

Persian Computing with Unicode

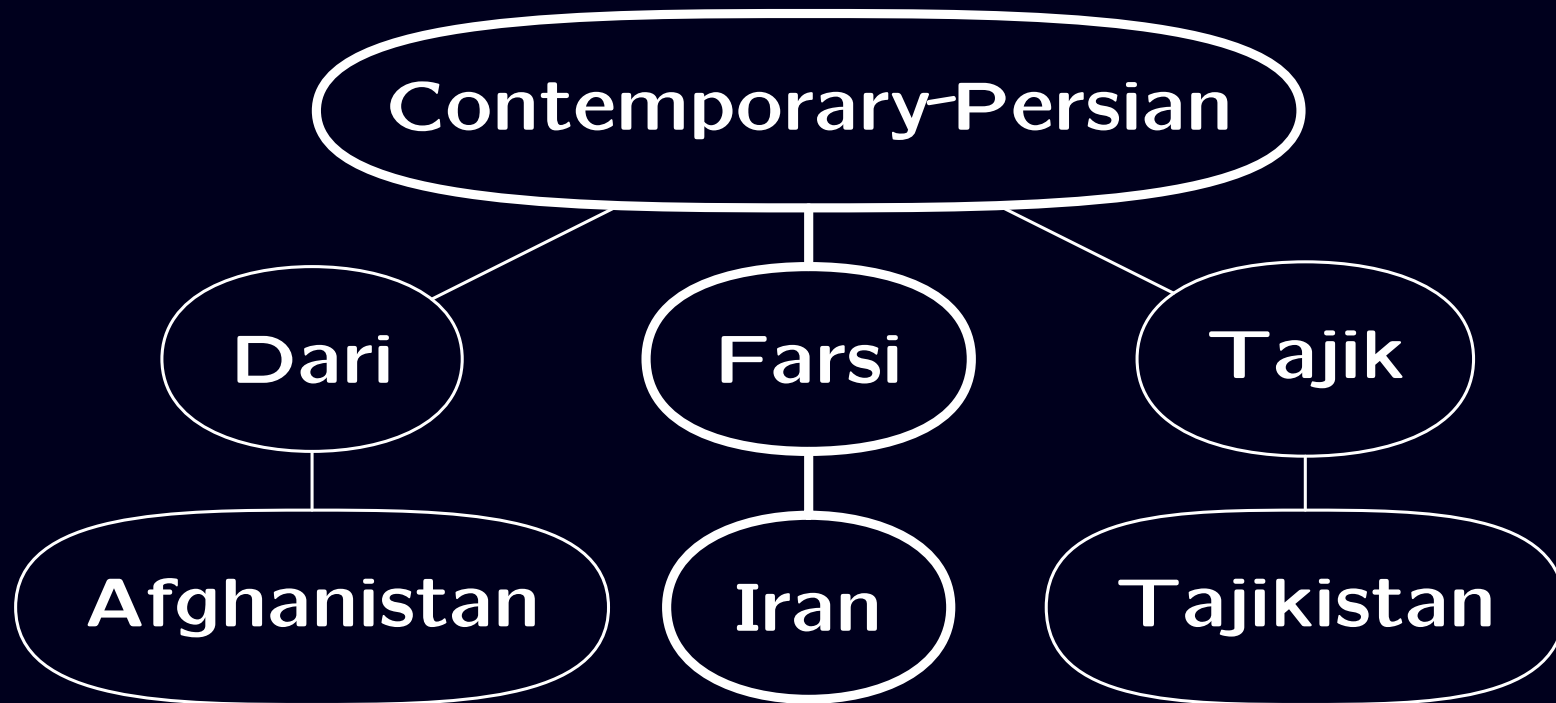
~~Behdad Esfahbod~~ no visa
unicode@behdad.org

Amir Youssefi
youssefi@cs.rpi.edu

The FarsiWeb Project
<http://www.farsiweb.info/>

March 31, 2004

What is Persian?



Persian in Computers

There are three relevant national standards:

- ISIRI 3342:1992 Farsi 8-bit Coded Character Set for Information Interchange (deprecated)
- ISIRI 2901:1994 Keyboard Layout for Farsi: Characters in Computer
- ISIRI 6219:2002 Information Technology -- Persian Information Interchange and Display Mechanism, using Unicode

Modern Persian Script

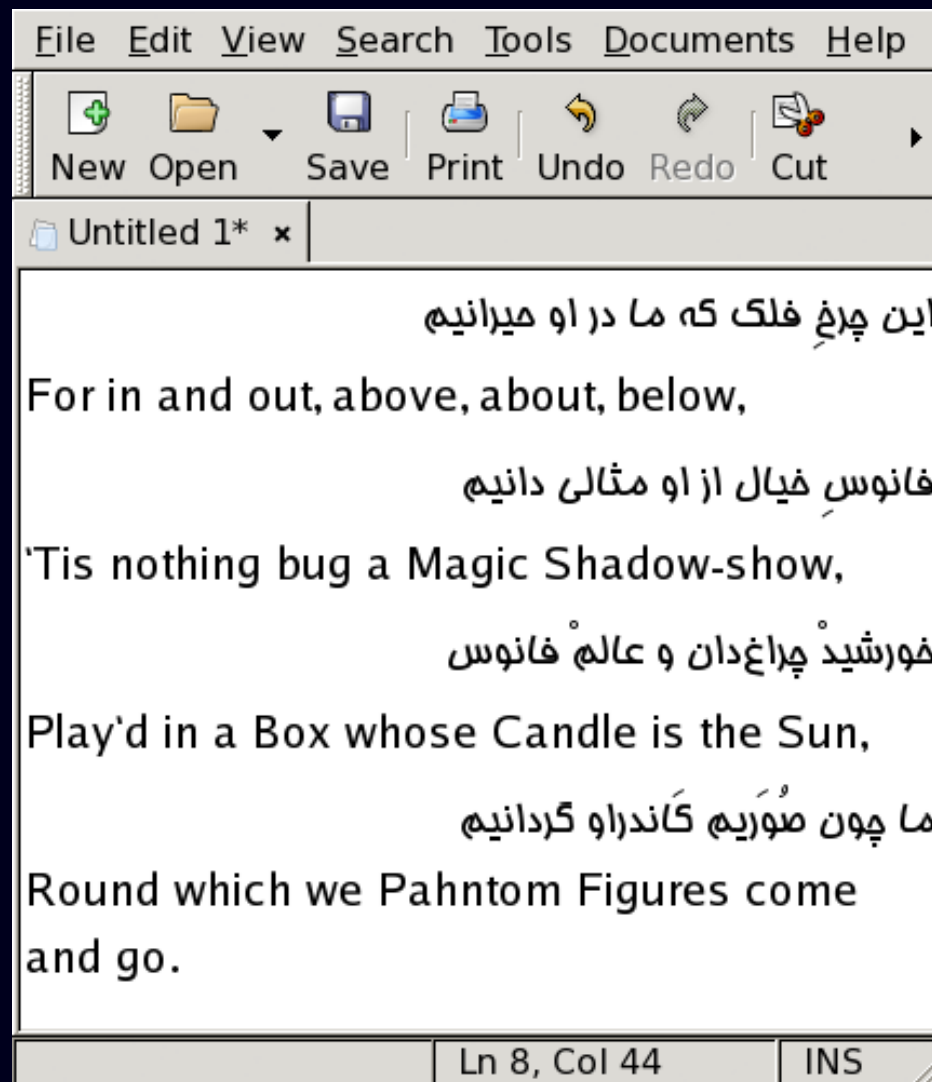
- Based on Arabic Script ($\langle U+0600..U+06FF \rangle$ block): With some extra letters, some modified letters
- But with completely different semantics and typographical habits
- Is a bidirectional script: Is written from right to left, except for numbers
- Needs cursive joining: Two adjacent letters may be *joined*, forming 1, 2, or 4 glyphs for each character: (for example س, س, سد, سد)

Arabic Script Rendering

Input text	Logical order	م ا ل س
After Bidirectional Algorithm	Visual order	س ل ا م
After Arabic Joining Algorithm	Glyph list	س ل ا م
After Ligation	Glyph list	س ل ا م
When Rendered	Output	سلام

With enough care, it is possible, to apply the above algorithms in a different order, and get the same result.

A Bidirectional Document (gedit/Pango)



Alphabet

- Extra letters:
U+067E *Peh* (پ), U+0686 *Tcheh* (چ),
U+0698 *Jeh* (ژ), U+06AF *Gaf* (گ)
- Modified letters:

Character	Isol	Fina	Medi	Init
U+0643 <i>Arabic Letter Kaf</i> (Arabic Kaf)	ك	ك	ك	ك
U+06A9 <i>Arabic Letter Keheh</i> (Persian Kaf)	ک	ک	ک	ک
U+064A <i>Arabic Letter Yeh</i> (Arabic Yeh)	ي	ي	ي	ي
U+06CC <i>Arabic Letter Farsi Yeh</i> (Persian Yeh)	ی	ی	ی	ی

Alphabet (continued)

- Three shapes of composed Hamza Above:
 - U+0623 *Alef with Hamza Above* (أ),
 - U+0624 *Waw with Hamza Above* (و),
 - U+0626 *Yeh with Hamza Above* (ي)
- Never used characters:
 - U+0649 *Alef Maksura* (ى): Like Yeh, but no dots at all
 - U+06C0 *Heh with Yeh Above* (هـ): Should **never** be used instead of ⟨*Heh, ZWNJ, Farsi Yeh*⟩ or ⟨*Heh, Hamza Above*⟩ sequence

Special Characters

- U+0640 *Arabic Tatweel*, for a longer joining stem

کتاب → کتاب

- U+200C *Zero Width Non-Joiner*, to prevent joining

کتابها → کتابها

- U+200D *Zero Width Joiner*, to choose a joined glyph when it would not join naturally

ه.ش → ه.ش

- U+200E *Left-to-Right Mark*, U+200F *Right-to-Left Mark*, and other bidirectional control chars ($\langle U+202A..U+202E \rangle$)

Numbers

- $\langle \text{U+06F0}..\text{U+06F9} \rangle$ *Extended Arabic-Indic Digits*:

٠ ١ ٢ ٣ ٤ ٥ ٦ ٧ ٨ ٩

instead of $\langle \text{U+0660}..\text{U+0669} \rangle$ *Arabic-Indic Digits*:


٠ ١ ٢ ٣ ٤ ٥ ٦ ٧ ٨ ٩

- U+066C *Arabic Thousands Separator* and U+066B *Arabic Decimal Separator*:

٩'٨٧٤'٥٤٣/٢١٠

- Western numerals in Latin context. Persian numerals everywhere else (page numbers, section numbers, ...)

Other Characters

- Harakat (Vowel) Non-spacing Marks:  . . .
- Arabic Punctuation Marks:
 - U+060C *Arabic Comma* (،),
 - U+061B *Arabic Semicolon* (;),
 - U+061F *Arabic Question Mark* (؟),
 - U+066A *Arabic Percent Sign* (٪),
- ⟨U+00AB, U+00BB⟩ *Double Angle Quotation Marks* (« »)
- Shared Punctuation Marks:
 - Latin full stop, exclamation mark, parenthesis, square brackets, . . .



												،	060D	060E	060F
0610	0611	0612	0613	0614	0615						:				؟
	ء	آ	أ	ؤ	إ	ئ	ا	ب	ة	ت	ث	ج	ح	خ	د
ذ	ر	ز	س	ش	ص	ض	ط	ظ	ع	غ					
-	ف	ق	ك	ل	م	ن	ه	و	ي	ي	="	"	"	"	"
-	ء	°	ˆ	ˆ	ˆ	ˆ	0657	0658							
•	۱	۲	۳	۴	۵	۶	۷	۸	۹	%	/	'	*	۰	۹
ا	أ	أ	إ	أ	أ	ؤ	ؤ	ئ	ئ	ئ	پ	پ	ت	پ	ث
پ	ح	ح	ج	ج	خ	چ	چ	ڈ	د	د	ڈ	ذ	ذ	ڈ	ذ
ڈ	ڑ	ڑ	ر	ر	ر	ر	ز	ژ	ژ	بب	پپ	پپ	ص	ض	ظ
غ	ف	ف	ف	ق	پ	ق	ف	ق	ک	ک	گ	ک	ک	گ	گ
گی	گی	گی	گ	گی	ل	ل	ل	پ	ب	ب	ٹ	ن	ٹ	ھ	چ
ہ	ہ	ہ	ہ	و	و	و	و	و	و	و	ی	ی	ی	ی	و
ی	ی	ے	ے	-	ہ	ط	ط	م	لا	ع	•	س		◉	•
•	ˆ	ˆ	س	•	•	ء	ء	ن	↑	•	•	•	م	06EE	06EF
•	۱	۲	۳	۴	۵	۶	۷	۸	۹	بب	ض	غ	ء	م	06FF

Keyboard Layout

ISIRI 2901:1994, features:

- All letters and punctuation marks
- Reasonable placement
- Persian digits
- Regular and shifted keys only
- Some empty slots

Iranian Standard Layout

	!	'	/	ٲ	%	x	،	*	()	-	+	BackSpC	
Tab	ّ	ّ	ّ	ّ	ّ	ّ	ّ	ّ	[]	{	}		
Caps	ؤ	ئ	ي	ا	آ	ة	«	»	:	؛	:	؛	Enter	
Shift	ك	ط	ژ	ر	ذ	د	ء	<	>	؟	Shift		\ -	
Ctrl	Win	Alt	 Space					Alt	Win	Menu	Ctrl			

Keyboard Layout (continued)

Proposed update, features:

- Fully backward-compatible with ISIRI 2901:1994
- Unicode 4.0 repertoire, and complete support for ISIRI 6219:2002
- Support quoting Arabic text
- Uses AltGr to add required but rarely used characters
- Adds all ASCII punctuation marks, useful for editing XML, ...
- Adds bidirectional control characters

Proposed Iranian Layout

÷	~	!	'	/	ج	%	x	‘	*	()	-	+	BackSpc													
ZWJ	~	ی	`	۲	@	۳	#	۴	\$	۵	%	۶	^	۷	&	۸	•	۹	LRM	•	RLM	-	_	=	-	BackSpc	
Tab	°	ص	ث	€	ق	ف	غ	ع	ه	LR O	خ	RL O	ح	PDF	ج	LR E	چ	RL E									
Caps	و	ئ	ي	ا	آ	ة	«	»	:	؛	”	Enter															
Shift	ك	ط	ژ	ا	ZWJ CNJ	ف	ء	<	>	؟	Shift		-														
Ctrl	Win	Alt	Space NBSP										Alt	Win	Menu	Ctrl											

Fonts

- Microsoft fonts are Arabic
- Tahoma is the best looking one
- Persian fonts are not Unicode compatible yet
- The only ligature: Lam-Alef
- Nastaliq is desired, but not possible yet

The Nastaliq Style

يا مقلب القلوب والابصار
يا مدبر الليل والنهار
يا محول الحول والاحوال
حول حالنا الى احسن الحال

Date and Time

- Three calendars in use!
 - Gregorian**, to synchronize with the rest of the world
 - Jalali**, the official calendar
 - Islamic**, for some holidays and ceremonies
- Islamic calendar depends on moon-sighting once a year
- Week starts Saturday
- Business weekdays from Saturday to Thursday
- 24-hour preferred in media
- No AM/PM equivalent

Collation

- Like Arabic basically

- $ي < ه < و < ن \rightarrow ن < و < ه < ي$

- Some L2 equal pairs:

ك < ₃ ك	ي < ₃ ي	ة < ₃ ت
ع < ₃ ع	ه < ₃ ه	ع < ₃ ط

- Traditional rules: Hamza variants are L2 equal:

ي <₃ و <₃ أ

- Modern rules: L2 equal with their base letter:

ي <₃ ي و <₃ و أ <₃ أ

Loose Searching

ZWNJ \simeq Space

حجة الاسلام \simeq حجت الاسلام

كتاب \simeq كتاب

ي \simeq ي

تأخير \simeq تاخير

مسألة \simeq مسأله

ملك \simeq ملك

٠ \simeq ٠ ...

٤ \simeq ٤

٥ \simeq ٥

٦ \simeq ٦

٩ \simeq ٩ ...

ZWNJ \simeq empty string

دايره المعارف \simeq دايرة المعارف

خانه \simeq خانه ي

ك \simeq ك

سؤال \simeq سوال

مسئول \simeq مسؤول

ملك \simeq ملك

Last Notes to Application Developers

- Typesetting Persian paragraphs:
 - Justified lines
 - No inter-letter spacing
 - No word hyphenation
 - Almost no inter-word spacing
 - Use Tatweel instead: کتاب من → کتاب من
- All text fields Right-to-Left
- Persian numbers
- Right-to-Left layout
- Beware: Right and Left are swapped!

An Old-style Persian Poem (MS Word)

زان پیش که حسبِ حال گویم
آن خالقِ ماه و خور که چون گوی
زو گوی سپهر مستدیر است
از حکمتِ اوست در زد و گیر
از ماه برین بلند ایوان
هر ذره ز ماه تا بماهی
صنعتش که ز مهر عالم افروخت
این گوی درست‌زر که مهرست
از شرق بغرب داده راهش
از صنایع ذوالجلال گویم
زو چرخ فتاده در تکاپوی
چوگانِ هلال گوشه‌گیر است
چوگانِ قضا و گوی تقدیر
گه گوی نموده گاه چوگان
بر وحدتِ او دهد گواهی
بر جیبِ سپهر گوی زر دوخت
در چرخ ز گردشِ سپهرست
کانجاست محلّ حالگاهش

A Right-to-Left Dialog (Gtk+/Pango)



A Persian Paragraph (Mozilla)

۲
این سفیدباشی از خیره‌سری است که
دم مار می‌گزد یا نشاطِ عیش کرده
است که گاه می‌نماید و گاه می‌رباید؟



چرا وقتی پشتِ ماه خمیده شد و تازی به ماری که
از درختِ عرعر بالا می‌رفت پارس کرد و او روی
پنجه‌ها گردن کشید تا ناله‌ی وصل کند خارگزی به
ساقم خلید تا حواسم پرتِ زقزق شود و نتوانم
شاهد نمه‌عرقی باشم که می‌گویند در ربع مسکون
تنها بر منخرین دختر شاهِ سمنگان می‌نشیند وقتی
شبِ چهارده از ایوان به کیوان می‌نگرد.

در این شب و مرتع و اسارت باقی به همین قیاس
گذشت تا سرانجام که دم به تله داد چنان شیپهای
کشید که حتی الاغ‌های آن اطراف هم دانستند که
اسبِ سفید ته‌مینه دیگر از خیره‌سری دم مار را
نخواهید گزید.

A Persian Table (Mozilla)

پورمقدم، یارعلی، ۱۳۳۰ -
یادداشت‌های یک اسب / یارعلی پورمقدم. - تهران: آرویج، ۱۳۸۰.
۴۵ص.

ISBN 964-7174-6-8

فهرست‌نویسی بر اساس اطلاعات فیپا.

۸فا۶۲/۳

PIR۷۹۹۲/۴۷ی۱۷

ی۷۵۶پ

۱۳۸۰

۱۳۸۰

۸۰-۵۶-۲۶۰م

کتابخانه ملی ایران

- یادداشت‌های یک اسب
- یارعلی پورمقدم
- طرح روی جلد: هومن خطیبی
- چاپ اول: ۱۳۸۰
- لیتوگرافی: پام مهر ۷۵۰۰۹۳۰
- چاپ: چکاد
- تیراژ: ۲۵۰۰
- قیمت: ۴۵۰ تومان
- شابک: ۹۶۴-۷۱۷۴-۶-۸
- انتشارات آرویج: خیابان شریعتی، بالاتر از سه‌راه طالقانی،
خیابان شهید کارگر، پلاک ۱۲، تلفن: ۷۵۲۵۱۶۵، تلفکس:
۷۵۳۷۰۷۶

A Bidirectional HTML Dialog (Mozilla)

The HTML/CSS markup generating the two forms below is exactly the same, except for:

- The Persian form has `lang=fa` once for the body tag.
- The Persian form has `direction: rtl` once for the body tag.
- The English form, the labels all have `align: right`, in the Persian one they have `align: left`.

Unfortunately there is no way in CSS to drop the annoying align difference.

User:	behdad	کاربر:	بهداد
Password:	***	گذرواژه:	***
Language:	English ▾	زبان:	فارسی ▾

Current Status -- Microsoft Windows

- Render correctly
- Shipped fonts work
- Keyboard layout is terrible
- No Persian digits support
- No Iranian calendar
- Locale data is wrong in places
- No interface translation
- Not trivial to enable Persian support

Current Status -- Linux

- Important systems support rendering
- No good fonts yet
- Standard keyboard layout
- No Persian digits support yet
- KDE claims Iranian calendar support
- Some interface translation done
- Not trivial to enable Persian support

Current Status -- MacOS

- Supports rendering
- No good fonts
- Legacy and standard keyboard layouts

References and Resources

- The Unicode Standard at <http://www.unicode.org/>
- Institute of Standards and Industrial Research of Iran at <http://www.isiri.com> (documents in Persian)
- The FarsiWeb Project at <http://www.farsiweb.info/>
- PersianComputing list at <http://lists.sharif.edu/mailman/listinfo/persiancomputing>
- Typing Persian Word Documents with Windows Tutorial at <http://students.washington.edu/irina/persianword/persianwp.htm>