

A Python Toolkit for Universal Transliteration

Ting Qian¹, Kristy Hollingshead², Su-youn Yoon³,
Kyoung-young Kim⁴, Richard Sproat⁵

University of Rochester¹, OHSU², ETS³, UIUC⁴, OHSU⁵
ting.qian@rochester.edu¹, hollingk@cslu.ogi.edu², syoon9@gmail.com³,
kkim36@illinois.edu⁴, rws@xoba.com⁵

LREC, Malta

May 21, 2010

Transliteration Examples from the Web

- कर्णम गुरु Guru गुरु गुरु गुरु गुरु गुरु
- 카메라 camera камеры Камер Camera Камера カメラ
- мелия Melia Μελία Мелиа મેલિયા Мелия मेलिया
- Хотэль хотел Хотел Хоутел Hotel отель
- Рома роми 口ウマ Roma Ρώμα ΡΩΜΗ
- Ванила ヴァニラ 𑂣𑂗𑂢𑂰 Ванили Vanilla Ванилия
バニラ वनिला
- Карма カルマ Karma 카르마 Кармы
- 브라운 Brown БРАУН Брайан Браун
- ๓๐๓๐ ๓๐๓๐ Тайм TIME
- सुनिल सुनिल सुनिल Sunil सुनिल

Basic Issues

- Cooccurrence - e.g. temporal correlation:
 - In parallel/comparable corpora we expect related concepts/terms to have similar distributions over space and time
- Edit distance:
 - Phonetic similarity
 - Graphical similarity
- Our goal: techniques for extracting plausible transliteration candidates for comparable corpora in n-tuples of languages that use different scripts.

Previous Work

- Transliteration: Knight & Graehl 1998; Meng et al. 2001; Gao et al. 2004; inter alia.
- Comparable corpora: Fung, 1995; Rapp 1995; Tanaka and Iwasaki, 1996; Franz et al., 1998; Ballesteros and Croft, 1998; Masuichi et al., 2000; Sadat et al., 2003; Tao and Zhai, 2005.
- Mining transliterations from multilingual web pages: Zhang & Vines, 2004
- Sproat, Tao & Zhai, ACL 2006:
 - Trained phonetic distance, similarity in temporal distribution

Previous Work

- Klementiev and Roth:
 - Discriminative model using letter n-gram features, and temporal distribution
- Tao et al, EMNLP 2006:
 - Untrained phonetic model and temporal distribution
- Yoon, Kim and Sproat, ACL 2007:
 - Untrained vs. discriminatively trained phonetic models
 - Unitran: Provides pronunciations for scripts in Basic Multilingual Plane
 - Hand-built phonetic model uses phonetic features as well as “pseudofeatures” derived from second-language learner errors
- Recent NEWS 2009 workshop (colocated with ACL in Singapore) highlighted a number of approaches to transliteration

Web Transliterations using Untran/Handbuilt Distance Model

- Find patterns of form $x_i x_{i+1} x_{i+2} \dots (y_i y_{i+1} y_{i+2} \dots)$ where at least some of $y_i y_{i+1} y_{i+2}$ are in a script different from $x_i x_{i+1} x_{i+2}$
- Use Untran to guess pronunciations for most strings:
- Festival for “English”
- Special tables for:
 - Chinese (Mandarin)
 - Kanji (kunyomi)
 - Extended Latin-1
- Rank by (untrained) phonetic edit distance

Web Transliterations using Unitran/Handbuilt Distance Model

ID	Score	Original	Cyrillic	Latin	Transliteration 1	Count 1	Transliteration 2	Count 2
1	1.50	2ni6k	2ni6ko	2pidk	2pidks	1	2	5 6
2	1.50	Adishi	адиш	A:diSi:	AdiSiS	1	23	5 6 6
3	1.50	Ardon	Ардон	A:rd>n	Ard>n_	1	3	5 6 6
4	1.50	Arquette	Аркетт	A:rKEt	Arkiett	1	3	5 6
5	1.50	Baby	бэйбэй	beibi:	beibi_	1	33	5 6
6	1.50	Bendz	Бендз	bEndz	biendz_	1	3	5 6 6
7	1.50	Brixlegg	Бриксlegt	brikslEg	briksliegg	1	3	8 9
8	1.50	Cartel	картель	kA:rTEl	kArtiel_	1	3	6 7
9	1.50	Casey	Кэйсэй	kelsi:	keisi_	1	33	5 6 6
10	1.50	Cepil	Чепиль	tSEpil	tSjiepil_	1	3	5 6
11	1.50	Christie	Кристис	krIsti:	kristis	1	3	6 7
12	1.50	Chubar	Чубар	tSubAr:	tSjubAr_	1	3	5 6 6
13	1.50	Crouzille	Крузиль	kruzil:	kruzil_	1	3	6 7
14	1.50	DIPOL	ДИПОЛЬ	dipol	dip>l_	1	3	5 6
15	1.50	Divigel	Дивигель	divigel	divigiel_	1	3	7 8
16	1.50	Effiel	Эффиль	effiel	EffiEl_	1	3	6 7
17	1.50	Elugel	Елюгель	elugel	ielugiel_	1	3	6 7
18	1.50	Endel	Эндель	endel	Endiel_	1	3	5 6
19	1.50	Filip	Филипп	fillp	filipp	1	3	5 6 6
20	1.50	Filipp	Филипп	fillp	filipp	1	3	5 6 6
21	1.50	Fitil	Фитиль	fitil:	fitil_	1	3	5 6 6
22	1.50	Football	футболь	fUtb>l	futb>l_	1	3	6 7
23	1.50	Fridrih	Фридрих	fri:dri:	fridrix	1	3	6 7
24	1.50	Gizelle	Гизель	glzEl	giziel_	1	3	5 6 6
25	1.50	Goygol	Гёйгёль	gojgol	giojgiol_	1	3	6 7
26	1.50	Isparih	Исперих	IspEri:	ispierix	1	3	6 7
27	1.50	Kamen	Камень	kA:mEn	kAmien_	1	3	5 6 6
28	1.50	Kamin	Каминь	kA:mi:n	kAmin_	1	3	5 6 6
29	1.50	Kamin	каминь	kA:mi:n	kAmin_	1	3	5 6 6

Web Transliterations using Unitran/Handbuilt Distance Model

Rank	Score	Original	Transliteration	Unitran	Handbuilt	Distance
1	1.50	うずまき	渦巻き	uzumAki	uzumAki_	32	35	7 8
2	1.65	관타나모	关塔那摩	kwanthanamo	kwanthanamo	24	35	10 10
3	1.69	민주당	民主党	mintSutaN	mintsrutaN	24	35	8 8
4	1.94	Setomaru	瀬戸丸	setomaru	setomAru	1	35	8 8
5	1.94	Toriimae	鳥居前	toriimae	toriimAe	1	35	8 8
6	2.00	NUMANOI	沼野井	numanoi	numAnoi	1	35	7 7
7	2.00	Takashimaya	高島屋	tAkAsimAja	tAkAsimAja	1	35	10 10
8	2.00	Yaneura	屋根裏	janeura	janeura	1	35	7 7
9	2.00	共產黨	공산당	koNtshraNtaN	koNshantaN	35	24	9 9
10	2.00	피용자	被用者	phijONtsa	phijONtsr&	24	35	7 7
11	2.06	Mikazuki	三日月	mlkA:zuki:	mikAtsuki	1	35	8 8
12	2.11	wenyanwen	文言文	wEnj@nw&n	w&njanw&n	1	35	9 9
13	2.12	Osezaki	大瀬崎	oUsEzA:ki:	oosesAki	1	35	8 8
14	2.12	Renminbi	人民币	rEnmlnbi:	r&nminpi	1	35	8 8
15	2.12	Renminbi	人民币	rEnmlnbi:	r&nminpi	1	35	8 8
16	2.12	renminbi	人民币	rEnmlnbi:	r&nminpi	1	35	8 8
17	2.12	renminbi	人民币	rEnmlnbi:	r&nminpi	1	35	8 8
18	2.20	ITALY	意大利	It&li:	itali	1	35	5 5
19	2.20	Italy	意大利	It&li:	itali	1	35	5 5
20	2.21	TaijiTu	太极图	taijitu	thaiocithu	1	35	7 7
21	2.21	Tulufan	吐鲁番	tuluf@n	thuluphan	1	35	7 7
22	2.21	동지사	同志社	toNtSisha	thoNtsrier&	24	35	7 7
23	2.31	mikazuki	三日月	mikazuki	mikAtsuki	1	35	8 8
24	2.33	Yaotouwan	摇头丸	jaotouwan	jauTh&uwan	1	35	9 9
25	2.33	巴基斯坦	파키스탄	paocisithan	phakhiSh4than	35	24	9 9
26	2.38	Mikadzuki	三日月	mlk@duki:	mikAtsuki	1	35	8 8
27	2.42	可 因	코카인	kh&khain	khokhain	35	24	6 6
28	2.43	BISHAMON	毘沙門	bIS&m&n	phisram&n	1	35	7 7
29	2.43	Bishamon	毘沙門	bIS&m&n	phisram&n	1	35	7 7

Web Transliterations using Unitran/Handbuilt Distance Model

66	2.80	Fritz	فريتس	frItS	frjts	1	6		5	5
67	2.80	MRirt	مريت	mrirt	mrjrt	1	6		5	5
68	2.80	Wilms	ويلمز	wjlmz	wilmz	6	1		5	5
69	2.80	Clint	كلينت	kljnt	klInt	6	1		5	5
70	2.81	zilatuun	زلاتوون	DlEtWUn	zilatuun	6	1		8	8
71	2.88	bird	بورڊ	b&rd	bord	1	6		4	4
72	2.88	black	بلاك	bl@k	bl&k	1	6		4	4
73	2.88	mustt	مست	m^st	mEst	1	6		4	4
74	2.88	Merge	مرج	mIrdZ	m&rdZ	6	1		4	4
75	2.90	shogol	شغول	S>g&l	SUGa~l	1	6		5	5
76	2.90	Muslimiyna	مُسلمين	mUslimIjnE	m^zliMIIn&	6	1		10	10
77	2.92	FRIENDS	فريينڊز	frEndz	frjndz	1	6		6	6
78	2.92	Friends	فريينڊز	frEndz	frjndz	1	6		6	6
79	2.92	Surendar	سرينڊر	srEndr	srjndr	1	6		6	6
80	2.92	Friends	فريينڊز	frjndz	frEndz	6	1		6	6
81	2.94	Välämäki	فيليمهكي	vhljmhkj	V_lim_ki	6	1		8	8
82	3.00	Gül	گُل	G_l	g_l	1	6		3	3

Web Transliterations using Unitran/Handbuilt Distance Model

1	2.25	YPL	gilisi	kilisi	gilisi	26	1		6	6
2	2.33	YAS	yonega	jonekA	jonega	26	1		6	6
3	2.56	TWFDZ	italiano	ithAliAno	ItAliAnoU	26	1		8	9
4	2.67	Kituwah	YSG	kItSuw&	kituWA	1	26		6	6
5	2.83	WHF	தமிழ்	thAmili	t{AmilA_	26	14		6	7
6	2.93	SPLI	dulisdi	tulisti	dul&sdi	26	1		7	7
7	3.08	SDH	Suomi	suomi	suoUmi	26	1		5	6
8	3.33	Kanuga	QNS	k@nug&	khAnukA	1	26		6	6
9	3.75	Keetoowah	YSG	kitu&	kituWA	1	26		5	6
10	3.88	adatasti	DLWQI	@d&t@sti	AtAthAsthi	1	26		8	8
11	4.04	dideloquasdi	ISQEI	dId&loUkw@zdi	titelokwAsti	1	26		13	11
12	4.04	ISQEI	dideloquasdi	titelokwAsti	dId&loUkw@zdi	26	1		11	13
13	4.08	limon	QNH	lim&n	lemAni	1	26		5	6
14	4.19	QWJ	unulahi	unulAhi	&njulAhi	26	1		7	8
15	4.28	ISLI	diyohidi	tijohiti	dIai&hidi	26	1		8	9
16	4.36	ugista	QKQW	uDZlst&	utsitsAthA	1	26		6	7
17	4.42	Cherokee	GWY	tSEr&ki	tsAlAki	1	26		6	6
18	4.42	QWW	Català	khAthAlA	Catal_	26	1		6	6
19	4.42	GWY	Cherokee	tsAlAki	tSEr&ki	26	1		6	6
20	4.44	QNSQ	nudadequa	nutAtekWA	n&dAdikw&	26	1		8	9

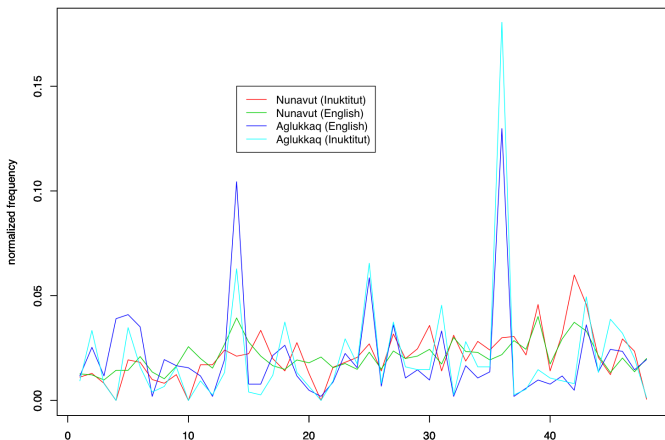
Web Transliterations using Untran/Handbuilt Distance Model

27	2.75	amaruq	◁Lᵖᵃᵃ	A:mA:r^k	AmAruq	1	27		6
28	2.75	nuliaq	ᵃᵃ◁ᵃᵃ	nuli@k	nuliAq	1	27		6
29	2.75	pualuuk	>◁ᵃᵃ	puA:l&k	puAlu:k	1	27		6
30	2.86	imiqtuq	ΔΓᵃᵃᵃᵃ	ImIkt^k	imiqtuq	1	27		7
31	2.90	Susan	ᵃᵃᵃ	suz&n	su:sAn	1	27		5
32	2.94	Qupanuq	ᵃᵃ◁ᵃᵃᵃᵃ	kupanuak	qupAnuAq	1	27		8
33	2.94	Sanirajak	ᵃᵃᵃᵃᵃᵃ	s&nIrA:dZ&k	sAnirAjAk	1	27		9
34	2.94	qikturiaq	ᵃᵃᵃᵃᵃᵃᵃᵃ	kIktSUrI@k	qikturiAq	1	27		9
35	3.00	Igluligaarjuk	Δᵃᵃᵃᵃᵃᵃᵃᵃ	Igl^IlgA:rdZ&k	iGlulIGArjuk	1	27		12
36	3.00	Inuktitut	Δᵃᵃᵃᵃᵃᵃᵃ	In^kt&t^t	inuktitut	1	27		9
37	3.00	arnaq	◁ᵃᵃᵃᵃ	A:rn@k	Arnaq	1	27		5
38	3.00	inuktitut	Δᵃᵃᵃᵃᵃᵃᵃᵃ	In^kt&t^t	inuktitut	1	27		9
39	3.00	panik	<ᵃᵃᵃᵃ	p@nIk	pAniq	1	27		5
40	3.00	qaniq	ᵃᵃᵃᵃᵃᵃ	kA:nIk	qAniq	1	27		5
41	3.00	sukak	ᵃᵃᵃᵃ	suk&k	sukAq	1	27		5

Temporal correlation: Nunavut Hansards

Nunavut
Aglukkaq

ᐅ ᐃ ᐅ ᐅ^c
ᐅ ᐅ ᐅ ᐅ ᐅ ᐅ^b



Synopsis

- ① Given comparable corpora, such as newswire text, in a pair of languages that use different scripts:
 - ScriptTranscriber provides an easy way to mine transliterations from comparable texts.
 - Particularly useful for underresourced languages
- ② ScriptTranscriber is an open source package that allows for ready incorporation of more sophisticated modules
- ③ Available as part of the nltk_contrib source tree at <http://code.google.com/p/nltk/>

Overview

- Approx. 7,500 lines of object-oriented Python
- Requires PySNoW
- Modules:
 - Document structure and XML representation.
 - Extractor: extracts terms from text. Specializations:
 - Capitalization-based extractor
 - Chinese foreign name extractor
 - Chinese personal name extractor
 - Thai extractor
 - Morph analyzer
 - Pronouncer. Specializations:
 - Unitran — UTF-8 pronouncer
 - English pronouncer
 - Hanzi (Chinese character) pronouncer
 - Comparator. Specializations:
 - Hand-built phonetic comparator
 - Time correlation comparator
 - Perceptron-based comparator

XML Fragment

埃及总统穆巴拉克、叙利亚总统阿萨德和沙特阿拉伯国王法赫德
28日和29日在埃及亚历山大市举行首脑会议。

Egyptian President Hosni Mubarak, Syrian president Hafez al-Assad and
King Fahd of Saudi Arabia held a meeting in the northern Egyptian port
city of Alexandria just before the end of last year.

```
<?xml version="1.0" encoding="UTF-8"?>
<doclist>
  <doc>
    <lang id="zh">
      ...
      <token count="1" morphs=""
        prons="sr a th & a m p ; a l a p o ;
        s u n A D U M g u m A k u d A k u D U M">沙特阿拉伯</token>
      <token count="1" morphs=""
        prons="f a x & a m p ; t & a m p ; ; n o r i A k a i o s i e">法赫德</token>
      <token count="1" morphs=""
        prons="m u p a l a kh & a m p ; ;
        j a w A r A g u d o m o e k u d A k u g A t s u">穆巴拉克</token>
      <token count="1" morphs="" prons="a s a t & a m p ;
        ; k u m A D U M o s i e">阿萨德</token>
    </lang>
    <lang id="en">
      <token count="1" morphs="" prons="@ l & a m p ; g z @ n d r i : & a m p ;">Alexandria</token>
      <token count="1" morphs="" prons="& a m p ; r e I b i : & a m p ;">Arabia</token>
      <token count="1" morphs="" prons="& a m p ; s A : d">Assad</token>
      <token count="1" morphs="" prons="I d Z I p S & a m p ; n">Egyptian</token>
      <token count="1" morphs="" prons="f A : d">Fahd</token>
      ...
      <token count="1" morphs="" prons="m u b A : r I k">Mubarak</token>
      ...
      <token count="1" morphs="" prons="s & g t ; d i :">Saudi</token>
      ...
    </lang>
  </doc>
</doclist>
```

Sample Program

```
#!/bin/env python
# -*- coding: utf-8 -*-

"""Sample transcription extractor based on the LCTL Thai parallel
data. Also tests Thai prons and alignment.
"""

__author__ = """
rws@uiuc.edu (Richard Sproat)
"""

import sys
import os
import documents
import tokens
import token_comp
import extractor
import thai_extractor
import pronouncer
from __init__ import BASE_

## A sample of 10,000 from each:

ENGLISH_      = '%s/testdata/thai_test_eng.txt' % BASE_
THAI_         = '%s/testdata/thai_test_thai.txt' % BASE_
XML_FILE_     = '%s/testdata/thai_test.xml' % BASE_
MATCH_FILE_   = '%s/testdata/thai_test.matches' % BASE_
```


Sample Program

```
BAD_COST_      = 6.0

def LoadData():
    t_extr = thai_extractor.ThaiExtractor()
    e_extr = extractor.NameExtractor()
    doclist = documents.Doclist()
    doc = documents.Doc()
    doclist.AddDoc(doc)
    ##### Thai
    lang = tokens.Lang()
    lang.SetId('th')
    doc.AddLang(lang)
    t_extr.FileExtract(THAI_)
    lang.SetTokens(t_extr.Tokens())
    lang.CompactTokens()
    for t in lang.Tokens():
        pronouncer_ = pronouncer.UnitranPronouncer(t)
        pronouncer_.Pronounce()
    ##### English
    lang = tokens.Lang()
    lang.SetId('en')
    doc.AddLang(lang)
    e_extr.FileExtract(ENGLISH_)
    lang.SetTokens(e_extr.Tokens())
    lang.CompactTokens()
    for t in lang.Tokens():
        pronouncer_ = pronouncer.EnglishPronouncer(t)
```

Sample Program

```
    pronouncer_.Pronounce()
return doclist

def ComputePhoneMatches(doclist):
    matches = {}
    for doc in doclist.Docs():
        lang1 = doc.Langs()[0]
        lang2 = doc.Langs()[1]
        for t1 in lang1.Tokens():
            hash1 = t1.EncodeForHash()
            for t2 in lang2.Tokens():
                hash2 = t2.EncodeForHash()
                try: result = matches[(hash1, hash2)] ## don't re-calc
                    except KeyError:
                        comparator = token_comp.OldPhoneticDistanceComparator(t1, t2)
                        comparator.ComputeDistance()
                        result = comparator.ComparisonResult()
                        matches[(hash1, hash2)] = result
    values = matches.values()
    values.sort(lambda x, y: cmp(x.Cost(), y.Cost()))
    p = open(MATCH_FILE_, 'w') ## zero out the file
    p.close()
    for v in values:
        if v.Cost() > BAD_COST_: break
        v.Print(MATCH_FILE_, 'a')
```

Sample Program

```
if __name__ == '__main__':  
    doclist = LoadData()  
    doclist.XmlDump(XML_FILE_, utf8 = True)  
    ComputePhoneMatches(doclist)
```

Interactive Use

```
>>> import pronouncer
>>> import tokens
>>> t1 = tokens.Token('WWJD')
>>> t2 = tokens.Token('拉拉瓜')
>>> p = pronouncer.UnitranPronouncer(t1)
>>> p.Pronounce()
>>> t1
#<WWJD 1 [] ['l A l A k u A'] >
>>> p = pronouncer.HanziPronouncer(t2)
>>> p.Pronounce()
>>> t2
#<拉拉瓜 1 [] ['l a l a k w a', 'k u d A k u k u d A k u u r i'] >
>>> import token_comp
>>> c = token_comp.OldPhoneticDistanceComparator(t1, t2)
>>> c.ComputeDistance()
>>> c.ComparisonResult()
#<comparator: WWJD <-> 拉拉瓜, 3.2857, "l A l A k u A <-> l a l a k w a">
>>> c.ComparisonResult().Cost()
3.2857142857142856
```

Summary

- ScriptTranscriber is a toolkit for extracting transliteration pairs from comparable corpora.
 - Works with any script in the Unicode Basic Multilingual Plane
 - Easy to extend the modules
- Available from the nltk_contrib source tree at <http://code.google.com/p/nltk/>.

Acknowledgments

Work reported here was partially funded by NBCHC040176 from the US Department of the Interior, a Google Research Award, and the National Science Foundation under grant #0705708 to the Center for Language and Speech Processing at the Johns Hopkins University.