

# The Tutorial Programme

**May 22, 2006**

<b>9:00 – 10:00</b>	<b>Introduction Arabic Orthography</b>
<b>10:00 – 11:00</b>	<b>Arabic Morphology</b>
<b>11:00 – 11:30</b>	<b>Break</b>
<b>11:30 – 12:30</b>	<b>Arabic Syntax</b>
<b>12:30 – 13:30</b>	<b>Arabic Dialects</b>

# **Tutorial Organiser**

**Nizar Habash**  
**Columbia University**  
**Center for Computational Learning Systems**

# **Tutorial Programme Committee**

**Nizar Habash**  
**Columbia University**  
**Center for Computational Learning Systems**

# Table of Contents

<b>Introduction</b>	<b>1</b>
<b>Focus of this Tutorial</b>	<b>1</b>
<b>Orthography</b>	<b>2</b>
<b>Arabic Script</b>	<b>3</b>
<b>MSA Phonology and Spelling</b>	<b>9</b>
<b>Recognizing Arabic vs. Persian/Urdu/...</b>	<b>12</b>
<b>Encoding Issues</b>	<b>14</b>
<b>Morphology</b>	<b>19</b>
<b>Derivational Morphology</b>	<b>21</b>
<b>Inflectional Morphology</b>	<b>24</b>
<b>Morphological Ambiguity</b>	<b>27</b>
<b>Arabic Computational Morphology</b>	<b>29</b>
<b>Syntax</b>	<b>31</b>
<b>Morphology and Syntax</b>	<b>31</b>
<b>Sentence Structure</b>	<b>32</b>
<b>Phrase Structure</b>	<b>35</b>
<b>Computational Resources</b>	<b>37</b>
<b>Machine Translation Issues</b>	<b>40</b>
<b>Morphology and Translation</b>	<b>40</b>
<b>Translation Divergences</b>	<b>42</b>
<b>Computational Resources</b>	<b>46</b>
<b>Dialects</b>	<b>47</b>
<b>General Definitions</b>	<b>48</b>
<b>Phonological Variation</b>	<b>49</b>
<b>Lexical Variation</b>	<b>50</b>
<b>Morphological Variation</b>	<b>51</b>
<b>Syntactic Variation</b>	<b>53</b>
<b>Code Switching</b>	<b>55</b>
<b>Computational Resources</b>	<b>56</b>
<b>Appendix</b>	<b>56</b>
<b>Resources</b>	<b>56</b>
<b>Conferences and Workshops</b>	<b>59</b>
<b>References</b>	<b>59</b>

# **Author Index**

**Nizar Habash**

## LREC 2006 Tutorial

Genoa, Italy

May 22, 2006

# Introduction to Arabic Natural Language Processing

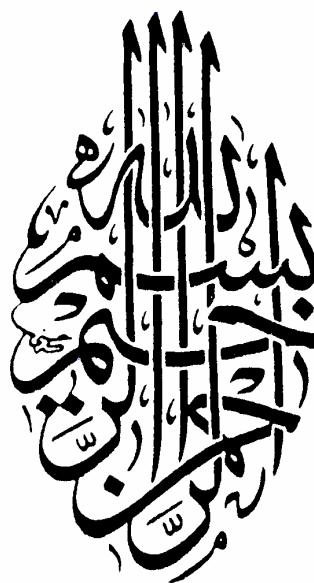
Nizar Habash

Columbia University

Center for Computational Learning Systems



- Focus of this tutorial
  - Phenomena
  - Concepts
  - Approaches & Resources
- What is 'Arabic'?
  - Arabic Script
  - Arabic Language
    - Modern Standard Arabic (MSA)
    - Arabic Dialects



## Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- Machine Translation Issues
- Dialects

3

## Road Map

- Introduction
- **Orthography**
  - Arabic Script
  - MSA Phonology and Spelling
  - Recognizing Arabic vs. Persian/Urdu/Pashto/Kurdish/Sindhi/...
  - Encoding Issues
- Morphology
- Syntax
- Machine Translation Issues
- Dialects

4

# Arabic Script

Modern Roman	A	B	G	D	E	F	Z	H	I	K	L	M	N	O	P	Q	R	S	T			
Early Latin	A	B	<	>	E	F	Z	H	z	k	l	m	n	o	p	q	r	s	t			
Greek	Α	Δ	Γ	Δ	Ξ	Α	Ζ	Β	Ζ	κ	λ	μ	ν	ο	π	φ	ρ	σ	τ			
Phoenician	𐤀	𐤁	𐤂	𐤃	𐤄	𐤅	𐤆	𐤇	𐤈	𐤉	𐤊	𐤋	𐤌	𐤍	𐤎	𐤏	𐤐	𐤑	𐤒	𐤓	𐤔	𐤕
Early Aramaic	𐤀	𐤁	𐤂	𐤃	𐤄	𐤅	𐤆	𐤇	𐤈	𐤉	𐤊	𐤋	𐤌	𐤍	𐤎	𐤏	𐤐	𐤑	𐤒	𐤓	𐤔	𐤕
Nabataean	𐤀	𐤁	𐤂	𐤃	𐤄	𐤅	𐤆	𐤇	𐤈	𐤉	𐤊	𐤋	𐤌	𐤍	𐤎	𐤏	𐤐	𐤑	𐤒	𐤓	𐤔	𐤕
Arabic	ا	ب	ت	ث	ج	ح	خ	د	ذ	ر	ز	س	ش	ص	ض	ط	ظ	ع	ف	ق	ك	گ

©Mamoun Sakka 1997

5

# Arabic Script

Arabic script is an alphabet with allographic variants, optional zero-width diacritics and common ligatures.

الخط العربي

Arabic script is used to write many languages: Arabic, Persian, Kurdish, Urdu, Pashto, etc.

6

## Arabic Script

### Alphabet

- letter forms

ع ط ص س ر د ح ب ا  
ء ي و ه ن م ل ل و

- letter marks

- Arabic only

•    •    •  
— — —    — — —    — — —    —

- Other languages

- Persian, Kurdish, Urdu, Pashto, etc.

•    •    •    •    •  
— — — — —    — — — — —    — — — — —    — — — — —

- *OCR output ambiguity*

7

## Arabic Script

### Alphabet (MSA)

- letters (form+mark)

- Distinctive

ب ب ت ت ش س ش  
/b/ /t/ /θ/ /t/ /s/ /ʃ/

- Non-distinctive

أ إ آ إ ئ و ء  
/?/  
*glottal stop aka hamza*

8



# Arabic Script

## Letter Shapes

- No distinction between print and handwriting
- No capitalization
- Right-to-left
- Ambiguous shapes
- Connective letters
- Disconnective letters

ز	د	ا	ن	ب	ك	م	ش	غ	Stand alone
			ز	ب	ك	م	شد	غ	initial
			ن	ب	ك	م	شد	غ	medial
ز	د	ا	ن	ب	ك	م	ش	غ	final

# Arabic Script

## Letter shaping

ك ت ب = ك ت ب ← ب ت ك  
 /katab/  
 to write

ك ت ا ب = ك ت ا ب ← ب ا ت ك  
 /kitāb/  
 book

## Arabic Script

### Diacritics

- Zero-width characters
- Used for short vowels

كَتَبَ /katab/ *to write*

- Nunation is used for nominal indefinite marker in MSA

كِتَابٌ /kitābun/ *a book*

Nunation	Vowel
بَ /ban/	بِ /ba/
بُ /bun/	بِ /bu/
بِ /bin/	بِ /bi/

11

## Arabic Script

### Diacritics

- No-vowel marker (*sukun*)

مَكْتَبٌ /maktab/ *office*

- Double consonant marker (*shadda*)

كَتَّبَ /kattab/ *to dictate*

- Combinable

بُّ    بَّ    بَّ  
/bbu/    /bbin/    /bban/

No Vowel
بْ /b/
Double Consonant
بَّبْ /bb/

12

## Arabic Script

### Putting it together

#### Simple combination

Arab /ʕarab/ ع ر ب ← ع ر ب = عرب

West /ɤarb/ غ ر ب ← غ ر ب = غرب

#### Ligatures

Peace /salām/ س ل ا م ← س ل ا م سلام 

13

## Arabic Script

### Tatweel

- 'elongation'
- aka kashida
- used for text highlight and justification

حقوق الانسان

حقوق الانسان

حقوق الانسان

حقوق الانسان

human rights /ħuqūq alʔinsān/

14

## Arabic Script

- Different styles
- High fluidity
- Optional ligatures
- Vertical arrangements

Arabic	Muhammad	algebra
عربي	محمد	الجبر
عربي	محمد	الجبر
عربي	محمد	الجبر
عربي	محمد	الجبر
/ʕarabi/	/muħammad /	/alʕabr/ <sub>15</sub>

## Arabic Script

### “Arabic” Numerals

- Decimal system
- Numbers written left-to-right in right-to-left text

استقلت الجزائر في سنة 1962 بعد 132 عاما من الاحتلال الفرنسي.

Algeria achieved its independence in 1962 after 132 years of French occupation.

- Three systems of enumeration symbols that vary by region

<b>Western Arabic</b> <i>Tunisia, Morocco, etc.</i>	0	1	2	3	4	5	6	7	8	9
<b>Indo-Arabic</b> <i>Middle East</i>	٠	١	٢	٣	٤	٥	٦	٧	٨	٩
<b>Eastern Indo-Arabic</b> <i>Iran, Pakistan, etc.</i>	۰	۱	۲	۳	۴	۵	۶	۷	۸	۹

## Road Map

- Introduction
- **Orthography**
  - Arabic Script
  - **MSA Phonology and Spelling**
  - Recognizing Arabic vs. Persian/Urdu/Pashto/Kurdish/Sindhi/...
  - Encoding Issues
- Morphology
- Syntax
- Machine Translation Issues
- Dialects

17

## MSA Phonology and Spelling

- Phonological profile of Standard Arabic
  - 28 Consonants
  - 3 short vowels, 3 long vowels, 2 diphthongs
- Arabic spelling is mostly phonemic ...
  - Letter-sound correspondence

ء آ أ إ و ئ ي ا ب ت ة ث ج ح خ د ذ ر ز س ش ص ض ط ظ ع غ ف ق ك ل م ن ه و ي  
ī j ū w h n m l k q f ʙ ʕ δ ʔ d ʒ ʃ s z r δ d x h ʔ θ t b ā ʔ

18

## MSA Phonology and Spelling

- Arabic spelling is mostly phonemic ...

### **Except for**

- Medial short vowels can only appear as diacritics
- Diacritics are optional in most written text
  - Except in holy scripture
  - Present diacritics mark syntactic/semantic distinctions
    - كَتَبَ /katab/ to write كُتِبَ /kutib/ to be written
    - حُبَّ /ḥubb/ love حَبَّ /ḥabb/ seed
- Dual use of ا, و, ي as consonant and long vowel
  - ا (/ā/, /ā/) و (/w/, /ū/) ي (/j/, /ī/)

19

## MSA Phonology and Spelling

- Arabic spelling is mostly phonemic ...

### **Except for (continued)**

- Morphophonemic characters
  - Feminine marker ة (*ta marbuta*)
    - كبير /kabīr/ (big ♂) كبيرة /kabīra/ (big ♀)
  - Derivation marker
    - /ʕaʕa/ (to disobey عصى) (a stick عصا)
- Hamza variants (6 characters for one phoneme!)
  - بهاء بهاءه بهائه (ء أأؤئ) /bahaʕ/ + 3MascSing (his glory)

20

## MSA Phonology and Spelling

- Arabic spelling can be ambiguous
  - optional diacritics and dual use of letter
- But how ambiguous? Really?
- Classic example
  - ths s wht n rbc txt lks lk wth n vwls
  - this is what an Arabic text looks like with no vowels
- Not exactly true
  - Long vowels are always written
  - Initial vowels are represented by an 'alef
  - Some final short vowels are represented

ths is wht an Arbc txt lks lik wth no vwls

*Will revisit ambiguity in more detail again under morphology discussion*

21

## Road Map

- Introduction
- **Orthography**
  - Arabic Script
  - MSA Phonology and Spelling
  - **Recognizing Arabic vs. Persian/Urdu/Pashto/Kurdish/Sindhi/...**
  - Encoding Issues
- Morphology
- Syntax
- Machine Translation Issues
- Dialects

22

# Arabic Script

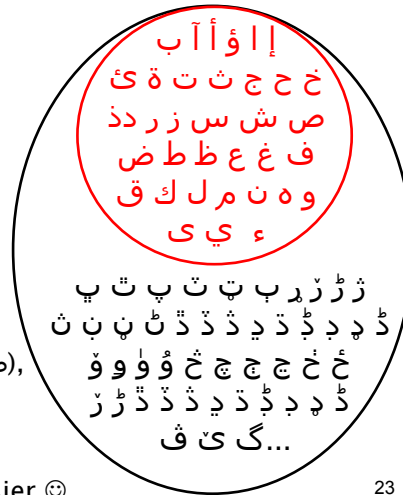
## Other languages

### Arabic

- No more than 3 dots
- Dots either above or below
- Marks are 1/2/3 dots, hamza (ء) or madda (~) only
- Rare borrowing for foreign words
  - پ/p/, ف/v/, گ/g/, چ/tʃ/
  - regionally variable

### Not Arabic

- Extra marks: haft (v), ring (o), taa (ط), four dots (::), vertical dots (:)
- Some Numerals (٤,٥,٦)



Once you learn the alphabet, it is easier ☺

23

Arabic  
 Not Arabic

بۆنە سووتی جگه رو بۆچی نه بی دل به که باب

بۆچی نه روا له ته نم روچی ره وان میسلی شه هاب (١)

بۆله سه ر چاوهیی چاو هه ئنه قوتی ره شحه یی خوین (٢)

بۆچ له فه ووارهیی موژگان نه تکی قه ترییی ناب

بۆله به ر نائه نه بی جه لقهی جه نغم به سروود

بۆله به ر گریه نه بی چه شمه ی چه شمم به سه راب

موونسی روژو شه ووم باعیسی نارامی دلیم (٤)

رویی وو من له غه می که وتمه نیو به حری عه زاب

به وقووعی سه فه ری قنادری نوستاد خدری (٥)

به جه فا عه یشمی تال کرد فه له کی خانه خه راب

چه نکا ونه ی لی مه ده موتریب که له به ر فیرقه تی نه و (٦)

رنه کی روچه له گویم نه غمه ی ناوازو روباب (٧)

ساغیری مه ی مه ده ساقی که له به ر دووری نه و (٨)

تاله و دک زه هری هه لایل له مه زاقم مه ی ناب (٩)



- Arabic
- Not Arabic

سجل... انا عربي...  
ورقم بطاقتي خمسون الف  
واطفالي ثمانية  
وتاسعهم سيأتي بعد صيف  
فهل تغضب  
سجل... انا عربي...  
واعمل مع رفاق الكدح في محجر  
واطفالي ثمانية  
اسلّ لهم رغيف الخبز والاثواب والدفتن  
من الصخر  
ولا اتوسل الصدقات من بابك  
ولا اصغر امام بلاط اعتابك  
فهل تغضب

25

سجل... انا عربي للشاعر: محمود درويش

- Arabic
- Not Arabic

شیلی بیٹی کے نام

تجھے جب بھی کوئی دکھ دے  
اس دکھ کا نام بیٹی رکھنا  
جب میرے سفید بال  
تیرے گالوں پر آن ہنسیں، رو لینا  
میرے خواب کے دکھ پہ سو لینا  
جن کھیتوں کو ابھی اگنا ہے  
ان کھیتوں میں<sup>3</sup>

# Road Map

- Introduction
- Orthography
  - Arabic Script
  - MSA Phonology and Spelling
  - Recognizing Arabic vs. Persian/Urdu/Pashto/Kurdish/Sindhi/...
  - Encoding Issues
- Morphology
- Syntax
- Machine Translation Issues
- Dialects

27

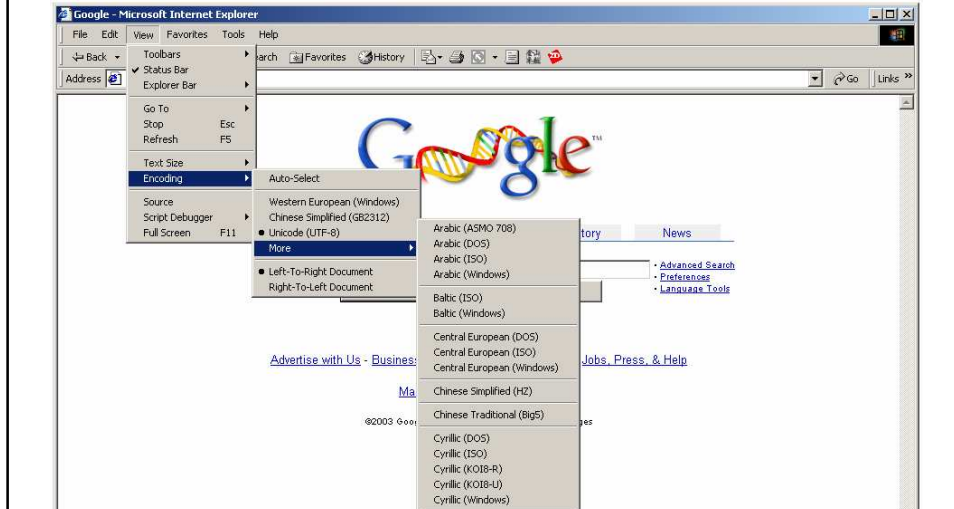
# Encoding Issues

- Encoding Arabic
  - Data entry, storage, and display
  - Ease of use for *Arabic-illiterate* users
  - Multi-script support
  - Multilingual support (extended Arabic characters)
- Types of Encoding
  - Machine character sets
    - Graphemic (shape insensitive, logical order)
    - Allographic (shape/direction sensitive) [obsolete]
  - Human accessible
    - Transliteration
    - Phonetic spelling (IPA)
    - Romanization

28

# Encoding Issues

- Many Conflicting Character Sets for Arabic



## Encodings

- CP-1256
  - Commonly used
  - 1-byte characters
  - Widely supported input/display
  - Minimal support for extended Arabic characters
  - bi-script support (Roman/Arabic)
  - Tri-lingual support: Arabic, French, English (ala ANSI)

Codepage 1256 - Arabic Windows

	-0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-A	-B	-C	-D	-E	-F
0-		0001	0002	0003	0004	0005	0006	0007	0008	0009	000A	000B	000C	000D	000E	000F
1-	0010	0011	0012	0013	0014	0015	0016	0017	0018	0019	001A	001B	001C	001D	001E	001F
2-		!	"	#	\$	%	&	'	(	)	*	+	,	-	.	/
3-	0020	0021	0022	0023	0024	0025	0026	0027	0028	0029	002A	002B	002C	002D	002E	002F
4-	0030	0031	0032	0033	0034	0035	0036	0037	0038	0039	003A	003B	003C	003D	003E	003F
5-	0040	0041	0042	0043	0044	0045	0046	0047	0048	0049	004A	004B	004C	004D	004E	004F
6-	0050	0051	0052	0053	0054	0055	0056	0057	0058	0059	005A	005B	005C	005D	005E	005F
7-	0060	0061	0062	0063	0064	0065	0066	0067	0068	0069	006A	006B	006C	006D	006E	006F
8-	0070	0071	0072	0073	0074	0075	0076	0077	0078	0079	007A	007B	007C	007D	007E	007F
9-	0080	0081	0082	0083	0084	0085	0086	0087	0088	0089	008A	008B	008C	008D	008E	008F
A-	0090	0091	0092	0093	0094	0095	0096	0097	0098	0099	009A	009B	009C	009D	009E	009F
B-	00A0	00A1	00A2	00A3	00A4	00A5	00A6	00A7	00A8	00A9	00AA	00AB	00AC	00AD	00AE	00AF
C-	00B0	00B1	00B2	00B3	00B4	00B5	00B6	00B7	00B8	00B9	00BA	00BB	00BC	00BD	00BE	00BF
D-	00C0	00C1	00C2	00C3	00C4	00C5	00C6	00C7	00C8	00C9	00CA	00CB	00CC	00CD	00CE	00CF
E-	00D0	00D1	00D2	00D3	00D4	00D5	00D6	00D7	00D8	00D9	00DA	00DB	00DC	00DD	00DE	00DF
F-	00E0	00E1	00E2	00E3	00E4	00E5	00E6	00E7	00E8	00E9	00EA	00EB	00EC	00ED	00EE	00EF



## Encoding Issues Arabic Display

- **Memory (logical order) →**

ÔÇÑßÊ ÝáÓØíä (Palestine) Ýí ÇæääÈíÇĬ (Olympics) 2000 æ 2004.  
نيطسلف تڪراش (Palestine) دايمل و ا ي ف (Olympics) 2000 و 2004.

*or this way for those with direction-bias*



.4002 æ 0002 )scipmylO( İÇiÈääæÇ íÝ )enitselaP( äiØÓáÝ ÊßÑÇÔ  
شاركف فلسطين (Palestine) في اولمبياد (Olympics) 2000 و 2004.

33

## Encoding Issues Arabic Display

- **Memory (logical order)**

ÔÇÑßÊ ÝáÓØíä (Palestine) Ýí ÇæääÈíÇĬ (Olympics) 2000 æ 2004.  
نيطسلف تڪراش (Palestine) دايمل و ا ي ف (Olympics) 2000 و 2004.

- **Display (visual order)**

- Bidirectional (BiDi) support

- Numbers and Roman script

.2004 æ 0002 )scipmylO( İÇiÈääæÇ íÝ )enitselaP( äiØÓáÝ ÊßÑÇÔ  
شاركف فلسطين (Palestine) في اولمبياد (Olympics) 2000 و 2004.

- Letter and ligature shaping

.2004 æ 0002 )scipmylO( İÇiÈääæÇ íÝ )enitselaP( äiØÓáÝ ÊßÑÇÔ  
شاركف فلسطين (Palestine) في اولمبياد (Olympics) 2000 و 2004.

34

# Display Problems

		Display Encoding			
		CP-1256	ISO-8859	Unicode	Western
Actual Encoding	CP-1256	تدشين منطقة حرة في دبي للتجارة الالكترونية	ة حرة تدشيل كلظ ترنبة دب ففتجارة افاف	γλ Ηγλ γ ψ Ωgāā Αλ	ÈÏÔia ääøþÉ ÎÑÉ Ýí ÎÊi ääËÏÑÉ ÇáÇääÈÑæäiÉ
	ISO-8859	ة حرة هو هتدش نتتجارة هدب هل هوهوانامتر	تدشين منطقة حرة في دبي للتجارة الالكترونية	γ 亲既 ηλ γ ψλ lÖgGG 親	ÈÏÔëæ äæ×âÉ ÎÑÉ áé ÎÊë ääËÏÇÑÉ ÇáÇääÈÑæëÉ
	Unicode	آ؟طهط طظظظ «آ» ظ...ظظظظظظظظظظ ظظظظظظظظظظ ظظظظظظظظظظظظ ظظظظظظظظظظظظ ظظظظظظظظظظظظ	لعللظظظظظظظظ ظظظظظظظظظظ لعللظظظظظظظظ ظظظظظظظظظظ ظظظظظظظظظظظظ ظظظظظظظظظظظظ ظظظظظظظظظظظظ ظظظظظظظظظظظظ	تدشين منطقة حرة في دبي للتجارة الالكترونية	i»¿øøøøùšù+ ù..ù+ø·ù,øø ø-ø±øø ùllùš øøøùš ù,,ù,,øøø-øøø±øøø øøù,,øøù,,ùføøø±ù ^ù+ùšøø

- Wrong encoding
- Partial support problems

# Encoding Issues Arabic Input

- Standard graphemic keyboard
- Logical order input

~	!	@	#	\$	%	^	&	*	(	)	-	+	
ذ	1	2	3	4	5	6	7	8	9	0	.	=	
	Q	W	E	R	T	Y	U	ا	و	ح	{	}	<
	ض	ص	ث	ق	ف	غ	ع	ه	خ	ج	[	]	د
	A	S	D	F	G	H	I	J	K	L	/	:	"
	ش	س	د	ف	ج	ه	ا	ي	ك	ل	;	;	ط
	Z	X	C	V	B	N	M	<	>	:	?	?	
	ذ	خ	ع	و	ر	ل	م	<	>	:	?	?	

↓

م ا ل س      →      س ل ا م

# Encodings

## Buckwalter Encoding

- Romanization
  - One-to-one mapping to Arabic script spelling
  - Left-to-right
  - Easy to learn/use
  - Human & machine compatible
- Commonly used in NLP
  - Penn Arabic Tree Bank
- Some characters can be modified to allow use with XML and regular expressions
- Roman input/display
- Monolingual encoding (can't do English and Arabic)
- Minimal support for extended Arabic characters

ء	ا	أ	ؤ	إ	ئ	أ	ب	ة	ت	ث	ج	ح	خ	د	ذ	ر	ز	س	ش	ص	ض	ط	ظ	ع	غ	ـ	ف	ق	ك	ل	م	ن	ه	و	ي	ـ	ن	ك	ا	ـ	ـ	ـ	ـ	ـ	ـ	ـ	ـ
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

37

# Road Map

- Introduction
- Orthography
- **Morphology**
  - Derivational Morphology
  - Inflectional Morphology
  - Morphological Ambiguity
  - Arabic Computational Morphology
- Syntax
- Machine Translation Issues
- Dialects

38

# Morphology

- Type
  - Concatenative: prefix, suffix, circumfix
  - Templatic: root+pattern
- Function
  - Derivational
    - Creating new words
    - *Mostly templatic*
  - Inflectional
    - Modifying features of words
      - Tense, number, person, mood, aspect
    - Mostly concatenative

39

# Road Map

- Introduction
- Orthography
- **Morphology**
  - Derivational Morphology
  - Inflectional Morphology
  - Morphological Ambiguity
  - Arabic Computational Morphology
- Syntax
- Machine Translation Issues
- Dialects

40



# Derivational Morphology

- Templatic Morphology

- Root

ك ت ب  
k t b

- Pattern



- Lexeme

مكتوب  
maktūb  
written

كاتب  
kātib  
writer

Lexeme.Meaning =  
(Root.Meaning+Pattern.Meaning)\*Idiosyncrasy.Random

41

# Derivational Morphology Root Meaning

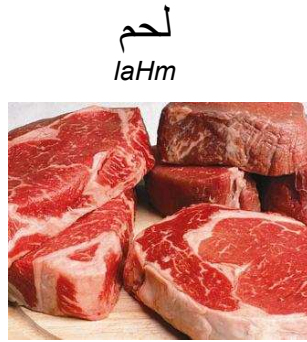
- ك ت ب KTB = notion of "writing"



42

## Derivational Morphology *Root Meaning*

- LHM-1
- Notion of “meat”
  - لحم /laḥm/
    - Meat
  - لحام /laḥḥām/
    - Butcher



43

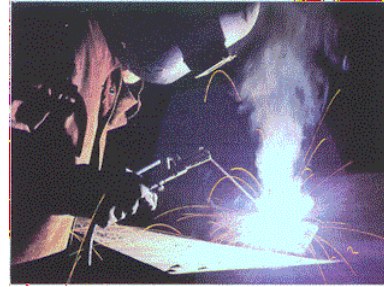
## Derivational Morphology *Root Meaning*

- LHM-2
- Notion of “battle”
  - ملحمة /malḥama/
    - Fierce battle
    - Massacre
    - Epic



## Derivational Morphology *Root Meaning*

- LHM-3
- Notion of “soldering”
  - لحم /laḥam/
    - Weld, solder, stick, cling
  - التحم /iltaḥam/
    - Be welded/soldered/fused
  - ملتحم /multaḥim/
    - Welded, soldered, fused



45

## Derivational Morphology *Pattern Meaning*

- Verb Pattern Meaning is hard to define

Pattern	Pattern Meaning	Example	Gloss
<b>I</b> 1a2a3	Basic sense of root	ktb → katab	write
<b>II</b> 1a22a3	Intensification, causation	ktb → kattab	dictate
<b>III</b> 1aA2a3	Interaction with others	ktb → kaAtab	correspond with
<b>IV</b> Aa12a3	Causation	jls → Ajlas	seat
<b>V</b> ta1a22a3	Reflexive of Pattern II	Elm → taEal~am	learn
<b>VI</b> ta1aA2a3	Reflexive of Pattern III	ktb → takaAtab	correspond
<b>VII</b> Ain1a2a3	Passive of Pattern I	ktb → Ainkatab	subscribe/enroll
<b>VIII</b> Ai1ta2a3	Acquiescence, exaggeration	ktb → Aiktatab	register
<b>IX</b> Ai12a33	Transformation	Hmr → AiHmarr	Turn red/blush
<b>X</b> Aista12a3	Requirement	ktb → Aistaktab	ask/make_write

46

# Road Map

- Introduction
- Orthography
- **Morphology**
  - Derivational Morphology
  - **Inflectional Morphology**
  - Morphological Ambiguity
  - Arabic Computational Morphology
- Syntax
- Machine Translation Issues
- Dialects

47

# Inflectional Morphology

- Derivational Morphology
  - Lexeme  $\approx$  Root + Pattern
- Inflectional Morphology
  - Word = Lexeme + Features
- Features
  - Part-of-speech
    - *Traditional*: Noun, Verb, Particle
    - *Computational*: N, PN, V, Adj, Adv, P, Pron, Num, Conj, Det, Aux, Pun, IJ, and others
  - Noun-specific
    - Number: singular, dual, plural, collective
    - Gender: masculine, feminine, Neutral
    - Definiteness: definite, indefinite
    - Case: nominative, accusative, genitive
    - Possessive clitic

48

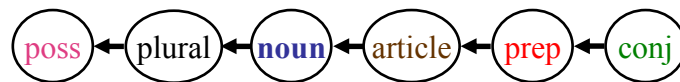
# Inflectional Morphology

- Features (continued)
  - Verb-specific
    - Aspect: perfective, imperfective, imperative
    - Voice: active, passive
    - Tense: past, present, future
    - Mood: indicative, subjunctive, jussive
    - Subject (Person, Number, Gender)
    - Object clitic
  - Others
    - Single-letter conjunctions
    - Single-letter prepositions

49

# Inflectional Morphology

## Nouns

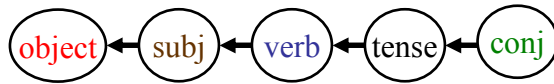


<p>وكبيوتنا /wakabiyūtinā/ و+ك+بيوت+نا wa+ka+biyūt+nā and+like+houses+our <i>And like our houses</i></p>	<p>وللمكتبات /walilmaktabāt/ و+ل+ال+مكتبة+ات wa+li+al+maktaba+āt and+for+the+library+plural <i>And for the libraries</i></p>
--	--

- Morphotactics (e.g. ل+ال → لل)
- Arabic *Broken Plurals* (templatic)

50

## Inflectional Morphology Verbs



فقلناها /faqlnāhā/ ف+قال+نا+ها fa+qul+na+hā so+said+we+it So we said it.	وسنقولها /wasanaqūluhā/ و+س+ن+قول+ها wa+sa+na+qūl+u+hā and+will+we+say+it And we will say it
---	---

- Morphotactics
- Subject conjugation (suffix or circumfix)

51

## Inflectional Morphology

- Perfect verb subject conjugation (*suffixes only*)

	Singular	Dual	Plural
1	كتبتُ katabtu	كتبنا katabnā	
2	كتبتَ katabta	كتبتما katabtumā	كتبتم katabtum
3	كتبَ kataba	كتبَا katabā	كتبوا katabtū

- Imperfect verb subject conjugation (*prefix+suffix*)

	Singular	Dual	Plural
1	أكتبُ aktubu	نكتبُ naktubu	
2	أكتبَ taktubu	نكتبان taktubān	نكتبون taktubūn
3	يكتبُ yaktubu	يكتبان yaktubān	يكتبون yaktubūn

52

*Feminine form and other verb moods not shown*

## Road Map

- Introduction
- Orthography
- **Morphology**
  - Derivational Morphology
  - Inflectional Morphology
  - **Morphological Ambiguity**
  - Arabic Computational Morphology
- Syntax
- Machine Translation Issues
- Dialects

53

## Morphological Ambiguity

- Derivational ambiguity
  - قاعدة: basis/principle/rule, military base, Qa'ida/Qaeda/Qaida
- Inflectional ambiguity
  - تكتب: you write, she writes
  - Segmentation ambiguity
    - وجد: he found; وجد: and+grandfather
    - ل:لغة: for a language; ل:اللغة: for the language
- Spelling ambiguity
  - Optional diacritics
    - كاتب: /kātib/ writer , /kātab/ to correspond
  - Suboptimal spelling
    - Hamza dropping: أ, إ → ا
    - Undotted ta-marbuta: ة → ه
    - Undotted final ya: ي → ى

54

# Morphological Ambiguity

- Multiple sources of ambiguity

بين

- /bayyana/ Verb *he declared/demonstrated*
- /bayyanna/ Verb *they [feminine] declared/demonstrated*
- /bayyin/ Adj *clear/evident/explicit*
- /bayna/ Prep *between/among*
- /biyin/ Proper Noun *in Yen*
- /biyn/ Proper Noun *Ben*

- Hard to measure specific causes of ambiguity

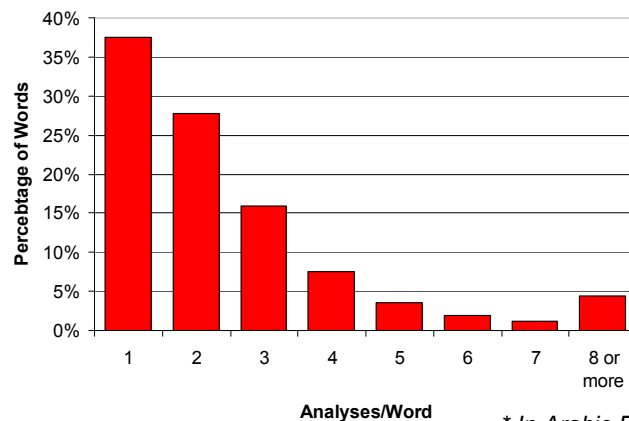
- Derivational ambiguity\* (diacritized tokens)
  - 1.09 entries/token
  - 1.01 entries/token (within same part-of-speech)
- Spelling ambiguity\* (undiacritized tokens)
  - 1.28 entries/token
  - 1.08 entries/token (within same part-of-speech)

55

\* in Buckwalter's Lexicon (~40,000 lexemes)

# Morphological Ambiguity

- Average overall ambiguity\* is 2.5 analyses/word
  - Compare to English ENGTWOL ambiguity (1.7-2.2 analyses/word)



56

\* In Arabic Penn Treebank 1



## Road Map

- Introduction
- Orthography
- **Morphology**
  - Derivational Morphology
  - Inflectional Morphology
  - Morphological Ambiguity
  - **Arabic Computational Morphology**
- Syntax
- Machine Translation Issues
- Dialects

57

## Arabic Computational Morphology

- Representation units
  - Natural token وللمكتبات
    - White space separated strings (as is)
    - Can include extra characters (e.g. tatweel/kashida)
  - Word وللمكتبات
  - Segmented word
    - Can include any degree of morphological analysis
    - Pure segmentation: و ل لمكتبات
    - Arabic Treebank tokens (with recovery of some deleted/modified letters): و ل المكتبات

58

## Arabic Computational Morphology

- Representation units (continued)
  - Prefix + Stem + Suffix
    - ولل+مكتب+ات
    - Can create more ambiguity
  - Lexeme + Features
    - [ل +و+ Def +Plural] مكتبة
  - Root + Pattern + Features
    - [و+ ل +Def +Plural] + م3ا21اة + كتب
    - Very abstract
  - Root + Pattern + Vocalism + Features
    - [و+ ل +Def +Plural] + م321ة + كتب
    - Very very abstract

59

## Arabic Computational Morphology

- Approaches
  - Finite state machines (Beesely,2001) (Kiraz,2001) (Habash et al, 2005b)
  - Concatenative analysis/generation (Buckwlater,2002) (Cavalli-Sforza et al, 2000)
  - Lexeme+Feature analysis/generation (Habash, 2004)
  - Shallow stemming (Darwish,2002) (Aljlayl and Frieder 2002)
  - Machine learning (Diab et al,2004) (Lee et al,2003) (Rogati et al, 2003) (Habash & Rambow 2005a) (Smith et al, 2005)
- Packages
  - AMIRA: Arabic SVM Toolkit (Diab et al, 2004)
  - MADA: Morphological Analysis and Disambiguation for Arabic (Habash and Rambow 2005a)
- Issues
  - Appropriateness of system representation for an application
    - Machine Translation vs. Information Retrieval
    - Arabic spelling vs. phonetic spelling
  - System coverage
  - System extensibility
  - Availability to researchers
  - Use for analysis and generation

60

# Road Map

- Introduction
- Orthography
- Morphology
- **Syntax**
  - **Morphology and Syntax**
  - Sentence Structure
  - Phrase Structure
  - Computational Resources
- Machine Translation Issues
- Dialects

61

# Morphology and Syntax

- Rich morphology crosses into syntax
  - Pro-drop / Subject conjugation
  - Verb subcategorization and object clitics
    - Verb<sub>transitive</sub>+subject+object
    - Verb<sub>intransitive</sub>+subject *but not* Verb<sub>intransitive</sub>+subject+object
    - Verb<sub>passive</sub>+subject *but not* Verb<sub>passive</sub>+subject+object
- Morphological interactions with syntax
  - Agreement
    - **Full**: e.g. Noun-Adjective on number, gender, and definiteness
    - **Partial**: e.g. Verb-Subject on gender (in VSO order)
  - Definiteness
    - Noun compound formation, copular sentences, etc.
    - Nouns+DefiniteArticle, Proper Nouns, Pronouns, etc.

62

## Morphology and Syntax

- Morphological interactions with syntax (continued)
    - Case
      - MSA is case marking: nominative, accusative, genitive
      - Almost-free word order
      - Case is often marked with optionally written short vowels
        - This effectively limits the word-order freedom in published text
    - Agglutination
      - Attached prepositions create words that cross phrase boundaries
- |                   |                         |
|-------------------|-------------------------|
| ل+المكتبات        | li+Almaktabāt           |
| for the-libraries | [PP li [NP Almaktabāt]] |
- Some morphological analysis (*minimally segmentation*) is necessary even for statistical approaches to parsing

63

## Road Map

- Introduction
- Orthography
- Morphology
- **Syntax**
  - Morphology and Syntax
  - **Sentence Structure**
  - Phrase Structure
  - Computational Resources
- Machine Translation Issues
- Dialects

64

## Sentence Structure

### *Two types of Arabic Sentences*

- Verbal sentences
  - [Verb Subject Object] (VSO)
    - كتب الاولاد الاشعار
    - Wrote the-boys the-poems
    - The boys wrote the poems
- Copular sentences
  - [Topic Complement]
  - الاولاد شعراء
  - the-boys poets
  - The boys are poets

65

## Sentence Structure

- Verbal sentences
  - Verb agreement with gender only
    - كتب الولد\الاولاد wrote<sub>3MascSing</sub> the-boy/the-boys
    - كتبت البنت\البنات wrote<sub>3FemSing</sub> the-girl/the-girls
  - Pronominal subjects are conjugated
    - كتبت wrote-you<sub>MascSing</sub>
    - كتبتُم wrote-you<sub>MascPlur</sub>
    - كتبوا wrote-they<sub>MascPlur</sub>
  - Passive verbs
    - Same structure: Verb<sub>passive</sub> Subject<sub>underlyingObject</sub>
    - Agreement with surface subject

66

## Sentence Structure

- Verbal sentences
  - Common structural ambiguity
    - *Third masculine/feminine singular are structurally ambiguous*
      - Verb<sub>3MascSingular</sub> Noun<sub>Masc</sub>  
Verb subject=he object=Noun  
Verb subject=Noun
    - Passive and active forms are often similar in standard orthography
      - كتب /kataba/ he wrote
      - كُتِبَ /kutiba/ it was written

67

## Sentence Structure

- Copular sentences
  - [Topic Complement]  
Definite Topic, Indefinite Complement
    - الولد شاعر  
the-boy poet  
*The boy is a poet*
  - [Auxiliary Topic Complement]  
Auxiliaries (*kāna and her sisters*)
    - Tense, Negation, Transformation, Persistence
    - كان الولد شاعرا *was* the-boy poet *The boy was a poet*
    - ليس الولد شاعرا *is-not* the-boy poet *The boy is not a poet*
  - Inverted order is expected in certain cases
    - Indefinite topic  
عندي كتاب /ʕandi kitābun/ at-me a-book *I have a book*

68

## Sentence Structure

- Copular sentences
  - Types of complements
    - Noun/Adjective/Adverb
      - الولد ذكي the-boy smart *The boy is smart*
    - Prepositional Phrase
      - الولد في المكتبة the-boy in the-library *The boy is in the library*
    - Copular-Sentence
      - الولد كتابه كبير [the-boy [book-his big]] *The boy, his book is big*
    - Verb-Sentence
      - الاولاد كتبوا الاشعار
      - [the-boys [wrote-they poems]] *The boys wrote the poems*
      - Full agreement in this order (SVO)
      - الاشعار كتبها الاولاد
      - [the-poems [wrote-it the boys]] *The poems, the boys wrote*

69

## Road Map

- Introduction
- Orthography
- Morphology
- **Syntax**
  - Morphology and Syntax
  - Sentence Structure
  - **Phrase Structure**
  - Computational Resources
- Machine Translation Issues
- Dialects

70

## Phrase Structure

- Noun Phrase

- Determiner Noun Adjective PostModifier

- هذا الكاتب الطموح القادم من اليابان  
this the-writer the-ambitious the-arriving from Japan  
*This ambitious writer from Japan*

- Noun-Adjective agreement

- number, gender, definiteness
      - الكاتبة الطموحة the-writer<sub>fem</sub> the-ambitious<sub>fem</sub>
      - الكاتبات الطموحات the-writer<sub>femPlur</sub> the-ambitious<sub>femPlur</sub>

71

## Phrase Structure

- Noun Phrase

- Idafa construction (إضافة)

- **Noun1 of Noun2** encoded structurally
    - Noun1-indefinite Noun2-definite
    - ملك الاردن  
king Jordan  
*the king of Jordan / Jordan's king*

- Noun1 becomes definite

- Agrees with definite adjectives

- Idafa chains

- $N^1_{indef} N^2_{indef} \dots N^{n-1}_{indef} N^n_{def}$   
ابن عم جار رئيس مجلس ادارة الشركة  
son uncle neighbor chief committee management the-  
company  
*The cousin of the CEO's neighbor*

72



# Phrase Structure

- Morphological *definiteness* interacts with syntactic structure

		Word 1 كاتب <i>writer</i>	
		definite	Indefinite
Word 2 فنان <i>artist</i>	definite	<b>Noun Phrase</b> الكاتب الفنان The <i>artist(ic)</i> writer	<b>Noun Compound</b> كاتب الفنان The writer of the artist
	indefinite	<b>Copular Sentence</b> الكاتب فنان The writer is an artist	<b>Noun Phrase</b> كاتب فنان An <i>artist(ic)</i> writer

73

# Road Map

- Introduction
- Orthography
- Morphology
- **Syntax**
  - Morphology and Syntax
  - Sentence Structure
  - Phrase Structure
  - **Computational Resources**
- Machine Translation Issues
- Dialects

74

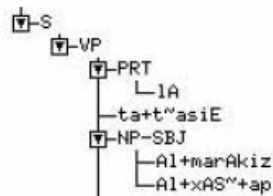
# Computational Resources

- Monolingual corpora for building language models
  - Arabic Gigaword
    - Agence France Presse
    - AlHayat News Agency
    - AnNahar News Agency
    - Xinhua News Agency
  - Arabic Newswire
  - United Nations Corpus (parallel with other UN languages)
  - Ummah Corpus (parallel with English)
- Distributors
  - Linguistic Data Consortium (LDC)
  - Evaluations and Language resources Distribution Agency (ELDA)

75

# Computational Resources

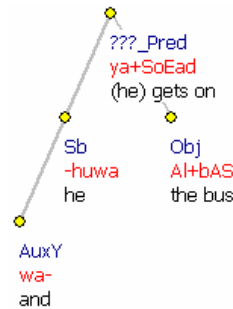
- Penn Arabic Treebank (PATB)
  - Started in 2001
  - Goal is 1 Million words
  - Currently 650K words
    - Agence France Presse , AlHayat newspaper, AnNahar newspaper
- POS tags
  - Buckwalter analyzer
  - Arabic-tailored POS list
- PATB constituency representation
  - Some modifications of Penn English Treebank
    - (e.g. Verb-phrase internal subjects)



76

## Computational Resources

- Prague Dependency Treebank
- Currently 100k words
- Partial overlap with PATB and Arabic Gigaword
  - Agence France Presse, AlHayat and Xinhua
- Morphological analysis
  - Similar to PATB
- Dependency representation



77

Graphic courtesy of Otakar Smrž: [http://ckl.mff.cuni.cz/padt/PADT\\_1.0/docs/slides/2003-eacl-trees.ppt](http://ckl.mff.cuni.cz/padt/PADT_1.0/docs/slides/2003-eacl-trees.ppt)

## Computational Resources

- Applications using Penn Arabic Treebank
  - Statistical parsing
    - Bikel's parser (Bikel 2003)
      - Same engine used with English, Chinese and Arabic
  - POS tagging and morphological disambiguation
    - (Diab et al, 2004) and (Habash and Rambow, 2005a)
- Arabic pos tagging (Khoja, 2001)
- Formalism conversion
  - Constituency to dependency (Žabokrtský and Smrž 2003)
  - Tree-adjoining grammar extraction (Habash and Rambow 2004)
- Automatic diacritization

78

## Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- **Machine Translation Issues**
  - Morphology and Translation
  - Translation Divergences
  - Computational Resources
- Dialects

79

## Morphology and Translation

*which level to go down to?*

- Natural token                      وللمكتبات
- Word                                      وللمكتبات
- Segmented Word                      و ل المكتبات
- Prefix + Stem + Suffix              ولل+مكتب+ات
- Lexeme + Features                  مكتبة [+Plural +Def +J +و]
- Root + Pattern + Features

و ل المكتبات [+Plural +Def +J +و] + م3ا21ا ة + ك ت ب

80

## Morphology and Translation

### *What approach?*

- Natural token Not Appropriate
- Word Statistical MT
- Segmented Word Statistical MT
- Prefix + Stem + Suffix Statistical/Symbolic
- Lexeme + Features Symbolic MT
- Root + Pattern + Features Too Abstract?

81

## Morphology and Translation

### *What resources?*

- Available resources may span different levels of representation!
- Most dictionaries are lexeme-based
- Buckwalter stem dictionary contains English glosses
- Statistical translation lexicons depend on the type of tokenization used before alignment
  - Word (no disambiguation necessary)
  - Segmented word (minimal disambiguation necessary)
  - Stem/Lexeme (machine/human disambiguation necessary)
- *Consistency is important*  
(Lee, 2004), (Habash, 2006), (Habash and Sadat, 2006), (Habash et al. 2006)

82

## Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- **Machine Translation Issues**
  - Morphology and Translation
  - **Translation Divergences**
  - Computational Resources
- Dialects

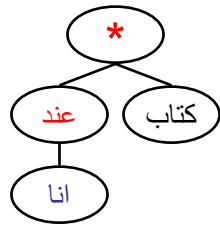
83

## Translation Divergences

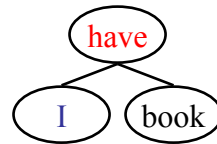
- Beyond word-order variation
  - Arabic VSO - English SVO
  - Arabic N Adj - English Adj N
- Meaning of two translationally equivalent constituents is distributed differently in two languages
- Divergence dimensions
  - Categorical Variation (*develop* → *development*)
  - Conflation (*become frozen* → *freeze*)
  - Inflation (*freeze* → *become frozen*)
  - Structural (*enter the room* → *enter into the room*)
  - Head Swap (*swim across the river* → *cross the river swimming*)
  - Thematic (*John likes Mary* → *Mary pleases John*)

84

## Translation Divergences *conflation*



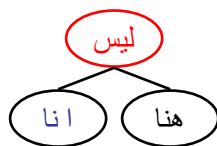
عندي كتاب  
at-me book



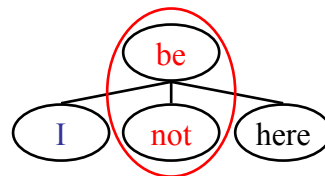
I have a book

85

## Translation Divergences *conflation*



ليست هنا  
I-am-not here



I am not here

86

## Translation Divergences *structural*

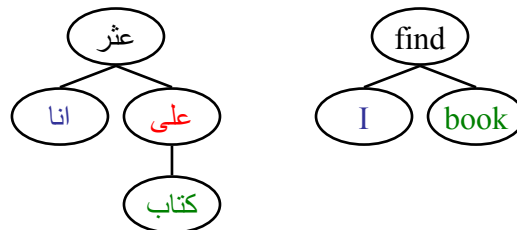


كتاب نزار  
book Nizar

Nizar's book  
Book of Nizar

87

## Translation Divergences *structural*



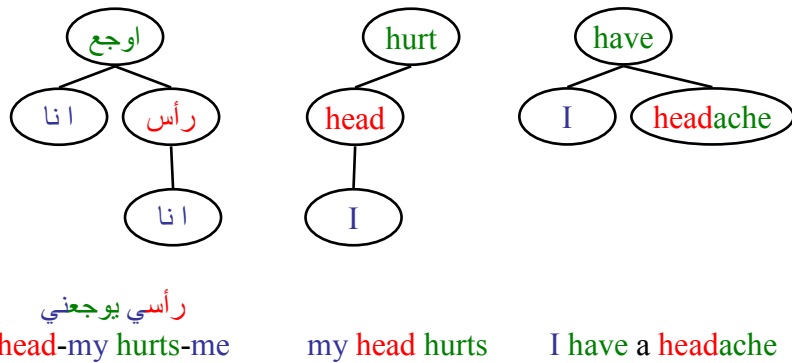
عثرت على الكتاب  
found-I upon the-book

I found the book

88

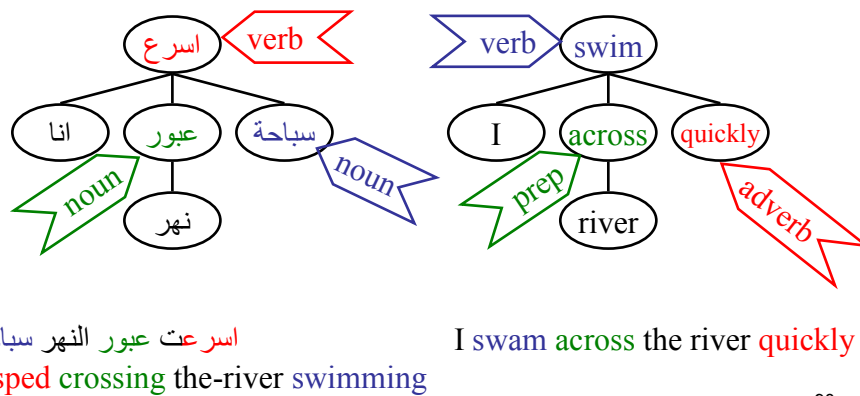


## Translation Divergences *thematic & conflational*



89

## Translation Divergences *head swap and categorial*



90

## Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- **Machine Translation Issues**
  - Morphology and Translation
  - Translation Divergences
  - **Computational Resources**
- Dialects

91

## Computational Resources

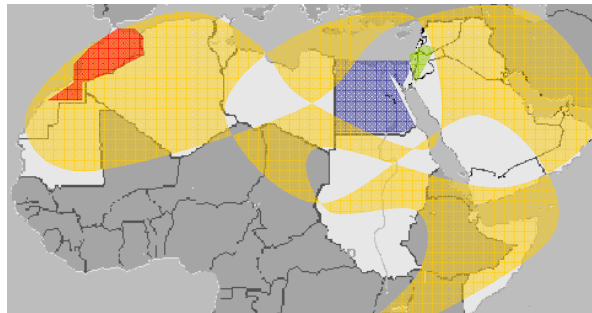
- Dictionaries
  - Buckwalter stem dictionary (LDC)
  - Salmone dictionary (Tufts university)
  - Online dictionaries – Ajeeb.com (Sakhr), Almisbar.com, Ectaco.com
- Parallel corpora (LDC)
  - United Nations Corpus (parallel with other UN languages)
  - Ummah Corpus (parallel with English)
  - Arabic News Translation Corpus
  - Arabic Treebank English Translation
  - *More on LDC webpage...*
- MT evaluation
  - Arabic-English Multi-translation Corpus (LDC)
  - NIST's MT-EVAL

92

# Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- Machine Translation Issues
- **Dialects**
  - General Definitions
  - Phonological & Lexical Variation
  - Morphological Variation
  - Syntactic Variation
  - Code Switching
  - Computational Resources

93



**lam jaʃtari nizār ʃawilatan ʒadīdatan** لم يشتري نزار طاولة جديدة

didn't buy Nizar table new

nizār maʃtarāʃ ʃarabēza gidīda ● نزار ماشراش طرييزة جديدة

nizār maʃtarāʃ ʃawile ʒdīde ● نزار ماشراش طاولة جديدة

nizar maʃrāʃ mida ʒdīda ● نزار ماشراش ميده جديدة

Nizar not-bought-not table new

94

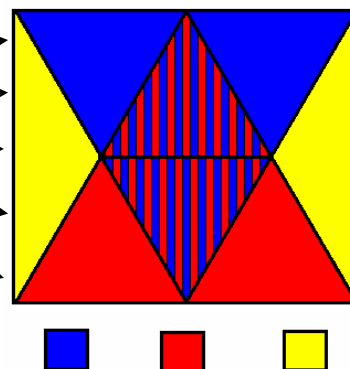
## General Definitions

- What is a 'dialect'?
  - Political and Religious factors
- Modern Standard Arabic
- Regional Dialects
  - Egyptian Arabic (EGY)
  - Levantine Arabic (LEV)
  - Gulf Arabic (GULF)
  - North African Arabic (NOR)
  - Iraqi, Yemenite, Sudanese, Maltese?
- Social dialects
  - City
  - Peasant
  - Bedouin

95

## General Definitions

- Diglossia
- Badawi's levels
  - Traditional Arabic
  - Modern Arabic
  - Educated Colloquial
  - Literate Colloquial
  - Illiterate Colloquial
- Polyglossia



Classical Dialect Foreign 96

# Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- Machine Translation Issues
- **Dialects**
  - General Definitions
  - **Phonological & Lexical Variation**
  - Morphological Variation
  - Syntactic Variation
  - Code Switching
  - Computational Resources

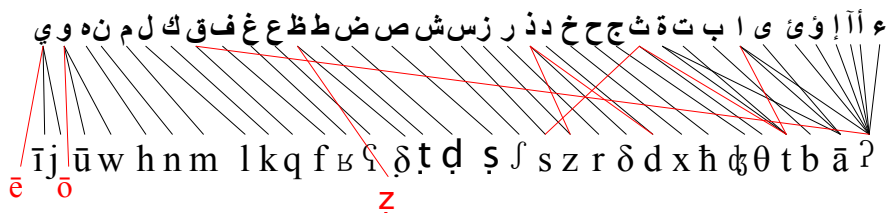
97

# Phonological Variation

## MSA



## LEV



• No dialect-specific standard orthography

98

## Lexical Variation

- Arabic Dialects vary widely lexically

English	table	cat	of	(I) want	there is	there isn't
MSA	Tāwila	qiTTa	<i>idafa</i>	'uridu	yūjadu	lā yujadu
Moroccan	mida	qeTTa	dyāl	bḡit	kāyn	mā kāynš
Egyptian	Tarabēza	'oTTa	bitā3	3āwez	fi	mafiš
Syrian	Tāwle	bisse	taba3	biddi	fi	mā fi
Iraqi	mēz	bazzūna	māl	'arid	aku	māku

- Arabic orthography allows consolidating some variations

99

## Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- Machine Translation Issues
- **Dialects**
  - General Definitions
  - Phonological & Lexical Variation
  - **Morphological Variation**
  - Syntactic Variation
  - Code Switching
  - Computational Resources

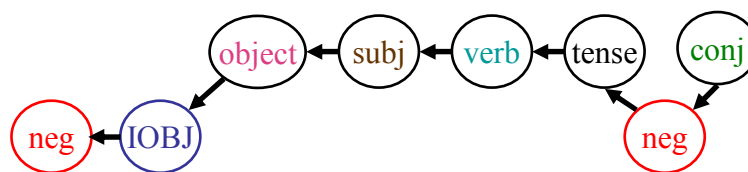
100

# Morphological Variation

- Nouns
  - No case marking
    - Word order implications
  - Paradigm reduction
    - Consolidating masculine & feminine plural
- Verbs
  - Paradigm reduction
    - Loss of dual forms
    - Consolidating masculine & feminine plural (2<sup>nd</sup>, 3<sup>rd</sup> person)
    - Loss of morphological moods
      - Subjunctive/jussive form dominates in some dialects
      - Indicative form dominates in others
  - Other aspects increase in complexity

101

## Morphological Variation Verb Morphology



MSA  
ولم تكتبوها له  
wa+lam taktubūhā lahu  
wa+lam taktubū+hā la+hu  
and+not\_past write\_you+it for+him

EGY  
وماكتبتهالوش  
wimakatabtuhalūš  
wi+ma+katab+tu+ha+lū+š  
and+not+wrote+you+it+for\_him+not

And you didn't write it for him

102

# Morphological Variation

## Verb conjugation

- Perfect verb derivation (*suffixes only*)

	1 <sup>st</sup> Person Singular	2 <sup>nd</sup> Person Singular ♂	2 <sup>nd</sup> Person Singular ♀
<b>MSA</b>	كُتِبْتُ katabtu	كُتِبْتَ katabta	كُتِبْتِ katabti
<b>LEV</b>	كُتِبْتَ katabt		كُتِبْتِي katabti

- Imperfect verb derivation (*prefix+suffix*)

	1 <sup>st</sup> Person Singular	2 <sup>nd</sup> Person Singular ♂	2 <sup>nd</sup> Person Singular ♀
<b>MSA</b>	أَكْتُبُ aktubu	تَكْتُبُ taktubu	تَكْتُبِينَ taktubīna تَكْتُبِي taktubī
<b>LEV</b>	أَكْتُبُ aktob	تَكْتُبُ toktob	تَكْتُبِي toktobi

# Morphological Variation

## Tense expression

	<b>Perfect</b>	<b>Imperfect</b>			
<b>M S A</b>	كُتِبَ kataba Past	يَكْتُبُ jaktubu Present			سَيَكْتُبُ sajaktubu Future
<b>L E V</b>	كُتِبَ katab Past	يَكْتُبُ jiktob 0-Tense	يَكْتُبُ bjoktob Present habitual	عَم يَكْتُبُ ʕam bjoktob Present progressive	حَيَكْتُبُ ħajiktob Future

104



# Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- Machine Translation Issues
- **Dialects**
  - General Definitions
  - Phonological & Lexical Variation
  - Morphological Variation
  - **Syntactic Variation**
  - Code Switching
  - Computational Resources

105

# Syntactic Variation

- Verbal sentences
  - The boys **wrote** the poems
  - MSA
    - **Verb** Subject Object (Partial agreement)  
كتب الاولاد الاشعار
    - **wrote<sub>masc</sub>** the-boys the-poems
    - Subject **Verb** Object (Full agreement)  
الاولاد كتبوا الاشعار
    - the-boys **wrote<sub>mascPlural</sub>** the-poems
  - LEV, EGY
    - Subject **Verb** Object  
الاولاد كتبوا الاشعار
    - The-boys **wrote<sub>mascPlural</sub>** the-poems
    - Less present: **Verb** Subject Object  
كتبوا الاولاد الاشعار
    - wrote<sub>mascPlural</sub>** the-boys the-poems
    - Full agreement in both order

	V-S <i>explicit subject</i>	V(S) <i>pro dropped subject</i>	S-V <i>explicit subject</i>
<b>MSA</b>	35%	30%	35%
<b>LEV</b>	10%	60%	30%

Verb-Subject distributions in the Levantine Arabic Treebank (Maamouri et al, 2006)

106

## Syntactic Variation

- Noun Phrase
  - Idafa construction
    - **Noun1 of Noun2** encoded structurally
    - ملك الاردن  
king Jordan  
*the king of Jordan / Jordan's king*
  - Dialects have an additional common construct
    - **Noun1 <particle> Noun2**
    - LEV: الملك تبع الاردن the-king *belonging-to* Jordan
    - <particle> differs widely among dialects
  - Pre/post-modifying demonstrative article
    - MSA: هذا الرجل this the-man *this man*
    - EGY: الرجل ده the-man this *this man*

107

## Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- Machine Translation Issues
- **Dialects**
  - General Definitions
  - Phonological & Lexical Variation
  - Morphological Variation
  - Syntactic Variation
  - **Code Switching**
  - Computational Resources

108

# Code Switching

MSA  
LEV

MSA and Dialect mixing in speech

- phonology, morphology and syntax

لا أنا ما يعتقد لأنه عملية اللي عم بيعارضوا اليوم تمديد للرئيس لحد هم اللي طالبوا بالتمديد للرئيس الهراوي وبالتالي موضوع منه موضوع مبدئي على الأرض أنا بحترم أنه يكون في نظرة ديمقراطية للأمور وأنه يكون في احترام للعبة الديمقراطية وأن يكون في ممارسة ديمقراطية ويعتقد إنه الكل في لبنان أو أكثرية ساحقة في لبنان تريد هذا الموضوع، بس بدي يرجع لحظة على موضوع إنجازات العهد يعني نعم نحكي عن إنجازات العهد لكن هل النظام في لبنان نظام رئاسي النظام في لبنان من بعد الطائف ليس نظام رئاسي وبالتالي السلطة هي عمليا بيد الحكومة مجتمعة والرئيس لحد أثبت خلال ممارسته الأخيرة بأنه لما بيكون في شخص مسؤول في منصب معين وأنا عشت هذا الموضوع شخصيا بممارستي في موضوع الاتصالات لما بياخذ مواقف صالحة ضمن خطاب ومبادئ خطاب القسم هو إلى جانبه إنما مش مطلوب من رئيس جمهورية هو يكون رئيس السلطة التنفيذية لأنه منه بقي في لبنان ما بعد اتفاق الطائف رئيس السلطة التنفيذية عليه التوجيه عليه إبداء الملاحظات عليه القول ما هو خطأ وما هو صح عليه تتمير جهود الوطنية الشاملة كي يظل في مصالحة وطنية كي يظل في توافق ما بين المسلم والمسيحي في لبنان يحتضن أبناء هذا البلد ما يترك المسار يروح باتجاه الخطأ نعم إنما خطاب القسم كان موضوع مبادئ طرحت هو ملتزم فيها اللي مشبوا معه وامنوا فيها التزموا فيها أنا أثبت خلال الأربع سنوات بالممارسة الحكومية أنني التزمت فيها ولما التزمنا بهذا الموضوع كان الرئيس لحد إلى جنبنا في هذا الموضوع، أما الموضوع الديمقراطي أنا بتفهم تماما هذا هالوجهة النظر بس ما ممكن نقول إنه الدستور أو تعديله هو أو إمكانية فتح إعادة انتخاب ديمقراطي ضمن المجلس والتصويت إلى ما هنالك لرئيس جمهورية بولاية ثانية هو مسح هيئة في جوهر الديمقراطية هذا بالأقل يعني قناعتي في هذا الموضوع.

109

Aljazeera Transcript [http://www.aljazeera.net/programs/op\\_direction/articles/2004/7/7-23-1.htm](http://www.aljazeera.net/programs/op_direction/articles/2004/7/7-23-1.htm)

# Road Map

- Introduction
- Orthography
- Morphology
- Syntax
- Machine Translation Issues
- **Dialects**
  - General Definitions
  - Phonological & Lexical Variation
  - Morphological Variation
  - Syntactic Variation
  - Code Switching
  - **Computational Resources**

110

# Computational Resources

- Most work on Arabic dialects focuses on Automatic Speech Recognition
- Speech/transcript corpora
  - Egyptian and Levantine Arabic (LDC)
  - Moroccan and Tunisian Arabic (ELDA)
  - Gulf Arabic (Appen)
  - Many other...
- Few lexicons/morphology resources
  - CallHome Egyptian Arabic monolingual lexicon (LDC)
  - CallHome Egyptian Verb transducer (LDC)
- Work on multi-dialectic resources
  - Linguistic Data Consortium
  - Columbia University Arabic Dialect Modeling (CADIM) Group
    - Pan-Arab lexicon and Pan-Arab Morphology
- Novel Approaches to Arabic Speech Recognition (JHU summer workshop 2002) (Kirchhoff et al, 2002)
- Parsing Arabic Dialects (JHU summer workshop 2005)  
(Rambow et al, 2005) , (Chiang et al., 2006)

# Resources

## Distributors

- [Linguistic Data Consortium](#)
- [NEMLAR \(Network for Euro-Mediterranean Language Resources\)](#)
- [ELSNET is the European Network of Excellence in Human Language Technologies](#)
- [ELDA Evaluation and Language resources Distribution Agency](#)

## Resources

### Reports

- Mohamed Maamouri and Christopher Cieri. 2002. [Resources for Natural Language Processing at the Linguistic Data Consortium](#). In Proceedings of the International Symposium on Processing of Arabic, pages 125--146, Manouba, Tunisia, April 2002.
- Mahtab Nikkhou and Khalid Choukri. [Survey on Arabic Language Resources and Tools in the Mediterranean Countries](#).
- [Arabic Information Retrieval and Computational Linguistics Resources](#) (thanks to Doug Oard)

113

## Resources

### Monolingual Corpora

- [Arabic Gigaword](#)
- [Arabic Newswire](#)

### Parallel Corpora

- [United Nations Parallel Corpus](#)
- [Ummah Parallel Corpus](#)
- [Arabic News Translation](#)
- [Multiple-Translation Arabic](#)

### Treebanks

- [Arabic Penn Treebank Webpage](#)
  - [Part 1 v 2.0](#), [Part 2 v 2.0](#), [Part 3 v 1.0](#), [10K-word English Translation](#)
- [Prague Arabic Dependency Treebank](#)

114

## Resources

### Morphology

- Buckwalter Arabic Morphological Analyzer
  - [Version 1.0](#), [Version 2.0](#)
- [Xerox Arabic Morphology](#) (online)

### Dialect Resources

- [CALLHOME Egyptian Arabic Transcripts](#)
- [CALLHOME Egyptian Arabic Speech](#)
- [Egyptian Colloquial Arabic Lexicon](#)
- [Levantine Arabic Resources](#)
- <http://www.orientel.org/>
- <http://www.appen.com.au/>
- CADIM: <http://www.ccls.columbia.edu/cadim>

115

## Resources

### Dictionaries

- [Buckwalter Stem Dictionary](#)
- H. Anthony Salmone. An Advanced Learner's Arabic-English Dictionary encoded by the Perseus Project, Tufts University (contact: David Smith [dasmith@perseus.tufts.edu](mailto:dasmith@perseus.tufts.edu))
- [Ajeeb Arabic-English Dictionary](#) (online)
- [Al-Misbar Dictionary](#) (online)
- [Ectaco Bilingual Dictionary](#) (online)

### Online MT systems

- [Ajeeb's Arabic-English Machine Translation](#) (online)
- [Al-Misbar English-Arabic Machine Translation](#) (online)

116

## Conferences and Workshops

*with some focus on Arabic*

- Parsing Arabic Dialects (JHU summer workshop 2005)
- ACL 2005 Workshop on Computational Approaches to Semitic Languages
- [Arabic Language Resources and Tools Conference 2004 Cairo, Egypt](#)
- [WORKSHOP Computational Approaches to Arabic Script-based Languages \(COLING 2004\)](#)
- [Traitement Automatique du Langage Naturel \(TALN ' 04\)](#)
- NIST MT EVAL (<http://www.nist.gov/speech/tests/mt/>)
- [MT Summit IX Workshop on Machine Translation for Semitic Languages in 2003](#)
- Novel Approaches to Arabic Speech Recognition (JHU summer workshop 2002)
- [LREC 2002 Arabic Language Resources and Evaluation Workshop](#)
- [ACL 2002 Workshop on Computational Approaches to Semitic Languages](#)
- International Symposium on Processing of Arabic 2002, Tunisia
- [Workshop on ARABIC Language Processing: Status and Prospects \(ACL/EACL 2001\)](#)
- [Arabic Translation and Localisation Symposium \(ATLAS 1999\)](#)
- [Computational Approaches to Semitic Languages \(COLING/ACL 1998\)](#)

117

## References

- Aljlal, M. and O. Frieder. 2002. [On Arabic Search: Improving the Retrieval Effectiveness via a Light Stemming Approach](#). ACM Conference on Information and Knowledge Management.
- Al-Sughaiyer, I. and I. Al-Kharashi. 2004. [Arabic Morphological Analysis Techniques: A Comprehensive Survey](#). Journal of the American Society for Information Science and Technology. Volume 55 , Issue 3.
- Beesley, K. 2001. [Finite-State Morphological Analysis and Generation of Arabic at Xerox Research: Status and Plans in 2001](#). EACL workshop on Arabic Language Processing: Status and Prospects.
- Bikel, D. 2002. [Design of a Multi-lingual, Parallel-processing Statistical Parsing Engine](#). HLT.
- Buckwalter, T. 2002. [Buckwalter Arabic Morphological Analyzer Version 1.0](#). LDC catalog number LDC2002L49.
- Cavalli-Sforza, V., A. Soudi, and T. Mitamura. 2000. [Arabic Morphology Generation Using a Concatenative Strategy](#). ANLP.
- Chiang, D., M. Diab, N. Habash, O. Rambow, and S. Shareef. 2006. [Arabic Dialect Parsing](#). EACL.
- Darwish, K. 2002. [Building a Shallow Morphological Analyzer in One Day](#). ACL workshop on Computational Approaches to Semitic Languages.
- Diab, M., K. Hacioglu and D. Jurafsky. 2004. [Automatic Tagging of Arabic Text: From raw text to Base Phrase Chunks](#). HLT-NAACL.
- Fischer, W. 2001. [A Grammar of Classical Arabic](#). Yale Language Series. Yale University Press. Translated by Jonathan Rodgers.
- Habash, N. and O. Rambow. 2004. [Extracting a Tree Adjoining Grammar from the Penn Arabic Treebank](#). TALN.
- Habash, N. and O. Rambow. 2005a. [Arabic Tokenization, Part-of-Speech Tagging in and Morphological Disambiguation One Fell Swoop](#). ACL.

118

## References

- Habash, N., O. Rambow and G. Kiraz. 2005b. Morphological Analysis and Generation for Arabic Dialects. *ACL workshop on Computational Approaches to Semitic Languages*.
- Habash, N. 2004. Large Scale Lexeme Based Arabic Morphological Generation. TALN.
- Habash, N. and F. Sadat. 2006. Arabic Preprocessing Schemes for Statistical Machine Translation. NAACL.
- Habash, N. 2006. "Arabic Morphological Representations for Machine Translation." Book Chapter. In Arabic Computational Morphology: Knowledge-based and Empirical Methods. Editors A. van den Bosch and A. Soudi.
- Habash, N., C. Mah, S. Imran, R. Calistri-Yeh, and P. Sheridan. 2006. The Design and Validation of an Arabic WordNet for Information Retrieval. LREC.
- Khoja, S. 2001. APT: Arabic Part-of-Speech Tagger. NAACL Student Research Workshop.
- Kiraz, G. 2001. Computational Nonlinear Morphology with Emphasis on Semitic Languages. Studies in Natural Language Processing. Cambridge University Press.
- Kirchhoff, K., J. Bilmes, S. Das, N. Duta, M. Egan, G. Ji, F. He, J. Henderson, D. Liu, M. Noamany, P. Schone, R. Schwartz and D. Vergyri. 2003. Novel Approaches to Arabic Speech Recognition: Report from the 2002 Johns-Hopkins Summer Workshop. IEEE Int. Conf. on Acoustics, Speech, and Signal Processing.
- Lee, Y., K. Papineni, S. Roukos, O. Emam and H. Hassan. 2003. Language Model Based Arabic Word Segmentation. ACL.
- Lee, Y. 2004. Morphological Analysis for Statistical Machine Translation. NAACL.
- Maamouri, M., A. Bies, T. Buckwalter, M. Diab, N. Habash, O. Rambow, D. Tabessi. 2006. Developing and Using a Pilot Dialectal Arabic Treebank. LREC.

119

## References

- Rambow, O., D. Chiang, M. Diab, N. Habash, R. Hwa, K. Sima'an, V. Lacey, R. Levy, C. Nichols, and S. Shareef. 2005. Parsing Arabic Dialects. Final Report, JHU Summer Workshop.
- Rogati, M., S. McCarley, and Y. Yang. 2003. Unsupervised Learning of Arabic Stemming Using a Parallel Corpus. ACL.
- Smrř, O. and P. Zemánek. 2002. Sherds from an Arabic Treebanking Mosaic. Prague Bulletin of Mathematical Linguistics, (78).
- Smith N., D. Smith, and R. Tromble. 2005. Context-Based Morphological Disambiguation with Random Fields. HLT-EMNLP.
- Soudi, A., V. Cavalli-Sforza, and A. Jamari. 2001. A Computational Lexeme-Based Treatment of Arabic Morphology. ACL workshop on Arabic Natural Language Processing.
- Xu J. 2002. UN Parallel Text (Arabic-English), LDC Catalog No.: LDC2002E15.
- Žabokrtský, Z. and O. Smrř. 2003. Arabic Syntactic Trees: from Constituency to Dependency. EACL.
- Zitouni, I., J. Olive, D. Iskra, K. Choukri, O. Emam, O. Gedge, M. Maragoudakis, H. Troupf, A. Moreno, A. Rodriguez, B. Heuft and R. Siemund. 2002. OrienTel: Speech-Based Interactive Communication Applications for the Mediterranean and the Middle East. ICSLP.

### Conference/Institution Name Abbreviations

**ANLP** = Applied Natural Language Processing

**ACL** = Association for Computational Linguistics

**ACM** = Association for Computing Machinery

**EMNLP** = Empirical Methods to Natural Language Processing

**EACL** = European ACL

**HLT** = Human Language Technology Conference

**ICSLP** = International Conference on Spoken Language Processing

**JHU** = Johns Hopkins University

**LREC** = Language Resources and Evaluation Conference

**LDC** = Linguistic Data Consortium, University of Pennsylvania

**NAACL** = North American ACL

**TALN** = Traitement Automatique du Langage Naturel

120