

Yet Another Intro to Arabic NLP

Arabic Natural Language Processing

Otakar Smrž

Institute of Formal and Applied Linguistics
Charles University in Prague

Department of Middle Eastern Studies
University of West Bohemia in Pilsen

Autumn 2005

Yet Another Introduction to Arabic NLP

Arabic Natural Language Processing

Otakar Smrž

Institute of Formal and Applied Linguistics
Charles University in Prague

Department of Middle Eastern Studies
University of West Bohemia in Pilsen

Autumn 2005

This is the series of lecture notes to the course on **Arabic Natural Language Processing** taught in the Winter Term of 2005 at the Department of Middle Eastern Studies, Faculty of Philosophy, University of West Bohemia in Pilsen.

Lecturer Otakar Smrž <otakar.smrz@mff.cuni.cz>

Website <http://ufal.mff.cuni.cz/~smrz/ANLP/>

Interest of the Course

Natural Language Processing application of engineering to the problems of human languages

Computational Linguistics study of linguistic problems using computational methods

:)

We *will not* do much of the crucial theories behind all that — theory of computation, logic, theoretical linguistics. We *will* explore how dealing with the natural language, esp. Arabic, is implemented today and what other solutions we can expect and contribute to.

Lecture Topics

1 Encodings, character sets, fonts, transliterations

- Unicode, UTF-8, CP-1256, ...
- Buckwalter transliteration
- Meta-encodings of ArabT_EX
- Encode::Arabic

2 Data formats, markup, text processing and rendering

- Plain text versus binary data
- XML, HTML and L^AT_EX documents
- Writing in ArabT_EX and MS Word
- Data re-use, transcription, advanced sorting

Lecture Topics

① Encodings, character sets, fonts, transliterations

- Unicode, UTF-8, CP-1256, ...
- Buckwalter transliteration
- Meta-encodings of ArabT_EX
- Encode::Arabic

② Data formats, markup, text processing and rendering

- Plain text versus binary data
- XML, HTML and L^AT_EX documents
- Writing in ArabT_EX and MS Word
- Data re-use, transcription, advanced sorting

Lecture Topics

1 Encodings, character sets, fonts, transliterations

- Unicode, UTF-8, CP-1256, ...
- **Buckwalter transliteration**
- Meta-encodings of ArabT_EX
- Encode::Arabic

2 Data formats, markup, text processing and rendering

- Plain text versus binary data
- XML, HTML and L^AT_EX documents
- Writing in ArabT_EX and MS Word
- Data re-use, transcription, advanced sorting

Lecture Topics

1 Encodings, character sets, fonts, transliterations

- Unicode, UTF-8, CP-1256, ...
- Buckwalter transliteration
- **Meta-encodings of ArabT_EX**
- Encode::Arabic

2 Data formats, markup, text processing and rendering

- Plain text versus binary data
- XML, HTML and L^AT_EX documents
- Writing in ArabT_EX and MS Word
- Data re-use, transcription, advanced sorting

Lecture Topics

1 Encodings, character sets, fonts, transliterations

- Unicode, UTF-8, CP-1256, ...
- Buckwalter transliteration
- Meta-encodings of ArabT_EX
- **Encode::Arabic**

2 Data formats, markup, text processing and rendering

- Plain text versus binary data
- XML, HTML and L^AT_EX documents
- Writing in ArabT_EX and MS Word
- Data re-use, transcription, advanced sorting

Lecture Topics

1 Encodings, character sets, fonts, transliterations

- Unicode, UTF-8, CP-1256, ...
- Buckwalter transliteration
- Meta-encodings of ArabT_EX
- Encode::Arabic

2 Data formats, markup, text processing and rendering

- Plain text versus binary data
- XML, HTML and L^AT_EX documents
- Writing in ArabT_EX and MS Word
- Data re-use, transcription, advanced sorting

Lecture Topics

- 1 Encodings, character sets, fonts, transliterations
 - Unicode, UTF-8, CP-1256, ...
 - Buckwalter transliteration
 - Meta-encodings of ArabT_EX
 - Encode::Arabic
- 2 Data formats, markup, text processing and rendering
 - Plain text versus binary data
 - XML, HTML and L^AT_EX documents
 - Writing in ArabT_EX and MS Word
 - Data re-use, transcription, advanced sorting

Lecture Topics

- 1 Encodings, character sets, fonts, transliterations
 - Unicode, UTF-8, CP-1256, ...
 - Buckwalter transliteration
 - Meta-encodings of ArabT_EX
 - Encode::Arabic
- 2 Data formats, markup, text processing and rendering
 - Plain text versus binary data
 - XML, HTML and L^AT_EX documents
 - Writing in ArabT_EX and MS Word
 - Data re-use, transcription, advanced sorting

Lecture Topics

- 1 Encodings, character sets, fonts, transliterations
 - Unicode, UTF-8, CP-1256, ...
 - Buckwalter transliteration
 - Meta-encodings of ArabT_EX
 - Encode::Arabic
- 2 Data formats, markup, text processing and rendering
 - Plain text versus binary data
 - XML, HTML and L^AT_EX documents
 - Writing in ArabT_EX and MS Word
 - Data re-use, transcription, advanced sorting

Lecture Topics

- 1 Encodings, character sets, fonts, transliterations
 - Unicode, UTF-8, CP-1256, ...
 - Buckwalter transliteration
 - Meta-encodings of ArabT_EX
 - Encode::Arabic
- 2 Data formats, markup, text processing and rendering
 - Plain text versus binary data
 - XML, HTML and L^AT_EX documents
 - Writing in ArabT_EX and MS Word
 - Data re-use, transcription, advanced sorting

Lecture Topics

- 3 Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - Xerox finite-state morphology
 - Functional morphology
- 4 Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- 5 Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Lecture Topics

- 3 Modeling of Arabic morphology, its various applications
 - **Buckwalter morphological analyzer**
 - Xerox finite-state morphology
 - Functional morphology
- 4 Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- 5 Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Lecture Topics

- 3 Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - **Xerox finite-state morphology**
 - Functional morphology
- 4 Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- 5 Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Lecture Topics

- 3 Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - Xerox finite-state morphology
 - **Functional morphology**
- 4 Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- 5 Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Lecture Topics

- 3 Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - Xerox finite-state morphology
 - Functional morphology
- 4 Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- 5 Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Lecture Topics

- ③ Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - Xerox finite-state morphology
 - Functional morphology
- ④ Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- ⑤ Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Lecture Topics

- ③ Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - Xerox finite-state morphology
 - Functional morphology
- ④ Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- ⑤ Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Lecture Topics

- 3 Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - Xerox finite-state morphology
 - Functional morphology
- 4 Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- 5 Morphological and syntactic annotation projects
 - **Prague Arabic Dependency Treebank**
 - Penn Arabic Treebank

Lecture Topics

- ③ Modeling of Arabic morphology, its various applications
 - Buckwalter morphological analyzer
 - Xerox finite-state morphology
 - Functional morphology
- ④ Lexicography and the other uses of linguistic corpora
 - Review of recent development and resources
- ⑤ Morphological and syntactic annotation projects
 - Prague Arabic Dependency Treebank
 - Penn Arabic Treebank

Other Information

ANLP Homepage > Software Installation Guide

<http://ufal.mff.cuni.cz/~smrz/ANLP/.php?html=tools>

ANLP Homepage > Bibliography and References

<http://ufal.mff.cuni.cz/~smrz/ANLP/.php?html=links>

Encodings and Transliterations

Arabic Natural Language Processing

Otakar Smrž

Institute of Formal and Applied Linguistics
Charles University in Prague

Department of Middle Eastern Studies
University of West Bohemia in Pilsen

October 6, 2005

Where to Begin . . .

Romanization, Transcription and Transliteration *by Ken Beesley*

<http://www.xrce.xerox.com/competencies/content-analysis/arabic/info/romanization.html>

orthography **conventions** for representing a language with some **set of symbols** or **character set**

transcription alternative (phonetic, morphological) representation of the language, possibly **romanization**

transliteration orthography with carefully **substituted symbols**, yet preserving the original conventions

encoding transliteration mapping the orthographical symbols into numbers interpreted as **characters** or **bytes**

Where to Begin . . .

Romanization, Transcription and Transliteration *by Ken Beesley*

<http://www.xrce.xerox.com/competencies/content-analysis/arabic/info/romanization.html>

orthography **conventions** for representing a language with some **set of symbols** or **character set**

transcription alternative (phonetic, morphological) representation of the language, possibly **romanization**

transliteration orthography with carefully **substituted symbols**, yet preserving the original conventions

encoding transliteration mapping the orthographical symbols into numbers interpreted as **characters** or **bytes**

Where to Begin . . .

Romanization, Transcription and Transliteration *by Ken Beesley*

<http://www.xrce.xerox.com/competencies/content-analysis/arabic/info/romanization.html>

orthography **conventions** for representing a language with some **set of symbols** or **character set**

transcription alternative (phonetic, morphological) representation of the language, possibly **romanization**

transliteration orthography with carefully **substituted symbols**, yet preserving the original conventions

encoding transliteration mapping the orthographical symbols into numbers interpreted as **characters** or **bytes**

Where to Begin . . .

Romanization, Transcription and Transliteration *by Ken Beesley*

<http://www.xrce.xerox.com/competencies/content-analysis/arabic/info/romanization.html>

- orthography **conventions** for representing a language with some **set of symbols** or **character set**
- transcription alternative (phonetic, morphological) representation of the language, possibly **romanization**
- transliteration orthography with carefully **substituted symbols**, yet preserving the original conventions
 - encoding transliteration mapping the orthographical symbols into numbers interpreted as **characters** or **bytes**

Bits of Information

Computers represent any information only through sequences of **bits**. Eight-tuples of bits form **bytes**. Every bit can only distinguish two values, symbols of a binary numeral system. Their series are identified as **numbers**, which may also encode, or be viewed as, **characters**.

Wikipedia > Binary numeral system

http://en.wikipedia.org/wiki/Binary_numeral_system

Wikipedia > ASCII

<http://en.wikipedia.org/wiki/ASCII>

Unicode in Its Words

Unicode provides a unique number for **every character**, no matter what the platform, no matter what the program, **no matter what the language**.

Před vznikem Unicode existovaly stovky rozdílných kódovacích systémů pro přiřazování těchto čísel.

وَلَمْ يُوجَدِ نِظَامٌ تَشْفِيرٍ وَاحِدٌ يَحْتَوِي عَلَى جَمِيعِ الْمَحَارِفِ الضَّرُورِيَّةِ.

Wa-lam yūğad nizāmu tašfirin wāḥidun yaḥṭawī ʿalā ġamīʿi 'l-maḥā-rifi 'd-ḍarūrīyati.

Unicode

<http://www.unicode.org/>

Unicode Arabic

The **identifiers** that are assigned to characters are called **code points**.

- Arabic Ux0600–Ux06FF
- Arabic Supplement Ux0750–Ux077F
- Arabic Presentation Forms-A UxFB50–UxFDFF
- Arabic Presentation Forms-B UxFE70–UxFEFF

Unicode also specifies **algorithms** for contextual and bidirectional rendering of the characters, while Universal Character Set does not.
(more on that later)

Wikipedia > Universal Character Set

http://en.wikipedia.org/wiki/Universal_Character_Set

Encodings of Unicode

Furthermore, the **numbers** of **code points** need to be mapped into some well-defined **sequences** of **bytes**. There are various ways to do that, and the mappings are called Unicode Transformation Formats.

Unicode Code Point	UTF-8 Byte Sequence	Character
Ux0639 0000 0 110 0011 1001	0xD8 0xB9 1101 1000 1011 1001	Arabic ع
Ux004C 0000 0000 0 100 1100	0x4C 0 100 1100	Latin L

Wikipedia > UTF-8

<http://en.wikipedia.org/wiki/UTF-8>

Encodings of Unicode

Furthermore, the **numbers** of **code points** need to be mapped into some well-defined **sequences** of **bytes**. There are various ways to do that, and the mappings are called Unicode Transformation Formats.

Unicode Code Point	UTF-8 Byte Sequence	Character
Ux0639 0000 0110 0011 1001	0xD8 0xB9 1101 1000 1011 1001	Arabic ع
Ux004C 0000 0000 0100 1100	0x4C 0100 1100	Latin L

Wikipedia > UTF-8

<http://en.wikipedia.org/wiki/UTF-8>

Buckwalter Transliteration

ﺹ	'	ء	<i>d</i>	d	د	<i>ḍ</i>	D	ض	<i>k</i>	k	ك
<i>b</i>	b	ب	<i>ḍ</i>	*	ذ	<i>ṭ</i>	T	ط	<i>l</i>	l	ل
<i>t</i>	t	ت	<i>r</i>	r	ر	<i>ẓ</i>	Z	ظ	<i>m</i>	m	م
<i>ṭ</i>	v	ث	<i>z</i>	z	ز	◌	E	ع	<i>n</i>	n	ن
<i>ǧ</i>	j	ج	<i>s</i>	s	س	<i>ǧ</i>	g	غ	<i>h</i>	h	ه
<i>ḥ</i>	H	ح	<i>š</i>	\$	ش	<i>f</i>	f	ف	<i>w</i>	w	و
<i>ḫ</i>	x	خ	<i>ṣ</i>	S	ص	<i>q</i>	q	ق	<i>y</i>	y	ي
<i>ā</i>	A	ا	<i>p</i>	P	پ	<i>v</i>	V	ف	◌	W	ؤ
◌	O	أ	<i>č</i>	J	چ	<i>g</i>	G	گ	◌	}	ئ
◌	I	إ	<i>ž</i>	R	ژ	<i>h</i>	p	ة	<i>ā</i>	Y	ى

Buckwalter Transliteration

ﺹ	’	ء	<i>d</i>	d	د	<i>ḍ</i>	D	ض	<i>k</i>	k	ك
<i>b</i>	b	ب	<i>ḍ</i>	*	ذ	<i>ṭ</i>	T	ط	<i>l</i>	l	ل
<i>t</i>	t	ت	<i>r</i>	r	ر	<i>ṣ</i>	Z	ظ	<i>m</i>	m	م
<i>ṭ</i>	v	ث	<i>z</i>	z	ز	◌	E	ع	<i>n</i>	n	ن
<i>ǧ</i>	j	ج	<i>s</i>	s	س	<i>ǧ</i>	g	غ	<i>h</i>	h	ه
<i>ḥ</i>	H	ح	<i>š</i>	\$	ش	<i>f</i>	f	ف	<i>w</i>	w	و
<i>ḫ</i>	x	خ	<i>ṣ</i>	S	ص	<i>q</i>	q	ق	<i>y</i>	y	ي
<i>ā</i>	A	ا	<i>p</i>	P	پ	<i>v</i>	V	ف	◌	W	ؤ
◌	O	أ	<i>č</i>	J	چ	<i>g</i>	G	گ	◌	}	ئ
◌	I	إ	<i>ž</i>	R	ژ	<i>h</i>	p	ة	<i>ā</i>	Y	ى

Buckwalter Transliteration

ﺹ	'	ء	<i>d</i>	d	د	<i>ḍ</i>	D	ض	<i>k</i>	k	ك
<i>b</i>	b	ب	<i>ḍ</i>	*	ذ	<i>ṭ</i>	T	ط	<i>l</i>	l	ل
<i>t</i>	t	ت	<i>r</i>	r	ر	<i>ẓ</i>	Z	ظ	<i>m</i>	m	م
<i>ṭ</i>	v	ث	<i>z</i>	z	ز	◌	E	ع	<i>n</i>	n	ن
<i>ǧ</i>	j	ج	<i>s</i>	s	س	<i>ǧ</i>	g	غ	<i>h</i>	h	ه
<i>ḥ</i>	H	ح	<i>š</i>	\$	ش	<i>f</i>	f	ف	<i>w</i>	w	و
<i>ḫ</i>	x	خ	<i>ṣ</i>	S	ص	<i>q</i>	q	ق	<i>y</i>	y	ي
<i>ā</i>	A	ا	<i>p</i>	P	پ	<i>v</i>	V	ف	◌	W	ؤ
◌	O	أ	<i>č</i>	J	چ	<i>g</i>	G	گ	◌	}	ئ
◌	I	إ	<i>ž</i>	R	ژ	<i>h</i>	p	ة	<i>ā</i>	Y	ى

Buckwalter Transliteration

›	’	ء	<i>d</i>	d	د	<i>ḍ</i>	D	ض	<i>k</i>	k	ك
<i>b</i>	b	ب	<i>ḍ</i>	*	ذ	<i>ṭ</i>	T	ط	<i>l</i>	l	ل
<i>t</i>	t	ت	<i>r</i>	r	ر	<i>ḟ</i>	Z	ظ	<i>m</i>	m	م
<i>ṭ</i>	v	ث	<i>z</i>	z	ز	◌	E	ع	<i>n</i>	n	ن
<i>ǧ</i>	j	ج	<i>s</i>	s	س	<i>ǧ</i>	g	غ	<i>h</i>	h	ه
<i>ḥ</i>	H	ح	<i>š</i>	\$	ش	<i>f</i>	f	ف	<i>w</i>	w	و
<i>ḫ</i>	x	خ	<i>ṣ</i>	S	ص	<i>q</i>	q	ق	<i>y</i>	y	ي
<i>ā</i>	A	ا	<i>p</i>	P	پ	<i>v</i>	V	ف	›	W	ؤ
›	O	أ	<i>č</i>	J	چ	<i>g</i>	G	گ	›	}	ئ
›	I	إ	<i>ž</i>	R	ژ	<i>h</i>	p	ة	<i>ā</i>	Y	ى

Buckwalter Transliteration

يُولَدُ جَمِيعُ النَّاسِ أَحْرَارًا مُتَسَاوِينَ فِي الْكِرَامَةِ وَالْحُقُوقِ. وَقَدْ وَهَبُوا عَقْلًا وَضَمِيرًا
وَعَلَيْهِمْ أَنْ يُعَامِلَ بَعْضُهُمْ بَعْضًا بِرُوحِ الْإِخَاءِ.

```
yuwladu jamiyEu {ln~aAsi OaHoraArFA mutasaAwiyna fiy
{lokaraAmapi wa{loHuquwqi. waqado wuhibuwA EaqlAF
waDamiyrFA waEalayohimo Oano yuEaAmila baEoDuhumo baEoDFA
biruwHi {loIixaA'i.
```

Buckwalter Transliteration

يولد جميع الناس أحرارا متساوين في الكرامة والحقوق. وقد وهبوا عقلا وضميرا
وعليهم أن يعامل بعضهم بعضا بروح الإخاء.

```
ywld jmyE AlnAs OHrArA mtsAwyn fy AlkrAmp wAlHqwq. wqd  
whbwA EqLA wDmyrA wElyhm On yEAml bEDhm bEDA brwH AlIxA'.
```


Buckwalter Transliteration

يُولَدُ جَمِيعُ النَّاسِ أَحْرَارًا مُتَسَاوِينَ فِي الْكِرَامَةِ وَالْحُقُوقِ. وَقَدْ وَهَبُوا عَقْلًا وَضَمِيرًا
وَعَلَيْهِمْ أَنْ يُعَامَلَ بَعْضُهُمْ بَعْضًا بِرُوحِ الْإِخَاءِ.

```
yuwladu jamiyEu {ln~aAsi OaHoraArFA mutasaAwiyna fiy
{lokaraAmapI wa{loHuquwqi. waqado wuhibuWA EaqlAF
waDamiyRFA waEalayohimo Oano yuEaAmila baEoDuhumo baEoDFA
biruwHi {loIixaA'i.
```

يولد جميع الناس أحرارا متساوين في الكرامة والحقوق. وقد وهبوا عقلا وضميرا
وعليهم أن يعامل بعضهم بعضا بروح الإخاء.

```
ywld jmyE AlnAs OHrArA mtsAwyn fy AlkrAmp wAlHqwq. wqd
whbWA EqLA wDmyrA wElyhm On yEAmI bEDhm bEDA brwH AlIxa'.
```

Notation of ArabTeX

›	’	ء	<i>d</i>	<i>d</i>	د	<i>ḍ</i>	<i>.d</i>	ض	<i>k</i>	<i>k</i>	ك
<i>b</i>	<i>b</i>	ب	<i>ḍ</i>	<i>_d</i>	ذ	<i>ṭ</i>	<i>.t</i>	ط	<i>l</i>	<i>l</i>	ل
<i>t</i>	<i>t</i>	ت	<i>r</i>	<i>r</i>	ر	<i>z</i>	<i>.z</i>	ظ	<i>m</i>	<i>m</i>	م
<i>t</i>	<i>_t</i>	ث	<i>z</i>	<i>z</i>	ز	‘	‘	ع	<i>n</i>	<i>n</i>	ن
<i>ǧ</i>	<i>^g</i>	ج	<i>s</i>	<i>s</i>	س	<i>ǧ</i>	<i>.g</i>	غ	<i>h</i>	<i>h</i>	ه
<i>ḥ</i>	<i>.h</i>	ح	<i>š</i>	<i>^s</i>	ش	<i>f</i>	<i>f</i>	ف	<i>w</i>	<i>w</i>	و
<i>ḥ</i>	<i>_h</i>	خ	<i>ṣ</i>	<i>.s</i>	ص	<i>q</i>	<i>q</i>	ق	<i>y</i>	<i>y</i>	ي
<i>ā</i>	<i>A</i>	ا	<i>p</i>	<i>p</i>	پ	<i>v</i>	<i>v</i>	ف	›	›	ؤ
›	›	أ	<i>č</i>	<i>^c</i>	چ	<i>g</i>	<i>g</i>	گ	›	›	ئ
›	›	إ	<i>ž</i>	<i>^z</i>	ژ	<i>h</i>	<i>T</i>	ة	<i>ā</i>	<i>Y</i>	ى

Notation of ArabTeX

›	’	ء	<i>d</i>	d	د	<i>ḍ</i>	.d	ض	<i>k</i>	k	ك
<i>b</i>	b	ب	<i>ḍ</i>	_d	ذ	<i>ṭ</i>	.t	ط	<i>l</i>	l	ل
<i>t</i>	t	ت	<i>r</i>	r	ر	<i>z</i>	.z	ظ	<i>m</i>	m	م
<i>t</i>	_t	ث	<i>z</i>	z	ز	‘	‘	ع	<i>n</i>	n	ن
<i>ǧ</i>	ˆg	ج	<i>s</i>	s	س	<i>ǧ</i>	.g	غ	<i>h</i>	h	ه
<i>ḥ</i>	.h	ح	<i>š</i>	ˆs	ش	<i>f</i>	f	ف	<i>w</i>	w	و
<i>ḥ</i>	_h	خ	<i>ṣ</i>	.s	ص	<i>q</i>	q	ق	<i>y</i>	y	ي
<i>ā</i>	A	ا	<i>p</i>	p	پ	<i>v</i>	v	ف	›	’	ؤ
›	’	أ	<i>č</i>	ˆc	چ	<i>g</i>	g	گ	›	’	ئ
›	’	إ	<i>ž</i>	ˆz	ژ	<i>h</i>	T	ة	<i>ā</i>	Y	ى

Notation of ArabTeX

›	'	ء	<i>d</i>	<i>d</i>	د	<i>d</i>	. <i>d</i>	ض	<i>k</i>	<i>k</i>	ك
<i>b</i>	<i>b</i>	ب	<i>d</i>	<i>_d</i>	ذ	<i>t</i>	. <i>t</i>	ط	<i>l</i>	<i>l</i>	ل
<i>t</i>	<i>t</i>	ت	<i>r</i>	<i>r</i>	ر	<i>z</i>	. <i>z</i>	ظ	<i>m</i>	<i>m</i>	م
<i>t</i>	<i>_t</i>	ث	<i>z</i>	<i>z</i>	ز	‘	‘	ع	<i>n</i>	<i>n</i>	ن
<i>ǧ</i>	<i>ˆg</i>	ج	<i>s</i>	<i>s</i>	س	<i>ǧ</i>	. <i>g</i>	غ	<i>h</i>	<i>h</i>	ه
<i>ḥ</i>	. <i>h</i>	ح	<i>š</i>	<i>ˆs</i>	ش	<i>f</i>	<i>f</i>	ف	<i>w</i>	<i>w</i>	و
<i>ḥ</i>	<i>_h</i>	خ	<i>ṣ</i>	. <i>s</i>	ص	<i>q</i>	<i>q</i>	ق	<i>y</i>	<i>y</i>	ي
<i>ā</i>	<i>A</i>	ا	<i>p</i>	<i>p</i>	پ	<i>v</i>	<i>v</i>	ف	›	' <i>w</i>	ؤ
›	' <i>a</i>	أ	<i>č</i>	<i>ˆc</i>	چ	<i>g</i>	<i>g</i>	گ	›	' <i>y</i>	ئ
›	' <i>i</i>	إ	<i>ž</i>	<i>ˆz</i>	ژ	<i>h</i>	<i>T</i>	ة	<i>ā</i>	<i>Y</i>	ى

Notation of ArabTeX

يُولَدُ جَمِيعُ النَّاسِ أَحْرَارًا مُتَسَاوِينَ فِي الْكِرَامَةِ وَالْحَقُوقِ. وَقَدْ وَهَبُوا عَقْلًا وَضَمِيرًا
وَعَلَيْهِمْ أَنْ يُعَامِلَ بَعْضُهُمْ بَعْضًا بِرُوحِ الْإِخَاءِ.

```
\cap yUladu ^gamI'u an-nAsi 'a.hrAraN mutasAwIna fI
al-karAmaTi wa-al-.huqUqi.
```

```
\cap wa-qad wuhibUA 'aqlaN wa-.damIraN wa-'alayhim 'an
yu'Amila ba'.duhum ba'.daN bi-rU.hi al-'i_hA'i.
```

Notation of ArabTeX

يولد جميع الناس أحرارا متساوين في الكرامة والحقوق. وقد وهبوا عقلا وضميرا
وعليهم أن يعامل بعضهم بعضا بروح الإخاء.

```
\cap yUladu ^gamI'u an-nAsi 'a.hrAraN mutasAwIna fI  
al-karAmaTi wa-al-.huqUqi.
```

```
\cap wa-qad wuhibUA 'aqlaN wa-.damIraN wa-'alayhim 'an  
yu'Amila ba'.duhum ba'.daN bi-rU.hi al-'i_hA'i.
```

Notation of ArabTeX

Yūladu ḡamīʿu ʿn-nāsi ʿaḥrāran mutasāwīna fī ʿl-karāmati wa-ʿl-ḥuqūqi. Wa-qad wuhibū ʿaqlan wa-ḍamīran wa-ʿalayhim ʿan yuʿāmila baḍduhum baḍan bi-rūḥi ʿl-ʿiḥāʿi.

```
\cap yUladu ^gamI' u an-nAsi 'a.hrAraN mutasAwIna fI
al-karAmaTi wa-al-.huqUqi.
```

```
\cap wa-qad wuhibUA 'aqlaN wa-.damIraN wa-'alayhim 'an
yu'Amila ba'.duhum ba'.daN bi-rU.hi al-'i_hA'i.
```

Notation of ArabTeX

يُولَدُ جَمِيعُ النَّاسِ أَحْرَارًا مُتَسَاوِينَ فِي الْكِرَامَةِ وَالْحَقُوقِ. وَقَدْ وَهَبُوا عَقْلًا وَضَمِيرًا
وَعَلَيْهِمْ أَنْ يُعَامِلَ بَعْضُهُمْ بَعْضًا بِرُوحِ الْإِخَاءِ.

يولد جميع الناس أحرارا متساوين في الكرامة والحقوق. وقد وهبوا عقلا وضميرا
وعليهم أن يعامل بعضهم بعضا بروح الإخاء.

*Yūladu ġamī'u 'n-nāsi 'aḥrāran mutasāwīna fī 'l-karāmati wa-'l-ḥuqūqi. Wa-
qad wuhibū 'aqlan wa-ḍamīran wa-ʿalayhim 'an yuʿāmila baḍuhum baḍan
bi-rūḥi 'l-iḥā'i.*

```
\cap yUladu ^gamI'u an-nAsi 'a.hrAraN mutasAwIna fI  
al-karAmaTi wa-al-.huqUqi.
```

```
\cap wa-qad wuhibUA 'aqlaN wa-.damIraN wa-'alayhim 'an  
yu'Amila ba'.duhum ba'.daN bi-rU.hi al-'i_hA'i.
```


Other Resources

Unicode > Code Charts

<http://www.unicode.org/charts/>

Buckwalter Transliteration

<http://www.qamus.org/transliteration.htm>

ArabTeX User Manual

<http://129.69.218.213/arabtex/html/arabdoc.pdf>

Encode::Arabic Online Interface

<http://ufal.mff.cuni.cz/~smrz/Encode/Arabic/>

Omniglot > Arabic Script

<http://www.omniglot.com/writing/arabic.htm>

Documents and Data — Using L^AT_EX

Arabic Natural Language Processing

Otakar Smrž

Institute of Formal and Applied Linguistics
Charles University in Prague

Department of Middle Eastern Studies
University of West Bohemia in Pilsen

October 13 & 20, 2005

Preliminaries

document **collection** of data, **logical** rather than physical (cf. file)
data **information** in the form suited for **computer processing**

We will explore how **re-usable** documents can be defined with the help of **markup** and a little bit of **programming**. We will separate the **information content** of our documents from their **visual form**. We will need \LaTeX and \ArabiTeX to create one from the other.

ArabiTeX User Manual Version 4.00

by Klaus Lagally

<http://129.69.218.213/arabtex/doc/arabdoc.pdf>

Expressing Meaning

Information requires **structure**. Explicit **annotation** of structure to a given **expression** in text or data is done through abstract **markup**. The **interpretation** of the markup in turn determines the **meaning** of the expression.

Example: structuring this very slide using the notation of \LaTeX

```
\begin{slide}
  \title{Expressing Meaning}

  Information requires \alert{structure}.
  Explicit \alert{annotation} of structure ...
\end{slide}
```

Writing in L^AT_EX

Note the **notation** for the markup, esp. `\ { }`, and the **pairing** and possible **nesting** of `\begin{...}` and `\end{...}`.

Example: elementary document in L^AT_EX

```
\documentclass{article}
```

```
\begin{document}
```

This is the first paragraph of my document.

This is the second paragraph, completing the document for now. There is not much to it.

```
\end{document}
```

Including ArabT_EX

We can use ready-made **packages** to extend or change the **interpretation** of our document. We can include الحظ العربيّ, of course.

Example: insertions of the right-to-left script

```
\documentclass{article}

\usepackage{arabtex}

\begin{document}

    This is the first and only paragraph of my document
    \<wa-fIhA _ha.t.tuN 'arabIyuN> \RL{'aw h_aka_dA}.

\end{document}
```

Flashcards!

`\documentclass{flashcards}` provides us with another functionality — it will typeset flashcards with vocabulary for us, with some **predefined** style, which we now modify explicitly to get the Arabic side both in the **orthography** and in the **transcription**.

Example: this is the document part of our document

```
\begin{flashcard}{sun beams}
  \RL{'a^si''aTu a^s-^samsi} % this text is a comment
  \\ \arabfalse\transtrue
  \RL{'a^si''aTu a^s-^samsi}
\end{flashcard}

\begin{flashcard}{official measures}
  \RL{'i^grA'AtuN rasmIyaTuN} \\ \arabfalse\transtrue
  \RL{'i^grA'AtuN rasmIyaTuN}
\end{flashcard}
```

Funny Morphology and MorphoTrees

Arabic Natural Language Processing

Otakar Smrž

Institute of Formal and Applied Linguistics
Charles University in Prague

Department of Middle Eastern Studies
University of West Bohemia in Pilsen

November 24, 2005

Morphology Disambiguation

- Arabic is a language of **rich morphology**, both derivational and inflectional, with **highly ambiguous** orthography
- Boundaries of syntactic units, **tokens**, are obscure in writing
- Orthographical **strings** consist of up to four syntactic tokens

Morphology Disambiguation

- Arabic is a language of **rich morphology**, both derivational and inflectional, with **highly ambiguous** orthography
- Boundaries of syntactic units, **tokens**, are obscure in writing
- Orthographical **strings** consist of up to four syntactic tokens

- **Disambiguation** encompasses subproblems like **tokenization**, **full morphological tagging** or its simplified 'part-of-speech' versions, **lemmatization**, **diacritization** or restoration of the structural components of words, **plus combinations** thereof

Functional Approximation

The underlying morphological engine for both the Penn Arabic Treebank and the Prague Arabic Dependency Treebank is the Buckwalter Arabic Morphological Analyzer. While PATB adopts the analyses in their original format, the PADT annotations take place on approximations of **Functional Arabic Morphology**, a novel theoretical and computational model, organized into MorphoTrees.

Functional Approximation

The underlying morphological engine for both the Penn Arabic Treebank and the Prague Arabic Dependency Treebank is the Buckwalter Arabic Morphological Analyzer. While PATB adopts the analyses in their original format, the PADT annotations take place on approximations of **Functional Arabic Morphology**, a novel theoretical and computational model, organized into MorphoTrees.

With respect to the linguistic view and the architecture of the software that we develop, we **unify the format of the morphological data** by converting all the Parts of PATB into the approximation, which is done in two steps: (a) the **morphs** of the original input strings are **re-grouped** to form **tokens** (b) the corresponding **sequences of tags** are **mapped into the fixed-width positional notation** of PADT.

He will notify them about that through SMS messages, the Internet, and other means. سَيُخَبِّرُهُمْ بِذَلِكَ عَنِ طَرِيقِ الرِّسَائِلِ الْقَصِيرَةِ وَالْإِنْتَرْنِتِ وَغَيْرِهَا.

He will notify them about that through SMS messages, the Internet, and other means. سيخبرهم بذلك عن طريق الرسائل القصيرة والإنترنت وغيرها.

String	Token	Token Tag	Buckwalter's M-Tags	Token Form	Token Gloss
		F-----	FUT	<i>sa-</i>	will
سيخبرهم		VIIA-3MS--	IV3MS+IV+IVSUFF_MOOD:I	<i>yu-ḥbir-u</i>	he-notify
		S----3MP4-	IVSUFF_DO:3MS	<i>-hum</i>	them
بذلك		P-----	PREP	<i>bi-</i>	about/by
		SD----MS--	DEM_PRON_MS	<i>dālika</i>	that
عن		P-----	PREP	<i>ʿan</i>	by/about
طريق		N-----2R	NOUN+CASE_DEF_GEN	<i>ṭarīq-i</i>	way-of
الرسائل		N-----2D	DET+NOUN+CASE_DEF_GEN	<i>ar-rasā'il-i</i>	the-messages
القصيرة		A-----FS2D	DET+ADJ+NSUFF_FEM_SG+ +CASE_DEF_GEN	<i>al-qaṣīr-at-i</i>	the-short
والإنترنت		C-----	CONJ	<i>wa-</i>	and
		Z-----2D	DET+NOUN_PROP+ +CASE_DEF_GEN	<i>al-ʾinternet-i</i>	the-internet
وغيرها		C-----	CONJ	<i>wa-</i>	and
		FN-----2R	NEG_PART+CASE_DEF_GEN	<i>ḡayr-i</i>	other/not-of
		S----3FS2-	POSS_PRON_3FS	<i>-hā</i>	them

He will notify them about that through SMS messages, the Internet, and other means. سيخبرهم بذلك عن طريق الرسائل القصيرة والإنترنت وغيرها.

String	Token	Token Tag	Buckwalter's M-Tags	Token Form	Token Gloss
		F-----	FUT	sa-	will
سيخبرهم		VIIA-3MS--	IV3MS+IV+IVSUFF_MOOD:I	yu-ḥbir-u	he-notify
		S----3MP4-	IVSUFF_DO:3MS	-hum	them
بذلك		P-----	PREP	bi-	about/by
		SD----MS--	DEM_PRON_MS	dālika	that
عن		P-----	PREP	ʿan	by/about
طريق		N-----2R	NOUN+CASE_DEF_GEN	ṭarīq-i	way-of
الرسائل		N-----2D	DET+NOUN+CASE_DEF_GEN	ar-rasā'il-i	the-messages
القصيرة		A-----FS2D	DET+ADJ+NSUFF_FEM_SG+ +CASE_DEF_GEN	al-qaṣīr-at-i	the-short
والإنترنت		C-----	CONJ	wa-	and
		Z-----2D	DET+NOUN_PROP+ +CASE_DEF_GEN	al-ʾinternet-i	the-internet
وغيرها		C-----	CONJ	wa-	and
		FN-----2R	NEG_PART+CASE_DEF_GEN	ḡayr-i	other/not-of
		S----3FS2-	POSS_PRON_3FS	-hā	them

He will notify them about that through SMS messages, the Internet, and other means. سيخبرهم بذلك عن طريق الرسائل القصيرة والإنترنت وغيرها.

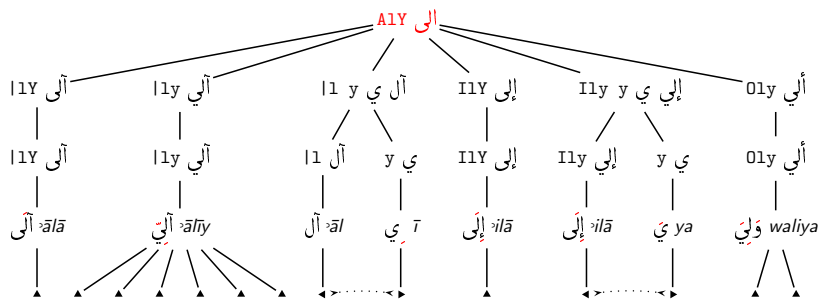
String	Token	Token Tag	Buckwalter's M-Tags	Token Form	Token Gloss
		F-----	FUT	<i>sa-</i>	will
سيخبرهم		VIIA-3MS--	IV3MS+IV+IVSUFF_MOOD:I	<i>yu-ḥbir-u</i>	he-notify
		S----3MP4-	IVSUFF_DO:3MS	<i>-hum</i>	them
بذلك		P-----	PREP	<i>bi-</i>	about/by
		SD----MS--	DEM_PRON_MS	<i>dālika</i>	that
عن		P-----	PREP	<i>ʿan</i>	by/about
طريق		N-----2R	NOUN+CASE_DEF_GEN	<i>ṭarīq-i</i>	way-of
الرسائل		N-----2D	DET+NOUN+CASE_DEF_GEN	<i>ar-rasā'il-i</i>	the-messages
القصيرة		A-----FS2D	DET+ADJ+NSUFF_FEM_SG+ +CASE_DEF_GEN	<i>al-qaṣīr-at-i</i>	the-short
والإنترنت		C-----	CONJ	<i>wa-</i>	and
		Z-----2D	DET+NOUN_PROP+ +CASE_DEF_GEN	<i>al-ʾinternet-i</i>	the-internet
		C-----	CONJ	<i>wa-</i>	and
وغیرها		FN-----2R	NEG_PART+CASE_DEF_GEN	<i>ḡayr-i</i>	other/not-of
		S----3FS2-	POSS_PRON_3FS	<i>-hā</i>	them

MorphoTrees vs. Lists

Suppose you have a list of morphological analyses [the next slide] for a given input string . . .

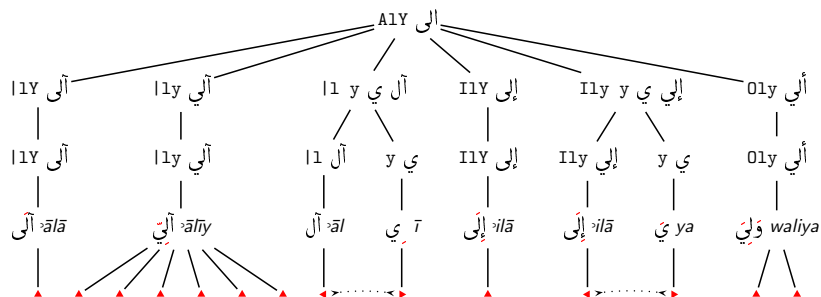
MorphoTrees vs. Lists

... organize them into a hierarchy with the string as its root



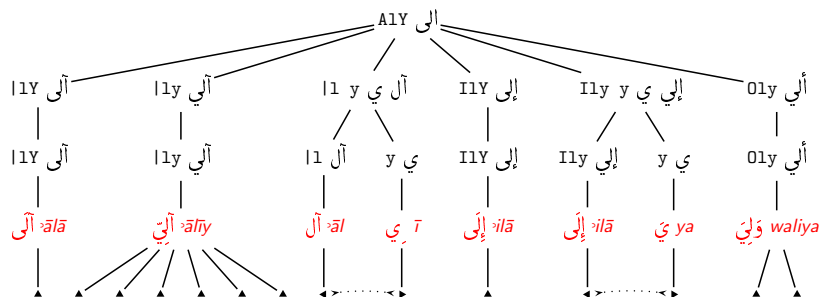
MorphoTrees vs. Lists

... organize them into a hierarchy with the **string** as its root and the **full tokens** as the leaves



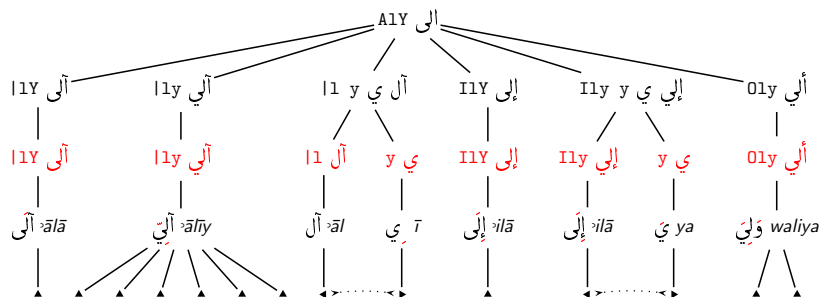
MorphoTrees vs. Lists

... organize them into a hierarchy with the **string** as its root and the **full tokens** as the leaves, grouped by their lemmas



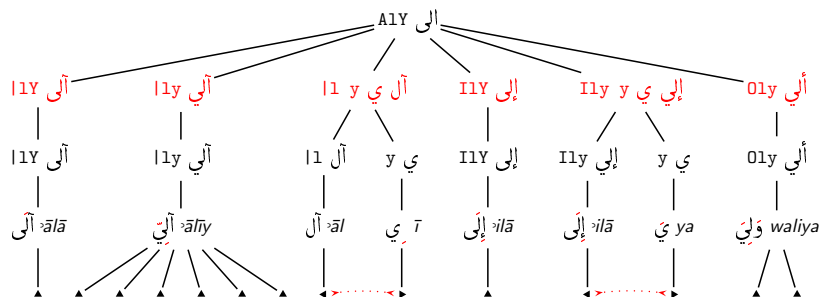
MorphoTrees vs. Lists

... organize them into a hierarchy with the **string** as its root and the **full tokens** as the leaves, grouped by their lemmas, canonical forms



MorphoTrees vs. Lists

... organize them into a hierarchy with the **string** as its root and the **full tokens** as the leaves, grouped by their lemmas, canonical forms and partitionings of the string into such forms:



Morphs	Form	Token Tag	Lemma	Glosses per Morph
laY+(null)	ʔālā	VP-A-3MS--	ʔālā	promise/take an oath + he/it
liy~	ʔālīy	A-----	ʔālīy	mechanical/automatic
liy~+u	ʔālīy-u	A-----1R	ʔālīy	mechanical ... + [def.nom.]
liy~+i	ʔālīy-i	A-----2R	ʔālīy	mechanical ... + [def.gen.]
liy~+a	ʔālīy-a	A-----4R	ʔālīy	mechanical ... + [def.acc.]
liy~+N	ʔālīy-un	A-----1I	ʔālīy	mechanical ... + [indef.nom.]
liy~+K	ʔālīy-in	A-----2I	ʔālīy	mechanical ... + [indef.gen.]
l +	ʔāl	N-----R	ʔāl	family/clan +
+ iy	-ī	S----1-S2-	ī	+ my
IilaY	ʔilā	P-----	ʔilā	to/towards
Iilay +	ʔilay	P-----	ʔilā	to/towards +
+ ya	-ya	S----1-S2-	ya	+ me
0a+liy+(null)	ʔa-lī	VIIA-1-S--	waliya	l + follow/come after + [ind.]
0a+liy+a	ʔa-liy-a	VISA-1-S--	waliya	l + follow/come after + [sub.]

List of the **morphological analyses** of the input string **AlY** الى with the quasi-functional token tags derived from the original Buckwalter's tags.

Morphs	Form	Token Tag	Lemma	Glosses per Morph
laY+(null)	ᵛāḷā	VP-A-3MS--	ᵛāḷā	promise/take an oath + he/it
liy~	ᵛālīy	A-----	ᵛālīy	mechanical/automatic
liy~+u	ᵛālīy-u	A-----1R	ᵛālīy	mechanical ... + [def.nom.]
liy~+i	ᵛālīy-i	A-----2R	ᵛālīy	mechanical ... + [def.gen.]
liy~+a	ᵛālīy-a	A-----4R	ᵛālīy	mechanical ... + [def.acc.]
liy~+N	ᵛālīy-un	A-----1I	ᵛālīy	mechanical ... + [indef.nom.]
liy~+K	ᵛālīy-in	A-----2I	ᵛālīy	mechanical ... + [indef.gen.]
l +	ᵛāl	N-----R	ᵛāl	family/clan +
+ iy	-ī	S----1-S2-	ī	+ my
IilaY	ᵛilā	P-----	ᵛilā	to/towards
Iilay +	ᵛilay	P-----	ᵛilā	to/towards +
+ ya	-ya	S----1-S2-	ya	+ me
0a+liy+(null)	ᵛa-lī	VIIA-1-S--	waliya	l + follow/come after + [ind.]
0a+liy+a	ᵛa-liy-a	VISA-1-S--	waliya	l + follow/come after + [sub.]

String Tag

VISA-1-S-----

Morphs	Form	Token Tag	Lemma	Glosses per Morph
laY+(null)	ᵛāḷā	VP-A-3MS--	ᵛāḷā	promise/take an oath + he/it
liy~	ᵛālīy	A-----	ᵛālīy	mechanical/automatic
liy~+u	ᵛālīy-u	A-----1R	ᵛālīy	mechanical ... + [def.nom.]
liy~+i	ᵛālīy-i	A-----2R	ᵛālīy	mechanical ... + [def.gen.]
liy~+a	ᵛālīy-a	A-----4R	ᵛālīy	mechanical ... + [def.acc.]
liy~+N	ᵛālīy-un	A-----1I	ᵛālīy	mechanical ... + [indef.nom.]
liy~+K	ᵛālīy-in	A-----2I	ᵛālīy	mechanical ... + [indef.gen.]
l +	ᵛāl	N-----R	ᵛāl	family/clan +
+ iy	-ī	S----1-S2-	ī	+ my
IilaY	ᵛilā	P-----	ᵛilā	to/towards
Iilay +	ᵛilay	P-----	ᵛilā	to/towards +
+ ya	-ya	S----1-S2-	ya	+ me
0a+liy+(null)	ᵛa-lī	VIIA-1-S--	waliya	I + follow/come after + [ind.]
0a+liy+a	ᵛa-līy-a	VISA-1-S--	waliya	I + follow/come after + [sub.]

String Tag

N-----RS----1-S2-----

Morphs	Form	Token Tag	Lemma	Glosses per Morph
laY+(null)	ᵛāḷā	VP-A-3MS--	ᵛāḷā	promise/take an oath + he/it
liy~	ᵛālīy	A-----	ᵛālīy	mechanical/automatic
liy~+u	ᵛālīy-u	A-----1R	ᵛālīy	mechanical ... + [def.nom.]
liy~+i	ᵛālīy-i	A-----2R	ᵛālīy	mechanical ... + [def.gen.]
liy~+a	ᵛālīy-a	A-----4R	ᵛālīy	mechanical ... + [def.acc.]
liy~+N	ᵛālīy-un	A-----1I	ᵛālīy	mechanical ... + [indef.nom.]
liy~+K	ᵛālīy-in	A-----2I	ᵛālīy	mechanical ... + [indef.gen.]
l +	ᵛāl	N-----R	ᵛāl	family/clan +
+ iy	-ī	S----1-S2-	ī	+ my
IilaY	ᵛilā	P-----	ᵛilā	to/towards
Iilay +	ᵛilay	P-----	ᵛilā	to/towards +
+ ya	-ya	S----1-S2-	ya	+ me
0a+liy+(null)	ᵛa-lī	VIIA-1-S--	waliya	l + follow/come after + [ind.]
0a+liy+a	ᵛa-liy-a	VISA-1-S--	waliya	l + follow/come after + [sub.]

Ambiguity Classes

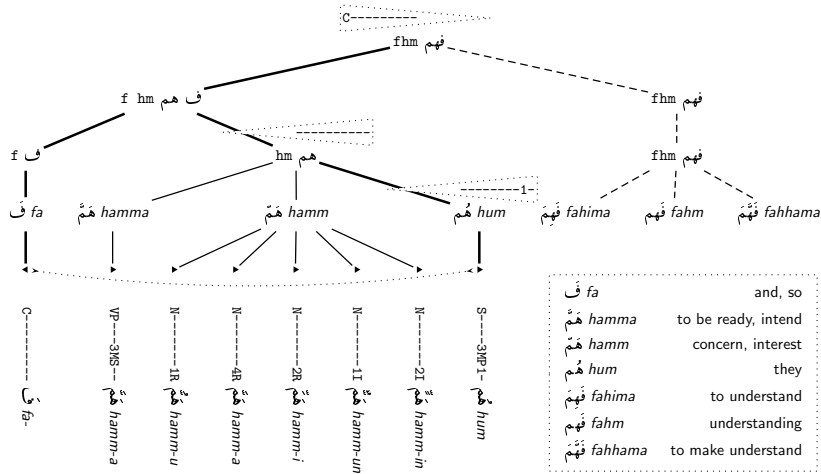
VP-A-3MS--
 A-I- ---1RS 1 S2
 NIS 1 2I
 P 4

Morphs	Form	Token Tag	Lemma	Glosses per Morph
laY+(null)	ᵛāḷā	VP-A-3MS--	ᵛāḷā	promise/take an oath + he/it
liy~	ᵛālīy	A-----	ᵛālīy	mechanical/automatic
liy~+u	ᵛālīy-u	A-----1R	ᵛālīy	mechanical ... + [def.nom.]
liy~+i	ᵛālīy-i	A-----2R	ᵛālīy	mechanical ... + [def.gen.]
liy~+a	ᵛālīy-a	A-----4R	ᵛālīy	mechanical ... + [def.acc.]
liy~+N	ᵛālīy-un	A-----1I	ᵛālīy	mechanical ... + [indef.nom.]
liy~+K	ᵛālīy-in	A-----2I	ᵛālīy	mechanical ... + [indef.gen.]
l +	ᵛāl	N-----R	ᵛāl	family/clan +
+ iy	-ī	S----1-S2-	ī	+ my
IilaY	ᵛilā	P-----	ᵛilā	to/towards
Iilay +	ᵛilay	P-----	ᵛilā	to/towards +
+ ya	-ya	S----1-S2-	ya	+ me
0a+liy+(null)	ᵛa-lī	VIIA-1-S--	waliya	l + follow/come after + [ind.]
0a+liy+a	ᵛa-liy-a	VISA-1-S--	waliya	l + follow/come after + [sub.]

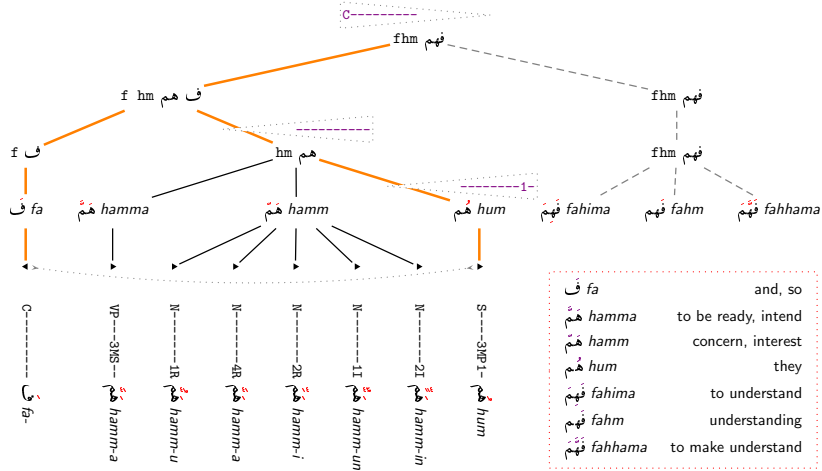
Ambiguity Classes

V -
 AP- 3 1-
 N-IA -MS2R- - --
 PIS--1--4IS-----1-S2-----

MorphoTrees with Restrictions

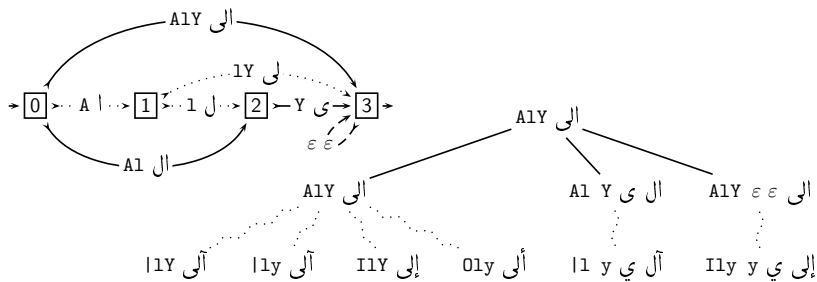


MorphoTrees with Restrictions



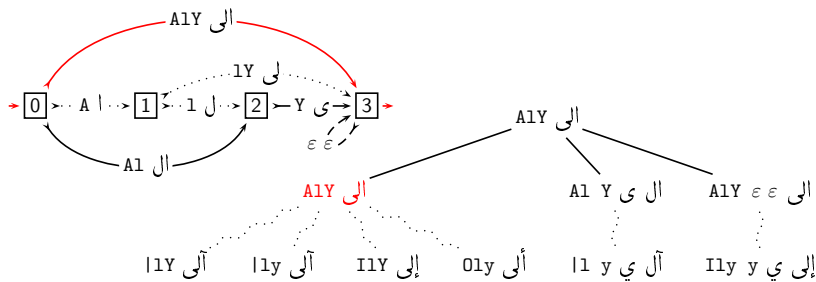
Tknz versus Tknz++

We introduce two measures for tokenization. **Tknz** is close to the evaluations that only check the partitioning determined by finding token boundaries between the **characters of the original string**, and do not, unlike **Tknz++**, require the tokenization to faithfully reconstruct the **canonical non-vocalized forms of tokens**, as is the case in MorphoTrees.



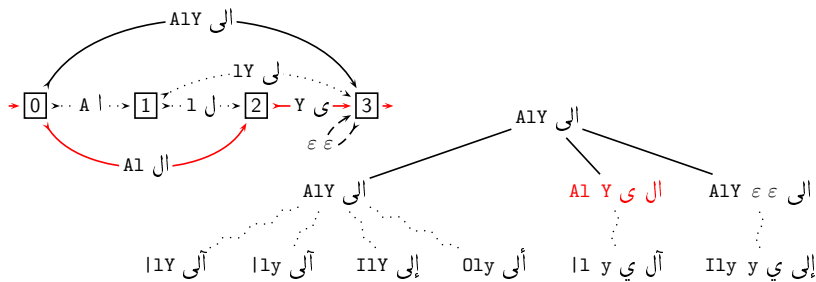
Tknz versus Tknz++

We introduce two measures for tokenization. **Tknz** is close to the evaluations that only check the partitioning determined by finding token boundaries between the **characters of the original string**, and do not, unlike **Tknz++**, require the tokenization to faithfully reconstruct the **canonical non-vocalized forms of tokens**, as is the case in MorphoTrees.



Tknz versus Tknz++

We introduce two measures for tokenization. **Tknz** is close to the evaluations that only check the partitioning determined by finding token boundaries between the **characters of the original string**, and do not, unlike **Tknz++**, require the tokenization to faithfully reconstruct the **canonical non-vocalized forms of tokens**, as is the case in MorphoTrees.



Tknz versus Tknz++

We introduce two measures for tokenization. **Tknz** is close to the evaluations that only check the partitioning determined by finding token boundaries between the **characters of the original string**, and do not, unlike **Tknz++**, require the tokenization to faithfully reconstruct the **canonical non-vocalized forms of tokens**, as is the case in MorphoTrees.

