

Linguistic Hacking

How to know what a text in an unknown language is about?

Martin.Haase@uni-bamberg.de

[maha@jabber.ccc.de](xmpp:maha@jabber.ccc.de)

Contents

- how to identify the language of a written text
 - in traditional ways,
 - with the help of computer technology.
- how to get at least some information out of an unknown text.

Hacking

“the intellectual challenge of creatively overcoming or circumventing limitations.”

Eric Raymond (1996):
*The New Hacker's
Dictionary.*

Spoken texts?

multi-language corpus of telephone calls

Writing Systems

- Roman (thousands of languages)
- Cyrillic (> 60 languages)
- Arabic (> 20 languages)
- Devanāgarī (> 10 languages, not counting derivative writing systems)
- Hebrew (~ 3 – 5 languages)

Devanāgarī

देवनागरी

DEVANAGARI.

(Ältere Druckschrift nach A. W. v. SCHLEGEL.)

Zeichen	Wert										
अ	a	आ	ा	ए	े	ए	े	ঁ	dha	ু	ু
आ	ā	ঁ	া	ঁ	া	ঁ	া	ন	na	ঁ	া
ই	i	ঈ	ৈ	ও	ো	ও	ো	ঁ	pa	ঁ	ো
ই	ī	ঈ	ী	ও	ৌ	ও	ৌ	ঁ	da	ঁ	ৌ
উ	u	আ	ো	চ	tša	ঁ	ঁ	ঁ	dha	ঁ	ব
উ	ū	আ	ু	ছ	tšha	ঁ	ঁ	ঁ	ঁ	ঁ	ব
ক	au	ঁ	ু	ঁ	dža	ঁ	ঁ	ঁ	ঁ	ঁ	ব
ক	ka	ঁ	ু	ঁ	džha	ঁ	ঁ	ঁ	ঁ	ঁ	ব
ক	kha	ঁ	ু	ঁ	ঁ	ঁ	ঁ	ঁ	ঁ	ঁ	ব
ক	্ৰ	ু	ু	ঁ	ঁ	ঁ	ঁ	ঁ	ঁ	ঁ	ব

DEVANAGARI.

Neuere Druckschrift nach alten Handschriften.

Zeichen	Wert	Zeichen	Wert	Zeichen	Wert	Zeichen	Wert	Zeichen	Wert	Zeichen	Wert
अ	a	ल	l	ग	ga	ट	ṭa	ध	dha	र	ra
आ	ā	ख	l̥	घ	gha	ठ	ṭha	न	na	ल्ल	la
इ	i	ए	e	ङ	ṅa	ड	da	प	pa	ङ्क	la
ई	ī	ऐ	ai	च	t̥sa	ढ	ḍha	फ	pha	व	va
उ	u	ओ	o	ञ	t̥sha.	ण	ṇa	ब	ba	स	sa
ऊ	ū	औ	au	ज	dža	त	ta	भ	bha	श्व	śa
कृ	r̥	कौ	ka	ঁ	džha	ঁ	tha	ম	ma	ষ	ša
কৃ	r̥	কৌ	kha	ঁ	ńa	ଦ	da	য	ya	ହ	ha

प	पा	पि	पी
pa	pā	pi	pī
पु	पु	पु	पु
pu	pū	pṛ	pī
प्प	प्पा	प्र	न्न्या
ppa	pta	pra	ṅkhyā

SAMPLE OF ARMENIAN

<i>Armenian:</i>	ի	սրտին	գործարանի	երեւութացեալ	հոգւոյն
<i>Transliteration:</i>	I	srtin	gorcarani	erewut'ac'eal	hogwoyn
<i>Transcription:</i>	i	sərtin	goṛsarani	eṛevutʰatsʰial	hogwoin
<i>Gloss:</i>	in	the.heart	of.the.organ	appearing	of.the.soul
<i>աչաց</i>	<i>թաթ</i>	<i>ձեռին</i>	<i>աջոյ</i>	<i>գրելով</i>	<i>ի վերայ վիմի.</i>
ač'ac'	t'at'	jerin	ajoy	grelov	i veray vimi.
atʰatsʰ	tʰatʰ	dzerin	adžo	gərelow	i vela vimi
of.the.eyes	the.palm	of.the.hand	of.the.right writing	on	the.rock
<i>զի</i>	<i>որպէս</i>	<i>ի</i>	<i>ձեան</i>	<i>վերջք դժին</i>	<i>կուտեալ ունէր քարն:</i>
zi	orpēs	i	jean	verjk'	kuteal unēr k'arn.
zi	voipes	i	dzian	verdžkʰ	kutial uneř kʰam
for	as	in	snow	traces	had.retained the.stone
<i>եւ</i>	<i>յարուցեալ</i>	<i>յաղօթիցն</i>	<i>եստեղձ</i>	<i>զնշանագիրս</i>	<i>մեր, հանդերձ</i>
Ew	yaruc'eal	yalot'ic'n	estełc	znšanagirs	mer, handerj
jev	harutsʰial	hayot'itsʰən	jesteyts	əznəʃanagiřəs	meř handeřz
and	arising	from.prayer	he.fashioned	the.letters	our together
<i>Հոռփիխնոսիւ կերպաձեւեալ</i>	<i>զգիրն</i>	<i>առ ձեռն</i>	<i>պատրաստ</i>	<i>Մեսրոպայ,</i>	
Hrop'ianosiw kerpajeweal	zgirn	air jern	patrast	Mesropay,	
hropʰianosiv keřpadzevial	əzgiřən	ar dzerən	patrast	mesropa	
with.Rufinus by.giving.shape	the.letters	by the.hand	prepared	of.Mesrop	

Hebrew

- Old and Modern Hebrew,
- Ladino (with different varieties),
- Judeo-Arabic,
- Yiddish.

א בְּרָאשִׁית בָּרָא אֱלֹהִים אֵת הַשָּׁמֶן וְאֵת הָאָרֶץ: ב וְהָאָרֶץ קִיְּתָה
תְּהוֹם וְבָהּוּ וְחַשְׁךְ עַל־פָּנָיו תְּהוּם וְנוּסָם אֱלֹהִים מְרֻחְכָּת עַל־פָּנָיו^ה
הַמְּפִימִים: ג וַיֹּאמֶר אֱלֹהִים יְהִי אֹור וַיְהִי אֹור: ד וַיֹּרֶא אֱלֹהִים אֲתָּה
הָאֹור כִּי־צָבוֹב וַיַּבְדֵּל אֱלֹהִים בּוּין הָאֹור וּבּוּין הַחַשָּׁךְ: ה וַיֹּקְרָא
אֱלֹהִים לְאֹור יוֹם וְלַחֲשָׁךְ קָרָא לִילָה וַיְהִי־עָרֵב וַיְהִי־בָּקָר יוֹם אֶחָד:

1 אַיְן אֲנָה יִבְחַט גַּאַט בְּאַשְׁאָפָּן דָּעַם הַיְמָל אֹוְן דֵּי עֶרֶד. 2 זָוְן דֵּי עֶרֶד אַיְז גַּעֲוָעָן וּוַיסְט
אוֹן לִיְדֵיק, אוֹן פְּינַצְטָעָרְנִיש אַיְז גַּעֲוָעָן אוַיְפָן גַּעַזְיִכְט פָּוּן טַהָּאָם, אוֹן דֵּר גִּיסְט פָּוּן גַּאַט
הַאַט גַּעַשְׂוָבָט אוַיְפָן גַּעַזְיִכְט פָּוּן דֵּי וּוַאֲסָעָרָן. 3 חַאַט גַּאַט גַּעַזְאָגָט: זָאַל וּוַעֲרָן לִיבְט. זָוְן עַס
אַיְז גַּעֲוָאָרָן לִיבְט. 4 זָוְן גַּאַט הַאַט גַּעֲזָעָן דָּאָס לִיבְט אַז עַס אַיְז גּוֹט; אוֹן גַּאַט הַאַט
פָּאָנָאָנְדָעָרְגָּעָשִׂידָט צְוֹוִישָׁן דָּעַם לִיבְט אוֹן צְוֹוִישָׁן דֵּר פְּינַצְטָעָרְנִיש. 5 זָוְן גַּאַט הַאַט גַּעַרְוָפָּן
דָּאָס לִיבְט טָאָג, אוֹן דֵּי פְּינַצְטָעָרְנִיש הַאַט עֶר גַּעַרְוָפָּן נָאָכָט. זָוְן עַס אַיְז גַּעֲוָעָן אָוּוֹנָט, אוֹן
עַס אַיְז גַּעֲוָעָן פְּרִימָאָרגָן, אַיְן טָאָג.

bu	3, 10, 11, 19, 23, 35, 40, 41, 42, 50
bú	27
bü	3
by	20, 22, 23, 24, 26, 28, 30, 32, 34, 35, 43
bý	27
bij	25
ć-	56
ca	10, 14, 15, 17a, 17c, 17d, 22, 23, 24, 40, 41, 53
cá	12, 41
cà	17c, 18
că	14
cã	12
ça	11, 15, 16, 53
çà	11
ce	11, 13, 14, 17c, 41, 42
cé	17c, 41
će	31, 17d, 35
će	2
če	33
ci	2, 11, 13, 14, 17c, 18, 30, 40, 41
ći	35
ći	2
či	32, 34
ćí	32
co	16, 17, 17a, 17b, 17c, 17d, 19, 30, 32, 41
có	41, 42
cò	16, 17
ço	16, 53
çò	16
čo	32, 34
cu	10, 14, 17, 17a, 17c, 17d, 40, 41, 53, 54
cú	12, 41
cù	18, 42
ća	31
ću	31
ću	2

1	Latin
2	Esperanto
3	Volapük
10	Spanish
11	French
12	Portuguese
13	Italian
14	Romanian
15	Catalan
16	Occitan
17	Ladin (Engadine)
17a	Sursilvan
17b	Surmiran
17c	Friulan
17d	Ladin (Dolomite)
18	Corsican
19	Jersey language
20	English
21	German
22	Danish
23	Swedish
24	Norwegian
25	Dutch
26	Afrikaans
27	Icelandic
28	Frisian
30	Polish
31	Serbo-Croat
32	Czech
33	Slovenian
34	Slovak
35	Sorbian (Wend)
40	Welsh
41	Irish Gaelic
42	Scottish Gaelic
43	Manx
44	Breton
50	Turkish
51	Finnish
52	Hungarian
53	Albanian (Tosk)
54	= 53 Albanian (Tosk) + some Geg words
55	Basque
56	Maltese
57	Estonian
58	Lithuanian
59	Latvian

Norman C. Ingle (1980): Language Identification Table. London: Technical Translation International.

* See Page 7

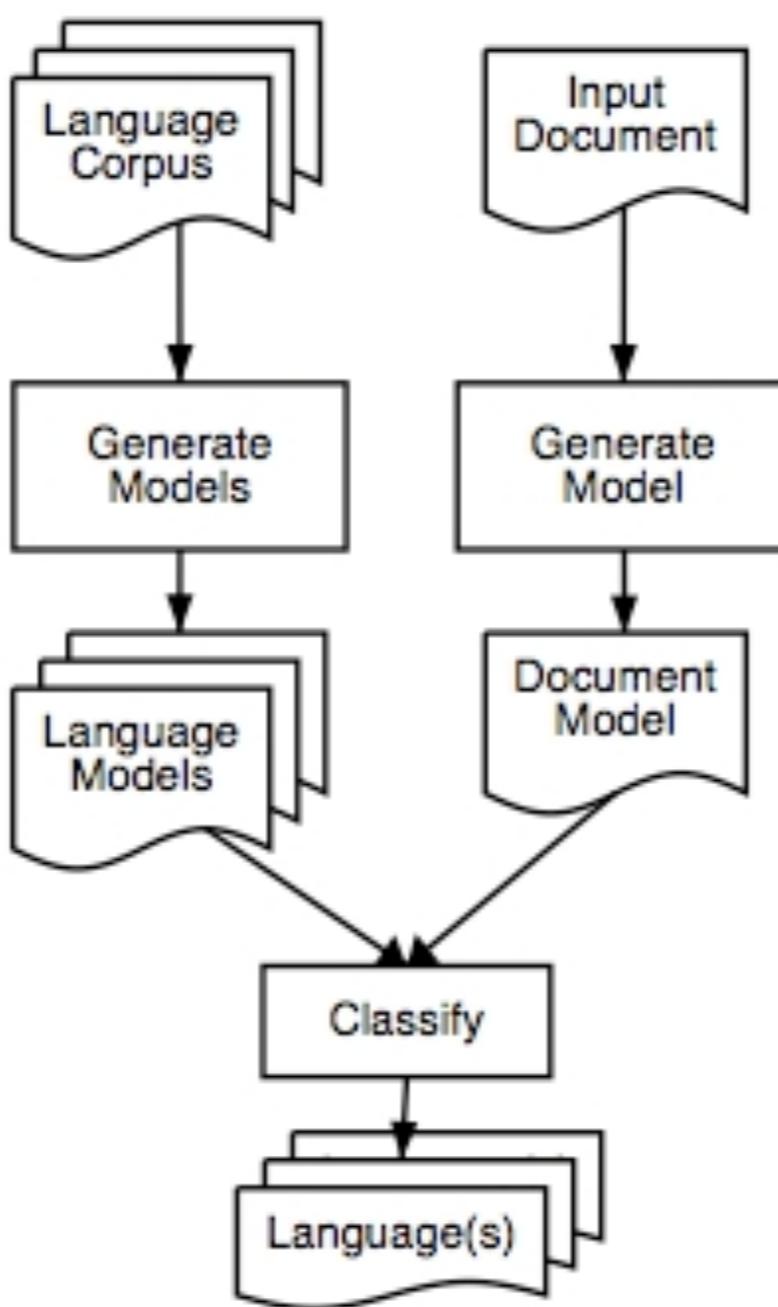
bu	3, 10, 11, 19, 23, 35, 40, 41, 42, 50	
bú	27	
bü	3	
by	20, 22, 23, 24, 26, 28, 30, 32, 34, 35, 43	
bý	27	
bij	25	
č-	56	
ca	10, 14, 15, 17a, 17c, 17d, 22, 23, 24, 40, 41, 53	
cá	12, 41	
cà	17c, 18	
că	14	
cã	12	
ça	11, 15, 16, 53	
çà	11	
ce	11, 13, 14, 17c, 41, 42	
cé	17c, 41	
će	31, 17d, 35	
	1 Latin	
	2 Esperanto	
	3 Volapük	
	10 Spanish	
	11 French	
	12 Portuguese	
	13 Italian	
	14 Romanian	
	15 Catalan	
	16 Occitan	
	17 Ladin (Engadine)	
	17a Sursilvan	
	17b Surmiran	
	17c Friulan	
	17d Ladin (Dolomite)	
	18 Corsican	
	19 Jersey language	
	20 English	
	21 German	
	22 Danish	
	23 Swedish	
	24 Norwegian	

Computer-aided identification

- frequencies of unique characters and character strings
- common words recognition
- *n*-gram analysis

“Text”

(TE), (TEX), (EXT), (XT)



variant of the unique character string approach

compression efficiency

reference
model

reference
model

+

text in language
to be identified

gzip



reference
model

+

text in language
to be identified

gzip



reference
model

+

text in language
to be identified

compression efficiency

Interesting applications

- measuring linguistic difference
 - > language families
- determining types of text
 - spam detection?

- TextCat (<http://odur.let.rug.nl/vannoord/TextCat/Demo/>), *n-gram* based, 76 languages, usable as a web application,
- Languid (<http://languid.cantbedone.org/>), downloadable program, web application not running,
- Langid (<http://complingone.georgetown.edu/~langid/>), *n-gram* based, 65 languages, web application,
- LanguageGuesser (<http://www.xrce.xerox.com/cgi-bin/mltt/LanguageGuesser>), frequency tests on characters and character sequences, about 40 languages, web application,
- Polyglot 3000 (<http://www.polyglot3000.com/>), corpora, method unknown, currently 441 languages, closed-source Windows freeware. :-(

approaching “content analysis”

Hacker's approach

- numbers, dates, words from another language
- typographic hints:
 - bold or italic print,
 - colored or underlined text chunks,
 - capital letters

Zipf's law

Very frequent words are shorter and contain less lexical information, whereas infrequent words are longer and contain more lexical information.

less lexical information implies more grammatical
information and vice versa

most interesting for us:
words with more specific lexical information

**Ignore all short words!
(even if they reiterate throughout the text)**

Ua salalau lenei gagana i le lalolagi atoa. ‘O lenei fo‘i gagana, ‘ua ‘avea ma gagana lona lua a le tele o tagata ‘o le vasa Pasefika, e pei ‘o Samoa. E iai le manatu, ‘o le gagana fa‘aperetania, ‘ua matuā talitonu i ai le tele o tagata Samoa e fa‘apea ‘o le gagana e maua ai le atamai ma le poto. ‘E talitonu fo‘i nisi o i latou, ‘e lē aoga la latou gagana. E lē sa‘o lea tāofi, ‘auā e ‘avatu le gagana fa‘aperetania i Samoa, ‘ua leva ona atamamai ma popoto tagata Samoa e fai lo latou soifua ma lo latou lalolagi.

Ua salalau lenei gagana i le lalolagi atoa. ‘O lenei fo‘i gagana, ‘ua ‘avea ma gagana lona lua a le tele o tagata ‘o le vasa **Pasefika**, e pei ‘o **Samoa**. E iai le manatu, ‘o le gagana fa‘aperetania, ‘ua matuā talitonu i ai le tele o tagata **Samoa** e fa‘apea ‘o le gagana e maua ai le atamai ma le poto. ‘E talitonu fo‘i nisi o i latou, ‘e lē aoga la latou gagana. E lē sa‘o lea tāofi, ‘auā e ‘avatu le gagana fa‘aperetania i **Samoa**, ‘ua leva ona atamamai ma popoto tagata **Samoa** e fai lo latou soifua ma lo latou lalolagi.

Ua salalau lenei **gagana** i le lalolagi atoa. ‘O lenei fo‘i **gagana**, ‘ua ‘avea ma **gagana** lona lua a le tele o tagata ‘o le vasa Pasefika, e pei ‘o Samoa. E iai le manatu, ‘o le **gagana fa‘aperetania**, ‘ua matuā talitonu i ai le tele o tagata Samoa e fa‘apea ‘o le **gagana** e maua ai le atamai ma le poto. ‘E talitonu fo‘i nisi o i latou, ‘e lē aoga la latou **gagana**. E lē sa‘o lea tāofi, ‘auā e ‘avatu le **gagana fa‘aperetania** i Samoa, ‘ua leva ona atamamai ma popoto tagata Samoa e fai lo latou soifua ma lo latou lalolagi.

Ua salalau lenei **gagana** i le lalolagi atoa. ‘O lenei fo‘i **gagana**, ‘ua ‘avea ma **gagana** lona lua a le tele o tagata ‘o le vasa Pasefika, e pei ‘o Samoa. E iai le manatu, ‘o le **gagana fa‘aperetania**, ‘ua matuā talitonu i ai le tele o tagata Samoa e fa‘apea ‘o le **gagana** e maua ai le atamai ma le poto. ‘E talitonu fo‘i nisi o i **latou**, ‘e lē aoga la **latou gagana**. E lē sa‘o lea tāofi, ‘auā e ‘avatu le **gagana fa‘aperetania** i Samoa, ‘ua leva ona atamamai ma popoto tagata Samoa e fai lo **latou soifua** ma lo **latou lalolagi**.