

# World languages

Languages on the Web

Written languages

Standardised Languages

Ethnologue lists **7,413** languages

ISO 639-3:2007 Codes: **7,589** entries

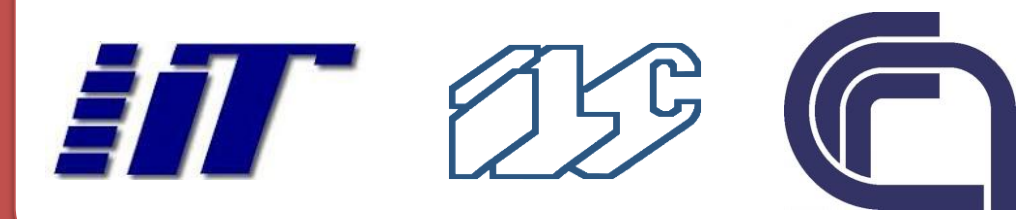
Books published in **480** languages (source: Google Books Project)

Official languages of country or region: **238** (source: Wikipedia)

Possibly **1,500** written languages on the **Web** (source: Scannel's study)

## Web Language Identification Tool - (IIT & ILC CNR)

matteo.abrate@iit.cnr.it clara.bacciu@iit.cnr.it francesca.frontini@ilc.cnr.it mariantonietta.lapolla@iit.cnr.it andrea.marchetti@iit.cnr.it monica.monachini@ilc.cnr.it



**Language identification or detection:**  
given a text, determine the language it is written in

### Language identification techniques

- Common words technique: (Grefenstette)
- Unique letter combinations: (Churcher)
- Statistical approach: (Dunning)
- N-grams approach: (Canvar et al.)
- Compression based approach: (Teahan et al.)

### Some open-source language identification tools

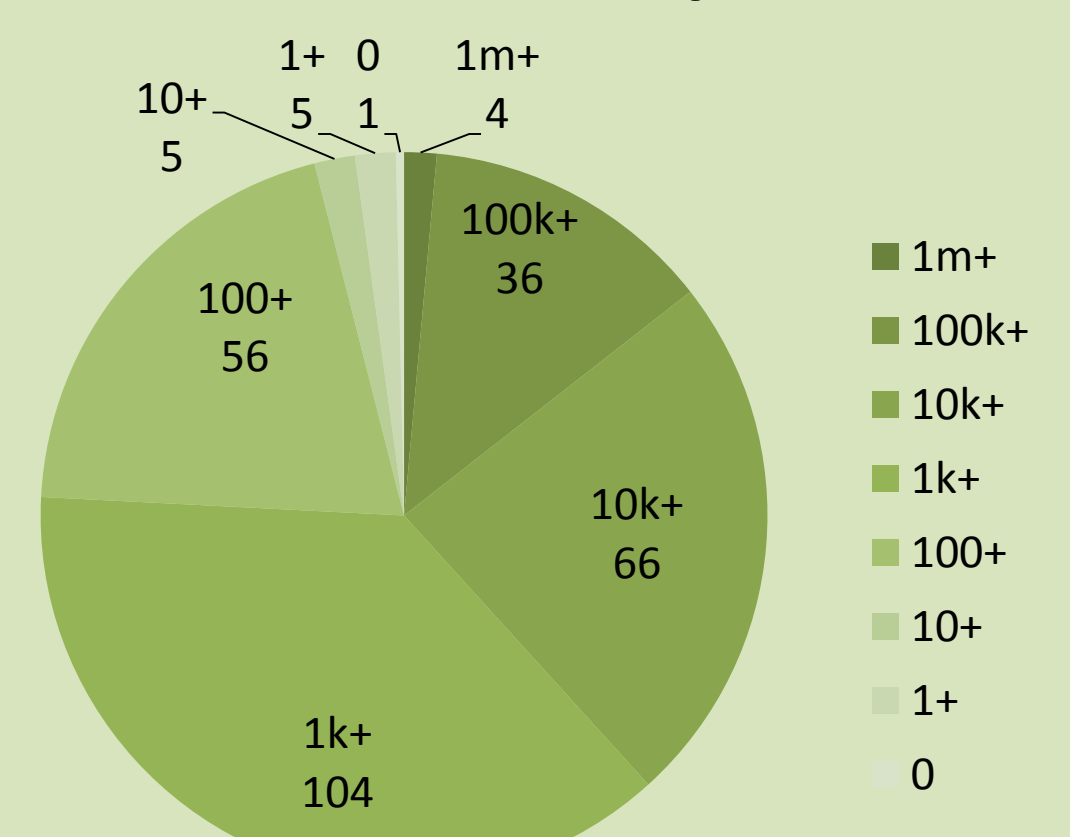
- Chromium Compact Language Detector (CLD):** declaring support for about **160 languages**. Method: n-grams. Language: C++ (we used an incomplete binding for Python)
- Lingua::Identify:** Supports **33 languages**. Methods: small words, prefix, suffix, n-grams. Language: Perl.

### Known problems:

- Underrepresented languages
- Short sentences
- Mix of languages in one page
- Unusual script/orthography
- Non-standard register

### Wikipedia:

a valuable source for underrepresented languages

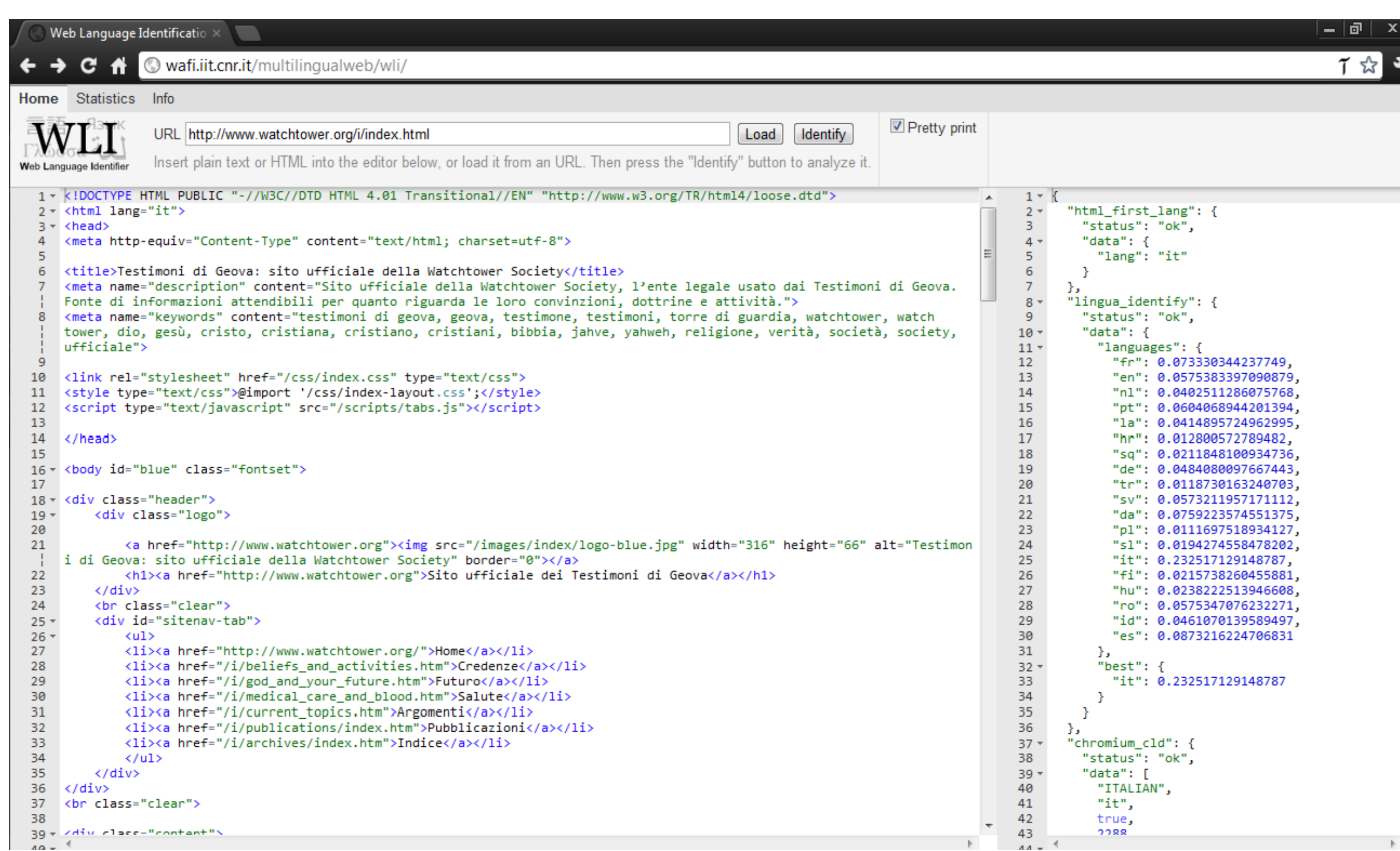
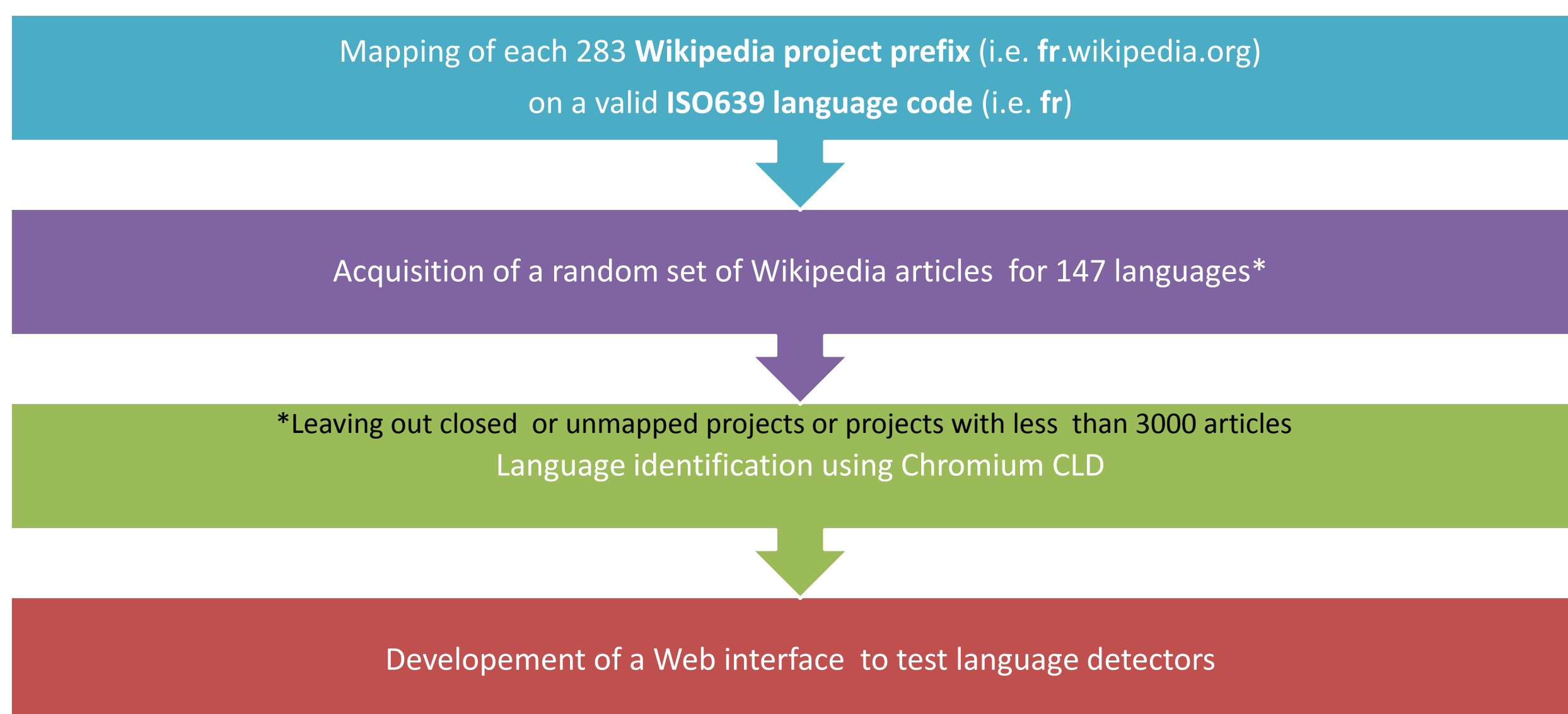


Projects by number of articles

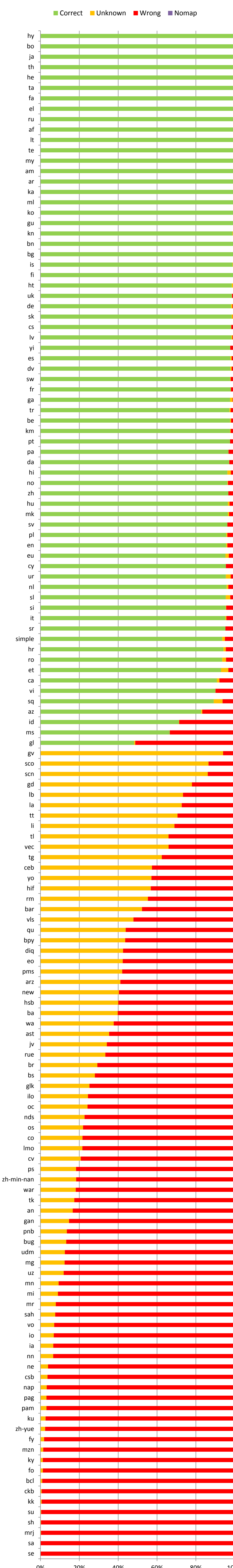
Test set: 1000 articles from 147 projects having more than 3000 articles

How does the best open source language identification library behave with respect to the languages on the Web?

### Testing Chromium CLD library



Testing result: a lot of data is still needed to build more language models



Testing Chromium Compact Language Detector library against 147.000 articles from 147 Wikipedia projects

### Language identification for the Wikipedia home pages

Prefix	ISO 639 code	Chromium CLD	Lingua: Identify	HTML lang
en	en	en	en	en
de	de	de	de	de
fr	fr	fr	fr	fr
it	it	it	it	it
es	es	es	es	es
ru	ru	ru	ru	ru
pt	pt	pt	pt	pt
nl	nl	nl	nl	nl
sv	sv	sv	sv	sv
pl	pl	pl	pl	pl
uk	uk	uk	uk	uk
ca	ca	ca	ca	ca
no	no	no	no	no
fi	fi	fi	fi	fi
cs	cs	cs	cs	cs
hr	hr	hr	hr	hr
ko	ko	ko	ko	ko
id	id	id	id	id
tr	tr	tr	tr	tr
ro	ro	ro	ro	ro
fa	fa	fa	fa	fa
ar	ar	ar	ar	ar
da	da	da	da	da
eo	eo	eo	eo	eo
sr	sr	sr	sr	sr
lt	lt	lt	lt	lt
sl	sl	sl	sl	sl
sk	sk	sk	sk	sk
ms	ms	ms	ms	ms
be	be	be	be	be
bg	bg	bg	bg	bg
kk	kk	kk	kk	kk
eu	eu	eu	eu	eu
vo	vo	vo	vo	vo
war	war	war	war	war
hr	hr	hr	hr	hr
hi	hi	hi	hi	hi
et	et	et	et	et
az	az	az	az	az
gl	gl	gl	gl	gl
an	an	an	an	an
simple	en	en	en	en