



Jacob
Lozano

안녕하세요



Language Profile

Southern Theory : Austronesian

- Fijian, Hawai'ian, Sunda, Tagalog ...1248 (311M)

Northern Theory : Altaic

- Japanese, Turkish, Uzbek ... 66 (250M)

National Language of North and South Korean

- Dialectal differences P'yöngyang & Seoul
 - China Japan, Russia, Thailand, Singapore ...
 - ~78 Million world wide (62 in Korea)

Original Scripts

- Hyangchal (향찰/鄉札), Gugyeol (구결/口訣) and Idu (이두/吏讀).

Current Scripts

- Invented 1444, promulgated 1446
 - Low status, women, children and uneducated
- 19th & 20th century Hanja and Hangeul popular up
- 1945 Chinese Characters insignificant
- N. 1949-1960 hanjaless
- S. debate of role, used in official & academic papers



Hangul (한글) Chosŏn'gŭl (조선글)

The Korean Alphabet System Including Diphthongs And Double Consonants

CONSONANTS

ㄱ	ㄴ	ㄷ	ㄹ	ㅁ	ㅂ	ㅅ	ㅇ	ㅈ	ㅊ
g, k	n	d, t	r, l	m	b, p	s	ng	j	ch

ㅋ	ㅌ	ㅍ	ㅎ	ㄱㄱ	ㄷㄷ	ㅍㅍ	ㅅㅅ	ㅈㅈ
k	t	p	h	kk	tt	pp	ss	jj

VOWELS

ㅏ	ㅑ	ㅓ	ㅕ	ㅗ	ㅛ	ㅜ	ㅠ	ㅡ	ㅣ
a	ya	eo	yeo	o	yo	u	yu	eu	i

ㅘ	ㅙ	ㅚ	ㅜ	ㅝ	ㅞ	ㅟ	ㅠ	ㅡ	ㅢ	
ae	yae	e	ye	wa	wae	oe	wo	we	wi	u

	g	n	d	r	m	b	s	o	l	ch	k	t	p	h
ae	gae	nae	dae	rae	mae	bae	saе	oe	jae	chae	kae	tae	pae	hae
yae	gyae	nyae	dyae	ryae	myae	byae	syae	yae	gyae	chyae	kyae	tyae	pyae	hyae
e	ge	ne	de	re	me	be	se	e	je	che	ke	te	pe	he
ye	gye	nye	dye	rye	mye	bye	syе	ye	gye	chye	kye	tye	pye	hye
wa	gwa	nwa	dwa	rwa	mwa	bwa	swa	wa	jwa	chwa	kwa	twa	pwa	hwa
oe	goe	noe	doe	roe	moe	boe	soe	oe	joe	choe	koе	toe	poe	hoe
wae	gwae	nwae	dwae	rwae	mwae	bwae	swae	wae	jwae	chwae	kwae	twae	pwae	hwae
wi	gwi	nwi	dwi	rwi	mwі	bwi	swi	wi	jwi	chwi	kwi	twi	pwi	hwi
wo	gwo	nwo	dwo	rwo	mwo	bwo	swo	wo	jwo	chwo	kwo	two	pwo	hwo
we	gwe	nwe	dwe	rwe	mwe	bwe	swe	we	jwe	chwe	kwe	twe	pwe	hwe
ui	gui	nui	dui	ruі	mui	bui	sui	ui	jui	chui	kui	tui	puі	hui

안녕

"Annyeong"

ng a n

안

n yeo ng

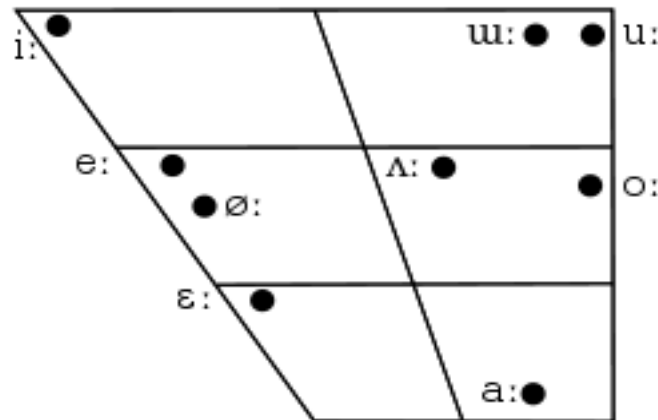
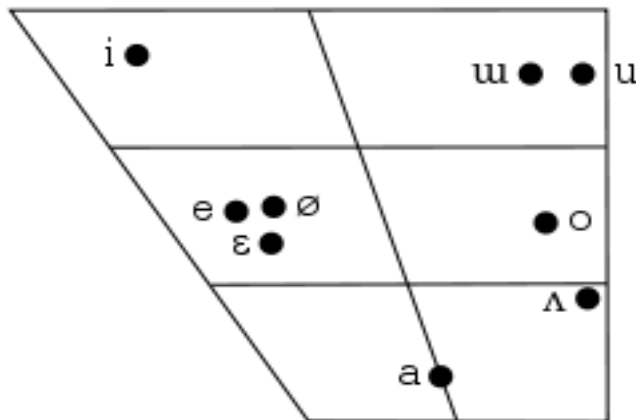
녕

Phonology

- Voiceless, unaspirated, aspirated, glottalized
- 100+ obsolete cons/vowels

		Bilabial	Alveolar	Post-alveolar	Velar	Glottal
Nasal		ㅁ /m/	ㄴ /n/		ㅇ /ŋ ^[19]	
Plosive and Affricate	plain	ㅍ /p/	ㅌ /t/	ㅊ /t͡ɕ/	ㅋ /k/	
	tense	ㅍㅍ /p̚/	ㅌㅌ /t̚/	ㅊㅊ /t͡ɕ̚/	ㅋㅋ /k̚/	
	aspirated	ㅍㅍ /pʰ/	ㅌㅌ /tʰ/	ㅊㅊ /t͡ɕʰ/	ㅋㅋ /kʰ/	
Fricative	plain		ㅅ /sʰ/			ㅎ /h/
	tense		ㅅㅅ /s̚/			
Liquid		/w/ ¹	ㄹ /l/	ㄹᵻ/ ¹		

Hangul	ㅁ	ㄴ	ㅊ	ㅋ	ㅍㅍ	ㅌㅌ	ㅊㅊ	ㅋㅋ	ㅍ	ㅌ	ㅊ	ㅅ	ㅎ	ㅅㅅ	ㅁ	ㄴ	ㅇ	ㄹ	
RR	<i>b</i>	<i>d</i>	<i>j</i>	<i>g</i>	<i>pp</i>	<i>tt</i>	<i>jj</i>	<i>kk</i>	<i>p</i>	<i>t</i>	<i>ch</i>	<i>k</i>	<i>s</i>	<i>h</i>	<i>ss</i>	<i>m</i>	<i>n</i>	<i>ng</i>	<i>r, l</i>
IPA	p	t	t͡ɕ	k	p̚	t̚	t͡ɕ̚	k̚	pʰ	tʰ	t͡ɕʰ	kʰ	s	h	s̚	m	n	ŋ	r, l



Hangul	ㅣ	ㅚ	ㅘ	ㅙ	ㅛ	ㅜ	ㅠ	ㅡ	ㅝ	ㅞ	ㅟ	ㅠ	ㅡ	ㅢ	ㅣ	ㅤ	ㅥ	ㅦ	ㅧ		
RR	<i>i</i>	<i>e</i>	<i>oe</i>	<i>ae</i>	<i>a</i>	<i>o</i>	<i>u</i>	<i>eo</i>	<i>eu</i>	<i>ui</i>	<i>ye</i>	<i>yae</i>	<i>ya</i>	<i>yo</i>	<i>yu</i>	<i>yeo</i>	<i>wi</i>	<i>we</i>	<i>wae</i>	<i>wa</i>	<i>wo</i>
IPA	i	e	ø	ɛ	a	o	u	ɯ	ɥi	je	je	ja	jo	ju	jɯ	ɥi	we	wɛ	wa	wɯ	

Morphology

Sentence : SOV

Nouns & Pronouns:

- no articles, gender/number marking
- Post-particles
 - Nominative, genitive, dative, accusative
 - Instrumental, locative, comitative
- Numeral Classifiers jang (장)
- I 'this', ku 'that', ce 'over there'

Verbs:

- No num/gen agree but polite
- Inflected base + ending (400+)
- Finite : honorific, tense, aspect, modal, formal, mood
- Tense
 - Marked: completed action
 - Unmarked: present
- Passive and Causative, +suffix
- Mood:
 - Declarative, interrogative, imperative, cohortative. Final of finite

a. 가 - 었 - 닷

ka-ess-ta

GO-past-decl

b. 가 - 시 - 었 - 닷

ka-si-ess-ta

GO-honorific-past-decl

c. 가 - 기 - 가

ka-ki-ka

GO-nominalizer-nom

d. 가 - 시 - 었 - 기 - 예 - 는

ka-si-ess-ki-ey-nun

GO-honorific-past-nominalizer-TO-topic

	Declarative	Ka-pni-ta	He is going
interrogative		ka-pni-k'a	Is he going
imperative		Ka-la	Go.!
coharative		Ka-ca	Let's go

MT Work : Morpheme tagger - spelling

Penn Tree Bank 54k

(3) 권한을 누가 갖고 있지?

kwenhan-ul-nwu-ka-kacko-iss-ci

AUTHORITY-acc WHO-nom HAVE BE-int

'Who has the authority?'

(4) (S (NP-OBJ-1 권한/NNC+을/PCA)
 (S (NP-SBJ 누구/NPN+가/PCA)
 (VP (VP (NP-OBJ *T*-1)
 가지/VV+고/EAU)
 있/VX+지/EFN))
 ?/SFN)

a. 동 해/NN ⇒ 동 해
tonghay tonghay
 'east sea'
 b. 동 해/VV ⇒ 동 학 어
tonghay tonghae
 'move'

Table II. Evaluation of spelling recovery

Process	Accuracy (%)
POS tagging for spelling recovery	92.03
Spelling recovery	96.33
Spelling recovery assuming 100% POS tagging	98.62

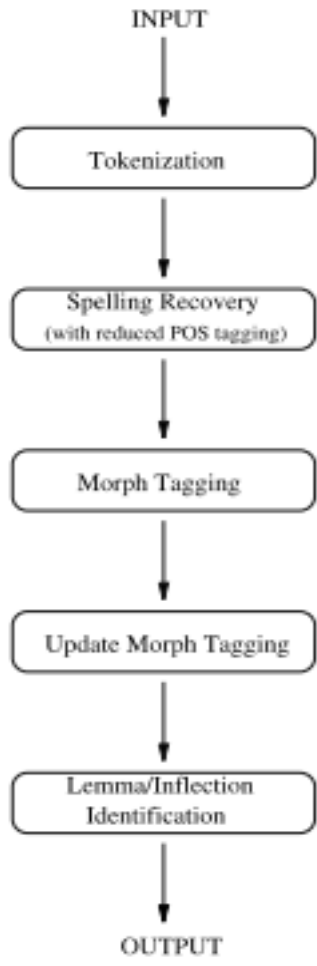


Figure 1. Overview of the tagger.

MT Work: Morpheme Tagger – 2/3 morph tags

165 complex tags 81%

Table III. Examples of complex tags

NNC
NNC+PCA
NNC+CO+EFN+PAD
NPR+PAD+PAU
NNU+PAU
VJ+EPF+EAN
VV+EPF+EFN
VV+EPF+EFN+PAU
VX+EPF+EFN+PAU

Few tag meaning

- NNC common
- CO copula
- EFN type
- PAD adverbial postposition
- EPF tense

594 templates 93%

Table IV. Examples of inflectional templates

있습 니 닷 <i>esssupnita</i>	VV+EPF+EFN
있 으 니 까 <i>essunikka</i>	VV+EPF+ECS
있 으 며 <i>essumye</i>	VV+EPF+ECS
있 으 므 로 <i>essumulo</i>	VV+EPF+ECS
있 음 에 <i>essumey</i>	VV+EPF+ENM+PAD
까 지 는 <i>kkacinun</i>	NNC+PAD+PAU
대 도 만 <i>tayloman</i>	NNC+PAD+PAU
로 부 터 의 <i>lopwuteuy</i>	NNC+PAD+PCA

a. Input:

제가 관측 사항을 보고 하였 습 니 닷 .
Cey-ka kwanchuk sahang-ul pokoha-ess-supnita.
 I-nom OBSERVATION ITEM-acc REPORT-past-decl
 ‘I reported the observation items.’

b. Output:

제 가 /NPN+PCA 관 측 /NNC 사 항 을 /NNC+PCA
 보 고 하 였 습 니 닷 /VV+EPF+EFN /SFN

a. Input:

제 가 /NPN+PCA 관 측 /NNC 사 항 을 /NNC+PCA
 잊 었 습 니 닷 /NNC
Cey-ka kwanchuk sahang-ul ic-ess-supnita.
 I-nom OBSERVATION ITEM-acc FORGET-past-decl
 ‘I forgot the observation items.’

b. Output:

제 가 /NPN+PCA 관 측 /NNC 사 항 을 /NNC+PCA
 잊 었 습 니 닷 /VV+EPF+EFN /SFN

MT Work: Morpheme Tagger – 4 Lemma/inflection identification

Table V. Example entries from inflection and stem dictionaries

Inflection dictionary		Stem dictionary	
ㅍ시타 (<i>psita</i>)	EFN	관객 (<i>kankyek</i>)	NN
ㅍ시오 (<i>psio</i>)	EFN	관단하 (<i>kantanha</i>)	VJ
가 (<i>ka</i>)	PCA	관략히 (<i>kanlyakhi</i>)	ADV
곶 (<i>keyss</i>)	EPF	관선 (<i>kansen</i>)	NNC
타 (<i>ta</i>)	EFN	관섭하 (<i>kansepha</i>)	VV
포타 (<i>pota</i>)	PAD	관섭 (<i>kansep</i>)	NNC
았 (<i>ass</i>)	EPF	관주하 (<i>kancwuha</i>)	VV
은 (<i>un</i>)	PAU	갈 (<i>kal</i>)	VV

(13) $abcd/p + x + y + z$ Longest, until pos



$$a/p + b/x + c/y + d/z$$

(14) a. Input:

제가/NPN+PCA 관측/NNC 사항/NNC+PCA
 보고/았/습/니/타/VV+EPF+EFN /SFN

b. Output:

제/NPN+가/PCA 관측/NNC 사항/NNC+을/PCA
 보고/았/VV+었/EPF+습/니/타/EFN /SFN

MT Work: Morpheme Tagger – Error and conclusion

Table VIII. Summary of error analysis

Type of error	Number	(%)	Gold	Test
Tagging	174	63	정직악 / VJ+지 / BCS <i>cengcikha-ci</i> HONEST-connective ‘be honest’	정직악 / VV+지 / BCS <i>cengikha-ci</i>
Spelling recovery	57	21	피우 / VV+리 / EAN <i>piwu-l</i> EMINATE-adnominal ‘which will emanate’	피울 / NNC <i>piwul</i>
Ambiguity	16	6	부대 / NNC+요 / PAD <i>pwutay-lo</i> SQUAD-TO ‘to the squad’	부 / NNC+태요 / PAD <i>pwu-taylo</i>
Missing NNC	15	5	학요 / NNC <i>hakto</i> CADET ‘cadet’	학 / NNC+요 / PAU <i>hak-to</i>
Missing template	11	4	훈련 / NNC+이 / CO+군 / EPN <i>hwunlyun-i-kwun</i> DRILL-cop-decl ‘be a drill’	훈련이군 / NNC <i>hwunlyunikwun</i>
Missing inflection	4	1	소대급 / NNC+에까지 / PAD <i>sotaykup-eykkaci</i> PLATOON LEVEL-TOWARDS ‘towards the platoon level’	소대급에 / NNC+까지 / PAD <i>sotaykupe-kkaci</i>

Table VI. Evaluation of the morphological tagger

Method	Precision (%)	Recall (%)
Treebank-trained	95.43	95.04
Treebank-trained ⁺	95.78	95.39

Table VII. Comparison with other morphological taggers

Reference	Accuracy (%)
Chan et al. (1998)	97.0
Yoon et al. (1999)	94.7
Lim et al. (1997)	94.8

MT Scores

Corpora

- Penn Tree Bank 34 files
- NIST 20 files
- SWRC KAIST corpus – korean-english-chinese 60k sent
- Speech ocean – korean-english-chinese-japanese 200k pairs

Some difficulties

- Borrowed words 트랙또르 ttŭrakttorŭ Ru. трактор (traktor)
트랙터 teuraekteo En. tractor tractor
- Difference in vocab 강냉이 kangnaeng-i 옥수수 oksusu corn

MT SCORE

English true

		BLEU-4 (mteval-v13a)			IBM BLEU (bleu-1.04)			NIST (mteval-v13a)			TER (tercom-0.7.25)			METEOR (meteor-0.7)		
SiteID	System	Overall	NW	WB	Overall	NW	WB	Overall	NW	WB	Overall	NW	WB	Overall	NW	WB
CMU	CMU_kor2eng_primary_un	0.1113	0.1172	0.1049	0.1110	0.1171	0.1044	5.4745	5.5402	5.1113	0.7856	0.7824	0.7890	0.3974	0.4122	0.3810
KUNLPL	KUNLPL_kor2eng_primary_un	0.1118	0.1173	0.1052	0.1119	0.1174	0.1054	5.3482	5.1900	5.2084	0.9031	0.9126	0.8929	0.4397	0.4434	0.4356
SAIC	SAIC_kor2eng_primary_un	0.0943	0.0962	0.0923	0.0942	0.0962	0.0921	5.2579	5.1684	5.0722	0.8606	0.8629	0.8582	0.4252	0.4315	0.4183
UVA	UVA_kor2eng_primary_un	0.0679	0.0665	0.0693	0.0679	0.0665	0.0693	2.2916	2.3803	2.0755	0.7438	0.7520	0.7351	0.2915	0.3022	0.2796

Korean true

		BLEU-4 (mteval-v13a)			IBM BLEU (bleu-1.04)			NIST (mteval-v13a)			TER (tercom-0.7.25)			METEOR (meteor-0.7)		
SiteID	System	Overall	NW	WB	Overall	NW	WB	Overall	NW	WB	Overall	NW	WB	Overall	NW	WB
CMU	CMU_kor2eng_primary_un	0.1047	0.1124	0.0963	0.1044	0.1124	0.0957	5.2021	5.3566	4.7619	0.7982	0.7881	0.8089	0.3815	0.3985	0.3626
KUNLPL	KUNLPL_kor2eng_primary_un	0.1014	0.1076	0.0943	0.1015	0.1077	0.0944	5.0712	5.0206	4.8328	0.9182	0.9208	0.9154	0.4189	0.4313	0.4052
SAIC	SAIC_kor2eng_primary_un	0.0881	0.0905	0.0854	0.0878	0.0902	0.0851	4.9513	4.9804	4.6584	0.8780	0.8720	0.8844	0.4049	0.4157	0.3929
UVA	UVA_kor2eng_primary_un	0.0608	0.0609	0.0605	0.0609	0.0610	0.0606	2.1462	2.2334	1.9402	0.7525	0.7602	0.7442	0.2775	0.2894	0.2643

Resources

Penn Korean Tree Bank

<http://www.newdesign.aclweb.org/anthology/Y/Y02/Y02-1007.pdf>

Morpheme tagger

<http://www.sfu.ca/~chunghye/papers/mt18-4-2.pdf>

KAIST Corpus

http://semanticweb.kaist.ac.kr/home/index.php/KAIST_Corpus

Speech Ocean

<http://www.speechocean.com/en-Text-Corpora/696.html>

Ethnologue

<http://www.ethnologue.com/language/kor>

Ominglot

<http://www.omniglot.com/writing/korean.htm>

About world languages

<http://aboutworldlanguages.com/korean>

감사합니다

gamsahabnida