

Local Language Support

AFGHANISTAN SCHOOL ON INTERNET GOVERNANCE

AFSIG 2017

KABUL, AFGHANISTAN

SAID ZAZAI

سعید خاکی

SAID@ZAZAI.CA

TWITTER: @SMZAZAI

APRIL 26, 2017

“Over six billion people live in over 200 countries spanning 24 time zones. These people use hundreds of currencies to conduct business in thousands of different languages and dialects. Their business practices are all over the map, ranging from simple barter to cash to electronic payment to sophisticated arbitrage.”

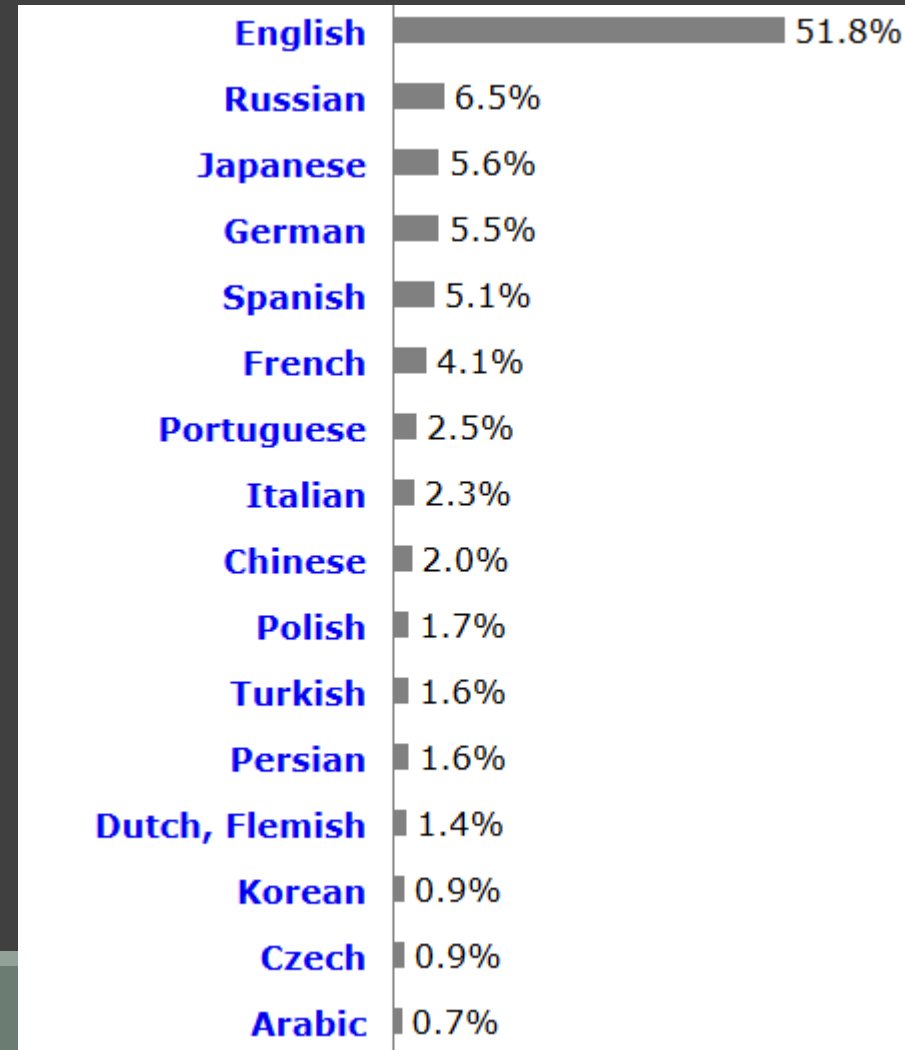
Donald DePalma, 2002

Usage of content languages for websites

❖ English is used by 51.8% of all the websites.

[Source: w3techs.com](http://w3techs.com)

❖ Less than 0.1% websites use Pashto language.
1.6% Persian (including Dari)



Communities advance when computers speak their language

IDRC reports

- A major obstacle to Internet use.
- With 3,500 local languages in the Asia-Pacific region, and fewer than 10% of people able to communicate in English, Internet use is typically restricted to urban areas.
- Highly skilled software engineers, linguists, and sociologists
 - work together to overcome the formidable technical obstacles to making local scripts compatible with computers, and to promote their use.
 - Requires a push from government, vendors, technical community, civil society and academia.

Local Language Support

Why need it?

- Read/Write (Community empowerment, intelligence/information society, knowledge sharing...)
- Research
- Language preservation
- Automate (Translation, text to speech, speech to text...)

Localization & Internationalization

- **Internationalization** is the process of designing a software application so that it can *potentially* be adapted to various languages and regions without engineering changes.
- **Localization** is the process of adapting internationalized software for a *specific* region or language by adding locale-specific components and translating text.
- A **locale** is a set of parameters that defines the user's language, region and any special variant preferences that the user wants to see in their user interface.
 - Numbers, calendar, currency, keyboard layout, paper size, date-time format, phone number format, system of measurement, postal address format etc.
- It may also consider differences in culture such as:
 - Local holidays, personal names and titles, local custom and conventions, architecture etc.

The Technology

Vendor/OS/Environment	Read (Fonts)	Write (Input Method Editor)
Microsoft Windows	Windows XP	Windows XP SP3
Google Android	Android 5	Android 7
MacOS		OS Capitan
Java	Unicode	Unicode
Oracle	Unicode	Unicode

Unicode

- Unicode is a universal encoded character set that enables information from any language to be stored using a single character set.
- Unicode provides a unique code value for every character, regardless of the platform, program, or language.

	FE7	FE8	FE9	FEA	FEB	FEC	FED	FEE	FEF
0	ا FE70	ء FE80	ب FE90	ج FEA0	ز FEB0	ض FEC0	غ FED0	ل FEE0	ى FEF0
1	ا FE71	آ FE81	ب FE91	ح FEA1	س FEB1	ط FEC1	ف FED1	م FEE1	ي FEF1
2	ا FE72	آ FE82	ب FE92	ح FEA2	س FEB2	ط FEC2	ف FED2	م FEE2	ي FEF2
		أ FE83	ت FE93	ث FEA3	د FEB3	ذ FEC3	ر FED3	ز FEE3	س FEF3

Standardization Issues

Calenders	
Hijri Solar (شمسي)	١٣٩٥
Hijri Lunar (قمري)	١٤١٠
Gregorian (ميلادي)	2017

Standardization Issues

ک ک

گ گ

ی ی

Solution: Normalization

ك	اَ	ب	اِ	ج	د	فب	گپ
0643	0653	0663	0673	0683	0693	06A3	06B3
ل	اَ	ع	ء	ج	ر	قا	گش
0644	0654	0664	0674	0684	0694	06A4	06B4
م	اَ	ه	اِ	ش	ر	پپ	ل
0645	0655	0665	0675	0685	0695	06A5	06B5
ن	اَ	و	و	چ	ر	قا	ن
0646	0656	0666	0676	0686	0696	06A6	06B6
ه	اَ	و	و	چ	ز	ف	ث
0647	0657	0667	0677	0687	0697	06A7	06B7
و	اَ	و	اِ	ط	ژ	ق	پ
0648	0658	0668	0678	0688	0698	06A8	06B8
ی	اَ	و	ط	د	ژ	ک	ن
0649	0659	0669	0679	0689	0699	06A9	06B9
ی	اَ	و	ن	د	بن	ک	ن
064A	065A	066A	067A	068A	069A	06AA	06BA

Languages of Afghanistan

Language	Alphabet Set (Written Language)	IME
Pashto	Yes	Yes
Dari	Yes	Yes
Uzbeki	Yes	No
Turkmeni	Yes	No
Pashayi	Yes (2003)	No
Sanglechi	No	No
Ishkashmi	No	No
Balochi	Yes	No
Qizilbashi	No	No
Kyrghyz	Yes	No

Q&A