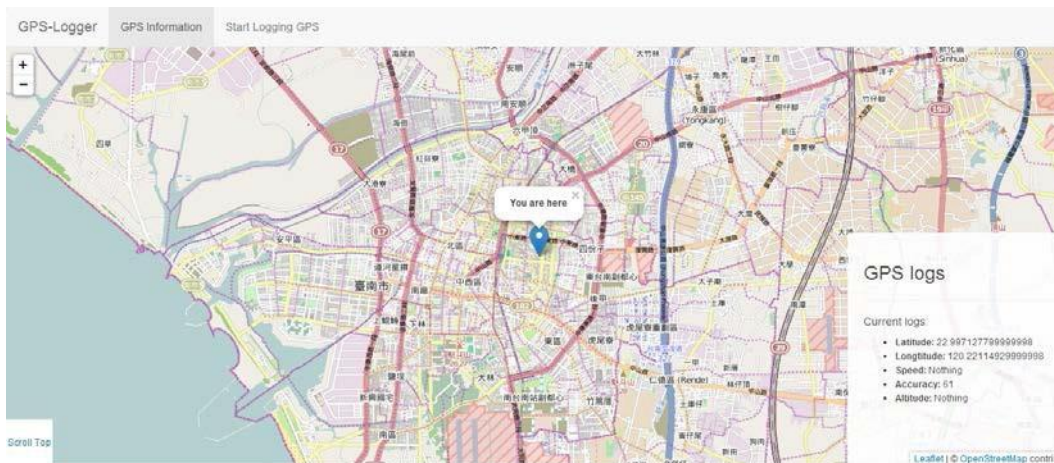
² GPS-Path-View : <http://chilijung.github.io/gps-path-view/>

者資訊更不容易被查詢。UUID 包含了 32 個十六進位數字，並且用 '-' 區分成 5 組，個數分別是 8, 4, 4, 4, 12，像是 523512e4-7b50-ee94-9475-f5dac6fe1cca，UUID 全部大約有 3.4×10^{38} 個，由於數量級龐大，所以有心人士並不容易透過我們的系統找到使用者的真實身份，即使每奈秒(10^{-9} 秒)生產一兆個 UUID，也要百億年才用的完。

當使用者取得 UUID 並且同意讓網站取得位置資料後，我們將會在地圖上標示出使用者的位置，如下圖(二)。如果點選了 GPS Information 則會顯示現在的經緯度、移動速度、位置誤差及高度資訊，其中經緯度的單位是 degree，位置誤差及高度是 meter，移動速度則是 meter/second。

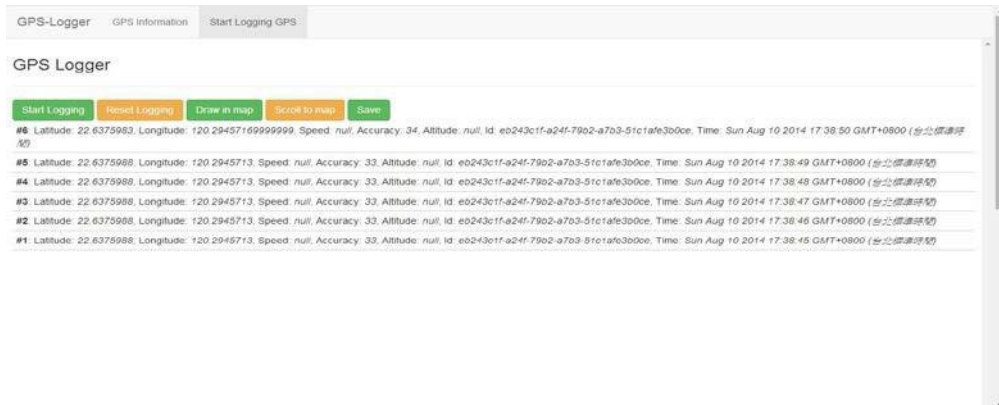


圖(二) GPS-Logger



圖(三) GPS-Information

當使用者按下 Start Logging GPS，會跳到另一頁面，如下圖四。按下 Start Logging 後，則會紀錄下每秒 GPS 的資訊，使用者也可以在頁面上看到自己的即時定位，此時資料都還沒上傳到 Database，所有記錄都還只是在本地端，當使用者按下 Save 後，資料才會上傳到我們的 Database；如果使用者不希望我們存取資料，也可以選擇 Draw in map，GPS-Logger 也會以動態的方式將資料呈現再本地端。

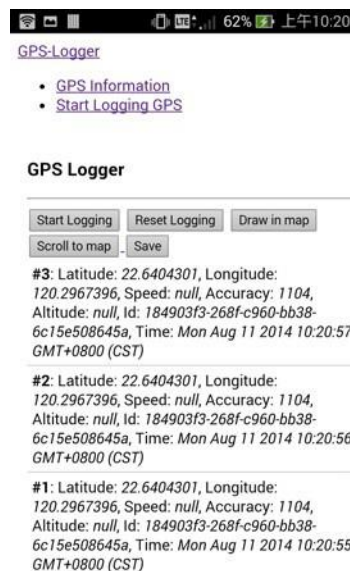


圖(四) Start Logging GPS

手機版的 app，我們是透過 phonegap [10]把 HTML、CSS 及 JavaScript 轉成 Android app，運作情況如下圖五及圖六。



圖(五) GPS-Logger app



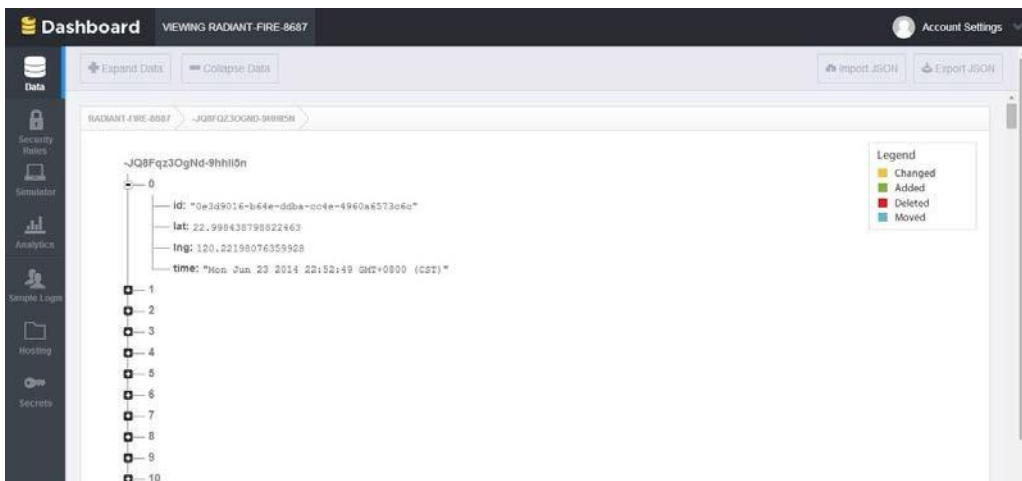
圖(六) GPS-Logger app Start Logging

4. Trajectory Database

Firebase 是目前矽谷最熱門雲端資料庫系統，有很多新興的即時(real-time)應用程式都以 Firebase 當資料庫，Firebase 讓開發者不用分心於 server 架設，就能同時容納數萬人上線，使開發人員可以專注於前端的開發，而不用擔心後端或是資料庫的狀況，即使是免費的方案，也能擁有 5GB 的流量及 100GB 的儲存空間，這項服務最好的地方就是，將資料存取和資料轉換在單一系統內合併，以前，資料傳送和存取原本是兩件分開的事，但現在兩者已結合成資料同步。

透過 Firebase，我們可以很方便的存取資料，GPS-Logger 每秒記錄一個點，當使用者按下 Save 後，我們會將本地端的資料都存入 Database，同時也將本地端的資料都清除，所以當使用者下次按下 Save 時，並不會存入一樣的資料，但是 UUID 是相同的，每次 Save 都會獨立成一條軌跡。

傳送到 Database 的資料是 array of json objects，json object 中包含了 id、lat、lng 及 time，id 是使用者的 UUID，lat 是緯度，lng 是經度，time 是時間資料，是由 javascript 的 Date object 轉成 string，當這些資料以 json object 儲存後，在往後的 Data Visualization 或是分析都可以很方便的擷取資料。當使用者每次儲存資料在 Firebase 時，Firebase 會亂數給一個編碼，像是-JQ8Fqz3OgNd-9hhli5n，而此編碼底下就是使用者存入的資料。在 Firebase 的網站，也可以很容易的察看資料庫中的資料，如下圖七。



圖(七) Firebase 資料庫查詢介面

5. Trajectory Visualization

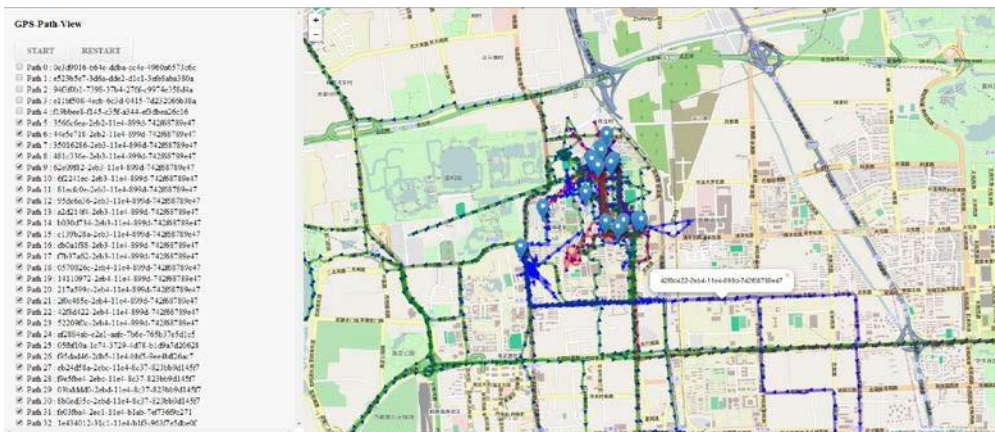
在 Data Visualization 中，我們使用的地圖是 Open Street Map，OSM 是目前圖資最豐富，

且更新最快的開放式地圖，顯示軌跡的圖層是 Leaflet，Leaflet 是很輕的 JavaScript Library，擁有豐富的畫圖 API，可以在地圖上畫 path、circle、polygon 等等，並且也有廣大的社群開發許多 Leaflet 套件，讓 Leaflet 功能更為廣泛，路徑上的箭頭與動畫即是 Leaflet 的插件，分別是 Leaflet.TextPath [11]及 Leaflet.AnimatedMarker [12]，點擊路徑時，會彈出一個對話框，告訴使用者路徑編號，而對話框也是 Leaflet 本身的 API。

在實作上，首先，我們會從 Firebase 取得全部資料，依序讀取 Firebase 隨機編號下的資料存入本地端的變數，再來將同一隨機編號底下的 Trajectory 集成成一條路徑，並加上包含 UUID 的對話框，最後會把路徑加上箭頭、做成動畫並產生 checkbox，讓使用者選擇是否把該路徑繪在自己的地圖上。Checkbox 的預設都是不勾選的，當使用者勾選某一條路徑時，系統會將地圖畫面移至該條路徑的起點。

在網頁的左上角有 START 跟 RESTART，按下 START，則地圖上的藍色標籤會開始沿著有勾選的路徑跑，當標籤在移動時，checkbox 會鎖起來，不能再新增或取消勾選，所以路徑無法從地圖上新增或移除；當使用者按下 RESTART，則會解鎖 checkbox，並將座標移回該條線的起點。

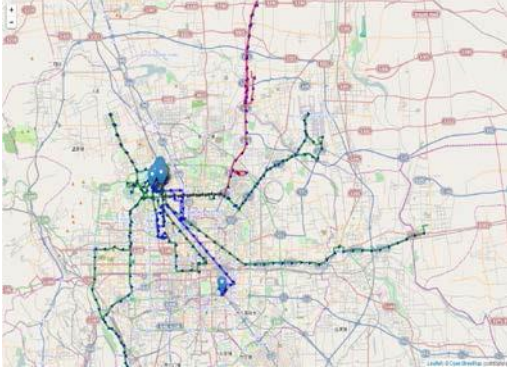
我們使用的 Trajectory 除了從 GPS-Logger 中收集的之外，還有從 Microsoft Research - GeoLife GPS Trajectories [13]取得軌跡，並且也能使用 My tracks³記錄軌跡，在透過我們的 kml 轉換器⁴，將軌跡上傳 Firebase。下圖八是 GPS-Path-View 的整個介面，下圖九跟十是部份 GeoLife GPS Trajectories 的 Data Visualization。



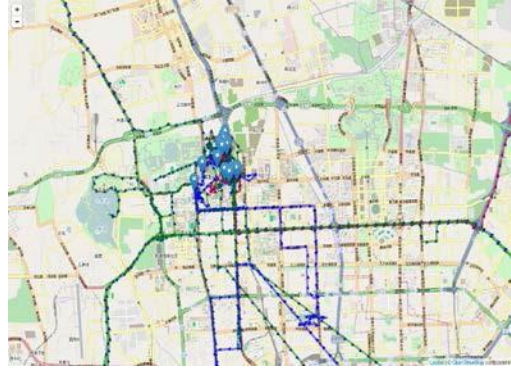
圖(八) GPS-Path-View with 170 Trajectories in Microsoft Research – GeoLife GPS Trajectories

³ My Tracks: <https://play.google.com/store/apps/details?id=com.google.android.maps.mytracks>

⁴ Kml 轉換器: <http://kml-parse.herokuapp.com/kmz/>



圖(九) GPS-Path-View Zoom out



圖(十) GPS-Path-View Zoom in

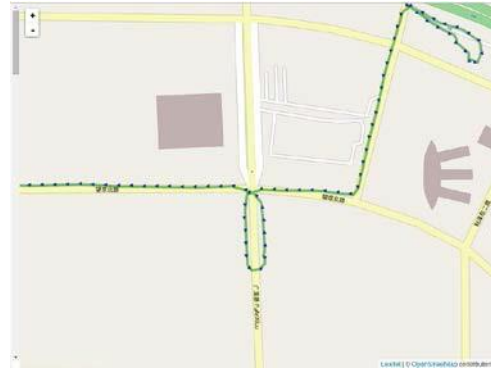
6. Future Work

使用者軌跡記錄及分析近年來已經普遍用於運動、健康等應用中。例如 NIKE+ RUNNING [14]，即為一個很熱門的運動路徑記錄 APP。然而目前這類軟體均為 APP-based 的架構而非 Web-based，在使用上多少會因為作業系統版本間的差異，而有系統更新的困難。我們認為，我們提出利用 Web-based 架構能克服平台差異造成的問題，並且我們使用的服務平台 Github Pages 及 Firebase 都是免費的，可以減輕開發者或研究員的負擔。後續利用我們所提出的軌跡記錄架構，除了能做為相關運動健康應用，我們正著手於以下應用上的開發。

首先，我們將針對 Open Street Map 的推廣跟便利度這方面來著手，加大 Open Source 平台的編輯便利度。OSM 是一個 Crowdsourcing 的平台，由於地圖的資料均由熱心的圖客們 (Mapper) 貢獻完成，如何讓圖客們很容易的了解某個地方可能有新的道路更新，包含道路新增，是一個很重要的課題。舉例來說，下圖十一是還沒加上路徑的 Open Street Map，下圖十二是標上路徑以後的 Open Street Map，我們發現本來地圖上沒有的路，但是卻有車子走過，表示這個地方可能有路，細察之後，我們發現這是北京市望京北路及廣澤路的交叉路口，配合比對 Google Map (下圖十三)，則會發現該地區路況已經變了，廣澤路在 Open Street Map 上只有一條道路，但在 Google Map 上卻有兩條道路，由此現象可以觀察到，未來我們將路徑資料與地圖資料做比對，當有一定數量的車子從原本地圖上沒有路的地方走過時，表示該地區有可能道路已經變更，如此可以增進 Open Street Map 改善的速度，讓 Open Street Map 的圖資更為完整。



圖(十一) Open Street Map



圖(十二) OSM with Trajectory



圖(十三) Google Map

此外，軌跡探勘技術已經被視為下一個階段的重要研究課題。研究上通常會透過評量軌跡之間的相似度來作為下一步應用的先備知識，然而要完成這項評量是非常困難的，因為軌跡的長度及空間上的關聯十分難用傳統的尤拉距離計算來評斷。我們正進行將軌跡轉換為使用者行動的 POI 的方式來提供研究人員進行 semantic 上的分析。例如一個使用者的軌跡，我們能轉換為”孔廟停留 30 分鐘”->”安平古堡停留 30 分鐘”。如此我們將更能提供研究技術人員提出更容易使用者個人化行銷或推荐等應用。此兩方面將是我們接下來著手的挑戰。

7. Conclusion

在本文中，我們提出了一個 Web-based 的軌跡擷取及資料即時呈現的雲端架構。透過現有的雲端服務，使用者可以不受限於手持式裝置的作業系統，均可便利的記錄其移動的細

節。我們認為，軌跡加值應用將是未來十分重要的一個 LBS 應用趨勢，透過資料的呈現及分析，可以提供個人化加值服務。我們也提出了一些未來的實作細節，包含將使用者的停留資訊擷取出來的的方法，以便後續的 data mining 分析過程，可以用 semantic 的角度來進行分析動作，我們預期能顯著的提高應用服務的精準度。

Acknowledgments

The authors especially thank the Taiwan Ministry of Economic Affairs and Institute for Information Industry for financially supporting this research : “Plan title : Fundamental Industrial Technology Development Program (2/4)”.

Reference

- [1] “Taipei | Food, Nightlife, Entertainment.” Available: <https://foursquare.com/>.
- [2] “Human - Activity Tracker: track walking, running, and biking.” Available: <http://human.co/>.
- [3] “世界迷霧 | 去實現你環遊世界的夢想吧!” Available: <http://zh-hant.fogofworld.com/>.
- [4] “GitHub Pages,” 15-Aug-2014. Available: <https://pages.github.com/>.
- [5] “Firebase - Build Realtime Apps.” Available: <https://www.firebase.com/>.
- [6] “OpenStreetMap.” Available: <http://www.openstreetmap.org/#map=5/51.509/-0.088>.
- [7] “Leaflet — an open-source JavaScript library for interactive maps.” Available: <http://leafletjs.com/>.
- [8] “Geolocation API Specification.” Available: <http://dev.w3.org/geo/api/spec-source.html>.
- [9] “World Wide Web Consortium (W3C).” Available: <http://www.w3.org/>.
- [10] “PhoneGap | Home.” Available: <http://phonegap.com/>.
- [11] “makinacorp/Leaflet.TextPath,” *GitHub*. Available: <https://github.com/makinacorp/Leaflet.TextPath>.
- [12] “openplans/Leaflet.AnimatedMarker,” *GitHub*. Available: <https://github.com/openplans/Leaflet.AnimatedMarker>.
- [13] “GeoLife GPS Trajectories.” Available: <http://research.microsoft.com/apps/mobile/download.aspx?p=b16d359d-d164-469e-9fd4-daa38f2b2e13>.
- [14] “Nike Running,” *Nike.com*. Available: http://www.nike.com/tw/zh_tw/c/running.



國立中央大學
National Central University