Minimum(ish) Working Example of Quarto with Arabic RTL text in an LTR document

The Authors

2024-11-16

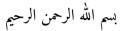
© 2023 Author Names

And Amirican and A

Table of contents

Pr	eface	1
1	Introduction	3
2	Arabic (العربية) support 2.1 In-line Arabic: spans 2.2 Arabic block text: divs 2.3 Pandoc Lua fiters 2.4 Arabic fonts 2.4.1 Specify Arabic font for HTML 2.4.2 Specify Arabic font for PDF	5 6 8 10 10
3	Transliteration of Arabic	13
4	Summary	17
Re	eferences	19

Preface



Quarto is a document publishing software. With it, you can write your document in Pandoc flavored markdown. Quarto will use Pandoc under the hood, and do a bunch of other fancy stuff, to output your markdown document in formats of your choice, like HTML for websites, and PDF (via Latex).

So far so good. But many of my documents are in English with Arabic content interspersed. Arabic is written right-to-left (RTL) whereas English is written left-to-right (LTR). The support of bidirectional (BiDi) text is a notoriously tricky problem. The cursive property of the Arabic script (with joining letters) compounds the issue.

In this write-up, I will describe how to configure Quarto to achieve this.

The code for this book can be used as a template for RTL document projects. Along with BiDi, I'll also discuss other aspects like fonts, figures, etc. إِنْ

The source code for this book can be found here: https://github.com/ada miturabi/quarto-arabic-mwe.

The rendered output is published here: https://adamiturabi.github.io/qu arto-arabic-mwe.

A PDF version of this document can be downloaded if you click on the PDF icon next to the title at the top left of this page.

1 Introduction

In order to create a Quarto book project (like this one), use this command:

quarto create project book mybook

It will create a bunch of files that will be needed for the book.

This procedure is described in detail here: https://quarto.org/docs/books/.

2 Arabic (العربية) support

There are two main ways to insert RTL Arabic text in an LTR document:

- 1. In-line Arabic within an LTR paragraph.
- 2. A block of Arabic text by itself

2.1 In-line Arabic: spans

For inline Arabic, we will use a Pandoc span. A span is written using this syntax:

```
[This is the span's *content text*]{.classname attributekey="attributeval"}
```

Within square brackets [] is the content of the span. This is what will be rendered in the output. Within the curly braces {} is a class name and some attributes that are needed by Quarto to properly process the span.

In order to render the Arabic content text correctly for HTML output, the span is typed thus in the .qmd source file. (Note that the following code listing does not render correctly in the PDF output, exemplifying how difficult BiDi support is. We haven't attempted to find a workaround for this.)

```
ar-txt dir="rtl" lang="ar"} نص [هذا
```

The class name is arbitrary. We suggest a descriptive name. We will be using it in the CSS for selecting the font later. The output HTML code will be something like:

```
<span class="reg-ar-txt" dir="rtl" lang="ar">>ه زا "span class="reg-ar-txt" dir="rtl" lang="ar
```

For PDF output, the dir="rtl" attribute is unneeded and actually clashes with the Xelatex PDF engine that Quarto mandates we use for documents with RTL text. So the span will need to be typed thus in the .gmd source SKUDY file:

```
ar-txt lang="ar"}-عربی.]{\reg.} نص [هذا
```

The output Latex code will be something like:

```
عربی.} نص هذا}{foreignlanguage{arabic}} نص هذا
```

Under the hood, \foreignlanguage is a command that is used by the Latex package babel that Pandoc specifies in its Latex template for handling multiple languages and their scripts.

Finally, this is an example of an English sentence with inline Arabic text within it. Locate this sentence in the source code file here to see نَصُّ عَرَبيُّ. how we wrote it.

2.2 Arabic block text: divs

In order to write a block (paragraph) of Arabic text within an LTR document we will use a Pandoc div. A div is written using this syntax:

```
:::{.classname attributekey="attributeval"}
This is the divs's *content text*.
:::
```

For HTML output, a div is typed thus in the .gmd source:

```
:::{.reg-ar-txt dir="rtl" lang="ar"}
سطرين. النص يبلغ حتى أكتب أن أريد طويل. عربي كلام هذا
الخارجي. الملف لإنتاج قوارطو برنامج أستعمل
قبل. من أستعمله كنت الذي بكداؤن البرنامج خلف قد جيد برنامج هو
```

§2.2 7

The class name reg-ar-txt is, again, arbitrary. The output HTML code will be:

For PDF output, a div is typed thus in the .qmd source:

```
:::(.otherlanguage data-latex="{arabic}" lang='ar'}
طويل. عربي كلام هذا
سطرين. النص يبلغ حتى أكتب أن أريد
الخارجي. الملف لإنتاج قوارطو برنامج أستعمل
قبل. من أستعمله كنت الذي بكداؤن البرنامج خلف قد جيد برنامج هو
```

In this case, the class name otherlanguage is not arbitrary. Furthermore another attribute data-latex="{arabic}" is also needed. And as with spans lang="ar" is needed but dir="rtl" should not be used. The output Latex code is:

\begin{otherlanguage}{arabic}

```
طويل. عربي كلام هذا
سطرين. النص يبلغ حتى أكتب أن أريد
الخارجي. الملف لإنتاج قوارطو برنامج أستعمل
قبل. من أستعمله كنت الذي بكداؤن البرنامج خلف قد جيد برنامج هو
```

\end{otherlanguage}

Finally, this is an example of an Arabic div. Locate it in the source code file here to see how we wrote it.

هذا كلام عربي طويل. أريد أن أكتب حتى يبلغ النص سطرين. أستعمل برنامج قوارطو لإنتاج الملف الخارجي. هو برنامج جيد قد خلف البرنامج بكداؤن الذي كنت أستعمله من قبل.

end

2.3 Pandoc Lua fiters

As you can see, the process for typing Arabic text is both lengthy, and different for HTML and PDF outputs. In order to simplify it, we can use Pandoc Lua filters.

We have created a Quarto filter extension (which is a grouping of Lua filters) to support Arabic divs and spans. The process for creating a Quarto filter extension is detailed here: https://quarto.org/docs/extensions/filters.html

This is the filter inline-arabic-span.lua that we wrote for handling Arabic spans:

```
-- Add attributes for Arabic text in a span
function Span (el)
 if el.classes:includes 'ar' or el.classes:includes 'aralt' then
    text = pandoc.utils.stringify(el)
    contents = {pandoc.Str(text)}
    if FORMAT:match 'latex' then
      -- for handling alternate Arabic font
      if el.classes:includes 'aralt' then
     -- can't seem to use string concatenate directly. Have to use RawInline
        table.insert(
          contents, 1,
          pandoc.RawInline('latex', '\\altfamily ')
        )
      end
    -- no dir needed for babel and throws error if it sees dir attribute. was previo
      return pandoc.Span(contents, {lang='ar'})
    elseif FORMAT: match 'html' then
      classval = 'reg-ar-text'
      if el.classes:includes 'aralt' then
       classval = 'alt-ar-text'
      end
     - dir needed for html otherwise punctuation gets messed up
    return pandoc.Span(contents, {class=classval, lang='ar', dir='rtl'})
  end
```

§2.3 9

This is the filter arabic-div.lua that we wrote for handling Arabic divs:

```
-- Add attributes for Arabic text in a div
function Div (el)
 if el.classes:includes 'ar' or el.classes:includes 'aralt' then
    text = pandoc.utils.stringify(el)
    contents = {pandoc.Str(text)}
    if FORMAT:match 'latex' then
      -- for handling alternate Arabic font
      if el.classes:includes 'aralt' then
     -- can't seem to use string concatenate directly. Have to use RawInline
        table.insert(
          contents, 1,
          pandoc.RawInline('latex', '\\altfamily ')
        )
    -- no dir needed for babel and throws error if it sees dir attribute. was p
    return pandoc.Div(contents, {class='otherlanguage', data latex="{arabic
    elseif FORMAT:match 'html' then
      classval = 'reg-ar-text'
      if el.classes:includes 'aralt' then
        classval = 'alt-ar-text'
      end
    -- dir needed for html otherwise punctuation gets messed up
    return pandoc.Div(contents, {class=classval, lang='ar', dir='rtl'})
  end
end
```

With activating these filters, now you can use Arabic divs and spans using a simplified syntax.

Input for an Arabic span:

```
عربی.]{.ar} نص[هذا
```

Input for an Arabic div:

```
:::{.ar}
```

```
طويل. عربي كلام هذا
سطرين. النص يبلغ حتى أكتب أن أريد
الخارجي. الملف لإنتاج قوارطو برنامج أستعمل
قبل. من أستعمله كنت الذي بكداؤن البرنامج خلف قد جيد برنامج هو
:::
```

The filters will process them correctly for HTML and PDF output. Note that the class name reg-ar-text is hardcoded in the filter. If you wish to modify it you can edit the Lua files directly.

2.4 Arabic fonts

You can use a specific font for the Arabic text which is different from the font used for the English text. This is usually desirable because the type-face design for the Latin font often does not optimize (or even sometimes support) an Arabic font.

For my project, I am using the Vazirmatn and Amiri fonts.

Both of these are well designed fonts. For me a major consideration is the correct typesetting of the hamza character (ε). See here for what I'm talking about: https://adamiturabi.github.io/hamza-rules/#typographical-limitations

To specify the Arabic fonts the process is different for HTML vs PDF output. We'll describe both below:

2.4.1 Specify Arabic font for HTML

For HTML output, the Arabic font is specified in the CSS file. The class name that we selected previously reg-ar-text is now assigned a font:

```
.reg-ar-text {
  font-family: Vazirmatn, serif;
  /* scaled up slightly w.r.t. the Latin font for readability */
  font-size: 1.2em;
  /* line spacing not scaled for visual congruence at the expense of clashes */
  line-height: 100%;
```

§2.4 11

```
}
.alt-ar-text {
  font-family: AmiriWeb, serif;
  font-size: 1.2em;
  line-height: 100%;
}
```

The font names Vazirmatn and AmiriWeb are defined in the same CSS file. See our CSS file for details.

Amiri is specified as an alternate font. It can be specified in the .qmd source with {.aralt} instead of {.ar}. You can also see how we handle it in the source code for the Lua filters above.

Here is an example of a div and a span in the alternate Arabic font. Span: هذا نص عربي،

Div:

By the way, I am not, by any means, at expert (or even proficient) in CSS, so if you see any problems with this method of specifying the font, feel free to let me know in the discussions or issues of the Github page for this book.

2.4.2 Specify Arabic font for PDF

As we mentioned earlier, Latex uses the babel package to handle multiple languages. In order to specify Arabic font(s), we need to add the following lines in the intermediate .tex file produced by Quarto:

```
\babelfont[arabic]{rm}[Language=Default]{Vazirmatn-Light}
\babelfont[arabic]{sf}[Language=Default]{Vazirmatn-Light}
\babelfont[arabic]{alt}[Language=Default]{Amiri}
```

Quarto provides hooks for inserting such additional code using includes and templates.

details, the property of the p The above lines of code need to be inserted at a specific point after the usepackage{babel} line. We found that replacing the partial template

3 Transliteration of Arabic

In my work, I frequently need to transliterate and transcribe Arabic text in Latin characters. There are various Romanization schemes in existence, using dots, macrons, etc. The Romanization scheme I am using for my work is tabulated below:

	Arabic letter	Transliterated output	ASCII input
	۶)	E
	I	ā	Α
	ب	b	b
	ت	t ,	t
	ث	t th	V
	ث ج ح د خ	j	j
	ح	ḥ kh	Н
	خ	kĥ	X
	2	d	d
	ذ	d dh	p
)	r	r
	j	Z	Z
	w O	$\frac{s}{sh}$	S
	ش	sh	С
	ص	Ş	S
A	ض	ģ	D
	ط	ţ	T
	ظ	ș d ț đh	Р
1	س ش ص ش ظ ط ض ص ق ك ق ق غ ع ظ	ε gh f	е
	غ	gĥ	g
407	فَ	f	f
	ق	q	q
	ك	k	k
▼	J	1	l
	م	m	m

Arabic letter	Transliterated output	ASCII input
Ů	n	n
٥	h	h
(C) و	W	W
(C) و (C) ي	y	У
。 (V)	ū	U
(V) ي	Ī	I

As you can see, I use digraphs \widehat{dh} , \widehat{gh} , etc. for some letters. This is because for my current work, I value readability over precision.

It is possible to input these special characters directly by modifying your keyboard layout or mapping, either at a operating system, or editor level. Andreas Hallberg has described a technique for inputing them in the vim editor here: https://andreasmhallberg.github.io/ergonomic-arabic-transcription/

For Quarto, I prefer to input the transliterated text as ASCII characters and write a Lua filter transliteration-span.lua to handle rendering them correctly. The ASCII to transliterated output mapping is shown in the table above. So if I type:

```
[pahabtu maphaban]{.trn}
```

It will be output as dhahabtu madh haban.

Note the dot character \cdot is automatically inserted by the filter between the digraph dh and the following h for helping in disambiguation.

With {.trn} the output is in italic (as above). But sometimes I need to have non-italic output, as in the case of names. For that I use {.trn2}. For example:

```
[#eAEicah]{.trn2} and [E#Adam]{.trn2} are studying
the [#qurEAn]{.trn2} and [#HadIv]{.trn2}.
```

This is rendered as:

Eā'ishah and 'Ādam are studying the Qur'ān and Hadīth.

Note how the hash character # is used to control capitalization.

ack a font alteration. I

4 Summary

riterate programs.

References

Moth in Progress. Not really for