DRAFT

ArabTeXa System for Typesetting Arabic

User Manual Version 3.09 $^{\rm 1}$ $^{\rm 2}$

Klaus Lagally

July 22, 1999

 $^{^1\}mathrm{Report~Nr.~1998/09},$ Universität Stuttgart, Fakultät Informatik, Breitwiesenstraße 20–22, 70565 Stuttgart, Germany $^2\mathrm{This~Report~supersedes~Reports~Nr.~1992/06}$ and 1993/11

Overview

ArabTEX is a package extending the capabilities of TEX/ETEX to generate the Perso-Arabic writing from an ASCII transliteration for texts in several languages using the Arabic script. It consists of a TEX macro package and an Arabic font in several sizes, presently only available in the Naskhi style. ArabTEX will run with Plain TEX and also with $\text{ETEX}2_{\varepsilon}$. It is compatible with Babel, CJK, the EDMAC package, and PicTEX (with some restrictions); other additions to TEX have not been tried.

ArabTEX is primarily intended for generating the Arabic writing, but the standard scientific transliteration can also be easily produced. For languages other than Arabic that are customarily written in extensions of the Perso-Arabic script some limited support is available.

ArabTEX defines its own input notation which is both machine, and human, readable, and suited for electronic transmission and E-Mail communication. However, texts in many of the Arabic standard encodings can also be processed.

Starting with Version 3.02, ArabTEX also provides support for fully vowelized Hebrew, both in its private ASCII input notation and in several other popular encodings.

ArabTEX is copyrighted, but free use for scientific, experimental and other strictly private, noncommercial purposes is granted. Offprints of scientific publications using ArabTEX are welcome. Using ArabTEX otherwise requires a license agreement. There is no warranty of any kind, either expressed or implied. The entire risk as to the quality and performance rests with the user.

Please send error reports, suggestions and inquiries to the author:

Prof. Klaus Lagally Institut für Informatik Universität Stuttgart Breitwiesenstraße 20-22 70565 Stuttgart GERMANY

lagally@informatik.uni-stuttgart.de

Copyright © 1992–1999, Klaus Lagally

Contents

| 1 | Inti | roduction to ArabT _E X | 7 | | | | | | | | | | | | |
|---|------|----------------------------------------------|-----------|--|--|--|--|--|--|--|--|--|--|--|--|
| 2 | Inp | nput to ArabT _E X | | | | | | | | | | | | | |
| | 2.1 | Arabic text elements | 12 | | | | | | | | | | | | |
| | 2.2 | Commands within an Arabic context | 13 | | | | | | | | | | | | |
| 3 | Rui | nning ArabT _E X | 15 | | | | | | | | | | | | |
| | 3.1 | Activating ArabTEX | 15 | | | | | | | | | | | | |
| | 3.2 | Language selection | 15 | | | | | | | | | | | | |
| | 3.3 | Font selection | 16 | | | | | | | | | | | | |
| 4 | Inp | ut encoding conventions | 18 | | | | | | | | | | | | |
| | 4.1 | ASCII Transliteration encoding | 18 | | | | | | | | | | | | |
| | | 4.1.1 Standard Arabic and Persian characters | 18 | | | | | | | | | | | | |
| | | 4.1.2 Vowelization | 21 | | | | | | | | | | | | |
| | | 4.1.3 Quoting | 22 | | | | | | | | | | | | |
| | | 4.1.4 Ligatures | 23 | | | | | | | | | | | | |
| | | 4.1.5 Coding examples for Arabic | 23 | | | | | | | | | | | | |
| | 4.2 | Verbatim input | 28 | | | | | | | | | | | | |
| | 4.3 | Alternate input encodings | 29 | | | | | | | | | | | | |
| | | 4.3.1 ASMO 449 = ISO 9036 | 29 | | | | | | | | | | | | |
| | | $4.3.2$ ASMO $449E = ISO 8859 - 6 \dots$ | 31 | | | | | | | | | | | | |
| | | 4.3.3 CP 1256 = Arabic Windows Encoding | 33 | | | | | | | | | | | | |
| | | 4.3.4 ISIRI 3342 | 33 | | | | | | | | | | | | |

CONTENTS 2

| | | 4.3.5 UNICODE Arabic | 36 |
|---|-----|-------------------------------------------------|-----------|
| 5 | Tra | nsliteration 3 | 39 |
| | 5.1 | ZDMG transliteration style | 39 |
| | 5.2 | Other transliteration styles | 40 |
| | 5.3 | Capitalization | 40 |
| 6 | Sup | oport for other languages | 41 |
| | 6.1 | Persian (Farsi, Dari), also Ottoman and Kurdish | 41 |
| | | 6.1.1 Coding examples for Persian | 42 |
| | 6.2 | Maghribi | 44 |
| | 6.3 | Urdu | 44 |
| | | 6.3.1 Coding examples for Urdū | 45 |
| | 6.4 | Pashto (Afghanic) | 48 |
| | 6.5 | Sindhi | 49 |
| | 6.6 | Kashmiri | 50 |
| | 6.7 | Uighuric | 50 |
| | 6.8 | Old Malay | 52 |
| | 6.9 | Other extensions of the Perso-Arabic script | 52 |
| 7 | Heb | orew mode | 53 |
| | 7.1 | Language switching | 53 |
| | 7.2 | Standard Hebrew encoding | 54 |
| | 7.3 | ISO 8859-8 and Hebrew MS-Windows | 55 |
| | 7.4 | HED, PC "oldcode" and "newcode" | 55 |
| | 7.5 | BHS encoding | 58 |
| | 7.6 | UNICODE Hebrew | 59 |
| | 7.7 | Hebrew fonts | 59 |
| | 7.8 | Hebrew transcription systems | 61 |
| 8 | Mis | cellaneous features | 62 |
| | 8.1 | Additional codings | 62 |
| | 8.2 | Dots on $y\bar{a}'$ | 63 |

CONTENTS 3

| | 8.3 | Vowel positioning | 63 | | | | | | | |
|----|-----------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------|--|--|--|--|--|--|--|
| | 8.4 | Abjad numerals | 63 | | | | | | | |
| | 8.5 | Automatic stretching | 63 | | | | | | | |
| | 8.6 | Uniform baselines | 64 | | | | | | | |
| | 8.7 | Verbatim copy of the input | 64 | | | | | | | |
| | 8.8 | Progress report | 64 | | | | | | | |
| | 8.9 | Module Reporting | 64 | | | | | | | |
| 9 | Con | mpatibility issues | 65 | | | | | | | |
| | 9.1 | Arabic document classes | 66 | | | | | | | |
| | 9.2 | Using ArabTeX with EDMAC | 66 | | | | | | | |
| | 9.3 | Using ArabTEX with Babel | 67 | | | | | | | |
| | 9.4 | Using ArabTEX with PicTEX \hdots | 67 | | | | | | | |
| | 9.5 | Using ArabTEX with CJK | 67 | | | | | | | |
| 10 | Ack | Abjad numerals | | | | | | | | |
| | | | | | | | | | | |
| Re | efere | nces | 68 | | | | | | | |
| | | | | | | | | | | |
| | | caining and installing ArabTeX | 71 | | | | | | | |
| | Obt | caining and installing ArabTEX Obtaining ArabTEX | 71 71 | | | | | | | |
| A | Obt A.1 A.2 | Caining and installing ArabTEX Obtaining ArabTEX | 71 71 72 | | | | | | | |
| A | Obt A.1 A.2 | Caining and installing ArabTEX Obtaining ArabTEX Installing ArabTEX ease history | 71 71 72 73 | | | | | | | |
| A | Obt A.1 A.2 Rele | Caining and installing ArabTEX Obtaining ArabTEX Installing ArabTEX ease history ArabTEX version 1.00 | 71 72 73 73 | | | | | | | |
| A | Obt A.1 A.2 Reld B.1 | Caining and installing ArabTEX Obtaining ArabTEX Installing ArabTEX ease history ArabTEX version 1.00 ArabTEX version 2.00 | 71 71 72 73 73 74 | | | | | | | |
| A | Obt A.1 A.2 Rele B.1 B.2 | Caining and installing ArabTEX Obtaining ArabTEX Installing ArabTEX ease history ArabTEX version 1.00 ArabTEX version 2.00 ArabTEX version 3.00 | 71 71 72 73 73 74 74 | | | | | | | |
| A | A.1 A.2 Rele B.1 B.2 B.3 B.4 | Caining and installing ArabTEX Obtaining ArabTEX Installing ArabTEX ease history ArabTEX version 1.00 ArabTEX version 2.00 ArabTEX version 3.00 ArabTEX version 4.00 | 71 71 72 73 73 74 74 75 | | | | | | | |
| В | A.1 A.2 Rele B.1 B.2 B.3 B.4 | Caining and installing ArabTEX Obtaining ArabTEX Installing ArabTEX ease history ArabTEX version 1.00 ArabTEX version 2.00 ArabTEX version 3.00 ArabTEX version 4.00 cellaneous utilities | 71 71 72 73 73 74 74 75 | | | | | | | |
| В | A.1 A.2 Reld B.1 B.2 B.3 B.4 Mis | Caining and installing ArabTEX Obtaining ArabTEX Installing ArabTEX ease history ArabTEX version 1.00 ArabTEX version 2.00 ArabTEX version 3.00 ArabTEX version 4.00 cellaneous utilities verses.sty | 71 71 72 73 73 74 74 75 76 | | | | | | | |

| CONTENTS | 4 |
|----------|---|
| | |

Index 78

List of Figures

| 1.1 | $Sample \ ArabT_{\hbox{\it E}}X \ input \ \ldots \ldots \ldots \ldots \ldots \ldots$ | 9 |
|-----|--------------------------------------------------------------------------------------|----|
| 1.2 | Sample ArabTEX output | 10 |
| 7.1 | Hebrew example | 58 |

List of Tables

| 4.1 | Standard encodings for Arabic and Persian consonants | 19 |
|-----|------------------------------------------------------|----|
| 4.2 | Additional encodings generally available | 20 |
| 4.3 | Verbatim encodings for the carrier of $hamza$ | 29 |
| 4.4 | ASMO 449 code table | 30 |
| 4.5 | ISO 8859-6 code table | 32 |
| 4.6 | Windows CP 1256 code table | 34 |
| 4.7 | ISIRI 3342 code table | 35 |
| 4.8 | UNICODE Arabic, Part 1 | 37 |
| 4.9 | UNICODE Arabic, Part 2 | 38 |
| 6.1 | The Urdu Alphabet | 46 |
| 6.2 | Additional codings for Pashto | 48 |
| 6.3 | The Sindhi Alphabet | 49 |
| 6.4 | The Kashmiri Alphabet | 50 |
| 6.5 | ArabTEX encoding of Uighuric | 51 |
| 7.1 | ISO 8859-8 and Windows CP 1255 code table | 56 |
| 7.2 | HED, CP 1255 and ISO 8859-8 code table | 57 |
| 7.3 | UNICODE Hebrew | 60 |
| 8.1 | Additional codings for special purposes. | 62 |

Chapter 1

Introduction to ArabT_EX

Note: This manual describes Version 4 of Arab T_EX , which in general is upwards compatible to earlier versions. Incompatible changes are flagged by an asterisk in the margin.

ArabTEX is a package extending the capabilities of TEX/IATEX to generate an extended Perso-Arabic and/or Hebrew script in addition to the customary left-to-right scripts (called *Roman* in the sequel). Besides Arabic and Hebrew there are provisions for many other languages that use some extensions of the Arabic script; we shall use the term *Arabic* generically to denote any supported right-to-left script, without any cultural or political connotation.

An ArabTeX document is usually multi-lingual and may contain right-to-left insertions within left-to-right paragraphs, and vice versa. There are various possibilities to encode right-to-left insertions: in addition to several standard encodings there are also 7-bit ASCII encodings modelled after various transliteration standards (there is more than one convention, so the intended language must be specified).

ArabTEX, like TEX and LATEX, is not a substitute for a word processor, and does not follow the WYSIWYG paradigm (What You See Is What You Get) where the author has all influence on, and therefore is also completely responsible for, every detail of the visual representation of his text. Instead it is an off-line system mirroring the classical publication process where an author, concentrating on the content and the structure of her paper only, will deliver a manuscript to the publisher, who will take care of a reasonable visual rendering of the text.

The publisher involved in our case is TEX, the famous typesetting program by D. E. Knuth; and the amount of technical typesetting knowledge embodied in the algorithms of TEX is such that, in general, the author will tacitly accept the resulting output if she had stated her preferences in a sufficiently precise way; if not, the formatting task has to be repeated under changed specifications, as

also happens in real life.

TEX is easily extensible by providing its own macro programming language, enabling the user to implement packages that contain algorithms embodying additional typesetting knowledge. A prominent example is LATEX by L. Lamport, a package providing standard formatting rules for several common document types. One of the basic concepts of LATEX, but already contained in TEX itself, is "structured markup": the user generally only indicates the logical structure of the document, whereas the formatting rules are laid down within descriptions of standard document classes.

ArabTEX, at the user interface, needs to add only very few features: we provide a way to indicate which parts of the input text are considered to be in Arabic (or some other supported language written from right to left) and thus have to be rendered accordingly, and a way of setting parameters influencing the rendering of Arabic text. Additionally we have to support a representation of the Arabic text that can be handled using the available standard equipment.

Internally, of course, there is much more to it. TeX may be very good at rendering left-to-right texts, but Arabic runs from right to left, so we have to instruct TeX to do it the other way; and if we want to mix Arabic and traditional left-to-right text within the same paragraph, things can get extremely complicated. This shows indeed: ArabTeX is very large and comparatively slow. Fortunately, computer technology is advancing at a high pace, so that the penalty of using very expensive algorithms will decrease just by simple waiting; and we thus never were too much concerned about efficiency (even if we tried to avoid the worst blunders).

In order to see how ArabTEX works, let us consider a very simple example: Figure 1.1 contains the complete TEX input text, Figure 1.2 the corresponding output. In this example (using LATEX) we use a standard "article" format, we activate the ArabTEX package, and we set a few options: our text is supposed to be in Arabic, we want to see all vowel marks (this is not always done in Arabic printing), we additionally want to get the scientific transcription, and, of course, the Arabic writing. The document proper consists of a centered headline and a sequence of short Arabic paragraphs, separated in the input by blank lines¹.

This example is *not* typical insofar as we produce both the scientific transcription of the text, and the Arabic writing, from the same input at the same time, interleaving them. This is rarely required; but in our example it allows us to demonstrate that the input notation is very closely related to the transcription. Details about this correspondence are covered in Chapter 4. What indeed *is* typical is the fact that, apart from the centered headline, we supplied no formatting information whatever; TEX will take care of that.

 $^{^1}$ If you happen to be curious and cannot read Arabic: the text contains a traditional story about a somewhat silly person named $\check{\text{G}}\text{u}\bar{\text{h}}\bar{\text{a}}$, trying *not* to lend his donkey to a friend, and failing.

```
\documentclass[12pt]{article}
\usepackage{arabtex}
\begin{document}
\setarab
           % choose the language specific conventions
\vocalize % switch diacritics for short vowels on
\ttranstrue \% additionally switch on the transliteration
\arabtrue % print arabic text ... is on by default anyway
\centerline {\RL{^gu.hA wa-.himAruhu}}
\begin{RLtext}
'at_A .sadIquN 'il_A ^gu.hA ya.tlubu minhu .himArahu li-yarkabahu
fI safraTiN qa.sIraTiN wa-qAla lahu:
sawfa 'u'Iduhu 'ilayka fI al-masA'i, wa-'adfa'u laka 'u^graTaN.
fa-qAla ^gu.hA:
'anA 'AsifuN ^giddaN 'annI lA 'asta.tI'u 'an 'u.haqqiqa
laka ra.gbataka, fa-al.himAru laysa hunA al-yawma.
wa-qabla 'an yutimmu ^gu.hA kalAmahu
bada'a al-.himAru yanhaqu fI i.s.tablihi.
fa-qAla lahu .sadIquhu:
'innI 'asma'u .himAraka yA ^gu.hA yanhaqu.
fa-qAla lahu ^gu.hA:
.garIbuN 'amruka yA .sadIqI!
'a-tu.saddiqu al-.himAra wa-tuka_d_dibunI?
\end{RLtext}
\end{document}
```

Figure 1.1: Sample ArabT_EX input

ğuḥā wa-ḥimāruhu تحكا وَحَمَارُهُ

Patā ṣadīqun Pilā ǧuḥā yaṭlubu minhu ḥimārahu li-yarkabahu fī-safratin qaṣīratin wa-qāla lahu:

أَتَى صَدِيقٌ إِلَى جُحَا يَطلُبُ مِنهُ حِمَارَهُ لِيَركَبَهُ فِي سَفرَةٍ قَصِيرَةٍ وَقَالَ لَهُ: sawfa 'u'īduhu 'ilayka fī 'l-maṣā'i, wa-'adfa'u laka 'uğratan.

سَوفَ أُعِيدُهُ إِلَيكَ فِي المَسَاءِ ، وَأَدفَعُ لَكَ أُجرَةً.

fa-qāla ğuḥā:

فَقَالَ مُجَا:

'anā 'āsifun ǧiddan 'annī lā 'astaṭī'u 'an 'uḥaqqiqa laka raġbataka, fa-'lḥimāru laysa hunā 'l-yawma.

أَنَا آسِفٌ جِدًّا أَنِّي لَا أَستَطِيعُ أَن أُحَقِّقَ لَكَ رَغَبَتَكَ ' فَالحِمَارُ لَيسَ هُنَا اليَومَ ' wa-qabla 'an yutimmu ğuḥā kalāmahu bada'a 'l-ḥimāru yanhaqu fī 'stablihi.

وَقَبَلَ أَن يُتِمُّ مُجَا كَلَامَهُ بَدَأً الحِمَارُ يَهَقُ في اصطلبهِ.

fa-qāla lahu ṣadīquhu:

فَقَالَ لَهُ صَديقُهُ:

innī iasma'u ḥimāraka yā ǧuḥā yanhaqu.

إِنِّي أَسْمَعُ حِمَارَكَ يَا نَجَا يَنْهَقُ.

fa-qāla lahu ǧuḥā:

فَقَالَ لَهُ حُحَا:

 $\dot{g}arar{\imath}bun$ 'amruka y $ar{a}$ ṣad $ar{\imath}qar{\imath}!$ 'a-tuṣaddiqu 'l-ḥim $ar{a}$ ra wa-tuka \underline{d} dibun $ar{\imath}?$

Figure 1.2: Sample ArabTFX output

Chapter 2

Input to ArabT_EX

ArabTEX is activated by the command \input arabtex (with Plain TEX) or \usepackage{arabtex} (with LATEX). After activating ArabTEX, select one of the language-specific Arabic writing styles, e.g., \setarab (see Chapter 3.2). Your modified TEX/LATEX system will recognize the following items:

- standard TEX/LATEX text and commands,
- Arabic quotations ¹ as arguments to the command \RL{ } (read: "right-to-left") within a Roman paragraph. A quotation may also be bracketed by < and > (or with \< and > except inside a IATEX {tabbing} environment)². An Arabic quotation forms a new group, so any assignments will be local by default.
- longer Arabic text segments called *Arabic environments*, which are bracketed by the commands \begin{RLtext} and \end{RLtext} (or also *\begin{arabtext} and \end{arabtext} with the same meaning) (even when using Plain TEX!), An *Arabic environment* consists of one or more paragraphs separated by blank lines or \par commands. It forms a group, so assignments will be local by default.

Arabic quotations and Arabic environments will be called Arabic contexts in the sequel.

 $^{^{1}}$ The former restriction that quotations must fit on the current line no more applies.

 $^{^2 \}mbox{Quotations}$ closed by > may not contain nested insertions; also observe that \< must be matched by >, not by \> !

2.1 Arabic text elements

Every *Arabic paragraph* and every *Arabic quotation* is a sequence of the following kinds of *Arabic items*, separated by blank spaces or newlines:

- isolated punctuation marks, interpreted as the corresponding Arabic punctuation mark;
- "numbers", i.e. character sequences starting with a digit, and possibly continued by digits, commas, hyphens, or slashes. A "number" will be processed using the normal writing sequence from left to right; however, a final punctuation mark will be split off and processed separately.
- "Arabic quotes" coded as two left quotes or two right quotes each, or as \lq and \rq. They should be written directly adjacent to a word.
- "words", i.e. character sequences starting with a letter or a special (non-digit) character followed by a letter. A final punctuation mark will be split off and processed separately. The (coded) characters of a word will in the output be arranged from right to left.
- a sequence of Arabic text elements (words, numbers, and special characters) enclosed in curly braces { and } . This introduces a new level of TeX grouping; otherwise the constituents are processed normally. This feature may be nested.

Output from all items will be arranged from right to left, lines will be broken as necessary.

Inside an Arabic Environment, or in an Arabic quotation, you may also have:

- ArabTeX commands with or without parameters. These will be executed immediately.
- Some, but not all, TeX/IATeX commands (see below). These will be executed immediately.
- Short mathematical insertions, bracketed by *single* \$ signs. They must fit on one output line and are processed as usual. TeX display mode (bracketed by \$\$) is not provided within an *Arabic environment*; if it is required, the user has to leave the *Arabic environment* temporarily.
- short left-to right ("Roman") quotations, containing text and possibly also TEX/IFTEX commands, as argument to \LR{ } (read: "left-to-right") or bracketed by < and >3. These introduce a new level of grouping, so if they contain any TEX/IFTEX assignments the effects of these will be local by default. When using "<" and ">" this feature is not available within an Arabic quotation. The alternate notation \< is not provided.

 $^{^3\}mathrm{Quotations}$ closed by > may not contain nested insertions.

2.2 Commands within an Arabic context

A control sequence inside an *Arabic context* should be separated from the preceding text item by white space or another control sequence, and may be of the following kinds:

- ArabTEX option changing commands. These may also be used outside an *Arabic Context*, and usually follow TEX's grouping rules.
- \\ for a line break; the current line will be padded out on the left with spaces.
- \| or \break for a line break; the current line will be spread out. If it comes out very badly spaced, automatic stretching might help (see Section 8).
- \indent or \par (or a blank line) for a new paragraph, \noindent for a new paragraph without indentation (not inside Arabic quotations).
- \emphasize { group_of_Arabic_items} will put a bar over the indicated group of Arabic items.
- \setnash, \setnashbf, \setnastaliq and other font selection commands, see Section 3.3.
- size changing LATEX commands like \large etc., only if LATEX is used!
- the following commands: \footnote (observe that the syntax for Plain TEX and LATEX is different!), \marginpar (also with Plain TEX, analogous to the LATEX usage).
- the TEX/IATEX commands \smallskip, \medskip, \bigskip, \input, \hfill, \hfill, \vfill, \u (for a space), \space, \, (small space), \newpage, \clearpage, \pagebreak with their usual meaning.
- \nospace will place the adjacent items in the output in direct contact, without any intervening space, except in case of a line break.
- \hspace{width} will introduce the indicated amount of spacing in the output. The same is true for \vspace, \hskip (observe TEX syntax!), and \vskip.
- \mbox{text} puts the text into a box that will not be split across a line break.
- \spreadbox{width}{text} spreads out the text to the indicated width.
 This may be useful e.g., when typesetting poetry.
 \spreadbox{width}{text\hfill } will inhibit the spreading,
 \spreadbox{width}{\hfill text\hfill } will center the text inside the box.

 $\spreadbox\{width\}\{\hfill \}$ or $\spreadbox\{width\}\{\label{preadbox}\}$ just introduces the indicated amount of horizontal space, as will $\hspace\{width\}$.

If two boxing commands follow each other without any intervening blank space in the input, there will also be no resulting space between the boxes in the output.

- \centerline{text} will start a new line whose contents are centered (not inside Arabic quotations).
- \spreadline{text} will start a new line whose contents are spread out over the whole width of the page (not inside Arabic quotations). It is approximately equivalent to \spreadbox{\hsize }{text}.
- User defined commands whose expansion produces legal ArabTEX input text may be called by \docommand{command_name and parameters}. The command is expanded exactly once, and the expansion text, after suitable substitution of parameters, will be processed by ArabTEX again.
- User defined commands may also be called directly within an Arabic context if they have been previously announced to the ArabTeX processor by \allowarab{command_name}. They are expanded exactly once, and the expansion text (after suitable substitution of parameters) will be processed by ArabTeX again, so it must be legal ArabTeX input text.
- Parameter assignments inside an *Arabic context* may be performed by \doassign{parameter}{value}. The effect is normally local except if the form \doassign{\global parameter}{value} is used.
- Any non-recognized command will generate an error message and will be echoed verbatim in the output. Even though ArabTeX tries hard to get into synchronization again, additional spurious errors may occur.
- inside an Arabic Context normally no further LATEX or ArabTEX environment may be nested; this restriction does not apply to the yet experimental LATEX document classes arabrep.cls, arabart.cls, arabbook.cls which are provided for right-to-left documents.

For a list of all available commands, consult the Index to this report. As a reminder, the command \arabstat will cause a list of all commands that are presently valid inside Arabic text to appear in the TEX log file.

⁴This is no strong restriction as the expansion may contain \docommand calls again.

Chapter 3

Running ArabT_EX

ArabTEX can be used both with Plain TEX and with IATEX, but is activated differently in either case.

3.1 Activating ArabTeX

With Plain T_EX, a small loader program is activated by the command \input arabtex at the beginning of an input text. It will define a default font, prepare a minimal environment simulating the (very few) L^AT_EX-like features needed, and load the ArabT_EX macro files.

With $\LaTeX 2\varepsilon$ the command \usepackage {arabtex} will do all the loading. Users still running $\LaTeX 2.09$ (horror!) should either add the option arabtex to the \documentstyle command, or upgrade to $\LaTeX 2\varepsilon$.

ArabTEX loads many internal files automatically, and defines a large numbers of internal commands. These all contain an "at"-sign (②) within their names and thus should not interfere with user defined commands. Collisions with other macro packages are possible, however, and may lead to surprises and interesting effects.

ArabTEX tries to diagnose the presence of some other packages with which it could run into conflicts, and sometimes locally modifies itself accordingly. For this to be possible, in case of doubt the ArabTEX package should be loaded last.

3.2 Language selection

The processing of input text in ASCII transliteration encoding is somewhat language dependent. Thus before the first *Arabic quotation* or *Arabic environment*

you have to indicate the desired processing mode by one of the language selection commands \setarab, \setfarsi, \seturdu, \setpashto, \setmaghribi, etc., ¹ or \setverb (no special processing; see however Section 4.2). The processing mode may be changed at any time, even inside an *Arabic environment* or an *Arabic quotation*.

Arabic insertions are generally included in \RL{ } or bracketed by \< and >. By selecting a language, the symbols < and > are also activated as shorthands to bracket short insertions in the chosen language. Whereas this is usually convenient, it also has some drawbacks: the angle brackets can thus no more be used for other purposes, except in mathematical mode where they retain their normal meaning as relational operators. To return them to their normal mode of operation, you can deselect them by \setnone.²

For further details on supported languages, see Section 6.

3.3 Font selection

For producing the extended Arabic script ArabTEX uses a special strategy to build up character shapes from a collection of fragments, which normally do *not* correspond to individual character glyphs. Therefore none of the available free or commercial Arabic fonts can be used; we provide our own "pseudo-fonts".

Presently the following pseudo-fonts are available:

- "nash14" is the default,
- "nash14bf" is a bold-face version of "nash14",
- "xnsh14" is an improved version of "nash14" providing additional shape * elements used for some exotic script extensions.
- "xnsh14bf" is a bold-face version of "xnsh14".

"nash14" and "xnsh14" are the default; use \setnashbf to switch to bold-face. \setnash will switch back.

With Plain TEX the fonts are available by default at 14 point size only, which cooperates well with the "cm" fonts at 10 points. Additional sizes are defined within the file "arabtex.tex"; they can be activated whenever needed by the command \setarabfont{font}.

¹We would prefer to use a single switching command like language, \setlanguage, or \selectlanguage, but these names have already been preempted by TEX3 and the Babel package.

²Note for advanced TEX users: All language selecting commands except \setnone set the character < to be active. If Arabic insertions are not needed, or are always started with \< or \RL, the user may reuse the command < for other purposes, or deactivate it by \catcode '\<=12 or \setnone to return it to its normal meaning.

With LATEX, the font size changing commands will also operate on the Arabic fonts.

We strongly recommend migrating to "xnsh14" and "xnsh14bf" by the command \newarabfont; \oldarabfont will switch back, if necessary. The old fonts "nash14" and "nash14bf" are of inferior quality and will be phased out gradually.

All fonts indicated presently are in the Naskhi style; we had started to write a Nastaliq font for Persian and Urdu, but ran into grave implementation problems, yet unsolved.

Due to a donation by Taco Hoekwater, the fonts "xnsh14" and "xnsh14" are also available in Postscript T1 format. Their use is highly recommended when using a Postscript interpreter, or converting the output to PDF format; readability is dramatically improved.

In Hebrew mode (see section 7) we can use the standard fonts available on CTAN (after installing them locally). As defaults the fonts "hclassic" and "hcaption" are provided with ArabTEX; these fonts support vowel points.

Chapter 4

Input encoding conventions

4.1 ASCII Transliteration encoding

The ASCII input notation for Arabic text has been modelled closely after the transliteration standards ISO/R 233 and DIN 31635. These standards do not guarantee unique re-transliteration and are also not 7-bit ASCII compatible, therefore some modifications were necessary. These follow the general rules:

- whenever the transliteration uses a single letter, code that letter;
- whenever the transliteration uses a letter with a diacritical mark, put the punctuation character most closely resembling the diacritical mark before the letter (and not behind it as in some other coding proposals, as otherwise the readability of the input would suffer, and the encoding could become ambiguous).
- use capital letters for writing variants.

4.1.1 Standard Arabic and Persian characters

The standard encodings for Arabic and Persian consonants are given in Table 4.1 and Table 4.2.

- For long vowels, we use the capital letters <A>, <I>, <U> or also <aa>, <iy>, <uw>, with the same meaning.
- To get the defective writing of long vowels, use <_a>, <_i>, <_u>.
- 'alif $maq s \bar{u} r a$ is <_A> or <Y>.

| a | ١ | a | 'alif | ъ | ب | b | $bar{a}'$ | р | پ | p | $par{a}$ |
|----|----|-----------------|------------------|----|---|-----------------|------------------------|----|----|----|--------------------------|
| t | Ç | t | $tar{a}$ ' | _t | ث | \underline{t} | $\underline{t}ar{a}$ ' | ^g | ج | ğ | $\check{g}\bar{\imath}m$ |
| .h | ح | ķ | $\dot{h}ar{a}$ ' | _h | خ | þ | $b ar{a}$ | d | د | d | $dar{a}l$ |
| _d | ٠ | \underline{d} | $d\bar{a}l$ | r | ر | r | $rar{a}$ ' | z | ز | z | $zar{a}y$ |
| s | س | s | $s\bar{\imath}n$ | ^s | ش | š | $\check{s}ar{\imath}n$ | .s | ص | s | $\dot{s}ar{a}d$ |
| .d | ض | ḍ | $d\bar{a}d$ | .t | ط | ţ | $t\bar{a}$ | .z | ظ | z. | $zar{a}$ |
| (| ں | ٧ | 'ayn | .g | غ | \dot{g} | $\dot{g}ayn$ | f | ف | f | $far{a}$ ' |
| q | و: | q | $qar{a}f$ | v | ڤ | v | $var{a}$ ' | k | اك | k | $k\bar{a}f$ |
| g | گ | g | $gar{a}f$ | 1 | J | l | $l\bar{a}m$ | m | م | m | $mar{\imath}m$ |
| n | ن | n | $nar{u}n$ | h | ٥ | h | $har{a}$ | W | و | w | $war{a}w$ |
| у | ي | y | $yar{a}$, | _A | ی | \bar{a} | 'alif | Т | ö | t | $tar{a}$ |
| | | | | | | | $maq s ar{u} r a$ | | | | $marbu \dot{t} a$ |

Table 4.1: Standard encodings for Arabic and Persian consonants.

- The short vowels *fatḥa*, *kasra*, *ḍamma* are coded <a>, <i>, <u> and need not normally be written except in the following cases:
 - at the beginning of a word where they generate 'alif,
 - adjacent to hamza where they will influence its carrier,
 - when the transliteration is required,
 - in the \vocalize and \fullvocalize modes.
- tanwīn is coded <aN>, <iN>, or <uN>. A silent 'alif, if required, is supplied automatically; it may also be explicitly written: <aNA>. Likewise, a silent wāw may be written <NU> as in <'amruNU>.
- hamza is denoted by a single right quote <'>. After selecting the language by \setarab the carrier of hamza will be determined from the context according to the rules for writing Arabic words; if that is not wanted, "quote" the hamza (see Section 4.1.3 below). In the \setverb mode, the carrier of hamza is determined by the following input character; see Section 4.2.

| С | ځ | c | $har{a}$ with $hamza$ |
|----|----|-------------|--------------------------------------------------|
| ^c | چ | č | $\check{g}\bar{\imath}m$ with three dots (below) |
| ,с | څ | ć | $h\bar{a}$ with three dots (above) |
| ^z | ژ | ž | $z\bar{a}y$ with three dots (above) |
| ^n | ڷۓ | \tilde{n} | $k\bar{a}f$ with three dots (Ottoman) |
| ^1 | Ŭ | ĩ | $l\bar{a}m$ with a bow accent (Kurdish) |
| .r | ړ | \dot{r} | $r\bar{a}'$ with a bow below (Kurdish) |

Table 4.2: Additional encodings generally available.

- madda on 'alif is generated by a right quote (hamza) before <A>: <'A>.

 It may also be written <^A>; likewise, <^I> and <^U> will produce madda on yā' and on wāw, as required in some older writing conventions.
- The coding <'> for 'ayn is a single left quote, beware of confusing it with hamza!
- The "invisible consonant" <|> may be inserted in order to break unwanted ligatures and to influence the *hamza* writing. It will not show in the Arabic output or in the transliteration. At the beginning of a word it will suppress a following short vowel; otherwise it acts like a consonant.
- The sequence <\,> will insert a small space, as does <"|> (see Section 4.1.3 * below). The adjacent characters will not be connected.
- *šadda* is indicated by doubling the appropriate letter coding. Therefore two equal consonants in sequence have to be separated by a short vowel indicator or <|> even in \novocalize mode.
- The definite article is separated from the following word by a hyphen. It may be written in the assimilated form (if it exists): as-salaamu, or always as al->; in that case a subsequent "sun letter" must be doubled: al-ssalaamu, to receive a šadda, and to prevent a $suk\bar{u}n$ on the $l\bar{u}m$. The transliteration in both cases is identical.
- Hyphens <-> are used for tying words together and for separating prefixes and the article; in these cases they start a new word. Hyphens can also be used to indicate inflectional endings, a connecting vowel in Arabic, or an

¹The former use of <~A>, <~I>, and <~U> has been discontinued in ArabTEX version 4.

izāfet connection in Persian. Hyphens will show up in the transliteration. Additionally, at the beginning and/or the end of an otherwise isolated word they enforce the use of the connecting form of the adjacent letter (if it exists), like e.g. in the date <1400 h->.

• A double hyphen <--> between two otherwise joining letters will break any ligature and will insert a horizontal stroke (tatwīl, kašīda) without appearing in the transliteration. It may be used repeatedly. See also Section 8.5: automatic stretching.

For special applications, it can also be coded $\langle B \rangle$; and $\langle |B \rangle$ will behave like any ordinary consonant and may carry vowel indicators, $tanw\bar{\imath}n$, $suk\bar{\imath}un$, and, in the combination $\langle |BB \rangle$: $\check{s}adda$.

4.1.2 Vowelization

There are three modes of rendering short vowels:

• \fullvocalize:

- Every short vowel written will generate the corresponding diacritical mark *fatha*, *kasra*, *damma*, except if quoted.
- If $< \mathbb{N} >$ follows a short vowel, the corresponding form of $tanw\bar{\imath}n$ is generated instead.
- Defective writing: The coding <_a> will produce a Qur'an 'alif accent (also called dagger 'alif) instead of an explicit 'alif character which would be coded <A> or <aa>. Likewise, <_i> will produce a small 'alif below the preceding consonant in place of <I> (<iy>), and <_u> will produce an inverted damma in place of <U> (<uw>).
- If a long vowel follows a consonant, the corresponding short vowel is implied. The long vowel itself carries no diacritical mark.
- If no vowel is given after a consonant, $suk\bar{u}n$ will be generated except if a double quote precedes the next consonant. The $l\bar{a}m$ of the definite article receives no $suk\bar{u}n$ if a doubled "sun letter" follows.
- 'alif at the beginning of a word carries waşla instead of the vowel indicator if the preceding word ended with a vowel.
- \vocalize: As above, but $suk\bar{u}n$ and waṣla will not be generated except if explicitly indicated by "quoting" (see section 4.1.3).
- \novocalize: No diacritics will be generated except if explicitly asked for by "quoting" (see section 4.1.3).

In all modes, a doubled consonant will generate $\check{s}adda$, and <'A> always generates madda on 'alif.

After <aN> the silent 'alif character is generated automatically if required. The silent 'alif may also be explicitly indicated by <aNA>, or coded literally as <A> in \novocalize mode. If a silent 'alif $maq \bar{s} \bar{u} ra$ is wanted instead, write <aN_A>, <aNY>, <_A> or <Y>.

The $tanw\bar{n}$ fatha is normally positioned on the last consonant of the word, even if a silent 'alif follows. If it is instead supposed to go onto the 'alif as required by some modern Arabic writing conventions, or in Persian, this behaviour can be achieved by the option \newtanwin. The option \oldrawin will restore the classical behaviour.

A silent 'alif after $w\bar{a}w$ is indicated by $\langle VA \rangle$ or $\langle WA \rangle$ (with a capital $\langle W \rangle$!).

4.1.3 Quoting

In \novocalize mode (see Section 4.1.2), a double quote <"> will modify the meaning of the following character as follows:

- if a short vowel follows, the appropriate diacritical mark *fatḥa*, *kasra*, *damma* will be put on the preceding character.
 - If <N> follows the short vowel, the appropriate form of tanwin will be generated instead.
 - At the beginning of a word, 'alif is assumed as the first character.
- if the following character is a single right quote, a hamza mark will be put on the preceding character even if in conflict with the hamza rules.

 At the beginning of a word, <"'> will generate an isolated hamza.
- if the following character is the "invisible consonant" <|>, the connection between the adjacent letters will be broken and a small space inserted. This can also be denoted <\,> instead of <"|>.
 - At the beginning of a word, 'alif with waşla will be generated.
- otherwise: a $suk\bar{u}n$ will be put on the preceding character. The following character will be processed again.

The double quote will not show up in the transliteration.

In \vocalize mode, (see Section 4.1.2), quoting will turn a short vowel off; likewise, in \fullvocalize mode, quoting will also turn a $suk\bar{u}n$ off. Put in other words: quoting will toggle the generation of short vowel indicators and $suk\bar{u}n$ on and off.

4.1.4 Ligatures

There is no way to explicitly enforce ligatures, as a large number of them are generated automatically. The results will not always look satisfactory, so we recommend inspecting the output after the first run. Any unwanted ligature can be suppressed by interposing the invisible consonant <|> between the two letters otherwise combined into a ligature. After \ligsfalse, in the middle of a word fewer ligatures will be produced; for some texts this looks better. You can return to the normal strategy by \ligstrue.

4.1.5 Coding examples for Arabic ²

The short vowels fatha, kasra, damma are denoted, as in the transliteration, by the small letters a, i, u:

mana 'a مَنَعَ mana'a, _dahaba ﴿ فَهَبَ dahaba, ^sariba مَنَعَ šariba, qabila مَنَع qabila, 'a.zuma عَظُم azuma, 'alu عُلُ 'alu, bal عَلُ bal, ni 'ma مَنَع ni'ma, yaktub ni'ma مَنَع ni'ma مَنَع ni'ma مَنَع ni

The long vowels \bar{a} , $\bar{\imath}$, \bar{u} are denoted by capitals A, I, U or by aa, iy, uw: qAtala لُومى $n\bar{u}zi^a$, 1UmI لُومى $l\bar{u}m\bar{\iota}$,

sIrI سِيرِي $sar{\imath}rar{\imath};$ lawm $ar{\imath}$ لُوْمی $sayrar{\imath}$ بسيرِي $sayrar{\imath}$

 ${\it Alif maqs\bar{u}ra}$ is coded as ${\tt _A}$ (or Y.)

ram_A رَمَى $ramar{a}$, _dikr_A دِكْرَى $\underline{d}ikrar{a}$, 'al_A عَلَى ' $alar{a}$, bal_A بَلَى $balar{a}$.

Silent 'alif: The plural suffixes $-\bar{u}$, -aw of the verb are denoted UA, aW or aWA: katabUA يَكْتُبُوا $yaktub\bar{u}$,

ramaWA رَمُوْا ramaw, yalqaW يَلْقَوْا yalqaw.

The defective notation of \bar{a} , $\bar{\imath}$, \bar{u} can be indicated by _a, _i, _u and leads to the appropriate spelling:

dAru-h_u هُارُهُ $d\bar{a}ru$ -h \bar{u} , ri^gli-h_i رجْلِهِ $rireve{g}li$ -h \bar{i} ,

however: ramA-hu رَمَاهُ $ramar{a}$ -hu, yarmI-hi يَرْمِيهِ $yarmar{i}$ -hi;

_dih_i وَهِ $dihar{\imath}$, h_a_dih_i هٰذِهِ $har{a}dihar{\imath}$, tih_i وَ $tihar{\imath}$, hAtih_i هُاتِهِ $har{a}tihar{\imath}$,

 $rabb_i$ ز ت $rabbar{\imath}$, .sAl_i صَالِ $sar{a}lar{\imath}$; hum_u مُمْ $humar{u}$;

qiy_amaTuN قِيْمَةٌ qiyāmatun, 'il_ahuN إِلٰهُ 'ilāhun,

²Most of the examples are taken from: Wolfdietrich Fischer, Grammatik des Klassischen Arabisch, 2. Auflage, Verlag Otto Harrassowitz, Wiesbaden 1987.

sam_awAtuN مَعْوَاتٌ samāwātun, _tal_a_tuN ثَلُثٌ talātun,

1_akin مُحُواتُ lākin, h_a_dA مُخَا hādā, 'al-11_ahu لُكِنْ 'al-lāhu,

'al-rra.hm_anu أُلَّةٌ عُنُ ar-raḥmānu, _d_alika لْلِكَ dālika.

To reproduce the historical writing correctly, a silent long vowel or 'alif $maqs\bar{u}ra$ after $_a$ receives no $suk\bar{u}n$ and is ignored in the transliteration:

.sal_aUTuN صَلُوةٌ ṣalātun, .hay_aUTuN عَيُوةٌ ḥayātun,
zak_aUTuN مِشْكُوةٌ zakātun, mi^sk_aUTuN وَكُوةٌ miškātun,
ar-rib_aU اَرِّبُو ar-ribā, tawr_aITuN تَوْزِيةٌ tawrātun,
ram_aYhu مِسْمُهُمْ ramāhu, sIm_aYhum رَمْهُ

The short vowel u can be written as a long vowel by _U:

$$^{\prime}$$
_Ul_A أُولُو $^{\prime}$ ى $^{\prime}$ ى $^{\prime}$ الَّوْلَاءِ $^{\prime}$ ى $^{\prime}$ ى

Tanwīn: The plural suffixes -un, -in, -an are written -uN, -iN, -aN or aNA. Silent 'alif in -an may be indicated by A or omitted; if necessary it is supplied from the context.

 ra^guluN رَجُلِ $ra\check{g}ulun$, ra^gulin , ra^gulan رَجُلِ $ra\check{g}ulin$, ra^gulan , $ra\check{g}ulan$, $ra\check{g}ula$

'i_daN إِذًا $i\underline{d}an$, samA'aN مَمَاءً $samar{a}$ an.

There is a special case:

ribaNU رِبُّو riban; 'amruNU عُمْرٍ و amrun, 'amriNU عُمْرٍ و amrun, 'amriNU عُمْرٍ و amrun, 'amraN عُمْرًا

Tanwīn fatḥa is traditionally put on the last consonant even if a silent 'alif follows. Some modern conventions, and also Persian practice, require to put it on the 'alif in this case. This behaviour may be switched on by \newtanwin, and off by \oldraddtanwin. \newtanwin mode is the default for Persian.

ra^gulaN رُجُلاً rağulan, 'i_daN وُرُجُلاً 'idan.

A silent 'alif maqsūra after tanwīn is written aNY or aN_A:

hudaNY هُدًى hudan, fataN_A هُدًى fatan; compare:

al-hudY اَلْهُدَى al-hud $ar{a}$, 'al-fat_A اَلْهُدَى 'al-fatar{a}.

 $T\bar{a}$ ' marbuta is denoted by T:

kalimaTun كُلِمَةٍ kalimatun, kalimaTin كُلِمَةً kalimatin, kalimaTin كُلِمَةً kalimaTan عُلَمَةً fatātun, fatATun فُتَاةً fatātun, fatATan فُتَاةً fatātan.

Hamza is indicated by '; the appropriate carrier is determined by the context:

'amrun, 'ibilun إِيلِّ 'ibilun, 'u_htun أَخْتُ 'uḥtun;
ra'sun أَوْأَسُ raʾsun, 'ar'asu أَوْأَسُ 'arʾasu, sa'ala وَأَنَّ raʾsun, 'ar'asu وَأَوْأَسُ saʾala,
qara'a أَوْقُ qaraʾa; bu'sun بُوُسٌ buʾsun, 'ab'usun قَرَأُ àabʾusun,
ra'ufa وَقُلَ raʾufa, ru'asA'u وُقَسَاءُ ruʾasāʾu; bi'run وُقُنَ biʾrun,
'as'ilaTun أَسُئِلُ ʾasʾilatun, ka'iba عَنْ kaʾiba, qA'imun وَقَاعُمْ saʾilatun,
ri'AsaTun الله عَنْ الله riʾāsatun, su'ila عَنْ suʾila; samA'un الله samāʾun,
barI'un بَرِيّ barīʾun, su'un عُنْ saʾuun, bad'un بَرِيّ badʾun,
'say'un بَرْيَ šayʾun, `say'in مَسْأَلَةٌ sayʾan;
sa'ala عَنْ saʾala, mas'alaTun مَسْأَلَةٌ masʾalatun,
saw'aTun أَوْ sawʾatun, _ha.tI'aTun خَطِيئَةٌ hatīʾatun.

Old *Hamza* convention: In an older writing style that is used, e.g., in some Qur'an editions, the *hamza* is sometimes put below its carrier or on the connecting line. This style may be switched on by \oldhamza (and off again by \newhamza):

'as'ilaTuN أُسْيِلُهُ 'as'ilatun, ka'iba كَبِيبَ ka'iba, qA'imuN قَامِم وَ $q\bar{a}$ 'imun, su'ila قَامِم su'ila, ^say'aN شَيْطًا šay'an, _ha.tI'aTuN خُطِيعَةٌ haṭratun.

Madda in the context ' \bar{a} is generated automatically:

'AkiluN وَأَنَّ ' $\bar{a}kilun$, qur'AnuN قُوْآنٌ qur' $\bar{a}nun$, ra'Ahu وَأَى ra' $\bar{a}hu$.

To reproduce the historic writing correctly, it can also be explicitly indicated by A , I , U in other contexts:

'a.sdiq^A'uh_u أَصْدِقَآؤُهُ 'a $sdiqar{a}$ uh $ar{u}$; yag°I'u عَجِيٓءُ ya $ar{y}ar{v}$ u, s $^{\circ}$ U'ila سُوۤئِلَ s $ar{u}$ vila.

Šadda: A double consonant must be written twice, even if it is coded by more than one character:

nazzala نَوَّرَ nazzala, ba^s^sAruN بَشَّارٌ baššārun, nawwara نَوَّرَ nawwara, sayyidun, sa', AluN سَأَّلُ sayidun, sa', AluN سَأَّلُ sayidun, sabiyyun, 'aduwwun, aduwwun. sabiyyun, 'aduwwun, aduwwun. Instead of iyy, uww one can also write Iy, Uw:
.sabIyuN صَبِيّ ṣabīyun, 'adUwuN عَدُوٌ adūwun.

Assimilation: the definite article may be always written al-; a following "sun letter" must be written twice like in the Arabic spelling. The transliteration and the use of $suk\bar{u}n$ are adjusted accordingly:

'al-ddAru أَلَّتُ ُعُلُ 'ad-dāru, 'al-rraˆgulu أَلَّتُ ُعُلُ 'ar-rağulu, 'al-ssanaTu أَلْسَنَهُ 'as-sanatu, 'al-nnAru أَلْسَنَهُ 'an-nāru; 'al-ˆgAru أُلْبَابُ 'al-ǧāru, 'al-bAbu أُلْبَابُ 'al-bābu; 'al-lisānu, أَلْبَسَانُ 'al-laylatu, 'al-llisAnu أَلْسَنَانُ 'al-lisānu,

'al-llisānu 'اللسَانُ al-laylatu, 'al-llisAnu 'الليلة 'al-laylaTu

'al-lāhu. أُللّٰهُ 'al-lāhu

The article may also be written in the assimilated form, with identical result:

'ad-dAru أَلَّدُّ أَرُّ جُلُ 'ad-dāru, 'ar-raʾgulu أَلَدُّارُ 'ar-rağulu, 'as-sanaTu أَلْنَاهُ 'as-sanaTu أَلْنَاهُ 'as-sanaTu أَلْنَاهُ 'as-sanaTu أَلْنَاهُ 'as-sanaTu أَلْنَاهُ 'as-sanaTu أَلْنَاهُ 'as-sanaTu مُعَامِعُهُ السَّنَةُ 'as-sanaTu أَلْنَاهُ 'as-sanaTu 'as-sanaT

In some special cases the literal spelling must be used:

'alla_dI ِ اَّلَّذِينَ alla_dIna ُ اَّلَذِينَ alla_dīna, 'allatī ِ اَلَّذِينَ allatī; 'alladīna, 'allatī ِ الله 'allatī;

'al-lla_dAni اَّللَّتَانِ 'al-ladāni, 'al-llatAni الَّلَّتَانِ 'al-latāni, 'al-llawAti الَّلَوَاتِي 'al-lawĀtī اللَّوَاتِي 'al-lawĀtī اللَّوَاتِي 'al-lawĀtī اللَّوَاتِي

Waṣla: an auxiliary vowel at the beginning of a word is always written, but in the middle of a sentence generally without *hamza*. If a vowel precedes the word, the auxiliary vowel will be omitted in the transcription, and the *waṣla* sign will be used in the spelling:

wa-ismuhu وَأَشْهُهُ wa-'smuhu, f--a-in.sarafa وَأَشْهُهُ fa-'nṣarafa. 3

This also works across word boundaries:

 $^{^{3}}$ In vowelized writing, it may sometimes be advisable to introduce a $ka\check{s}\bar{\imath}da$ to prevent the vowel marks from bumping into each other.

yA ibnI يَا آبْنِي $yar{a}$ ' $bnar{\imath}$, h_a_dA ibnuh_u هٰذُا آبْنُهُ $har{a}dar{a}$ ' $bnuhar{u}$, qAla u_hru^g قَالَ آخْرُجْ $qar{a}la$ ' $hru\check{g}$.

An auxiliary vowel at the end of the preceding word may be separated by a hyphen:

qad-i in.sarafa قَدِ آنْصَرَفَ qad-i 'nṣarafa,

ra'aW-u al-bAba رَأَوُا آلْبَابَ ra'aw-u 'l-bāba,

min-i ibnih_i مِن آئِنِه min-i 'bnih $\bar{\imath}$.

This also works for the article preceding 'alif al-wasl:

'al-i-ismu أَلاِّتْمَرُ 'al-i-'smu, 'al-i-i^stirA'u أَلِاَّمْرُاءُ 'al-i-'štirā'u,

and even if the auxiliary vowel is omitted in the spelling:

ra^guluN-i ibnatuh_u ^gamIlaTuN رُجُلٌ آبْنَتُهُ جَمِيلَةٌ rağulun-i 'bnatuhū ğamīlatun,

mu.hammaduN-i al-quraˆsIyu مُحَمَّدٌ ٱلْفُرَشِيُّ muḥammadun-i ʾl-qurašīyu.

The particles li- and la- must be combined with the article except before $l\bar{a}m$:

lil-rra^guli لِلرَّ جُل lir-rağuli, lal-ma^gdu لِلرَّ جُل lal-mağdu;

however:

li-llaylati, li-ll_ahi لِلَّيْلَةِ li-llaylati, li-ll_ahi.

The Name of God is written with a special ligature if it is recognized from the input sequence ll_ah:

'al-ll_ahu الله 'al-lāhu, ta-al-ll_ahi الله 'al-lāhu, ta-al-ll_ahi الله 'ta-'l-lāhi.

Increased spacing (*Tatwīl*) between adjoining characters may be produced by a double hyphen --;

qabila قَبِلَ qabila, qa--bi--la قَبِلَ qabila, q--ab--ila قَبِلَ qabila, q-ab--ila قَبِلَ qabila

This feature should be used with discretion; automatic spreading usually leads to a better result.

Ties between words are indicated by a single hyphen:

bi-baladin بِبَلَدٍ bi-baladin, ta-al-ll_ahi بِبَلَدٍ ta-'l-lāhi,

sa-ya'tI لِيَفْرَحَ sa-ya'tī, li-yafra.ha لِيَفْرَحَ li-yafraḥa,

wa-iswadda وَآَسُودٌ
$$wa$$
-'s $wadda$, ba'da-mA بَعْدُمَا ba 'da- $mar{a}$, .tAla-mA وَعَلَامُ $tar{a}la$ - $mar{a}$, fI-ma فيم $far{i}$ - ma , 'alA-ma عَلَامُ عَلَامُ

A single hyphen at the beginning or end of a word will enforce the use of the joining form of the first resp. the last character, if that form exists (for special uses only):

Digit sequences are written in the natural order:

Hyphen and comma as a decimal separator do not terminate the number:

 $\label{ligatures} \textbf{Ligatures} \ \ \text{are generated automatically; they can be suppressed by } \ | \ :$

Abbreviations and emphasis are indicated by \emphasize:

\emphasize {\.sl'm} صلعم
$$slm$$
 \emphasize {\ab\,^g} \overline{l} $ab\check{g}$ \emphasize {\alayhi as-salAmu} عليه السّلام $alayhi$'s-salāmu

4.2 Verbatim input

After disabling language specific processing by \setverb, ArabTEX will not use any context information to determine the carrier of hamza. Instead the user has to supply this information himself by the next character typed after <'>. Generally this character will be used as the carrier; for examples and some exceptions see Table 4.3. A short vowel indicator may follow.

To ease automatic conversion, an initial 'alif may also be coded <A>.

| 'a | ١ | hamza on 'alif | 'i | 1 2 | hamza below 'alif |
|-----|---|-------------------------|-----|-----|-------------------|
| , M | ؤ | hamza on wāw | 'у | ده | hamza on a tooth |
| 'h | ٥ | $hamza$ on $h\bar{a}$ ' | 'В | ٤_ | hamza on the line |
| , | £ | isolated hamza | , A | Ĩ | madda on 'alif |

Table 4.3: Verbatim encodings for the carrier of hamza

4.3 Alternate input encodings

The ArabTEX input notation has been very carefully designed for flexibility, readability, and ease of use for linguists confined to standard 7-bit ASCII equipment for processing and transmitting data. However, it does not make much sense re-coding existing machine-readable text files that have been encoded according to other standards. Thus, some alternate reading modules have been written (as there are more than 10 different codings in current use, this is an open-ended activity), and a general code switching procedure has been provided.

An alternate reading module, e.g. asmo449.sty for the ASMO 449 code, is installed by \usepackage{asmo449} or by \input asmo449.sty. Afterwards, a code_name (in this case asmo449) is defined. Input encoding is switched by the command \setcode{code_name} that changes the coding for Arabic text globally. Encoding may be switched several times in the same document, provided the appropriate reading modules are installed; \setcode{arabtex} or \setcode{standard} returns to the standard ArabTEX notation.

As texts coded in an alternate encoding are always rendered verbatim, the commands \novocalize, \vocalize, \fullvocalize and the language selection commands \setarab etc. generally make no sense and are temporarily disabled.

4.3.1 ASMO 449 = ISO 9036

ASMO 449 (see Table 4.4) is a 7-bit code, differing from ASCII (ISO 646) mainly by replacing the Roman letters by the Arabic letter characters and diacritical marks; the Arabic digits share their positions with the ASCII digits. The positions of special and control characters in both codes are identical. ASMO 449 is supported by Arabic MS-DOS.

The file asmo449.sty contains a reading module for the ASMO 449 code (identical to ISO 9036). It is installed by the IATEX command usepackage {asmo449} or by \input asmo449.sty. The module is activated by \setcode {asmo449} or \setcode {iso9036}; all following Arabic text will be considered to be coded

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|----|-----|-----|----|----|-----|----------|-----|------|
| 00 | NUL | DLE | SP | ٠ | @ | ذ | _ | = |
| 01 | SOH | DC1 | ! | ١ | ٤ | ر | ف | 3 |
| 02 | STX | DC2 | " | ۲ | 1 | ز | و: | • |
| 03 | ETX | DC3 | # | ٣ | , | س | ك | w3 I |
| 04 | ЕОТ | DC4 | \$ | ٤ | ؤ | ش | J | 116 |
| 05 | ENQ | NAK | % | 0 | - 4 | ص | م | 311 |
| 06 | ACK | SYN | & | ٦ | ۱ ه | ض | ن | 131 |
| 07 | BEL | ЕТВ | , | Υ | 1 | ط | ٥ | 331 |
| 08 | BS | CAN |) | ٨ | J. | 畄 | و | 311 |
| 09 | НТ | EM | (| ď | :0 | ع | ی | |
| 10 | LF | SUB | * | | ij | ن. | ي | |
| 11 | VT | ESC | + | ٠. | ث |] | w.I | } |
| 12 | FF | IS4 | 4 | ^ | ن | / | * - | I |
| 13 | CR | IS3 | | = | ح | С | I. | { |
| 14 | SO | IS2 | • | \ | خ | ' | 11 | ~ |
| 15 | SI | IS1 | / | ? | د | - | 2 | DEL |

Table 4.4: ASMO 449 code table

according to the ASMO 449 standard.

Texts in ASMO 449 are usually not fully vowelized; thus the transliteration cannot be expected to be correct. This is especially true for Egyptian texts which commonly do not differentiate between $y\bar{a}'$ and 'alif magsūra.

A minimal driver file for processing existing ASMO 449 text, e.g. in a file asmotext.dat, could look as follows:

```
\documentclass {article}
\usepackage{arabtex}
\usepackage{asmo449}
\begin {document}
\setcode {asmo449}
\begin {RLtext}
\input asmotext.dat
\end {RLtext}
\end {document}
```

4.3.2 ASMO 449E = ISO 8859 - 6

The file iso88596.sty contains a reading module for the ISO 8859-6 code (extended ASMO 449 = ASMO 449E). It is installed by the IATEX command \usepackage{iso88596} or by \input iso88596.sty. The module is activated by \setcode{iso8859-6}; all following Arabic text will be considered to be coded according to the ISO 8859-6 standard. The ArabTEX notation may be reactivated by \setcode{arabtex}.

ISO 8859-6 (see Table 4.5) is an 8-bit code closely related both to 7-bit ASCII and to ASMO 449; whereas the lower 128 positions are identical to ASCII (ISO 646), the upper 128 positions contain the Arabic characters of ASMO 449 in the analogous places, plus a few additional graphic and control characters.

We exploit the close relationship of these codes by reusing the ASMO 449 reading routines, after suitable modification of the input. This only works correctly if the input text does not contain genuine ASCII letters, as we project the Arabic characters onto their locations in ASMO 449. Some of the code switching messages in the log file are spurious; do not worry.

The notes on vowelization and transliteration of ASMO 449 apply also.

The driver file indicated for ASMO 449 will be usable after the obvious modifications; however, your TEX installation must be capable of processing 8-bit data input. This is nowadays usually the case; otherwise you can try to locally find some utility program that will strip the highest order bit off the characters in your file, and process the result via ASMO 449.

| | 00 | 01 | 02 | (|)3 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 |
|----|-------------|-----|----|---|----|----|----|----|-----|----|----|-----|----|-----|-----|-----|-------------|
| 00 | NUL | DLE | SP | 0 | ٠ | @ | Р | | p | | | NSP | | ASP | ذ | _ | I. |
| 01 | SOH | DC1 | ! | 1 | ١ | A | Q | a | q | | | | | ء | ر | ف | 3.1 |
| 02 | STX | DC2 | " | 2 | ۲ | В | R | b | r | | | | | Ĩ | ز. | ق | • |
| 03 | ETX | DC3 | # | 3 | ٣ | С | S | c | S | | | | | ٲ | س | ك | W3 |
| 04 | ЕОТ | DC4 | \$ | 4 | ٤ | D | Т | d | t | | | Ħ | | ؤ | ش | J | 231 |
| 05 | ENQ | NAK | % | 5 | ٥ | E | U | e | u | | | | | إ | و | م | 311 1 |
| 06 | ACK | SYN | & | 6 | 7 | F | V | f | v | | | | | رء | ص | ن | 131 |
| 07 | $_{ m BEL}$ | ETB | , | 7 | ٧ | G | W | g | w | | | | | 1 | 4 | ٥ | 73 I |
| 08 | BS | CAN | (| 8 | ٨ | Н | X | h | x | | | | | ب | 冯 | و | 31.1 |
| 09 | НТ | EM |) | 9 | ď | Ι | Y | i | у | | | | | ö | ل | ی | |
| 10 | LF | SUB | * | | : | J | Z | j | z | | | | | ت | ره. | ي | |
| 11 | VT | ESC | + | | ; | K | Г | k | { | | | | 4 | ث | | w.1 | |
| 12 | FF | IS4 | , | | < | L | / | 1 | _ | | | Ĺ | | ج | | 5.1 | |
| 13 | CR | IS3 | _ | = | = | M |] | m | } | | | SHY | | ح | | 1, | |
| 14 | so | IS2 | • | 7 | > | N | ^ | n | ~ | | | | | خ | | 1 | |
| 15 | SI | IS1 | / | | ? | О | - | О | DEL | | | | ? | د | | 2 | |

Table 4.5: ISO 8859-6 code table

4.3.3 CP 1256 = Arabic Windows Encoding

The file arabwin.sty contains a reading module for the Arabic part of the code page 1256 used within MS Arabic Windows. It is installed by the LATEX command \usepackage{arabwin} or by \input arabwin.sty. The module is activated by \setcode{arabwin} or \setcode{cp1256}; all following Arabic text will be considered to be coded according to the MS Arabic Windows standard. The ArabTeX notation may be reactivated by \setcode{arabtex}.

The code page 1256 used in MS Arabic Windows (see Table 4.6) is an 8-bit code closely related to 7-bit ASCII; whereas the lower 128 positions are identical to ASCII (ISO 646), some of the upper 128 positions contain the Arabic characters plus additional graphic and control characters.

We reuse the ASMO 449 reading routines, after suitable modification of the input. This only works correctly if the input text does not contain genuine ASCII letters, as we project the Arabic characters onto their locations in ASMO 449. Please note that only the characters that appear in Table 4.6 are processed correctly. Some of the code switching messages in the log file may be spurious; do not worry.

The notes on vowelization and transliteration of ASMO 449 apply also.

The driver file indicated for ASMO 449 will be usable after the obvious modifications; however, your TeX installation must be capable of processing 8-bit data input.

4.3.4 ISIRI 3342

The file isiri.sty contains a reading module for the ISIRI 3342 Persian Standard Code. It is installed by the LATEX command \usepackage{isiri} or by \input isiri.sty. The module is activated by \setcode{isiri}; all following Arabic text will be considered to be coded according to the ISIRI 3342 standard. The ArabTeX notation may be reactivated by \setcode{arabtex}.

The ISIRI 3342 code (see Table 4.7) is an 8-bit code closely related to 7-bit ASCII; whereas the lower 128 positions are identical to ASCII (ISO 646), some of the upper 128 positions contain the Arabic/Persian characters plus additional graphic and control characters.

The notes on vowelization and transliteration of ASMO 449 apply also.

The driver file indicated for ASMO 449 will be usable after the obvious modifications; however, your TEX installation must be capable of processing 8-bit data input.

| | 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 |
|----|----------------------|-----|----|----|----|----|----|-----|--------|----|-----|----|-----|--------------|----|----------|
| 00 | NUL | DLE | SP | 0 | @ | Р | , | p | پ | گ | NSP | | ٥ | ذ | | <u>"</u> |
| 01 | SOH | DC1 | ! | 1 | A | Q | a | q | | | 4 | | ۽ | ر | J | 2 |
| 02 | STX | DC2 | " | 2 | В | R | b | r | | | | | Ī | ز: | | I w |
| 03 | ETX | DC3 | # | 3 | С | S | c | s | | | | | ٲ | س | م | 1 |
| 04 | ЕОТ | DC4 | \$ | 4 | D | Т | d | t | | | ಜ | | ۇ | ش | ن | |
| 05 | ENQ | NAK | % | 5 | Ε | U | e | u | | | | | 1 4 | و | ٥ | 2 |
| 06 | ACK | SYN | & | 6 | F | V | f | v | | | | | ئ | ض | و | 11 |
| 07 | BEL | ETB | , | 7 | G | W | g | w | | | | | 1 | | | |
| 08 | BS | CAN | (| 8 | Н | X | h | x | | ک | | | ب | ا | | 31 |
| 09 | НТ | EM |) | 9 | Ι | Y | i | у | | | | | 10 | 诌 | | |
| 10 | LF | SUB | * | | J | Z | j | z | 4) | ל | a | ٠. | ij | ل | | •1 |
| 11 | VT | ESC | + | ; | K | Г | k | { | | | | | Ů | ن. | | |
| 12 | FF | IS4 | , | \ | L | / | 1 | ı | | | | | ن | 1 | ی | |
| 13 | CR | IS3 | _ | | М |] | m | } | ر ا | | SHY | | ح | ف | ي | LRO |
| 14 | SO | IS2 | | > | N | ۲ | n | ~ | ژ | | | | خ | ق | | RLO |
| 15 | SI | IS1 | / | ? | О | - | О | DEL | ځ | J | | ? | د | ك | | ال |

Table 4.6: Windows CP 1256 code table

| | 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 |
|----|-----|-----|----|----|----|----|----|-----|----|----|-----|----------|--------|----|----|-----------|
| 00 | NUL | DLE | SP | 0 | @ | Р | | p | | | SP | • | Ĩ | س | ٥ | 1 |
| 01 | SOH | DC1 | ! | 1 | A | Q | a | q | | | PSP | ١ | 1 | ش | ی | 7 |
| 02 | STX | DC2 | " | 2 | В | R | b | r | | | PCN | ۲ | u | ص |] | 2 |
| 03 | ETX | DC3 | # | 3 | С | S | с | s | | | ! | ٣ |). | ض | [| "- |
| 04 | ЕОТ | DC4 | \$ | 4 | D | Т | d | t | | | Ħ | ¥ |)» | ط | { | 1 % |
| 05 | ENQ | NAK | % | 5 | E | U | e | u | | | % | Q | 9 | ظ | } | 2 |
| 06 | ACK | SYN | & | 6 | F | V | f | v | | | • | ۶ | Ċ | ع | " | 31 |
| 07 | BEL | ЕТВ | , | 7 | G | W | g | w | | | , | Υ | ج | غ | " | - |
| 08 | BS | CAN | (| 8 | Н | X | h | x | | | (| ٨ | ا ا | ف | * | ٲ |
| 09 | нт | EM |) | 9 | Ι | Y | i | у | | |) | ď | ح | ق | - | ؤ |
| 10 | LF | SUB | * | •• | J | Z | j | z | | | × | | خ. | ک | ı | ١ |
| 11 | VT | ESC | + | ; | K | [| k | { | | | + | ٠. | 1 | گ | ١ | ئ |
| 12 | FF | IS4 | , | \ | L | / | 1 | _ | | | Ĺ | ^ | ٠, | J | | ö |
| 13 | CR | IS3 | _ | = | M |] | m | } | | | _ | II | ر | م | | 실 |
| 14 | SO | IS2 | | > | N | ` | n | ~ | | | / | ~ | ز | ن | | ي |
| 15 | SI | IS1 | / | ? | О | - | О | DEL | | | / | ? | ڗٛ | و | | |

Table 4.7: ISIRI 3342 code table

4.3.5 UNICODE Arabic

The file utf8.sty contains a reading module for the Arabic and the Hebrew segment of UNICODE in UTF-8 encoding. It is installed by the LATEX command \usepackage{utf8} or by \input utf8.sty.

UTF-8 (UNICODE Transmission Format, see tables 4.8, 4.9, and 7.3) is a multibyte encoding which, for Arabic and Hebrew, uses two bytes per character whereas ASCII characters use a single byte. Far-eastern languages are encoded in three bytes per character. This is in contrast to UNICODE itself which always uses two bytes per character.

The module is activated by \setcode{utf8}; all following Arabic and Hebrew text will be considered to be coded according to the UTF-8 encoding standard. To use the correct font, select the appropriate language. The ArabTEX notation may be reactivated by \setcode{arabtex}.

| | 060 | 061 | 062 | 063 | 064 | 065 | 066 | 067 |
|---|-----|-----|---------------|------|-----------|------|-----|----------------|
| 0 | | | | ذ | - | | | <u>,</u> |
| 1 | | | ٤ | ر | ف | 1 80 | ١ | Ī |
| 2 | | | Ĩ | ز | ق | • | ۲ | Í |
| 3 | | | ١ | س | ك | 17 | ٣ | ١ |
| 4 | | | ۇ، | ش | J | u l | ٤ | ٤ |
| 5 | | | - u | و | 4 | ١u | 0 | ١٤ |
| 6 | | | ئ | .બ | ن | | 7 | ^ئ و |
| 7 | | | - | 4 | 0 | | Υ | ء ۋ |
| 8 | | |). | 当 | و | | 人 | عى |
| 9 | | | ;0 | ع | ی | | ď | ڻ |
| A | | | ŋ | ر ه. | ي | | % | ٿ |
| В | | ٠. | ڽ | | "- | | , | ٻ |
| С | (| | ج | | × - | | , | ټ |
| D | | | ح | | I۱۱ | | * | ٽ |
| Е | | | <u>ح</u> خ | | 1 | | | پ |
| F | | ? | د | | 2 | | | ٿ |

Table 4.8: UNICODE Arabic, Part 1 $\,$

| | 068 | 069 | 06A | 06B | 06C | 06D | 06E | 06F |
|---|---------------|-------------|-----|--------|-----|-----|-----|------------|
| 0 | ڀ | ڐ | | گ | ő | ې | | • |
| 1 | ځ | رد | و | گ"، | î | ۑ | | ١ |
| 2 | ڂ | رر | ڢ | ڲ | ٥, | _ | | ۲ |
| 3 | ڃ | ړ | ڣ | ڳ گ | ~ | ک | | ٣ |
| 4 | | ڔ | ڤ | ڰ | ٩ | - | | þ |
| 5 | <u>ڄ</u> څ | ڕ | ڥ | Ŭ | و | 0 | | ۵ |
| 6 | چ | <i>ن</i> ا. | ڦ | ڶ | ۆ | | | ۶ |
| 7 | = | ڗ | ڧ | ڷ | ۇ | | | Υ |
| 8 | ڈ | ڗٛ | ڨ | ڸ | ۈ | | | Л |
| 9 | ٦ | ڙ | ک | ڹ | ۉ | | | ٩ |
| A | ڊ | ښ | J | J | ۊ | | | ۺ |
| В | ڋ | پس | سک | ڻ | ۋ | | | ۻ |
| С | ڌ | ڜ | نے | ڼ | ى | | | . ف |
| D | ڌ | ڝ | ڷۓ | ڽ | ی | | | <u>=</u> ۵ |
| Е | ڎ | ڞ | پ | æ | ێ | | | Q = |
| F | ڏ | ظ | گ | ڿ | ۏ | | | |

Table 4.9: UNICODE Arabic, Part 2 $\,$

Transliteration

In addition to the arabic writing, the standard scientific transliteration may also be obtained from a fully vowelized input text. This mode is activated by \transfuse. If only the transliteration is wanted, you can deactivate the arabic writing by \arabfalse; it can be reactivated by \arabfurue. If both modes are active their output will be interleaved line by line. The font used for the transliteration is normally *italic*, it can be changed by \settransfont{font}.

5.1 ZDMG transliteration style

The "ZDMG transliteration" is in fact a family of closely related, but slightly different, transliteration conventions for several languages using the Perso-Arabic script. Therefore for producing it correctly, the appropriate language mode must have been selected.

For Arabic text, the following special cases are handled:

- after the definite article, a double consonant will be assimilated;
- an initial vowel will be replaced by an apostrophe whenever the preceding word ended with a vowel (in this case a waṣla appears in the Arabic writing). If that is not wanted, start with hamza.
- a silent 'alif or 'alif $maqs\bar{u}ra$ after $<\mathbb{N}>$ $(tanw\bar{u}n)$ and $<\mathbb{U}>$ is omitted in the transliteration. The same happens after $w\bar{u}w$ if it is written as a capital $<\mathbb{W}>$.
- To correctly reproduce some historical writings, a silent long vowel after <_a> is omitted in the transliteration. For examples, see section 4.1.5.

¹The former option "atrans" is no more necessary.

For Persian texts, the Izafet connection is handled specially, and a final silent h will be omitted in the transliteration.

5.2 Other transliteration styles

Since there is no general agreement on transcriptions, a number of variants have been provided:

- \settrans{english} will switch to the style of the Encyclopedia of Islam which is close to the conventions of the Library of Congress.
- \settrans{iranica} will produce the style used in the Encyclopedia Iranica.
- \settrans{farsi} produces a variant of the style used in the Encyclopedia Iranica.
- \settrans{lazard} switches to the conventions of Gilbert Lazard: "La langue des plus anciens documents de la prose persane".
- \settrans{urdu} switches to the conventions for Urdu used in the ALA-LC tables.
- \settrans{kashmiri} switches to the conventions for Kashmiri used in the ALA-LC tables; this mode is also chosen automatically by \setkashmiri.
- \settrans{turk} will produce a style similar to modern Turkish; it only makes sense for Ottoman texts.
- \settrans{standard} or \settrans{zdmg} will revert to the standard ZDMG mode.
- Transcription conventions for Hebrew are given in section 7.8.

The transliteration mode may be switched at any time. If the input text is not fully vowelized, the transcription cannot be expected to be correct.

5.3 Capitalization

If transcription output is used as part of a Roman text, it may be desirable to have some words start with a capital letter. This can be achieved by prefixing the command \cap to the word in question. If the first letter is hamza or 'ayn, the next letter will be capitalized. This feature may also be used after the article or a prefix, and even in other arbitrary positions; \cap will only influence the following letter. The Arabic writing is not affected.

Support for other languages using Perso-Arabic script

ArabTEX is primarily intended for typesetting texts in classical and modern Arabic, but it also provides some support for several other languages that are customarily written using the Arabic alphabet or some extension of it.

In order to switch to the conventions for one of these languages, say \setfarsi, \seturdu, \setpashto, \setmaghribi, etc.; \setverb will switch off any language specific processing. \setarab can be used to switch back to the Arabic conventions.

After selecting the language, < and > are active as delimiters for quotations; \setnone will return < and > to their normal TEX meaning. Quotations still can be bracketed by \< and > or by using \RL{}.

This part of ArabTEX relies heavily on contributions from the user community; we want to especially mention Ivan Derzhanski who completely reimplemented the routines for processing Persian. As we extensively modified these contributions again while integrating the system, we are solely responsible for any remaining, or newly introduced, errors.

6.1 Persian (Farsi, Dari), also Ottoman and Kurdish

The Persian mode is activated by \setfarsi.

 All characters needed for writing Farsi are available by default. The short vowels <e> and <o> are mapped to <i> and <u>, the long vowels <E> and <0> to <1> and <0> without a vowel indicator. <H> denotes final silent $h\bar{a}$. This $h\bar{a}$ receives no $suk\bar{u}n$ even in fully vowelized mode.

- For fatha or kasra followed by a final silent $h\bar{a}$ you can also write <,a> or <,e> in place of <aH> and <eH> (deprecated).
- The *izāfet* connection may always be written <-i> or <-e> (with hyphen); then ArabTEX tries to determine the correct spelling from the context. Likewise the *yā'-i-wahdat* can always be written <-I> or <-E>.
- The present tense forms of the copula are coded <-am>, <-I>, <-ast>, <-Im>, <-Id>, <-and>. In the output they are written as separate words after a little space.
- The final $y\bar{a}'$ carries no dots. Farsi uses the Nasta'liq font if available, otherwise Naskh.

6.1.1 Coding examples for Persian¹

The short vowels α (\check{a}), e (\check{i}), o (\check{u}) are denoted by the lowercase letters a, e or i, o or u:

bar بُرْ
$$bar$$
, beh بهٔ beh , bon بُرْ bon .

The long vowels a (\bar{a}), i (\bar{i} , \bar{e}), u (\bar{u} , \bar{o}) are denoted by the capital letters A, I or E, U or O. \cancel{Elef} mxdde is automatically generated for word-initial a:

Ab بُودٌ
$$ar{a}b$$
, bAd بُادٌ $bar{a}d$, bId بيدُ $bar{t}d$, bUd بُودُ $bar{u}d$.

Note that I yields a ya-ye mæ'ruf (with $z\bar{\imath}r$), whilst E yields a ya-ye mæjhul (without $z\bar{\imath}r$). Similarly, U yields a waw-e mæ'ruf (with $pi\bar{s}$), whilst O yields a waw-e mæjhul (without $pi\bar{s}$):

The diphthongs \hat{ei} and \hat{ou} are written ay and aw:

pay ئۇ
$$pay$$
, naw ئۇ naw .

Intervocalic *hæmze* is written ':

$$pA'Iz$$
 پَائِيزْ $par{a}$ γiz ; miy $A'I$ مِيكُو ئِي $miyar{a}$ γi , mIg $U'I$ مِيكُو ئِي $mar{c}gar{w}$ i ; tawAn $A'I$ زَنَاشُو ئَی $tawar{a}nar{a}$ γi , zan $A^*sU'I$ زَنَاشُو ئَی $zanar{a}$ $sar{w}$ i .

¹We gratefully acknowledge the voluntary help by Ivan Derzhanski who wrote this chapter, *and* implemented the language-specific processing. As we extensively modified his routines during system integration, all responsibility for any remaining, or new, errors rests with us.

Silent word-final $w\bar{a}w$ is generated by _U or 0:

$$t_U$$
 تُو tu , d_U دُو tu , d_U تُو tu , d_U تُو $t\bar{o}$, d_U

Waw-e mx'dul is written w; it is omitted in the transliteration and the preceding xe receives no jxezm:

_hwAb خۇدْ
$$har{a}b$$
, _hwI^s خويْش $har{a}b$, _hwod خۇدْ $har{a}b$

Ha-ye hæwwæz-e mæxfi is generated by H, or optionally by ,e, ,a or ,A. It does not receive a jæzm even in fully vocalised mode and is not joined to a following letter:

_hAneH خَانِه
$$har{a}neh$$
, ^c,e چِه $\check{c}eh$, naH خَانِه nah , yal_aH يُله $yalar{a}h$, yal_ ah

_hAneHhA خَانِه هَا $b\bar{a}nehh\bar{a}$, _hAneH-hA خَانِه هَا $b\bar{a}neh-h\bar{a}$.

Short edafe is written -e or -i:

ketAb-e U رَاهِ تُو
$$ket\bar{a}b$$
-e \bar{u} , rAh-e t_U رَاهِ تُو $r\bar{a}h$ -e tu , nAmeH-i man نَامِهِ مَنْ $n\bar{a}meh$ - i man , bInI-e An mard بِينِيِّ آَنْ مَرُدْ $b\bar{n}n\bar{i}$ -e $\bar{a}n$ $mard$, pA-i In zan پَايِ اِينْ زَنْ $p\bar{a}$ - i $\bar{i}n$ zan , bAzU-i In zan بَازُ و ي اِينْ زَنْ $b\bar{a}z\bar{u}$ - i $\bar{i}n$ zan .

Long edafe is written -_i:

dAr-_i man دَارِ مَنْ
$$dar{a}rar{\imath}$$
 man, _hU-_i t_U خُوي تُو $\hbarar{u}ar{\imath}$ tu.

Hæmze as ya-ye wæhdæt/nesbæt/xeṭab is likewise written -_i:

nAmeH-_i نَامِهِ
$$nar{a}meh-ar{\imath},$$
 sormeH-_i نُامِهِ $sormeh-ar{\imath},$ gofteH-_i نُامِهِ $gofteh-ar{\imath}.$

Ye-ye wæḥdæt is written -I or -E:

The present tense forms of the verb buden and the pronominal clitics are written as they are spoken:

rafteH-I رَفْتِه اِيدُ $rafteh-\bar{\imath}d$, rafteH-Id رَفْتِه اِيدُ $rafteh-\bar{\imath}d$, rafteH-ast رَفْتِه اَسْتُ rafteh-ast, rafteH-and رَفْتِه اَسْتُ rafteh-and; mard-Id مُرْدِيدُ $mard-\bar{\imath}d$, asb-etĀn اَسْبِتَانْ $asb-et\bar{\imath}an$; An^gA-st أُوسْتُ $\bar{\imath}an\check{\jmath}a\bar{\imath}-st$, U-st أُوسْتُ $\bar{\imath}an\check{\jmath}a\bar{\imath}-st$, t_U-st تُوسْتُ $\bar{\imath}an\check{\jmath}a\bar{\imath}-st$, t_U-st تُوسْتُ $\bar{\imath}an\check{\jmath}a\bar{\imath}-st$, nAmeH-I-st تَابِيسْتُ $\bar{\imath}an\bar{\imath}aeh-\bar{\imath}-st$.

The preposition be- can be written with or without a hyphen:

be-man بِتُو
$$be$$
- man , be-t_U بِتُو be - tu ; be-An بِآنْ be - $ar{a}n$, be-In بِأُو be - $ar{a}n$, be-In بِأُو be - $ar{a}n$, be-In باكو be - an

The components of compounds can be separated by $\,$, or "|:

. sA.heb\,_hAneH صَاحِب خَانِه
$$sar{a}hebhar{a}neh,$$
 ta_ht-e-"|_hwAb تَخْتِ خوَابْ $taht$ -e- $har{a}b;$ pas\,andAz يَس اَنْدَازْ $pasandar{a}z,$ naw"|AmUz نَو آمُوزْ $nawar{a}mar{u}z,$ bI\,_hwod بى خۇدْ $bar{v}hod.$

Digit sequences are written in their natural order:

1234567890 \ \ \ \ \ \ \ \ \ \ \ \ \ \ 1234567890

6.2 Maghribi

This works nearly like Arabic, but using a different writing convention. $f\bar{a}'$ is written with one dot below the letter, $q\bar{a}f$ with one dot above the normal letter form of $f\bar{a}'$. The three dots of $v\bar{a}'$ are put below the letter.

Switch to this mode by \setmaghribi.

6.3 Urdu

The Urdu mode is activated by \seturdu.

- For Urdu, additional codings are available, see Table 6.1. Some of the given codings also occur in Pashto but with a different meaning, see Section 6.4.
- Urdu uses the Nasta'liq font if available, otherwise Naskh.

6.3.1 Coding examples for $Urd\bar{u}^2$

The short vowels \check{a} , \check{i} , and \check{u} are encoded by the lowercase letters a, i, and u, and are marked respectively by zabar, $z\bar{e}r$, and $p\bar{e}\check{s}$:

par پُر
$$par$$
 , dam $/$ fir , din پُر din / sukh , dukh چُک $dukh$

The long vowels \bar{a} , $\bar{\imath}$, \bar{u} , \bar{e} , and \bar{o} are encoded by the capital letters A, I, U, E, and 0:

 \circ Note: 'alif madda is automatically generated for word-initial \bar{a} :

Ap آم
$$\bar{a}p$$
 / Am آم $\bar{a}m$

tIn يَين
$$t ar{\imath} n$$
 , la,rkI رُور $t a r k ar{\imath} / d U r$ دير $t ar{\imath} n$, cer ,

ba,rE بَرْ
$$ba\acute{r}ar{e}\ /\ {
m mOr}$$
 مور $mar{o}r$

 \circ Note: I yields a ya-ye ma'ruf (with $z\bar{e}r$), while E yields a ya-ye mağhul (without $z\bar{e}r$).

$$tIn$$
 بین $tar{t}in$, $rItI$ بین $rar{t}iar{t}$ / $mErE$ میرے $mar{e}rar{e}$, la,rkE کرکے $la\acute{r}kar{e}$

The diphthongs ae and ao are encoded ae, and ao:

kaesA يُودُا
$$kaesar{a} / paodA$$
 يُودُا $paodar{a}$

 \circ Note: U yields a $w\bar{a}w$ -e-ma'ruf (with $p\bar{e}\check{s}$), while 0 yields a $w\bar{a}w$ -e- $ma\check{g}hul$ (without $p\bar{e}\check{s}$), and **ao** is indicated by a zabar. Compare:

pUr پُور
$$par{u}r$$
 / pOtA پُور $par{o}tar{a}$ / paodA پُور $paodar{a}$

Intervocalic hamza is written ':

^cA'E چَائِے
$$ar car aar e$$
 / ma'I مَئِی $ma^{ar i}$ / kO'I کو ئِی $kar oar a$

Aspiration is produced by coding a **h** after the consonant to be aspirated. Aspiration in Urdū is produced by adding $d\bar{o}$ $\check{c}a\check{s}m\bar{\imath}$ he after the consonant:

khEt گہر
$$khar{e}t$$
 / ghar b گہر $ghar$ / mujhE کھے $har{e}t$ / dharm دھرم $har{e}t$ / dharm

Nasalization is indicated by $n\bar{u}n$ -e-junnah coded as .n. Note that the nuqta for $n\bar{u}n$ is not written when it is used to represent nasalization:

$$mae.n$$
 مَیں $mae.n$ مَیں $mae.n$ مَیں $ahinsar{a}$

 $^{^2 \}mbox{Please}$ contact Anshuman Pandey [apandey@u.washington.edu] with questions or comments regarding this section.

| $\mathrm{URD}\bar{\mathrm{U}}$ | NAME | CODE |
|--------------------------------|----------------------------------------|----------|
| ĺa | 'alif | a |
| <i>b</i> ب | be | b |
| ⊯ bh | bhe | bh |
| ج. p پ | pe | р |
| $\not = ph$ | phe | ph |
| $\ddot{}$ t | te | t |
| th تيم | the | th |
| ڻ ث | $\acute{t}e$ | ,t |
| ź th | $\acute{t}he$ | ,th |
| <u>t</u> ث | $\underline{t}e$ | _t |
| ۆ ج | $\check{g}\bar{\imath}m$ | j / ^g |
| خ ğh | $\check{g}h\bar{\imath}m$ | jh / ^gh |
| č | $\check{c}e$ | ^c |
| čh چِھ | $\check{c}he$ | ^ch |
| ب _ب | $bacute{r}ar{\imath}\ he\ /\ \dot{h}e$ | .h |
| ن ئ خ | he / khe | _h |
| ک d | $dar{a}l$ | d |
| dh دھ | $dhar{a}l$ | dh |
| d ڈ | $dar{a}l$ | ,d |
| d́h ڈھ | $\acute{d}har{a}l$ | ,dh |
| <u>ن</u> <u>d</u> | $ar{d}ar{a}l$ | _d |
| <i>r</i> ر | re | r |
| rh رھ | rhe | rh |
| ŕ ڑ | $\acute{r}e$ | ,r |
| <i>ŕh</i> ڑھ | $\acute{r}he$ | ,rh |
| z | ze | z |
| ž ژ | $\check{z}e$ | ^z |
| <i>s</i> س | $sar{\imath}n$ | s |
| ة ش š | $\check{s}ar{\imath}n$ | ^s |
| ج ص ج | $arsigmaar{a}d$ | .s |
| ض d | $\dot{q}ar{a}d$ | .d |
| <i>ب</i> ط | $\dot{t}oi$ | .t |
| <i>z</i> ظ | <i>zoi</i> | .z |

| r. | | | | | |
|-------------------------------------------------------------------------------|------------------------------------------|----------------------------------|--|--|--|
| $URD\bar{U}$ | NAME | CODE | | | |
| ، ع | `ain" | ć | | | |
| غ غ | $\dot{g}ain$ | · g | | | |
| \check{b} f | fe | f | | | |
| $egin{array}{c} \dot{g} \ \dot{g} \ \dot{f} \ \dot{g} \ \ddot{q} \end{array}$ | $qar{a}f$ | q | | | |
| $\searrow k$ | $kar{a}f$ | k | | | |
| ∡ kh | $khar{a}f$ | kh | | | |
| <i>و</i> گ | $gar{a}f$ | g | | | |
| <i>gh</i> گ∡ | $ghar{a}f$ | gh | | | |
| $\bigcup l$ | $l\bar{a}m$ | 1 | | | |
| ⊿ lh | $lhar{a}m$ | lh | | | |
| <i>m</i> م | $mar{\imath}m$ | m | | | |
| ∞ mh | $mhar{\imath}m$ | mh | | | |
| رن n | $nar{u}n$ | n | | | |
| , nh | $nhar{u}n$ | nh | | | |
| w | $war{a}w$ | $w \ / \ U, \ O, \ ao$ | | | |
| \circ h | $\check{c}har{o}\acute{t}ar{\imath}\;he$ | ,h | | | |
| <i>h</i> ه | $dar{o}$ čaš $mar{\imath}$ he | h | | | |
| <i>y</i> ى | $\check{c}har{o}\acute{t}ar{\imath}\ ye$ | ${	t y} \; / \; {	t I}, {	t E}$ | | | |
| ē | $bacute{r}ar{\imath}\ ye$ | E / ae | | | |
| ب ن | $nar{u}n$ - e - $\dot{g}unnah$ | .n | | | |
| (ء | ḥamza | , | | | |
| ö t | $te\ marbu ta$ | T | | | |
| <i>≟</i> a | a | a | | | |
| $\vdash \bar{a}$ | $ar{a}$ | A | | | |
| = i | i | i | | | |
| يی $ar{\imath}$ | $ar{\imath}$ | I | | | |
| 2 u | u | u | | | |
| <i>ū</i> ـُو | $ar{u}$ | U | | | |
| ے ē | $ar{e}$ | E | | | |
| ≟ ae | ae | ae | | | |
| ō ـو | $ar{o}$ | 0 | | | |
| ao <u>ئو</u> | ao | ao | | | |

Table 6.1: The Urdu Alphabet

The "hanging he" or Ha-ye hawwaz-e-mahfī is generated by H. It does not receive a jazm even in fully vocalised mode and is not joined to a following letter:

 $Tanw\bar{\imath}n$ is coded by aN:

taqr
Iban أَقْوِراً
$$taqr\bar{\imath}ban$$
 / faoran فُوراً $faoran$

 $Ta\check{s}d\bar{\imath}d$ is produced by coding the consonant twice:

 \circ Note that double consonants in Urdū verbs are written without $ta\check{s}d\bar{\imath}d$. $Ta\check{s}d\bar{\imath}d$ is disabled by inserting | between the double consonants:

$$ban|nA$$
 لَنْنَا $bannar{a}$ / $gin|nA$ گُنْنَا $ginnar{a}$ / $jAn|nA$ بَنْنَا $\check{g}ar{a}nnar{a}$

 $Id\bar{a}fat$ is written -e or -i:

Ab-e .hayAt آب حَيَات
$$\bar{a}b$$
-e $hayar{a}t$

 $W\bar{a}w$ -e $ma'd\bar{u}la$, or the " $w\bar{a}w$ which is passed over", is written w; it is omitted in the transliteration and the preceding he receives no jazm:

_hwAb خۇاب
$$har{a}$$
 / _hwAja خۇاج $har{a}$ خۇاب $har{a}$

 $W\bar{a}w$ -e 'atf, or the " $w\bar{a}w$ of conjunction" is coded as -0

sarw-0 san0bar سَروو سَنو بَر sarw-ō sanōbar

tar-O tAzaH تُرو تَازَه
$$tar-ar{o}\ tar{a}zah$$

'Alif $mags\bar{u}ra$ is encoded as $_A$ or Y:

$$fatw_A$$
 وَعَوَى $fatwar{a}$ / da 'wY دُعوَى da 'w $ar{a}$

Digit sequences are written in the natural order:

| ,t | ټ | ť | $t\bar{a}'$ with a small loop |
|----|-----|-------------|---------------------------------------------------|
| ,d | 7 | \acute{d} | $d\bar{a}l$ with a small loop |
| ,r | ړ | ŕ | $r\bar{a}'$ with a small loop |
| .n | ړ | \dot{v} | $n\bar{u}n$ with a small loop |
| g | گ | g | $g\bar{a}f$ with a small loop instead of a bar |
| ,z |). | ź | $r\bar{a}$ ' with one dot above and one below |
| ,s | بوز | ś | $s\bar{\imath}n$ with one dot above and one below |
| ae | ئئ | ae | the diphtong ae |
| Ee | نی | ey | the diphtong ey |
| ee | ئئ | ey | the diphtong ey |
| E | ئې | \bar{e} | the long vowel \bar{e} |
| 0 | ئو | \bar{o} | the long vowel \bar{o} |
| U | بُو | \bar{u} | the long vowel \bar{u} |

Table 6.2: Additional codings for Pashto.

6.4 Pashto (Afghanic)

Switch to this mode by \setpashto. For writing some Pashto words in the Urdu style, write the command \seturdu and afterwards switch back.

- For Pashto, additional codings are available, see Table 6.2. Some of the given codings also occur in Urdu but with a different meaning.
- The codings $\langle H \rangle$, $\langle ,a \rangle$ and $\langle ,e \rangle$ are used as in Persian. The rules for $iz\bar{a}fet$ and $y\bar{a}$ '-i-wahdat apply.
- The short vowel <e> is indicated by a zwarakay, <o> by an inverted damma. Observe also the following codings:
 - $\langle w"' \rangle$ و $hamza ext{ on } w\bar{a}w$

6.5 Sindhi

| a | ĺ | a | ~n | ڃ | \tilde{n} | z | ز | z | kh | ک | kh |
|-----|----------|-----------------|-----|-----------|-----------------|-----|--------|------------------------|----|-----|----------------|
| Ъ | ب | b | ^c | چ | č | s | س | s | g | گ | g |
| :b | ٻ | \ddot{b} | ^ch | | čh | ^s | ش | \check{s} | :g | يگ | g |
| bh | ڀ | bh | .h | ح | ķ | . s | ص | \dot{s} | gh | گھ | gh |
| t | ت | t | _h | خ | \dot{b} | .d | ض | \dot{q} | :n | گ. | \ddot{n} |
| th | ٿ | th | d | د | d | .t | ط | ţ | 1 | J | l |
| ,t | ٽ | \acute{t} | dh | ڌ | dh | .z | ظ | z | m | A | m |
| ,th | ٺ | ťh | :d | ڏ | \dot{a} | ć | ع | c | n | ن | n |
| _s | ث | \underline{t} | ,d | ڊ | \widetilde{d} | .g | ره. ره | \dot{g} | ,n | ڻ | \acute{n} |
| р | پ | p | ,dh | ڍ | $\acute{d}h$ | f | ف | f | W | و | w |
| j | ج | ğ | _d | ذ | \underline{d} | ph | ڦ | ph | ,h | ٥ | h |
| :j | ٦ | ij | r | ر | r | q | ق | q | h | A | h |
| jh | جھ | ğh | ,r | ڙ | \acute{r} | k | J | k | У | ي | y |
| a | <u></u> | a | е | ŀ | e | i | 4 | i | 0 | , | 0 |
| u | <u>,</u> | u | A | ٢ | \bar{a} | E | Œ. | \bar{e} | I | یي | $\bar{\imath}$ |
| 0 | _و | \bar{o} | U | <u>رُ</u> | \bar{u} | ae | Ĺ | ae | ao | ـُو | ao |
| i | Ţ | i | _A | ئى | \bar{a} | ' A | Ĩ | ${}^{\mathbf{j}}ar{a}$ | 'a | ١ | a |
| 'i | } | i | 'у | ئ | \mathbf{y} | ' W | ۋ | \dot{w} | ' | s. |) |

Table 6.3: The Sindhi Alphabet

To activate the Sindhi mode, select the language by \setsindhi. Sindhi input texts are encoded in a modification of the standard ArabTEX encoding. The alphabet is given in Table 6.3 on page 49.

- Use hyphens to resolve ambiguities with aspired consonants.
- There are two special codings: \MIN_{2} , \IN_{2} .
- The user might want to break some ligatures by inserting a vertical bar to get the correct writing, or just for a better appearance of the script.

6.6 Kashmiri

| a | ĺ | a | d | د | d | .d | ض | z | m | م | m |
|-----|----------|-----------|-----|-------------|----------------|----|----------|-----------------|-----|--------------|--------------|
| Ъ | ب | b | ,d | ۲ | d | .t | ط | \underline{t} | n | ن | n |
| р | ڕ | p | _d | ذ | <u>z</u> | .z | ظ | z | W | و | w |
| t | ت | t | r | ر | r | r | ع | c | ,h | ٥ | h |
| ,t | ڻ | t | ,r | ڑھ | ŗ | .g | غ | gh | У | ى | y |
| _t | ڽ | <u>s</u> | z | ز | z | f | ف | f | h | A | h |
| j | ج | j | ^z | ڗٛ | ts | q | ق | q | Е | _ | \bar{e} |
| ^c | چ | c | s | س | s | k | ک | k | , | ۽ | > |
| .h | ح | \dot{h} | ^s | ش | ś | g | گ | g | Т | ö | t |
| _h | خ | <u>kh</u> | . ន | ص | ş | 1 | J | l | . у | ~ | \dot{y} |
| a | ŀ | a | i | 4 | i | u | <u>'</u> | u | .0 | لم | _o |
| A | تا | \bar{a} | I | } :- | $\bar{\imath}$ | U | ئ | $ar{u}$ | .0 | ۔ وا | $ar{o}$ |
| .a | <u> </u> | ạ | .u | - u | u' | 0 | ـۆ | 0 | е | <u>`</u> | e |
| . A | ڵ | \bar{a} | .U | ۶ | \bar{u} | 0 | ۔و | \bar{o} | E | <u> </u> | \bar{e} |

Table 6.4: The Kashmiri Alphabet

Select Kashmiri by \setkashmiri. The input codes are given in Table 6.4. The transcription follows the ALA-LC romanization conventions.

6.7 Uighuric

Switch to this mode by \setuighur.

Uighuric input texts are encoded in a modification of the standard ArabTEX encoding, see column 5 of Table 6.5. Please observe that in Uighuric all characters are coded verbatim.

| | 1 | 2 | 3 | 4 = 5 | (6)7 |
|----|----|----------------|------------|---------------------------------|---------|
| 01 | | | L | 1 = a | (01) a |
| 02 | | | هـ | o = :a | (02)ä |
| 03 | | | ٨ | c = d | (09) de |
| 04 | | | ٠ | r = ر | (10) re |
| 05 | | | بز | z = ز | (11) ze |
| 06 | | | ىژ | z ^ c | (12) že |
| 07 | | | _و | ه = و | (25) o |
| 08 | | | <u>-</u> ۆ | ە: = ۆ | (27)ö |
| 09 | | | _ۇ | u = ۇ | (26) u |
| 10 | | | -ۈ | u: e و | (28)ü |
| 11 | | | ـۋ | w = ۋ | (29) we |
| 12 | بـ | - - | ـب | b = ب | (03) be |
| 13 | پر | - - | ڀ | p = پ | (04) pe |
| 14 | تـ | ت | ـت | = t | (05) te |
| 15 | نہ | نـ | -ن | $\dot{\mathtt{o}} = \mathtt{n}$ | (23) ne |
| 16 | جـ | ج | ے | j = 5 | (06) je |
| 17 | چـ | چ | ے | c = چ | (07) če |

| | 1 | 2 | 3 | 4 = 5 | (6)7 |
|----|-----|------------------|------------|--------------------------|---------------|
| 18 | خـ | خ | خ | x = خ | (08) xe |
| 19 | ئ | ئ | <u>-</u> | i = <u>ځ</u> | (31) i |
| 20 | ېـ | - - - | - ې | e = ې | (30) e |
| 21 | یہ | - - | ي | و = ي | (32) y |
| 22 | سـ | | _س | s = س | (13) se |
| 23 | شــ | | ـش | s^ = ش | (14) še |
| 24 | غـ | خـ | ڂ | غ = ^ g | (15)ğe |
| 25 | ف | ف_ | ف | = ف | (16) fe |
| 26 | ق | ـقـ | ـق | p = g | (17) qe |
| 27 | 2 | بک | ک | $\mathcal{L}= k$ | (18) ke |
| 28 | څ | څ | ٿ | n^ = ثُد | (20) ηe |
| 29 | گ | گ | گ | $\mathcal{J}=\mathrm{g}$ | (19) ge |
| 30 | ٦ | 7 | ـل | J = 1 | (21) le |
| 31 | م | ۔ | _م | $_{ m f}$ = m | (22) me |
| 32 | هـ | + | AT. | h = ه | (24) he |
| 33 | | _\$_ | | = ' | () |
| 34 | | | X | $rac{1}{2}=1$ a | () |

- 1. initial shape
- 2. medial shape
- 3. final shape
- 4. isolated shape
- 5. external encoding
- 6. sorting position
- 7. name

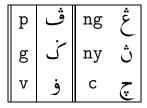
Table 6.5: Arab T_{EX} encoding of Uighuric

6.8 Old Malay

The preliminary ArabTEX language mode \setmalay is provided for processing Old Malay texts in the extended Arabic writing.

Old Malay (Jawi) input texts are encoded in a modification of the standard ArabTEX encoding, see below.

Additional encodings:



This language mode is strictly experimental and expected to contain many errors; it will be adapted to the users' requirements. Please report your experience and suggestions for changes and improvements to the author.

6.9 Other extensions of the Perso-Arabic script

This is up to experimentation by the user. If \setarab or \setfarsi will not produce the desired result, try \setverb for verbatim mode.

The vowelization and the transliteration cannot generally be expected to be correct, but might work by accident.

In case some character variants not yet provided are needed, feel free to ask the author for help. There is no simple way for the user to modify the script directly.

Hebrew mode

On the request of some users, starting with Version 3.02 ArabTEX has been extended by some modules adding support for Hebrew. Whereas the initial applications only called for short Hebrew quotations within Roman texts, possibly containing Arabic insertions too, adding "Hebrew environments" proved comparatively easy. We also added most commands provided by the HebrewTEX package (an alternative TEX extension developed in Israel, that requires TEX--XET). The Hebrew date quite probably will not work correctly.

To process Hebrew input with ArabT_EX, proceed as follows:

- for use with Plain TEX, say \input hebtex; a small loader module will load both ArabTEX and the Hebrew extension.
- with $\LaTeX 2\varepsilon$ say: \usepackage{hebtex}.

The extension provides a language mode \sethebrew, and several common encodings of texts in Hebrew, that may be switched by the \setcode command. One (nameless) encoding is compatible with Dov Grobgeld's editor HED, so files prepared for HebrewTEX are supposed to be compatible. In addition, the standard ArabTEX encoding has been extended to cater for Hebrew too.

7.1 Language switching

\sethebrew switches to Hebrew mode, \setarab back to Arabic.

Remember to switch the encoding and the vowelization mode too!

7.2 Standard Hebrew encoding

\setcode{standard} or \setcode{arabtex} will switch to the standard ArabTEX Hebrew encoding, defining the consonants as follows:

| , | Z | aleph | b | ı | beth | g | ג | gimel | d | J | daleth |
|----|--------|--------|----|---|------|---|---|--------|----|----|-----------------------|
| h | П | heh | W | 1 | waw | z | 1 | zayin | _h | П | chet |
| _t | D G | teth | у | , | yod | k | ⊐ | kaph | 1 | 5 | lamed |
| m | מ | mem | n | נ | nun | s | D | samekh | ć | ン | ayin |
| р | Ð | peh | .s | Z | sade | q | ק | qof | r | Ţ | resh |
| ,s | 27 | \sin | ^s | W | shin | S | W | s(h)in | t | IJ | taw |

Note: without punctuation, the characters sin, shin and s(h)in look identical; otherwise sin $\overset{\bullet}{\mathbf{w}}$ has a dot to the left, shin $\overset{\bullet}{\mathbf{w}}$ has a dot to the right, s(h)in $\overset{\bullet}{\mathbf{w}}$ is the form without a dot.

There are alternative encodings for soft consonants: $\langle v \rangle$ for $\langle b \rangle$, $\langle f \rangle$ for $\langle p \rangle$.

Vowels are encoded as follows:

| s | hor | t vowels | | long | vowels | | def | ective | | half | f vowels | |
|---|-----|------------------|---|------|---------------|----|-----|--------|----|------|-------------------|--|
| a | ī | pathach | A | Ŧ | qames | | | | .a | ï | chateph patach | |
| е | Ÿ | segol | Е | , i | sere yod | _e | ı | sere | .e | ¥ | chateph segol | |
| i | ٠ | chireq | I | | chireq yod | | | | .i | : | shewa | |
| 0 | T | qames chatuph | 0 | • | cholem waw | _0 | ٠ | cholem | .0 | T: | chateph qames | |
| u | \ | qibbus | U | ٦ | shureq | | | _ | .u | | no vowel mark | |

The matres lection is can also be written explicitly, e.g., <_ey> for <E>, <iy> for <I>, <_ow> for <O>.

- \vocalize (default) activates vowel points and special punctuation; \novocalize switches them off again.
- patach furtivum is written <.a> before its carrier: <rU.a_h> ☐☐.
- dagesh lene with <b g d k p t> and mappiq with <h> is expressed by prefixing a dot: <.b>, <.g>, <.d>, <.k>, <.p>, <.t>; <.h>

- dagesh forte is expressed by doubling the consonant; thus two equal consonants in sequence (even in \novocalize mode) must be separated by some short vowel indicator (or <.u>), if the standard encoding is used.
- dagesh orthophonicum is coded like dagesh forte.
- meteg is indicated by <|> after the vowel.
- maqqef is <--> (en-dash; a single hyphen will be ignored)
- Prefixes may be separated by a single hyphen, which appears in the transcription without changing the Hebrew writing.
- For those rare cases where a consonant is missing, input <| "> (bar quote);
 this may also carry vowels.
- raphe, accents, and cantillation marks are not supported.

Abbreviations may not be used in this mode as we know of no obvious way of denotating them. Suggestions are welcome.

7.3 ISO 8859-8 and Hebrew MS-Windows

ISO 8859-8 is an 8-bit encoding extending 7-bit ASCII with Hebrew letters. Within Hebrew MS-Windows, the code page CP 1255 provides a superset of ISO 8859-8, containing a full complement of Hebrew vowels, and a host of miscellaneous special characters.

ArabTeX provides within the package cp1255.sty a reading module for a subset of CP1255, containing all Hebrew characters from ISO 8859-8 and all Hebrew vowels of CP1255, but omitting the extra special characters (see Table 7.1). This encoding is activated by the switching command \setcode{cp1255} or \setcode{hwin}.

7.4 HED, PC "oldcode" and "newcode"

There is a default reading module for the Hebrew characters in code positions $96\cdots 122$ (HebrewTeX "pccode"), in code positions $128\cdots 154$ as generated by the editor HED, and also in code positions $224\cdots 250$ (HebrewTeX "newcode", ISO 8859-8). In fact these are three different encodings catered for by a single reading module; see the code assignments in Table 7.2. Observe that the Hebrew characters in ISO 8859-8 are supported both by this reading module and by the Hebrew MS-Windows encoding mentioned in Section 7.3. The code switching commands \setcode{hed}, \setcode{newcode}, \setcode{pccode}, and \setcode{iso8859-8} all activate the default reading module.

| | 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 |
|----|-----|-----|----|-----|----|--------------|----|-----|----|----|-----|----|----------|-----|----|-----|
| 00 | NUL | DLE | SP | 0 | @ | Р | , | p | | | NSP | | - | | × | נ |
| 01 | SOH | DC1 | ! | 1 | A | Q | a | q | | | | | □ ;: | i i | п | П |
| 02 | STX | DC2 | " | 2 | В | R | b | r | | | | | ï. | | ג | ע |
| 03 | ETX | DC3 | # | 3 | С | S | c | s | | | | | 1: | : | Г | J. |
| 04 | ЕОТ | DC4 | \$ | 4 | D | Т | d | t | | | ¤ | | · | Ċ | Π | U |
| 05 | ENQ | NAK | % | 5 | Е | U | e | u | | | | | ı. | וו | ١ | r |
| 06 | ACK | SYN | & | 6 | F | V | f | v | | | | | u, | ۱, | ١ | Z. |
| 07 | BEL | ЕТВ | , | 7 | G | W | g | w | | | | | ū. | 11 | П | ק |
| 08 | BS | CAN | (| 8 | Н | X | h | х | | | | | | | छ | J |
| 09 | нт | EM |) | 9 | Ι | Y | i | у | | | | | .0 | | , | u |
| 10 | LF | SUB | * | : | J | \mathbf{Z} | j | z | | | | | • | | 7 | ם |
| 11 | VT | ESC | + | ; | K | Г | k | { | | | | | | | n | |
| 12 | FF | IS4 | , | \ | L | \ | 1 | 1 | | | | | · | | 5 | |
| 13 | CR | IS3 | _ | III | M |] | m | } | | | SHY | | - | | П | LRO |
| 14 | SO | IS2 | | > | N | , | n | ~ | | | | | - | | מ | RLO |
| 15 | SI | IS1 | / | ? | О | - | О | DEL | | | | | - | | Ĭ | |

Table 7.1: ISO 8859-8 and Windows CP 1255 code table

| | 00 | 01 | 02 | 03 | 04 | 05 | 06 | 07 | 08 | 09 | 10 | 11 | 12 | 13 | 14 | 15 |
|----|-----|-----|----|----|----|----|----|----|----|----|----|----|----|----|----|-----|
| 00 | NUL | DLE | SP | 0 | @ | Р | × | נ | z | נ | | | | | × | נ |
| 01 | SOH | DC1 | ! | 1 | A | Q | П | ם | п | ם | | | | | П | ם |
| 02 | STX | DC2 | " | 2 | В | R | ג | ソ | ן | ン | | | | | ג | ソ |
| 03 | ЕТХ | DC3 | # | 3 | С | S | Г | ſ | Г | Ĺ. | | | | | Г | Ľ. |
| 04 | ЕОТ | DC4 | \$ | 4 | D | Т | Ε | n | Ε | n | | | | | Ε | ណ |
| 05 | ENQ | NAK | % | 5 | Е | U | 7 | ۲ | • | ٤ | | | | | 7 | ٢ |
| 06 | ACK | SYN | & | 6 | F | V | 1 | Y | 1 | צ | | | | | 1 | Z |
| 07 | BEL | ЕТВ | , | 7 | G | W | П | ף | ⊐ | ٦ | | | | | П | ק |
| 08 | BS | CAN |) | 8 | Н | X | Я | ſ | Я | ſ | | | | | Я | ſ |
| 09 | НТ | EM | (| 9 | Ι | Y | , | Э | , | Ð | | | | | , | a |
| 10 | LF | SUB | * | | J | Z | Г | G | Γ | G | | | | | Г | L |
| 11 | VT | ESC | + | ; | K |] | n | } | n | | | | | | n | |
| 12 | FF | IS4 | , | ^ | L | / | ſ | _ | ſ | | | | | | ſ | |
| 13 | CR | IS3 | - | II | M | [| Д | { | П | | | | | | Д | |
| 14 | SO | IS2 | | < | N | ^ | വ | ~ | വ | | | | | | വ | |
| 15 | SI | IS1 | / | ? | О | - | Ť | | 1 | | | | | | Ť | DEL |

Table 7.2: HED, CP 1255 and ISO 8859-8 code table

In this encoding vowel points, dagesh and meteg cannot be used, as they cannot be represented in the input. Abbreviations may be expressed by a single or double apostrophe (right quote). The final and the medial forms of characters are equivalent; ArabTeX will choose the appropriate shape automatically.

7.5 BHS encoding

a"MT"b"001"c"Gen"x1

בראשית בַּרַא אֵלהִים אֵת הַשָּׁמַיִם וְאֵת היתה תהו ובהו וחשך על־פני תהום ורוח אלהים מרחפת ניאמר אלהים יהי אור ווהי־אור: 3 ובין האור בין אלהים קָרָא אַלהִים אַת־הַרַקִיעַ וַיַּבְדֵּל ויעש ויהי־ערב יקוו הַמַּים מִתַּחַת הַשַּׁמִים אֵל־מַקוֹם אחד ויקרא אלהים למינו עשה פרי אצותו 12 הַאָרֵץ רַשָּא עשב מוריע וַרַע לִמִינָהוּ למינהו וירא אלהים כיישוב: אלהים לָהַבְרִיל בֵּין הַיוֹם ובֵין הַלַּיִלָה וְהַיוֹ לאתת ולמוערים ולימים ושנים:

Figure 7.1: Hebrew example

The package bhs.sty provides support for the encoding that is used in the machine-readable version of BHS (Biblia Hebraica Stuttgartensia). After loading the package, you can switch to this encoding by the command \setcode{bhs}.

The line-breaks of the source are (usually) respected. BHS line numbers and comments are only partially supported. For an example, see Figure 7.1.

\setcode{witbhs} switches to a variant of the BHS encoding that was developed by the Werkgroep Informatica of the Vrije Universiteit Amsterdam. Activate it by \usepackage{witbhs}.

7.6 UNICODE Hebrew

The file utf8.sty contains a reading module for the Arabic and the Hebrew segment of UNICODE in UTF-8 encoding. It is installed by the LATEX command \usepackage{utf8} or by \input utf8.sty.

UTF-8 (UNICODE Transmission Format, see table 7.3) is a multi-byte encoding which, for Arabic and Hebrew, uses two bytes per character whereas ASCII characters use a single byte. Far-eastern languages are encoded in three bytes per character. This is in contrast to UNICODE itself which always uses two bytes per character.

The module is activated by \setcode{utf8}; all following Arabic and Hebrew text will be considered to be coded according to the UTF-8 encoding standard. To use the correct font, select the appropriate language. The ArabTEX notation may be reactivated by \setcode{arabtex}.

7.7 Hebrew fonts

- As a default, the fonts "hclassic" (default) and "hcaption" are distributed with ArabTEX. Switch to "hcaption" by \hp, and back by \hc. These fonts have been designed and donated by Joel Hoffman, who also wrote several macro packages from which we took a few ideas for positioning punctuation. There is a variety of other usable fonts on the CTAN archives.
- If no vowel points are required, all the standard fonts "DeadSea", "OldJaffa", "TelAviv", and "Jerusalem" can be used if locally available. They are activated by the commands \ds, \oj, \ta, \jm; the command \hc switches back to the default "hclassic" font.
- The "Shalom" family of fonts, if available, can be activated by \shlmold, \shlmscr, and \shlmstk. Their vowel points presently do *not* work since they must be handled differently from the ArabTEX strategy.
- In case a font appears to be locally available but is not found, check and, if required, correct the exact spelling of the font name within the file "uheb.fd". We have seen various incompatible variants on CTAN and on the InterNet.

| | 058 | 059 | 05A | 05B | 05C | 05D | 05E | 05F |
|---|-----|-----|-----|-----------|-----|-----|----------|-----|
| 0 | | | | | | Z | | וו |
| 1 | | | | | | П | U | יי |
| 2 | | | | : | . 🗆 | ג | K | לל |
| 3 | | | | <u></u> : | ** | Г | Ĺ | , |
| 4 | | | | · | · 🔲 | Γ | A | " |
| 5 | | | | | | 1 | 7 | |
| 6 | | | | | | • | И | |
| 7 | | | | П | | Γ | Γ | |
| 8 | | | | Пь | | ß | J | |
| 9 | | | | . 🗆 | | , | ह | |
| A | | | | • | | 7 | บ | |
| В | | | | | | n | | |
| С | | | | | | J | | |
| D | | | | | | П | | |
| Е | | | | _ | | വ | | |
| F | | | | Ō | | Ī | | |

Table 7.3: UNICODE Hebrew

 Activating other Hebrew fonts by the command \sethebfont{font} might work.

Note: We recommend to set \lineskiplimit -20pt whenever Hebrew and Roman script are used within the same paragraph; this will lead to uniform line spacing. The value of \baselineskip may have to be adjusted.

7.8 Hebrew transcription systems

\transtrue activates the standard ZDMG transcription, and there are provisions for additional transcription systems:

- \settrans{zaw} switches to the conventions of "Zeitschrift für die Alttestamentliche Wissenschaft" (recommended);
- \settrans{gesenius} activates the system used in W. Gesenius' Hebrew Grammar, 26th edition (deprecated).
- \settrans{standard} restores the standard ZDMG transcription.

Miscellaneous features

| .k | <u>.s</u>]. | .k | $k\bar{a}f$ in the final position without a mark |
|----|--------------|-----------|---------------------------------------------------------------------|
| ^d | וֹי | d | $d\bar{a}l$ with a dot below |
| .f | ڣ | .f | $f\bar{a}'$ without a dot |
| .b | ب. | .b | $b\bar{a}'$ without a dot |
| .n | ڹ | .n | $n\bar{u}n$ without a dot (not available in Pashto mode) |
| Y | ی | \bar{a} | 'alif $maq s \bar{u} ra; y \bar{u}$ ' without dots in all positions |

Table 8.1: Additional codings for special purposes.

8.1 Additional codings

To reproduce exotic, erroneous or archaic texts exactly as they are written, some additional codings are available, see Table 8.1.

If further variants are needed, write to the author and indicate:

- the required shape,
- the assumed transliteration,
- a suggestion for the input coding,
- some information on the intended use.

We are willing to consider any suggestion. Adding a new character might be easy, or else it might be outright impossible. ArabTEX is rather flexible, but there are also some technical limitations.

8.2 Dots on $y\bar{a}$

Whether $y\bar{a}'$ in the final position carries dots or not is controlled by the chosen language convention. You can override this, after selecting the language, by γ and γ and γ .

8.3 Vowel positioning

In vowelized Arabic text, the short vowel marks are by default positioned close to the basic character glyphs. If this is not wanted, they may be raised to approximately uniform height by the command \accentshigh. You may revert to the default strategy by \accentslow.

8.4 Abjad numerals

The command \abjad {#1} will convert its argument, which has to be a legal representation of a number between 1 and 1999, to the Arabic 'abjad notation used in some mediaeval manuscripts. The result of the conversion will not look perfect, and the legal 'abjad number 0 can presently not be generated. The command \abjad{#1} can be used inside and outside of an Arabic context.

This routine profited greatly from suggestions by Dr. Benno van Dalen (Utrecht University).

8.5 Automatic stretching

For special purposes, e.g. for headlines and for Arabic paragraphs containing long mathematical or non-Arabic insertions, the connection between adjacent Arabic letters may be made "elastic", if they form no ligature. Thus a $ka\check{s}\bar{\imath}da$ is inserted whose length will be adjusted automatically to uniformly fill the output line.

This feature increases the already high storage demands of ArabTeX, and should therefore be used sparingly. It can be switched on with \spreadtrue and switched off again with \spreadfalse. Inside an Arabic Environment, it will also be switched off automatically at the end of every paragraph.

8.6 Uniform baselines

The Arabic and Hebrew fonts are optically compatible with the standard Roman fonts, but have larger ascenders and descenders; this will lead to unequal distances between the baselines of consecutive lines, especially if Roman and non-Roman text are mixed within the same Roman paragraph.

Typesetting on a grid will improve line spacing. We recommend to set \lineskiplimit -20pt whenever Roman script and Arabic and/or Hebrew are used within the same Roman paragraph. The value of \baselineskip may have to be adjusted; with LATEX use \baselinestretch.

Also within an *Arabic environment* typesetting on a grid may lead to a better result.

8.7 Verbatim copy of the input

For testing purposes, the Arabic input may be reproduced verbatim after \showtrue in addition to the normal output; \showfalse switches this feature off again. Commands will not usually be shown. The output will generally not look pleasant, and this feature is only provided in order to trace down errors, or to demonstrate the operation of ArabTeX as in the examples above.

8.8 Progress report

Since ArabT_EX is still rather slow (due to evolving technology it is getting faster every year), it will produce some terminal output while running to indicate it is still alive. If that is not wanted, e.g., on a very fast computer system, or while running a batch job, say \quiet or \tracingarab = 0 (outside an Arabic Environment; otherwise say \doassign {\tracingarab } {0}). The setting \tracingarab = 1 will only report Arabic paragraphs, a value of 2: Arabic lines and insertions, a value of 3 or more: individual Arabic items.

8.9 Module Reporting

A complete list of the modules loaded in a particular run will be put into the TEX log file (before the run statistics), if LATEX is used. This is believed to be useful when tracing down errors. This list is also available to the user, even with Plain TEX, as the contents of the control sequence \arabtexconfig.

Compatibility issues

ArabTEX relies only on part of the powerful features of the TEX typesetting engine (neither mathematical mode nor the alignment mechanism are used), and few of the features provided by the Plain TEX package and none of LATEX are required, but may be necessary in other parts of a multi-lingual document. Of course, TEX's macro processor is very heavily used.

It turned out that ArabTEX could be made to cooperate with a number of other macro packages, sometimes after some adjustments to ArabTEX when detecting the presence of another system, and sometimes by compatible adjustments to the other system. However there are some problem areas:

- The resource requirements of ArabTEX and usually also of the other packages are very high, and might reach the limits of a small TEX system. Fortunately, nowadays very large TEX implementations are available.
- The running time is not negligible (however, computers are still becoming faster, and typesetting this very document takes only about 20 seconds on a Pentium 233 PC running emT_EX).
- Tracking down errors in a combination of several large macro packages might be difficult and time consuming.
- There might be conflicts between the names of internal commands of several packages. The resulting effects can be very obscure; there seems to be no easy solution.
- ArabTEX assumes that the special and punctuation characters have their original category codes both when it is loaded, and when Arabic processing begins. If some macro package changes these codes, Arabic processing will usually be broken. This does not apply to Babel nor to "german.sty"; these packages are specially handled.

• Conversely ArabTEX changes the category code of < which might break other packages. Loading ArabTEX as the last module usually helps, and enables ArabTEX to detect the presence of other packages.

9.1 Arabic document classes

The experimental LaTeX2 $_{\mathcal{E}}$ classes "arabart", "arabbook", "arabrep" extend the standard classes "article", "book", and "report" in several respects: The overall document layout has been "arabized": page numbers are in Indic numerals, and columns run from right to left. The format of running heads depends on the context of the corresponding sectioning commands.

Within Arabic environments which are bracketed by \begin{RLtext} and \end{RLtext} most IATEX commands and environments are allowed, including all sectioning commands, \tabular, \tabbing, even \tableofcontents, and use an "Arabic looking" format. All arguments that denote text to be typeset are interpreted according to the currently activated Arabic encoding. Other arguments keep their IATEX standard meaning, including the preamble of \tabular, whose columns are processed from left to right (visual formatting). Generally only the basic functionality is available; optional arguments in brackets are not yet supported.

The commands \pagenumbering{abj} and \abj{ctr} generate "'abjad" numerals for page numbers and/or arbitrary LaTeX counters.

The document will start out in Roman mode, but may even be made into a single *Arabic environment*. Outside of Arabic environments the L^AT_EX commands revert to their standard meaning. The picture environment and mathematical displays presently only work in Roman mode, but may contain *Arabic insertions*.

9.2 Using ArabTEX with EDMAC

ArabTEX will cooperate with EDMAC, a Plain TEX macro package for critical editions, written by John Lavagnino and Dominik Wujastyk. If EDMAC is already present when ArabTEX is loaded, the EDMAC commands will, after suitable local modifications, be available inside an *Arabic environment*. Their arguments are considered Roman text but may contain *Arabic quotations*.

For further details, see the EDMAC documentation.

EDMAC has been extended to work with LATEX too, and ArabTEX still cooperates most of the time. However the three macro packages involved are very complicated and interact in very subtle ways, so the user may sometimes get a surprise. In this case, please contact the author.

9.3 Using ArabT_FX with Babel

The Babel package by Johannes Braams provides support for multi-lingual texts in a large number of, mostly European, languages. ArabTEX does not use the language-switching facilities provided, but is otherwise compatible.

If ArabTEX is used in a Babel document, "Roman insertions" within an *Arabic context* are interpreted according to the presently active Babel language mode. Conversely, a "Roman paragraph" in a Babel document may contain *Arabic insertions*.

9.4 Using ArabT_EX with PicT_EX

With some caution, ArabTEX can be used together with PicTEX. However, PicTEX uses the angle brackets < and > for labeling diagrams, and this requires switching off their special meaning within ArabTEX by the command \setnone. Therefore short Arabic insertions must be included as arguments of \RL{} or bracketed with \< and >.

9.5 Using ArabTEX with CJK

The CJK package by Werner Lemberg, supporting typesetting of texts in Chinese, Japanese, and Korean, to our surprise proved to be compatible with ArabTEX (after a very small adjustment). Due to the high resource requirements of both packages, a *Very Big TEX* may be required for processing texts of substantial size.

Acknowledgments

The development of ArabTEX would not have been possible without the assistance of many people, and it is impossible to acknowledge every individual contribution. Besides our local team, i.e. Udo Merkel and Heribert Schlebbe, helpful advice came, among others, from Chahriar Assad, Benno van Dalen, Ivan Derzhanski, Wolfdietrich Fischer, Ahmed El-Hadi, Yannis Haralambous, Abdelsalam Heddaya, Nicholas Heer, Taco Hoekwater, Yussuf Jabri, Iqbal Khan, Tom Koornwinder, Eberhard Krüger, Asif Lakehsar, Jan Lodder, Richard Lorch, Pierre MacKay, Eberhard Mattes, Fathy Neamat-Allah, Anshuman Pandey, Bernd Raichle, Ulrich Rebstock, Adrian Rezus, Paul Roochnik, Mohamed Saba, Waheed Samy, Annemarie Schimmel, Nariman Shehab, Arian Verheij, Dominik Wujastyk, and Michio Yano. We also have to thank all users who sent error reports, comments, and suggestions.

References

- B. Alavi, M. Lorenz: Lehrbuch der persischen Sprache.
- 5. Auflage 1988. VEB Verlag Enzyklopädie, Leipzig.
- A. A. Ambros: Einführung in die moderne arabische Schriftsprache.
- 1. Auflage 1969. Max Hueber Verlag, München.

ASMO 449: 7-bit coded Arabic character set for information interchange. Arabic Standards and Measurements Organization, 1982.

- J. D. Becker: *Arabic Word Processing*. Comm. ACM **30/7**, 600-610 (1987).
- T. Borg: Arabisch für Ausländer. Ein Lehrbuch für modernes Hocharabisch.
- 2. Auflage 1979. Verlag Borg GmbH, Hamburg.

J. A. Boyle: Grammar of Modern Persian.

Wiesbaden: Otto Harrassowitz, 1966.

B. Comrie (ed.): The World's Major Languages.

Croom Helm, London 1987.

DIN 31 635: Umschrift des Arabischen Alphabets.

Deutsches Institut für Normung e.V., 1982.

J. Lavagnino and D. Wujastyk: An Overview of EDMAC: A plain TeX format for critical editions.

TUGboat 11/4, 623-643 (1990).

L. P. Elwell-Sutton: Elementary Persian Grammar.

Cambridge University Press, 1963.

C. Faulmann: Das Buch der Schrift, enthaltend die Schriften und Alphabete aller Zeiten und aller Völker des gesammten (sic!) Erdkreises.

K. K. Hof- und Staatsdruckerei, Wien 1878.

W.D. Fischer: Grammatik des Klassischen Arabisch.

2. Auflage 1987. Verlag Otto Harrassowitz, Wiesbaden.

A. Grohmann: Arabische Paläographie (Teil I und II).

Österreichische Akademie der Wissenschaften, Philosophisch-historische Klasse, Denkschriften 94, 1. Wien 1967.

E. Harder, A. Schimmel: Arabische Sprachlehre.

15. Auflage 1983. Julius Groos Verlag, Heidelberg.

هاشم محمّد الخطّاط، قواعد الخطّ العربيّ

Hāšim Muḥammad al-Ḥaṭṭāṭ: Qawāʻid al-Ḥaṭṭi al-ʻArabī.

Maktaba an-Nahda, Baghdad; Dār al-Qalam, Beirut, 1400/1980.

 ${\rm ISO/R}$ 233 - 1961: International System for the Transliteration of Arabic Characters.

International Standards Institution, 1961.

ISO 8859 - 6: Information processing — 8-bit single-byte coded graphic character sets — Part 6: Latin/Arabic alphabet.

International Organization for Standardization, 1987.

ISO 9036: Information processing — Arabic 7-bit coded character set for information interchange.

International Organization for Standardization, 1987.

D. E. Knuth: The METAFONTbook.

Addison Wesley Publishing Comp., Reading, Mass., 1986.

D. E. Knuth: The TEXbook.

Sixth printing. Addison Wesley Publishing Comp., Reading, Mass., 1986.

D. E. Knuth and P. MacKay: Mixing right-to-left texts with left-to-right texts.

TUGboat 8/1, 14-25 (1987).

Ann K. S. Lambton: *Persian Grammar*. Cambridge University Press, 1953.

L. Lamport: Lambort: Lambort:

M. Lorenz: Lehrbuch des Pashto (Afghanisch).

2. Auflage 1982. VEB Verlag Enzyklopädie, Leipzig.

P. A. MacKay: Typesetting Problem Scripts.

BYTE **11/2**, 201-216 (1986).

H. Ritter: Über einige Regeln, die beim Drucken mit arabischen Typen zu beachten sind.

ZDMG **100/2**, 577-580 (1951).

Friedrich Rückert: Grammatik, Poetik und Rhetorik der Perser.

Wiesbaden: Otto Harrassowitz, 1966.

C. Salemann, V. Shukovski: Persische Grammatik.

4. Auflage 1947. Verlag Otto Harrassowitz, Leipzig.

A. Schimmel: *Islamic Calligraphy*. E.J.Brill, Leiden, Netherlands 1970.

H.J. Vermeer, W. Akhtar, A. Akhtar: Urdu-Lautlehre und Urdu-Schrift.

3. Auflage 1985. Julius Groos Verlag, Heidelberg.

Appendix A

Obtaining and installing ArabT_EX

A.1 Obtaining ArabT_EX

The ArabTEX system is available from the author's institution (by anonymous FTP from ftp.informatik.uni-stuttgart.de (129.69.211.2), in the directory pub/arabtex) and from many other common servers, e.g. the CTAN network

- ftp.dante.de/tex-archive/language/arabtex
- ftp.tex.ac.uk/tex-archive/language/arabtex
- ctan.tug.org/tex-archive/language/arabtex

The files may be transferred individually or as a package: arabtex.zip for PC systems, arabtex.tar.Z for U*IX systems; we recommend to get and inspect the file arabtex.htm or readme.txt first. Successfull operation on the Apple Macintosh in conjunction with OzTEX has also been reported.

At the time of this writing, version 4.00 is current. The Nasta'liq font is still under development; Naskh will be substituted automatically.

ArabTEX is being maintained and extended, if required, on a regular schedule. The current status can be found in the file arabtex.htm or by sending an EMail message (with arbitrary content) to

arabtex@informatik.uni-stuttgart.de

ArabTEX is copyrighted, but free use for scientific, experimental and other strictly private, noncommercial purposes is granted. Offprints of any publications using ArabTEX are welcome. Using ArabTEX otherwise requires a license agreement.

A.2 Installing ArabTeX

The installation procedure is strongly system dependent, and we recommend securing the assistance of a local TeXpert. You have to install the fonts provided ("nash14", "nash14bf", "xnsh14bf", "hclassic", "hcaption") with their "*.pk" and "*.tfm" files on the font search path of your TeX system, and the "*.sty" files, "arabtex.tex", and "hebtex.tex" on the source search path (usually TEXINPUT) of your system. Possibly you will also have to rename the "*.pk" files according to local conventions, and as a last resort you can try to recreate the fonts from the "*.mf" METAFONT sources. Additional fonts, whenever available, are installed analogously.

ArabTEX has been found to cooperate well with TEX versions 3.xxx, \LaTeX versions 2.09 of 1991 or later, MlTêX, NFSS and NFSS2 (not required), and previewers that can handle fonts of more than 128 characters. TEX-XET or TEX-XET are not required, and their additional features are presently not exploited.

The T_EX "hash size" should be at least 3000 to 3500, especially when using ArabT_EX in conjunction with I^AT_EX, and if the transliteration module is used. Use of a BIG T_EX may be necessary when using the NFSS2 due to the latter's high demand on string storage. Space and time requirements are not negligible, and have increased during development; however, ArabT_EX currently still runs, albeit slowly, even on a PC XT standard configuration.

Appendix B

Release history

The development of the ArabTEX system began around 1991 as a private project of the author, for his personal use. However it turned out soon that the package, if at all feasible, could be of use for others also who see the need of printing Arabic text without involving a special publishing agency. As prospective users we mainly considered Orientalists which we believed very short on funding (this proved to be a drastic understatement). There was no Arabic word processing available at that time, and using TEX as a platform looked like the only remaining possibility except perhaps implementing some complete system from scratch, which probably would necessitate building multiple versions for various computer platforms and operating systems, which we did neither dare nor could afford.

Basing the design on T_EX required the minimal user interface to become extremely lean, to facilitate the use by non-programmers. In fact, only three commands and the input notation conventions have to be learned, once the user is familiar with T_EX or I^AT_EX, to be able to use ArabT_EX for a standard Arabic document. Additional features can be looked up as required.

Using the TEX typesetting engine as a machine independent platform suggested implementing the internal algorithms in TEX's powerful internal macro language. Unfortunately, it is not easy to use, and errors can be very hard to find and eliminate.

B.1 ArabT_EX version 1.00

ArabTEX version 1.00 was a prototype, to check the basic feasibility of our approach, and to get some operating experience. Many of its features were only available in a very primitive form. It is no more supported.

B.2 ArabT_EX version 2.00

ArabTEX Version 2 was the first stable version of ArabTEX. It was not fully compatible with Version 1; however, moving to the new version usually caused little problems. Apart from some extensions, most changes were introduced in order to better conform to the transliteration standards, and to have less compatibility problems with TEX and LATEX.

The main differences between versions 1 and 2 were:

- The font size was increased, so the document layout changed. The old font "nash10" was abolished and replaced by "nash14"; the character locations have been assigned differently.
- Some Arabic characters were now coded differently: 'ayn is denoted by a left quote, and <c>, <^z>, <^t>, and <.n> have been assigned new meanings in order to better conform to the standard transliteration.
- Many more ligatures than before were supplied. This normally did not concern the user.
- \vocalize no more generated sukūn and waṣla except if explicitly indicated by quoting. See \fullvocalize.
- Arabic Environments are now always bracketed by the new control sequences \begin{arabtext} and \end{arabtext} even if only the transliteration is wanted.

B.3 ArabT_EX version 3.00

The changes introduced in Version 3.00 fall into one of two categories: error corrections, and upward compatible extensions. Details are not given here, but are documented in the text file "changes.txt" that is part of the distribution package of ArabTeX. The earlier change history up to Version 3.00 is described in the text file "changes2.txt".

Version 3 is upwards compatible with version 2. However, many new features were introduced gradually, among them support for additional input encodings and a multitude of languages that use the Perso-Arabic script. We gratefully acknowledge the cooperation of several users who contributed information, documentation, and even helped with the coding.

On some users' request, a Hebrew mode was added, as well as support for nearly all the Hebrew TEX fonts that are available on the CTAN server network.

B.4 ArabT_EX version 4.00

Version 4 is an upwards compatible extension of version 3, and many modules have been rewritten. All presently supported features are documented in this manual. It proved impractical to indicate every extension explicitly; the basic user interface is still the same.

In a few instances we had to abolish certain old features that, to our knowledge, were rarely used or not at all, because of ambiguities or conflicts with the extensions. These places are indicated in the manual, and are flagged by an asterisk in the margin.

The most important incompatible changes are:

- Tilde (<~>) is no more used as a prefix in the transliteration encoding, because of conflicts with TEX's use of tilde for a stable space. Caret (<^>) is now used instead in all cases.
- The double bar (<||>) for indicating a small unbreakable space has been * replaced by <\,>.

Users who still need the old features should contact the author; there might be a workaround.

Appendix C

Miscellaneous utilities

The following packages are not part of ArabTEX proper, and are not supported in any way, but are distributed along with ArabTEX as possibly a convenience to the users. There is no warranty whatsoever.

C.1 verses.sty

This is a small utility for typesetting classical Arabic poetry in two parallel blocks, such that every line contains two half-verses. For its use, see the file itself.

C.2 twoblks.sty

This LATEX option will define a command \twoblocks {#1}{#2} which will place the two parameters #1 and #2, usually two paragraphs, into two boxes side by side, separated by space of length \colsep. If necessary, the resulting boxes will be split across a page boundary.

This feature is useful if two versions of a text are to be contrasted. They may be in different languages, and one of them might be in Arabic (if enclosed in \begin {arabtext}\...\end {arabtext}).

This sentence has been written twice: in the English language and in the Arabic language.

Otherwise this command does not depend on ArabTEX in any way, and indeed originated in a completely different context.

Beware that the two "blocks" should each not contain much more than one, not too long, paragraph of text, otherwise TEX's main storage might overflow. There must be no \verbatim text inside the parameters of \twoblocks, nor any \catcode changes; and all TEX groups and \if \cdots \fi sequences must be properly nested.

C.3 raw.sty

This is a small utility to ease the processing of input files that have been produced by some OCR reading program. It will deactivate most of TEX's special characters.

This package depends strongly on the special application; if you need it or a variant of it, enquire with the author.

Index

| " (quoting), 22 | \hfil, 13 |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------|
| "1, 20, 22 | \hfill, 13, 14 |
| \$, 12 | \hp, 59 |
| , 21, 55 | \hskip, 13 |
| _, 13 | \hspace, 13, 14 |
| , 20, 22 | \indent, 13 |
| \<, 11, 12, 16, 41 | \input, 13 |
| \LR, 12 | \ligsfalse, 23 |
| \RL, 11, 16, 41 | \ligstrue, 23 |
| \ 13 | \lineskiplimit, 64 |
| \abjad, 63 | \lq, 12 |
| \abj{ctr}, 66 | \marginpar, 13 |
| \accentshigh, 63 | \mbox, 13 |
| \accentslow, 63 | \medskip, 13 |
| \allowarab, 14 | \newarabfont, 17 |
| \arabfalse, 39 | $\new hamza, 25$ |
| \arabstat, 14 | $\new page, 13$ |
| \arrange | $\newtanwin, 22, 24$ |
| \arabtrue, 39 | $\noindent, 13$ |
| \baselineskip, 64 | $\noindent $ |
| $\begin{tabular}{l} \begin{tabular}{l} tabu$ | $\novemath{\mathtt{novocalize}},21,22,54$ |
| \begin{RLtext}, 11 | \oj, 59 |
| $\begin{arabtext}, 11, 74$ | $\oldsymbol{\colored}$ |
| \bigskip, 13 | $\oldsymbol{\colored}$ |
| \cap, 40 | $\texttt{\oldtanwin}, 22, 24$ |
| \centerline, 14 | $\parbox{13}$ |
| \clearpage, 13 | $\parbox{pagenumbering{abj}, 66}$ |
| $\colsep, 76$ | \par, 11, 13 |
| \doassign, 14 | $\quiet, 64$ |
| \document | \rq, 12 |
| \ds, 59 | \starab , 11, 15, 19, 41, 53 |
| \emphasize, 13 | $\strut 16$ |
| $\ensuremath{\mbox{\ensuremath{\mbox{RLtext}}}\xspace, 11}$ | $\strut 9$ |
| $\ensuremath{\mbox{\mbox{end}\{arabtext\}, 11, 74}}$ | $\verb \setcode{arabtex} , 29, 31, 33, 36,$ |
| \footnote, 13 | 54, 59 |
| $\verb \fullvocalize , 19, 21, 22, 74 $ | $\strut = 13$ |
| \hc, 59 | $\strut_{asmo449}, 31$ |
| | |

| \+d-[hh-] 50 | \1 |
|----------------------------------------------------------|---------------------------|
| \setcode{bhs}, 58 | \spreadline, 14 |
| \setcode{cp1255}, 55 | \spreadtrue, 63 |
| \setcode{cp1256}, 33 | \ta, 59 |
| \setcode{hed}, 55 | \tabular environment, 66 |
| \setcode{hwin}, 55 | \tracingarab, 64 |
| \setcode{isiri}, 33 | \transfalse, 39 |
| \setcode{iso8859-6}, 31 | \transtrue, 39 |
| \setcode{iso8859-8}, 55 | \twoblocks, 76 |
| \setcode{iso9036}, 31 | \usepackage{hebtex}, 53 |
| \setcode{newcode}, 55 | \vfil, 13 |
| \setcode{pccode}, 55 | \vfill, 13 |
| \setcode{standard}, 29, 54 | \vocalize, 21, 22, 54, 74 |
| \setcode{utf8}, 36, 59 | \vskip, 13 |
| \setcode{witbhs}, 59 | \vspace, 13 |
| \setfarsi, 15, 41 | \yahdots, 63 |
| \sethebfont, 61 | \yahnodots, 63 |
| \sethebrew, 53 | <, 11, 12, 16, 41, 66 |
| \setmaghribi, 15, 41 | >, 11, 12, 16, 41 |
| \setnash, 13, 16 | \ , 13 |
| \setnashbf, 13, 16 | "hanging he", 47 |
| \setnastaliq, 13 | 1, 20, 22, 23 |
| \setnone, 16, 41 | l", 55 |
| \setpashto, 15 , 41 , 48 | IB, 21 |
| \settransfont, 39 | BB, 21 |
| $\operatorname{settrans}\{\operatorname{english}\}, 40$ | 11, 20, 22 |
| $\texttt{settrans{farsi}}, 40$ | ((()) 20 |
| $\operatorname{settrans}\{\operatorname{gesenius}\}, 61$ | ' ('ayn), 20 |
| $\operatorname{settrans\{iranica\}}, 40$ | '(hamza), 19 |
| $\operatorname{settrans}\{\operatorname{kashmiri}\}, 40$ | A 18 22 28 |
| $\texttt{\scale}$, 40 | A, 18, 23, 28 |
| $\texttt{settrans}\{\texttt{standard}\}, 40, 61$ | 'A, 20, 21, 25 |
| $\operatorname{turk}, 40$ | ,A, 43 |
| $\strans{urdu}, 40$ | ^A, 20 |
| $\operatorname{settrans{zaw}}, 61$ | _A, 18, 22, 23 |
| $\operatorname{settrans}\{\operatorname{zdmg}\}, 40$ | ,a, 42, 43, 48 |
| \sturdu , 15 , 41 , 48 | _a, 18, 21, 23 |
| $\texttt{\sc tverb}, 15, 41, 52$ | a (fatḥa), 19, 23 |
| $\shlmold, 59$ | aa, 18, 23 |
| \shlmscr, 59 | abbreviation, 28 |
| \shlmstk, 59 | 'abjadnumbers, 66 |
| $\showfalse, 64$ | abjad.sty, 63 |
| \showtrue, 64 | 'abjad numbers, 63 |
| \slash smallskip, 13 | ae, 45 |
| \space, 13 | Afghanic, 48 |
| \spreadbox, 14 | aH, 42 |
| \spreadfalse, 63 | 'ayn, 20 |
| | |

| al-, 20, 39 | aWA, 23 |
|--------------------------------------|----------------------------|
| 'alif, 28 | ay, 42 |
| dagger, $18, 21, 23$ | D 01 |
| initial, 28 | B, 21 |
| $maq s \bar{u} ra,\ 18,\ 22-24,\ 47$ | Babel, 67 |
| silent, 22, 24 | baselines |
| Qur'an, 21, 23 | uniform, 64 |
| silent, 22–24, 39 | be-, 44 |
| small, $21, 23$ | bgdkpt, 54 |
| below, 21, 23 | boxing commands, 14 |
| 'Allah (spelling), 27 | Braams, Johannes, 67 |
| aN, 19, 22, 24 | breaking connections, 22 |
| $aN_A, 22, 24$ | |
| aNA, 19, 22, 24 | cantillation, 55 |
| aNY, 24 | capital letter, 40 |
| ao, 45 | category codes, 65, 66 |
| arabart.cls, 14, 66 | CJK, 67 |
| arabbook.cls, 14, 66 | code |
| Arabic, 7 | 7-bit, 29 |
| generic term, 7 | 8-bit, 31, 33 |
| Arabic LATEX classes, 66 | arabtex, 29 |
| Arabic context, 11–13 | ASCII, 29, 31, 33 |
| Arabic environment, 11 | ASMO 449, 29, 31 |
| Arabic fonts, 13, 16 | ISIRI 3342, 29, 33 |
| Arabic group, 12 | ISO 646, 29, 31 |
| Arabic item, 12 | ISO 8859-6, 29, 31 |
| Arabic MS Windows, 33 | ISO 8859-8, 55 |
| Arabic MS-DOS, 29 | ISO 9036, 29, 31 |
| Arabic number, 12 | MS-Windows, 33 |
| Arabic quotation, 11 | UNICODE, 36, 59 |
| Arabic quotes, 12 | UTF-8, 36 , 59 |
| Arabic script, 7 | coding conventions, 18, 74 |
| Arabic word, 12 | commands |
| arabrep.cls, 14, 66 | $ArabT_{EX}$, 12, 13 |
| ArabTeX commands, 12, 13 | boxing, 14 |
| archaic text, 62 | illegal, 14 |
| ASCII, 29, 31, 33 | IAT _E X, 12, 13 |
| ASMO 449, 29, 31 | overview, 14 |
| | size changing, 13, 17 |
| aspiration, 45 | T _E X, 12, 13 |
| Assad, Chahriar, 68 | user defined, 14 |
| assignment, 14 | compatibility, 65 |
| global, 14 | Babel, 67 |
| assimilation, 20, 21, 26, 39 | CJK, 67 |
| automatic stretching, 63 | EDMAC, 66 |
| aW, 23 | PicT _E X, 67 |
| aw, 42 | - 10 - E-1, 0 · |

| compounds, 44 | tabular, 66 |
|-------------------------------|-----------------------------------|
| connecting form, 21 | extra characters, 52, 63 |
| copyright, 1, 72 | , , |
| CP 1255, 55 | Farsi, 41 |
| CP 1256, 33 | fatha, 19, 21, 22 |
| CTAN, 17 | Fischer, Wolfdietrich, 23, 68 |
| - , : | font |
| dō čašmī he, 45 | additional, 16 |
| dagesh, 54, 58 | Arabic, 13, 16 |
| forte, 55 | nash10, 74 |
| lene, 54 | nash14, 16, 71, 72, 74 |
| orthophonicum, 55 | nash14bf, 16, 72 |
| dagger 'alif, 18, 21 | Naskh, 16, 71, 72 |
| damma, 19, 21, 22 | Nasta'liq, 16, 71 |
| inverted, 21, 23, 48 | xnsh14, 16, 72 |
| Dari, 41 | xnsh14bf, 16, 72 |
| date, 21 | bold, 16 |
| default font, 16 | commercial, 16 |
| defective writing, 18, 21, 23 | default, 16 |
| definite article, 20, 26, 39 | Hebrew, 17, 59 |
| Derzhanski, Ivan, 42, 68 | DeadSea, 59 |
| diacritics, 21 | hcaption, 59, 72 |
| diphthongs, 42, 45 | hclassic, 59, 72 |
| display mode, 12 | Jerusalem, 59 |
| document classes, 66 | OldJaffa, 59 |
| dots on $y\bar{a}'$, 42, 63 | Shalom, 59 |
| 0 , , | standard, 59 |
| E, 42 | TelAviv, 59 |
| -E, 42 | installation, 72 |
| ,e, 42, 43, 48 | nastaʻliq, 42, 44 |
| -e, 42 | selection, 13 |
| EDMAC, 66 | spelling, 59 |
| eH, 42 | standard, 17 |
| El-Hadi, Ahmed, 68 | transliteration, 39 |
| emphasis, 28 | , |
| Encyclopedia Iranica, 40 | grid, 64 |
| Encyclopedia of Islam, 40 | Grobgeld, Dov, 53 |
| ending, 21 | grouping, 12, 28 |
| environment | |
| Arabic, 11 | H, 42, 43, 47, 48 |
| arabtext, 11 | h |
| IATEX, 66 | silent, 40 |
| picture, 66 | h-, 21 |
| RLtext, 11 | hamza, 19, 22, 25, 42, 43, 45, 48 |
| Roman, 11 | carrier, 25, 28 |
| tabbing, 11 | old style, 25 |
| | |

| harakāt, 19, 21–23, 42 | $\mathrm{Urd}ar{\mathrm{u}},47$ |
|----------------------------|--------------------------------------|
| on $tatwil, 21$ | Jabri, Yussuf, 68 |
| Haralambous, Yannis, 68 | jazm, 43 |
| hcaption, 17 | Juzini, 49 |
| hclassic, 17 | $ka\check{s}\bar{\imath}da,21,27,63$ |
| Hebrew consonants, 54 | kasra, 19, 21, 22 |
| Hebrew fonts, 59 | Khan, Iqbal, 68 |
| Hebrew mode, 53 | Knuth, Donald E., 7 |
| Hebrew script, 7 | Koornwinder, Tom, 68 |
| Hebrew vowels, 54 | Krüger, Eberhard, 68 |
| HebrewT _E X, 53 | Kurdish, 41 |
| hebtex.tex, 53 | Ruidisii, 41 |
| HED, 53 | la-, 27 |
| Heddaya, Abdelsalam, 68 | Lakehsar, Asif, 68 |
| Heer, Nicholas, 68 | Lamport, Leslie, 8 |
| Hoekwater, Taco, 68 | language selection, 11, 15 |
| hyphen, 21, 27, 28 | LATEX commands, 12, 13, 66 |
| | LATEX environment, 66 |
| I, 18, 23 | Lavagnino, John, 66 |
| -I, 42 | Lemberg, Werner, 67 |
| ^I, 20 | li-, 27 |
| -i, 42 | Library of Congress, 40 |
| _i, 18, 21, 23 | ligature, 23, 28, 74 |
| i (kasra), 19, 23 | breaking, 20, 21, 23, 28 |
| implementation | Lodder, Jan, 68 |
| Mac, 71 | long vowels, 18, 21 |
| PC, 71 | Lorch, Richard, 68 |
| U*IX, 71 | Loren, Richard, 00 |
| iN, 19, 24 | Macintosh, 71 |
| input switching, 29 | MacKay, Pierre, 68 |
| insertion | madda, 20, 21, 42, 45 |
| mathematical, 12 | Maghribi, 44 |
| non-Arabic, 12 | mappiq, 54 |
| Roman, 12 | maqqef, 55 |
| installation, 72 | mathematical insertion, 12 |
| inverted damma, 21, 48 | matres lectionis, 54 |
| invisible consonant, 20 | Mattes, Eberhard, 68 |
| ISIRI 3342, 33 | Merkel, Udo, 68 |
| ISO 646, 29, 31, 33 | METAFONT, 72 |
| ISO 8859-6, 31 | meteg, 55, 58 |
| ISO 8859-8, 55 | MlTêX, 72 |
| ISO 9036, 31 | • |
| item | Module list, 64 MS Arabic Windows 33 |
| Arabic, 12 | MS Arabic Windows, 33 |
| iy, 18, 23 | MS-Windows, 33 |
| izāfet, 21, 40, 42, 43, 48 | N, 21, 22, 39 |
| • • • • • • • • | ,,, |

| nūn-e ġunnah, 45 | raphe, 55 |
|---------------------------|---------------------------------------------------|
| naming conflict, 65 | raw.sty, 77 |
| nasalization, 45 | reading module, 29 |
| Naskh, 16, 17, 71, 72 | Rebstock, Ulrich, 68 |
| Nasta'liq, 17, 42, 44, 71 | Rezus, Adrian, 68 |
| Neamat-Allah, Fathy, 68 | Roman, 7 |
| nesting, 12, 14 | generic term, 7 |
| NFSS, 72 | script, 7 |
| nikudot, 17 | Roman environment, 11 |
| no vowel, 55 | Roman insertion, 12 |
| non-Arabic insertion, 12 | Roochnik, Paul, 68 |
| NU, 19, 24 | , , |
| numbers, 28, 44 | Saba, Mohamed, 68 |
| 'abjad, 63 | Samy, Waheed, 68 |
| Arabic, 12 | Schimmel, Annemarie, 68 |
| | Schlebbe, Heribert, 68 |
| O, 42, 43 | script |
| -0, 47 | Arabic, 7 |
| option | Hebrew, 7 |
| abjad, 63 | Roman, 7 |
| asmo449, 29 | šadda, 20, 21, 25 |
| iso88596, 29 | on $tatw\bar{\imath}l$, 21 |
| twoblks, 76 | Shehab, Nariman, 68 |
| Ottoman, 41 | short vowels, 19 |
| , | silent 'alif, 22, 39 |
| Pandey, Anshuman, 45, 68 | silent h, 40 |
| Pashto, 44, 48 | size changing, 13, 17 |
| patach furtivum, 54 | space |
| PC implementation, 71 | small, 20 |
| Persian, 41 | unbreakable, 20 |
| Persian copula, 42 | special codings, 62 |
| PicT _E X, 67 | stretching, 13, 21, 63 |
| picture environment, 66 | automatic, 63 |
| $pi\check{s}, 42$ | $suk\bar{u}n,\ 21,\ 22,\ 43,\ 47,\ 74$ |
| Progress report, 64 | on $l\bar{a}m$, 20 |
| pseudo fonts, 16 | on $tatw\bar{\imath}l$, 21 |
| punctuation, 12 | sun letter, 20 |
| | |
| quotation | T, 25 |
| Arabic, 11 | tabbing environment, 11 |
| non-Arabic, 12 | $tar{a}'\ marbuta,\ 25$ |
| Roman, 12 | $tanw\bar{\imath}n,\ 19,\ 21,\ 22,\ 24,\ 39,\ 47$ |
| quoting, 19, 21, 22 | $fatha,\ 24$ |
| Qur'an 'alif, 21 | on $tatw\bar{\imath}l$, 21 |
| | $ta\check{s}d\bar{\imath}d,20,21,47$ |
| Raichle, Bernd, 68 | disabling, 47 |
| | 97 |

| Urdū verbs, 47 | vowel points, 17 |
|-----------------------------|---------------------------------------|
| $tatw\bar{\imath}l,21,27$ | vowels |
| TeX commands, 12, 13 | Hebrew, 54 |
| TEX hash size, 72 | invisible, 55 |
| text | long, 18, 21, 23, 42, 45 |
| archaic, 62 | positioning, 63 |
| erroneous, 62 | short, 19, 23, 42, 45 |
| $T_{E}X - X_{E}T$, 53, 72 | silent, 39 |
| transcription, 40 | , |
| transliteration, 18, 39, 74 | W, 39 |
| Encyclopedia Iranica, 40 | wāw-e 'aṭf, 47 |
| Encyclopedia of Islam, 40 | WA, 22 |
| Farsi, 40 | wasla, 21, 22, 26, 39, 74 |
| Kashmiri, 40 | Wujastyk, Dominik, 66, 68 |
| Lazard, 40 | 3 0 , , , , |
| Library of Congress, 40 | Y, 18, 23 |
| standard, 40 | $y\bar{a}$, |
| | dots, 42, 63 |
| turkish, 40 | $y\bar{a}$ '-i-waḥ dat , 42, 43, 48 |
| Urdu, 40 | Yano, Michio, 68 |
| ZDMG, 39, 40 | , |
| twoblks.sty, 76 | $z\bar{e}r,45$ |
| U, 18, 23, 39 | $z\bar{\imath}r,~42$ |
| ^U, 20 | zwarakay, 48 |
| _U, 24, 43 | · · |
| | |
| _u, 18, 21, 23 | |
| u (damma), 19, 23 | |
| U*IX implementation, 71 | |
| UA, 22, 23 | |
| uheb.fd, 59 | |
| uN, 19, 24 | |
| UNICODE, 36, 59 | |
| Arabic, 37, 38 | |
| Hebrew, 60 | |
| uniform baselines, 64 | |
| Urdu, 44, 48 | |
| user defined commands, 14 | |
| UTF-8, 36, 59 | |
| uw, 18, 23 | |
| van Dalan Ronna 62 69 | |
| van Dalen, Benno, 63, 68 | |
| verbatim, 28 | |
| Verheij, Arian, 68 | |
| verses.sty, 76 | |
| visual formatting, 66 | |
| vowel marks, 21 | |
| | |