

**MEADR 2009**

**Cairo, Egypt**  
**April 21, 2009**

# Introduction to Arabic Natural Language Processing

**Nizar Habash**

Columbia University

Center for Computational Learning Systems

*habash@ccls.columbia.edu*



- 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

# Road Map

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

# Arabic Script

# 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.

# Arabic Script

## Alphabet

- letter forms

ع ط ص س ر د ح ب ا  
ء ى و ه ب م ل ف

- 
- letter marks

- Arabic only

. .. ^ n ~  
— — — — — — — —  
— .. — — — — — —

- Other languages

- Persian, Kurdish,  
Urdu, Pashto, etc.

.. : ; v b / i  
— — — — — — — —  
; .. : — o — — — —

- *OCR output ambiguity; common spelling errors*

# Arabic Script

## Alphabet (MSA)

- letters (form+mark)

- Distinctive

ب ت ث س ش

/ʃ/ /s/ /θ/ /t/ /b/

- 
- Non-distinctive

ا ا ا ا ا و و و

/ʔ/

*glottal stop aka hamza*

# Arabic Script

## Letter Shapes

- No distinction between print and handwriting
- No capitalization
- Right-to-left
- Ambiguous shapes
- Connective letters
- Disconnective letters
- (ة ر ز د ذ و ؤ ا )  
cause word-internal visual spacing

Stand alone	initial	medial	final
ب	ب	ب	ب
ن	ر	ن	ن
ذ	د	ذ	ذ

# Arabic Script

## Letter shaping

ك ت ب = كتب ←

/katab/      b    t    k

*to write*

ك ت ا ب = كتاب

/kitāb/      b    ā    t    k

*book*

# Arabic Script

## Diacritics

- Zero-width characters
- Used for short vowels

كتب /katab/ *to write*

- Nunation is used for nominal indefinite marker in MSA

كتاب /kitābun/ *a book*

Nunation
بَنٌونٍ /ban/
بَنُونٍ /bun/
بَنِينٍ /bin/

Vowel
بَأْ /ba/
بَعْ /bu/
بَيْ /bi/

# Arabic Script

## Diacritics

- No-vowel marker (*sukun*)

مَكْتَب /maktab/ *office*

- Double consonant marker (*shadda*)

كَتَب /kattab/ *to dictate*

- Combinable

بُ

بُّ

بُّّ

/bbu/

/bbin/

/bban/

### No Vowel

بْ

/b/

### Double Consonant

بُّ

/bb/

# Arabic Script

## Putting it together

### *Simple combination*

Arab /ʕarab/   عَرَبْ ← عَرَب = عَرَب

West /karb/   غَرْبْ ← غَرَب = غَرَب

### *Ligatures*

Peace /salām/ سلام ← سلام X

# Arabic Script

## Tatweel

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

حقوق الانسان

حقـوقـاـنـسـانـاـنـ

حقـوـقـاـنـسـانـاـنـ

حقـوـقـاـنـسـانـاـنـ

human rights /ħuqūq alʔinsān/

# Arabic Script

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

Arabic	Muhammad	algebra
عربی	محمد	الجبر

/arabi/   /muhammad /   /alqabir/

# 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>	٠	١	٢	٣	٤	٥	٦	٧	٨	٩
<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

# 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



# 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
    - كتب /katab/ to write    كُتُب /kutib/ to be written
    - حُب /ħubb/ love    حَب /ħabb/ seed
- Dual use of ا, و, ي as consonant and long vowel
  - ا (/'/, /ā/) و (/w/, /ū/) ي (/j/, /ī/)

# MSA Phonology and Spelling

- Arabic spelling is mostly phonemic ...

***Except for (continued)***

- Morphophonemic characters
  - Feminine marker ة (ta marbuta)
    - /kabīr/ (big ♂) كِبِيرٌ /kabīra/ (big ♀) كِبِيرَةٌ
  - Derivation marker
    - /ʕaṣa/ (to disobey) عصى (a stick عصا)
- Hamza variants (6 characters for one phoneme!)
  - بَهَاءٌ بَهَاءٌ بَهَاءٌ (ءُ آءُ ؤُ ئُ ) /baha'/ + 3MascSing (his glory)

# 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*

# Proper Name Spelling

- The Qaddafi-Schwarzenegger problem
  - Foreign Proper name spelling is often ad hoc
  - Multiplicity of spellings causes increased sparsity

قذافي	→	Gadafi Gaddafi Gaddfi Gadhafi Ghaddafi Kadaffy Qaddafi Qadhafi ...
شوارزنیغر		
شوارزنغر		
شوارزنیجر	←	Schwarzenegger
شوارترنجر		

# Road Map

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

# Arabic Script

# Other languages

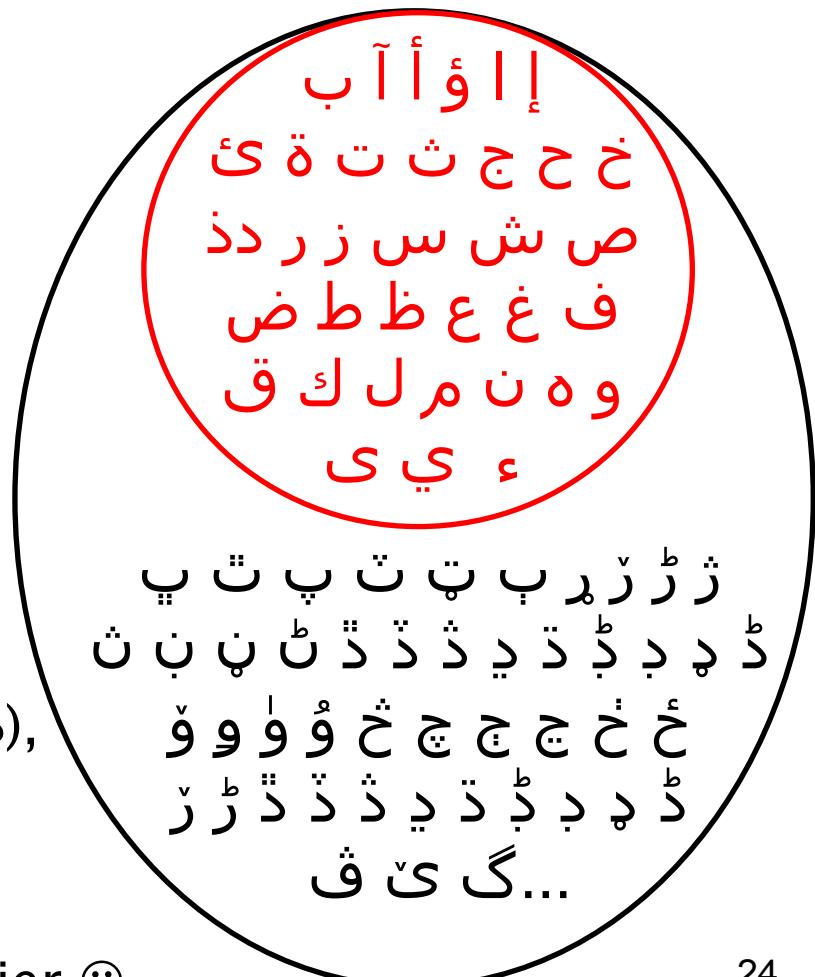
# Arabic

- 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/, ڦ /g/, ڙ /tʃ/
    - regionally variable

# Not Arabic

- Extra marks: haft (v), ring (o), taa (b), four dots (::), vertical dots (:)
  - Some Numerals (೨, ೫, ೬)

Once you learn the alphabet, it is easier 😊



بُونه سووتی جَگه رو بُوچی نه بی دل به که باب

(۱) بُوچی نه روا له ته نم روحی رهوان میسلی شه هاب

(۲) بوله سه ر چاوه یی چاو هه لنه قولی ره شجه یی خوین

بُوچ له فه وواره یی موژگان نه تکی قه تره یی ناب

بوله به ر ناله نه بی حه لقهی حه لقم به سروود

بوله به ر گریه نه بی چه شمهی چه شمم به سه راب

(۴) موونسی روزو شه ووم باعیسی نارامی دلم

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

(۵) به وقووعی سه فه ری قادری نوستاد خدری

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

(۶) چه نک ونهی لی مه ده موتریب که له به ر فیرقه تی ئه و

(۷) رنه کی روحه له گوئیم نه غمهی ئاوازو رو باب

(۸) ساغیری مهی مه ده ساقی که له به ر دووری نه و

(۹) تاله وه ک زه هری هه لایل له مه زاقم مهی ناب

## □ Not Arabic

## □ Arabic

سجل... انا عربي...  
ورقم بطاقي خمسون الف  
واطفالي ثمانية  
وتاسعهم سياتي بعد صيف  
فهل تغضب  
سجل... انا عربي...  
واعمل مع رفاق الكدح في مجر  
واطفالي ثمانية  
اسل لهم رغيف الخبز والاثواب والدفتر  
من الصخر  
ولا اتوسل الصدقات من بابك  
ولا اصغر امام بلاط اعتابك  
فهل تغضب

# شلی بیٹی کے نام

Not Arabic

تجھے جب بھی کوئی دکھدے

اس دکھکا کا نام بیٹی رکھنا

جب میرے سفیدے بال

تیرے گالوں پر آن ہنسیں، رو لینا

میرے خواب کے دکھ پہ سولینا

جن کھیتوں کو ابھی اگنا ہے

ان کھیتوں میں<sup>7</sup>

# Road Map

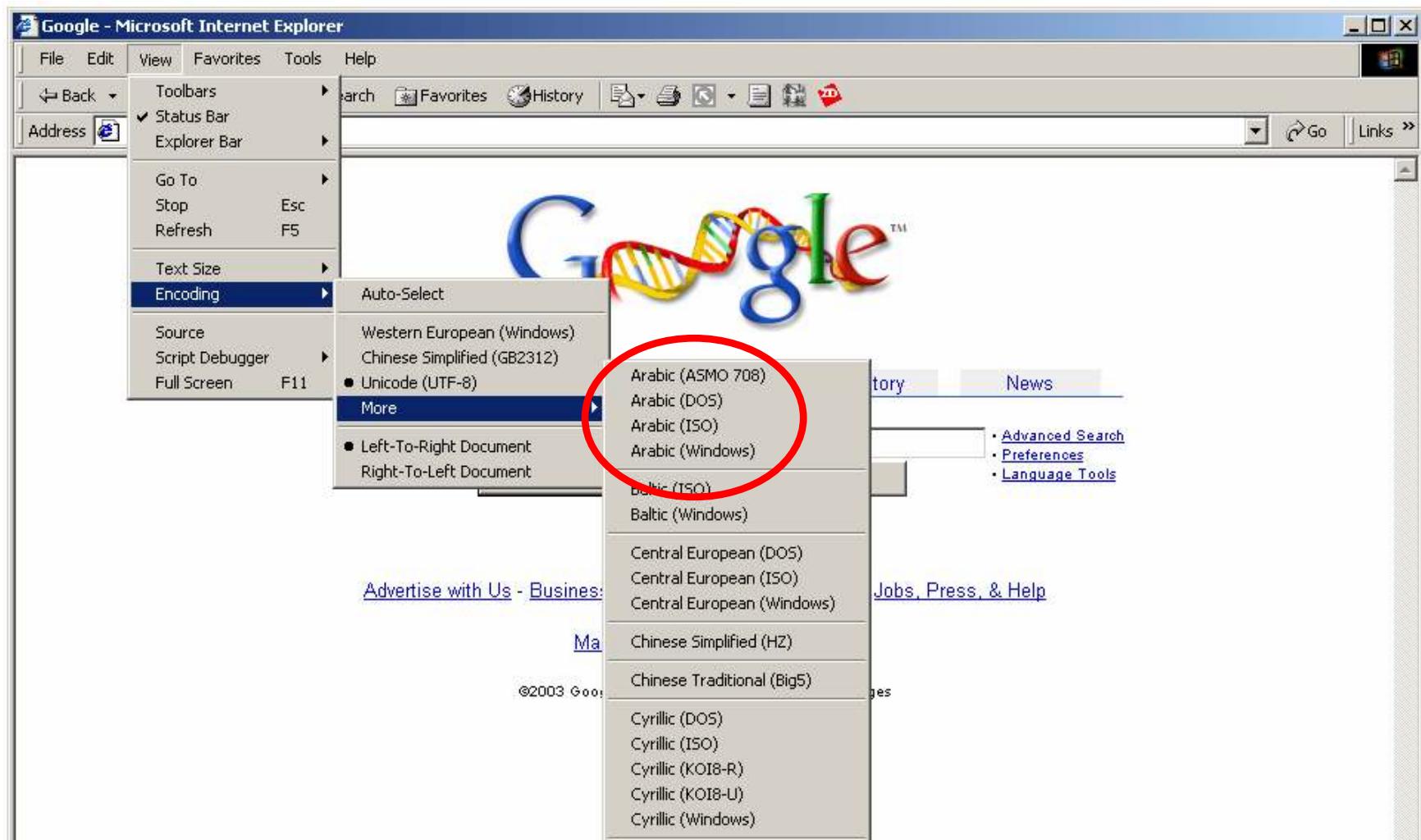
- Introduction
- Orthography
  - Arabic Script
  - MSA Phonology and Spelling
  - Recognizing Arabic vs. Persian/Urdu/Pashto/Kurdish/Sindhi/...
  - Encoding Issues
- Morphology
- Syntax

# 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

# 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-	0020	!	"	#	\$	%	&	'	(	)	*	+	,	-	.	/
3-	0030	0031	0032	0033	0034	0035	0036	0037	0038	0039	003A	003B	003C	003D	003E	003F
4-	0040	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N
5-	0050	P	Q	R	S	T	U	V	W	X	Y	Z	[	\	]	^
6-	0060	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n
7-	0070	p	q	r	s	t	u	v	w	x	y	z	{		}	~
8-	20AC	€	،	f	،،	...	†	‡	^	%eo	„	œ	Œ	߻	߼	߾
9-	06AF	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ
A-	00A0	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ
B-	00B0	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ
C-	0621	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ
D-	0630	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ
E-	00E0	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ
F-	064B	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ	ݚ

# Encodings

- Unicode
  - Becoming the standard more and more
  - 2-byte characters
  - Widely supported input/display
  - Supports extended Arabic characters
  - Multi-script representation

	060	061	062	063	064	065	066	067	068	069	06A	06B	06C	06D	06E	06F
0			ذ	-	ڦ	·	ڻ	ڦ	ڦ	ڦ	ڳ	ڳ	ڙ	ڙ	ڙ	·
1			ء	ڦ	ڦ	ا	ا	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ا
2			آ	ڦ	ڦ	ُ	ُ	ا	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ُ
3			أ	ڦ	ڦ	س	س	إ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ُ
4				ڦ	ڦ	ل	ش	ؤ	ڦ	ڦ	ڦ	ڦ	ڦ	و	-	ـ
5					إ	م	ص	!	ڦ	ڦ	ڦ	ڦ	ڦ	و	ه	ـ
6						ن	ض	ئ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ـ
7						ه	ط	ا	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ـ
8							و	ظ	ب	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ـ
9							ي	ع	ة	ڦ	ڦ	ڦ	ڦ	ڦ	ڦ	ـ
A								غ	ت	٪	٪	٪	٪	٪	٪	ـ
B		:						ي	غ	ت	ـ	ـ	ـ	ـ	ـ	ـ
C	,								ث	ـ	ـ	ـ	ـ	ـ	ـ	ـ
D									ـ	ـ	ـ	ـ	ـ	ـ	ـ	ـ
E									ـ	ـ	ـ	ـ	ـ	ـ	ـ	ـ
F										ـ	ـ	ـ	ـ	ـ	ـ	ـ

# Encodings

- Unicode
  - Supports presentation forms (shapes and ligatures)

FE70

Arabic Presentation Forms-B

FEFF

	FE7	FE8	FE9	FEA	FEB	FEC	FED	FEE	FEF
0	؀	؁	؂	؃	؄	؅	؆	؈	؉
1	؀	؁	؂	؃	؄	؅	؆	؈	؉
2	؀	؁	؂	؃	؄	؅	؆	؈	؉
3	؀	؁	؂	؃	؄	؅	؆	؈	؉
4	؀	؁	؂	؃	؄	؅	؆	؈	؉

FC40

Arabic Presentation Forms-A

FD1F

	FC4	FC5	FC6	FC7	FC8	FC9	FCA	FCB	FCC	FCD	FCE	FCF	FD0	FD1
0	،	ؔ	ؕ	ؖ	ؗ	ؘ	ؙ	ؚ	؛	؜	؝	؞	؟	ؠ
1	،	ؔ	ؕ	ؖ	ؗ	ؘ	ؙ	ؚ	؛	؜	؝	؞	؟	ؠ
2	،	ؔ	ؕ	ؖ	ؗ	ؘ	ؙ	ؚ	؛	؜	؝	؞	؟	ؠ
3	،	ؔ	ؕ	ؖ	ؗ	ؘ	ؙ	ؚ	؛	؜	؝	؞	؟	ؠ

# Encoding Issues

## Arabic Display

- Memory (logical order) →

ÔÇÑËÊ ÝáÓØíä (Palestine) Ýí ÇæáâÈíÇÏ (Olympics) 2000 æ 2004.  
نیطسلف تکر اش (Palestine) دایبم لو ایف (Olympics) 2000 و 2004.

*or this way for those with direction-bias*



.4002 æ 0002 ) scipmylO( ïÇíÈääæç íÝ ) enitselaP( äíØÓáÝ ÈßÑÇÔ  
شاركت فلسطين ( في اولمبیاد ) enitselaP( 0002 و .4002

# Encoding Issues

## Arabic Display

- **Memory (logical order)**

ÔÇÑËÊ ÝáÓØíä (Palestine) Ýí ÇæáãÈíÇÏ (Olympics) 2000 æ 2004.  
شراكت فلسطين (Palestine) في أولمبياد (Olympics) 2000 و 2004.

- **Display (visual order)**

- Bidirectional (BiDi) support

- Numbers and Roman script

.2004 و 2000 (Olympics) في اولمبياد (Palestine) شراكت فلسطين

- Letter and ligature shaping

.2004 و 2000 (Olympics) في اولمبياد (Palestine) شراكت فلسطين

# Display Problems

Display Encoding				
	CP-1256	ISO-8859	Unicode	Western
Actual Encoding	تدشين منطقة حرة في دبي للتجارة الالكترونية	ة حرة تدشيل كلظ ترنلة دب ففتجارة افاف	Yí Ígá gá ΨÓgáá Áí	ÊÍÔíä ääøþé íñé Ýí ïèí ááêíçñé Çáçáßêñæäíé
CP-1256	ة حرة xâ و هو تدش ن التجارية ê دب ل ة و è انامتر	تدشين منطقة حرة في دبي للتجارة الالكترونية	Y 染既 gí gá ψí lÖgGG 株親g	ÊÍÔéæ åæxâé íñé áé ïèé ääêíçñé Çäçääßêñææé
ISO-8859	« طهط طظظ + ئ ظ...ظ ط ط+ ط ط ظظظ ط ط- ط ط ظ، ط طظظ ط ط ظ، ط طظظ، ط ط ظظظ ط طظظ	ع ع ظ ظ ظ ظ ظ ظ ظ ع ع ظ ظ ظ ظ ظ ظ	تدشين منطقة حرة في دبي للتجارة الالكترونية	» خ Ø^ Ø Ø Ø Ø Ø Ù...Ù Ø Ù Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ù Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø Ø
Unicode	؟ طهط طظظ + ئ ظ...ظ ط ط+ ط ط ظظظ ط ط- ط ط ظ، ط طظظ ط ط ظ، ط طظظ، ط ط ظظظ ط طظظ	؟ ظ ظ ظ ظ ظ ظ ظ ظ ظ ظ ظ ظ ظ ظ		

- Wrong encoding
- Partial support problems

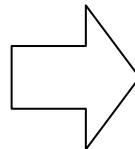
# Encoding Issues

## Arabic Input

- Standard graphemic keyboard
  - Logical order input



م ا ل س



سلام

# 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	ِ	ذ	*	ڙ	ِ
ِ	ِ	ر	r	ِ	m
ِ	>	ز	z	ِ	n
ِ	&	س	s	ِ	h
ِ	<	ش	\$	ِ	w
ِ	}	ص	s	ِ	y
ِ	A	ض	d	ِ	y
ِ	b	ٻ	t	ِ	F
ِ	p	ڦ	z	ِ	N
ِ	t	ڻ	E	ِ	K
ِ	v	ڻ	g	ِ	a
ِ	j	ـ	ـ	ِ	u
ِ	H	ڻ	f	ِ	i
ِ	x	ڻ	q	ِ	~
ِ	d	ڻ	k	ِ	o

# Road Map

- Introduction
- Orthography
- Morphology
- Syntax

# Morphology

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

# Road Map

- Introduction
- Orthography
- Morphology
  - Derivational Morphology
  - Inflectional Morphology
  - Morphological Ambiguity
  - Arabic Computational Morphology
- Syntax

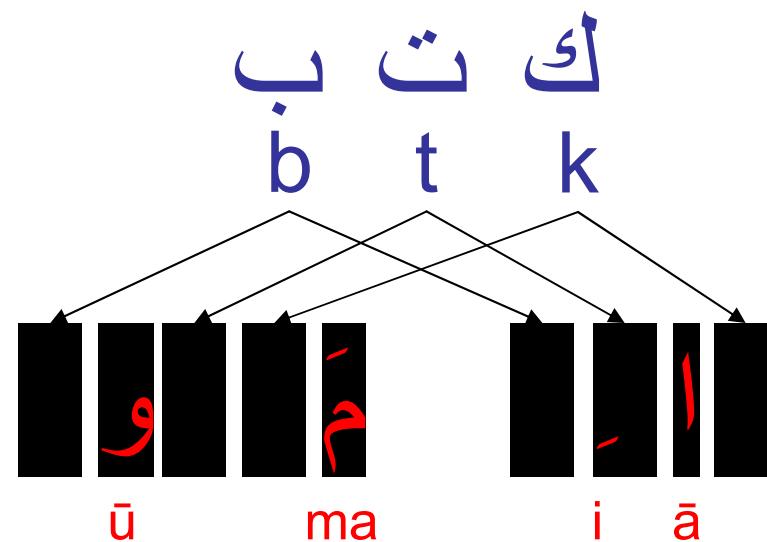
# Derivational Morphology

- Templetic Morphology

- Root

- Pattern

- Lexeme



مکتب

# maktūb

*written*

کانٹب

## kātib

## *writer*

## *Lexeme.Meaning* =

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

# Derivational Morphology

## *Root Meaning*

- كتب KTB = notion of “writing”



# Derivational Morphology

## *Root Meaning*

- LHM-1
- Notion of “meat”
  - لحم /laħm/
    - Meat
  - لحام /laħħām/
    - Butcher

لحم  
laHm



# Derivational Morphology

## *Root Meaning*

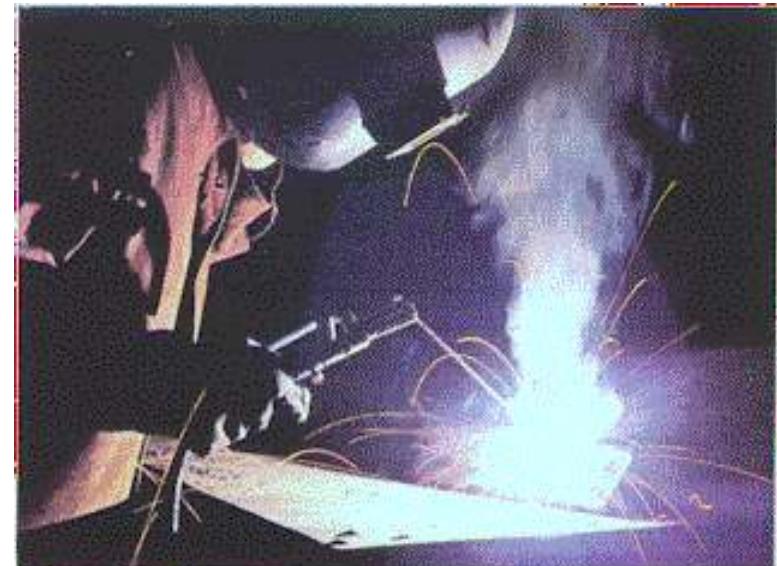
- LHM-2
- Notion of “battle”
  - ملحمة /malħama/
    - Fierce battle
    - Massacre
    - Epic



# Derivational Morphology

## *Root Meaning*

- LHM-3
- Notion of “soldering”
  - لحم /laħam/
    - Weld, solder, stick, cling
  - التحام /iltaham/
    - Be welded/soldered/fused
  - ملتحم /multaħim/
    - Welded, soldered, fused



# Derivational Morphology

## *Pattern Meaning*

- Verb Pattern Meaning is hard to define

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

# Road Map

- Introduction
- Orthography
- Morphology
  - Derivational Morphology
  - Inflectional Morphology
  - Morphological Ambiguity
  - Arabic Computational Morphology
- Syntax

# Inflectional Morphology

- Derivational Morphology
  - Lexeme ≈ Root + Pattern
- Inflectional Morphology
  - Word = Lexeme + Features
- Part-of-speech
  - *Traditional*: Noun, Verb, Particle
  - *Computational*: N, PN, V, Adj, Adv, P, Pron, Num, Conj, Det, Aux, Pun, IJ, and others
  - Many tag sets exist ranging from 3 to over 22K tags
- Noun-specific Features
- Verb-specific Features
- Other Features

# Inflectional Morphology

- **Noun-specific Features**

- Number: singular, dual, plural, collective
- Gender: masculine, feminine
- Definiteness: definite, indefinite
- Case: nominative, accusative, genitive
- Possessive clitic

- **Verb-specific Features**

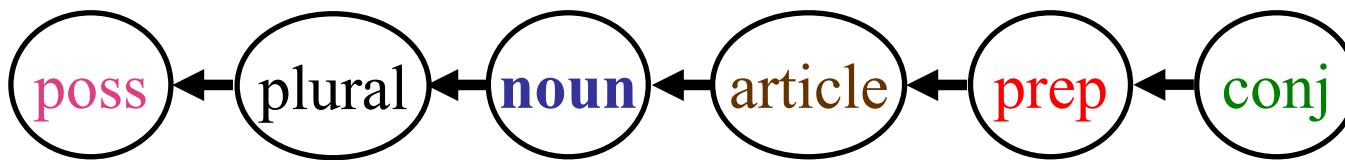
- Aspect: perfective, imperfective, imperative
- Voice: active, passive
- Tense: past, present, future
- Mood: indicative, subjunctive, jussive
- Subject (Person, Number, Gender)
- Object clitic

- **Other Features**

- Single-letter conjunctions
- Single-letter prepositions

# Inflectional Morphology

## Nouns

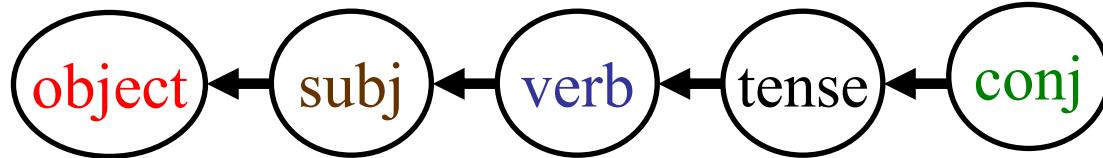


و كبيوتنا	وللمكتبات
/wakabiyūtinā/	/walilmaktabāt/
و + <u>ك</u> + بيوت + نا	و + ل + مكتبة + ات
wa+ka+biyūt+nā	wa+li+al+maktaba+āt
and+like+houses+our	and+for+the+library+plural
<i>And like our houses</i>	<i>And for the libraries</i>

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

# Inflectional Morphology

## Verbs



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

- Morphotactics
- Subject conjugation (suffix or circumfix)

# Inflectional Morphology

*katab ‘to write’*

- **Perfect** verb subject conjugation (*suffixes only*)

	Singular	Dual	Plural
1	كتب katabtu	كتبنا katabnā	
2	كتبت katabta	كتبتما katabtumā	كتبتم katabtum
3	كتب kataba	كتبنا katabā	كتبوا katabtū

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

	Singular	Dual	Plural
1	اكتب aktubu	نكتبnaktabu	
2	تكتب taktabu	تكتبان taktabān	تكتبون taktabūn
3	يكتب yaktubu	يكتبان yaktabān	يتكتبون yaktabūn

# Road Map

- Introduction
- Orthography
- Morphology
  - Derivational Morphology
  - Inflectional Morphology
  - Morphological Ambiguity
  - Arabic Computational Morphology
- Syntax

# Morphological Ambiguity

- Derivational ambiguity
  - قاعدة: basis/principle/rule, military base, Qa'ida/Qaeda/Qaida
- Inflectional ambiguity
  - تكتب /taktub/: you write, she writes
  - Segmentation ambiguity
    - وجده: he found; وجد: and+grandfather
    - لغة اللغة: for a language; لغة: for the language

# Morphological Ambiguity

- 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 → ي
- Multiple sources of ambiguity

بین

– /bayyana/	Verb	<i>he demonstrated</i>
– /bayyanna/	Verb	<i>they [feminine] demonstrated</i>
– /bayyin/	Adj	<i>clear/evident/explicit</i>
– /bayna/	Prep	<i>between/among</i>
– /biyin/	Proper Noun	<i>in Yen</i>
– /biyn/	Proper Noun	<i>Ben</i>

# Morphological Disambiguation

## *in English*

- Select a morphological tag that fully describes the morphology of a word
- Complete English morphological tag set (Penn Treebank): 48 tags

Verb:

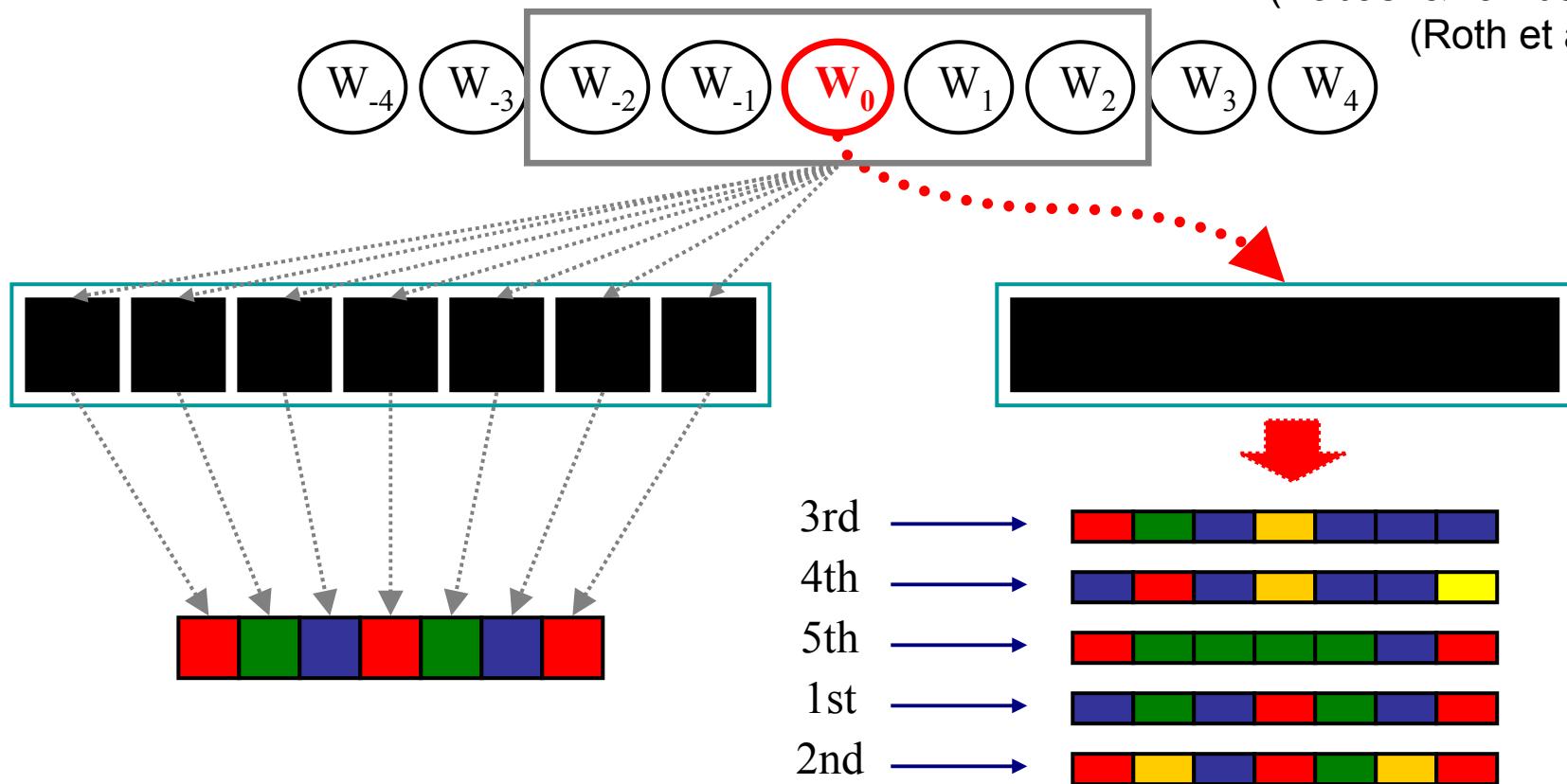
VB	VBD	VBG	VBN	VBP	VBZ
go	went	going	gone	go	goes

- Same as “POS Tagging” in English

# Morphological Disambiguation *in Arabic*

- Morphological tag has 14 subtags corresponding to different linguistic categories
  - Example: Verb  
Gender(2), Number(3), Person(3), Aspect(3),  
Mood(3), Voice(2), Pronominal clitic(12),  
Conjunction clitic(3)
- 22,400 possible tags
  - Different possible subsets
- 2,200 appear in Penn Arabic Tree Bank Part 1 (140K words)
- Example solution: MADA (Habash&Rambow 2005)

MADA (Habash&Rambow 2005)  
(Habash&Rambow 2007)  
(Roth et al. 2008)



## MORPHOLOGICAL CLASSIFIERS

- Multiple independent classifiers
- Corpus-trained

## RANKER

- Heuristic or corpus-trained

## MORPHOLOGICAL ANALYZER

- Rule-based
- Human-created

(Habash&Rambow 2005)  
(Habash&Rambow 2007)  
(Roth et al. 2008)

# MADA

## Morphological Analysis and Disambiguation for Arabic

::: SENTENCE AsbAnyA tnfy tjmyd AlmsAEdp AlmmnwHp llmgrb

::WORD AsbAnyA

::MADA: AsbAnyA art-NO aspect-NA case-NOCASE clitic-NO conj-NO def-DEF mood-NA n  
um-SG part-NO per-3 pos-PN voice-NA

\*0.78571 <isobAniyA=[<isobAniyA\_1 POS:PN BW:+<isobAniyA/NOUN\_PROP+]=Spain

^0.71429 >asobAniyA=[<isobAniyA\_1 POS:PN BW:+>asobAniyA/NOUN\_PROP+]=Spain

\_0.50000 <isobAniy~A=[<isobAniy~\_2 POS:AJ +MASC +DU +NOM +POSS BW:+<isobAniy~/ADJ+A/NSUFF\_MASC\_DU\_

\_0.50000 <isobAniy~AF=[<isobAniy~\_2 POS:AJ +ACC +INDEF BW:+<isobAniy~/ADJ+AF/CASE\_INDEF\_ACC]=Spanish/S

\_0.57143 <isobAniy~A=[<isobAniy~\_1 POS:N +MASC +DU +NOM +POSS BW:+<isobAniy~/NOUN+A/NSUFF\_MASC\_DU\_

\_0.57143 <isobAniy~AF=[<isobAniy~\_1 POS:N +ACC +INDEF BW:+<isobAniy~/NOUN+AF/CASE\_INDEF\_ACC]=Spanish.

-----  
::WORD tnfy

::MADA: tnfy art-NA aspect-IV case-NA clitic-NO conj-NO def-NA mood-I num-SG par  
t-NO per-3 pos-V voice-ACT

\*1.00000 tanofiy=[nafaY\_1 POS:V +IV MOOD:I +S:3FS BW:ta/IV3FS+nofiy/IV+(null)/IVSUFF\_MOOD:I]=disavow/deny/reje

\_0.76923 tunofayo=[nafA-u\_1 POS:V +IV +PASS MOOD:SJ +S:2FS

BW:tu/IV2FS+nof/IV\_PASS+ayo/IVSUFF\_SUBJ:2FS\_MOOD:SJ]=be\_rejected/be\_refuted/be\_denied

\_0.84615 tanif~iy=[naf~i\_1 POS:V +IV MOOD:SJ +S:2FS BW:ta/IV2FS+nif~/IV+iy/IVSUFF\_SUBJ:2FS\_MOOD:SJ]=blow\_1

\_0.84615 tanofiy=[nafA-u\_1 POS:V +IV MOOD:SJ +S:2FS BW:ta/IV2FS+nof/IV+iy/IVSUFF\_SUBJ:2FS\_MOOD:SJ]=refute/

\_0.84615 tanofiy=[nafaY\_1 POS:V +IV MOOD:SJ +S:2FS BW:ta/IV2FS+nof/IV+iy/IVSUFF\_SUBJ:2FS\_MOOD:SJ]=disavow

\_0.84615 tanofiya=[nafaY\_1 POS:V +IV MOOD:S +S:2MS BW:ta/IV2MS+nofiy/IV+a/IVSUFF\_MOOD:S]=disavow/deny/reje

\_0.92308 tanofiy=[nafaY\_1 POS:V +IV MOOD:I +S:2MS BW:ta/IV2MS+nofiy/IV+(null)/IVSUFF\_MOOD:I]=disavow/deny/rej

\_0.92308 tanofiya=[nafaY\_1 POS:V +IV MOOD:S +S:3FS BW:ta/IV3FS+nofiy/IV+a/IVSUFF\_MOOD:S]=disavow/deny/rejec

# Road Map

- Introduction
- Orthography
- Morphology
  - Derivational Morphology
  - Inflectional Morphology
  - Morphological Ambiguity
  - Arabic Computational Morphology
- Syntax

# Arabic Computational Morphology

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

# Arabic Computational Morphology

- Representation units (continued)
  - Prefix + Stem + Suffix
    - مكتب+ات+ولل
    - Can create more ambiguity
  - Lexeme + Features
    - مكتبة [+Plural +Def +و +ج]
  - Root + Pattern + Features
    - كتب + ةا3ا21ا + [+Plural +Def +ج +و]
    - Very abstract
  - Root + Pattern + Vocalism + Features
    - كتب + ة321م + a.a.a + [+Plural +Def +ج +و]
    - Very very abstract

# TOKAN

- A generalized tokenizer
- Assumes disambiguated morphological analysis
  - a la MADA
- Declarative specification of tokenization scheme

**wsyktbhA=[katab\_1 POS:V +IV w+ s+ +S:3MS +O:3FS]**

Example	Scheme	Specification
w+ syktbhA	D1	w+ f+ REST
w+ s+ yktbhA	D2	w+ f+ b+ k+ l+ s+ REST
w+ s+ yktb +hA	D3	w+ f+ b+ k+ l+ s+ AI+ REST +P: +O:
w+ syktb +hA	TB	w+ f+ b+ k+ l+ REST +P: +O:
w+ s+ ktb/VBZ S:3MS +hA	EN	w+ f+ b+ k+ l+ s+ AI+ LEXEME + BIESPOS +S:

- Uses generator (Habash 2004)

# Arabic Computational Morphology

- Approaches
  - Finite state machines (Beesely,2001) (Kiraz,2001) (Habash&Rambow 2006)
  - Concatenative analysis/generation (Smrz, 2007) (Buckwlatr,2002) (Cavalli-Sforza et al, 2000)
  - Lexeme+Feature analysis/generation (Habash, 2004) (Habash&Rambow 2006)
  - 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)
  - Survey article: (Al-Sughaiyer&Al-Kharashi, 2004)
- Issues
  - Appropriateness of system representation for an application
    - Machine Translation vs. Information Retrieval
    - Arabic spelling vs. phonetic spelling
  - System coverage
  - System extendibility
  - Availability to researchers
  - Use for analysis and generation

# Road Map

- Introduction
- Orthography
- Morphology
- Syntax
  - Morphology and Syntax
  - Sentence Structure
  - Phrase Structure
  - Computational Resources

# Morphology and Syntax

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

# 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

# Road Map

- Introduction
- Orthography
- Morphology
- **Syntax**
  - Morphology and Syntax
  - **Sentence Structure**
  - Phrase Structure
  - Computational Resources

# Sentence Structure

## *Two types of Arabic Sentences*

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

# Sentence Structure

- Verbal sentences
  - Verb agreement with gender only
    - Default singular number
    - كتب الولد \ الولاد wrote<sub>3MascSing</sub> the-boy/the-boys
    - كتبت البنات \ البنات wrote<sub>3FemSing</sub> the-girl/the-girls
  - Pronominal subjects are conjugated
    - كتبت wrote-you<sub>MascSing</sub>
    - كتبتم 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

# Sentence Structure

- Verbal sentences
  - Common structural ambiguity
    - *Third masculine/feminine singular is 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
      - كتب /kutiba/ it was written

# 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  
عندی كتاب / **Indi kitābun/** at-me a-book *I have a book*

# 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<sub>3rdMascPlur</sub> poems]] *The boys wrote the poems*
      - Full agreement in this order (SVO)
      - الاشعار كتبها الاولاد – [the-poems [wrote<sub>3rdMascSing</sub>-them the boys]] *The poems, the boys wrote*

# Road Map

- Introduction
- Orthography
- Morphology
- Syntax
  - Morphology and Syntax
  - Sentence Structure
  - Phrase Structure
  - Computational Resources

# 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>FemSing</sub> the-ambitious<sub>FemSing</sub>
    - الكاتبات الطموحات – the-writer<sub>FemPlur</sub> the-ambitious<sub>FemPlur</sub>
  - Exception: Plural non-persons
    - definiteness agreement; feminine singular default
    - المكتب الجديد – the-office<sub>MascSing</sub> the-new<sub>MascSing</sub>
    - المكتبة الجديدة – the-library<sub>FemSing</sub> the-new<sub>FemSing</sub>
    - المكاتب الجديدة – the-offices<sub>MascBPlur</sub> the-new<sub>FemSing</sub>
    - المكتبات الجديدة – the-libraries<sub>FemPlur</sub> the-new<sub>FemSing</sub>

# Phrase Structure

- Noun Phrase

- Idafa construction (اضافة)

- 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*

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

# Agreement in Arabic

- Verb-Subject agreement
  - Verb agrees with subject in full (gender,number)
    - Exception: partial agreement (number=singular) in VSO order
    - Exception: partial agreement (number=singular; gender=feminine) for non-person plural subjects regardless of order
- Noun-Adjective
  - Adjective agrees with noun in full (gender, number, definiteness and case)
    - Exception: partial agreement (number=singular; gender=feminine) for non-person plural nouns
- Noun-Number
  - Number is the syntactic-case head
  - for numbers [3..10]: Noun is plural+genitive (idafa); number gender is inverted gender of noun!
  - for numbers [11..99]: Noun is singular+accusative (tamyizz/specification); number gender is even more complicated ☺
  - for numbers [100,1K,1M]: Noun is singular+genitive (idafa)

bnyt ‘was built’	vIAv ‘three’	jAmEAt ‘universities’	jdydp ‘new’
Fem+Sg	Masc+Sg+Nom	Fem+PL+Gen	Fem+Sg+Gen
Verbs in VSO order are always Sg and agree in gender only	Numbers agrees by gender inversion		Adjectives of plural non-person nouns are Fem+Sg

# Road Map

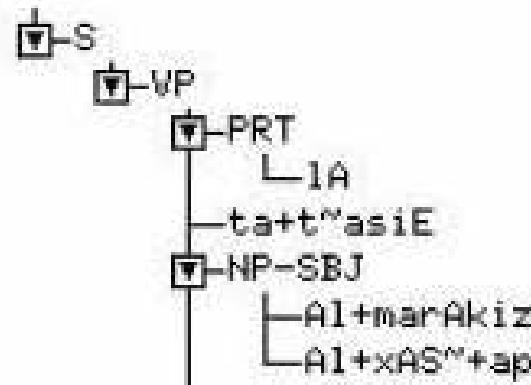
- Introduction
- Orthography
- Morphology
- **Syntax**
  - Morphology and Syntax
  - Sentence Structure
  - Phrase Structure
  - Computational Resources

# 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)

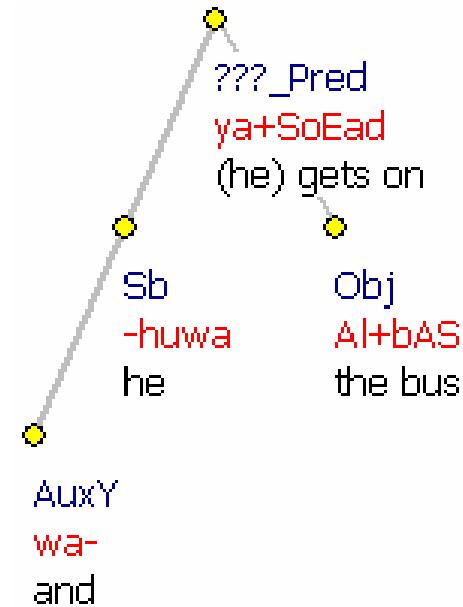
# Computational Resources

- 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)



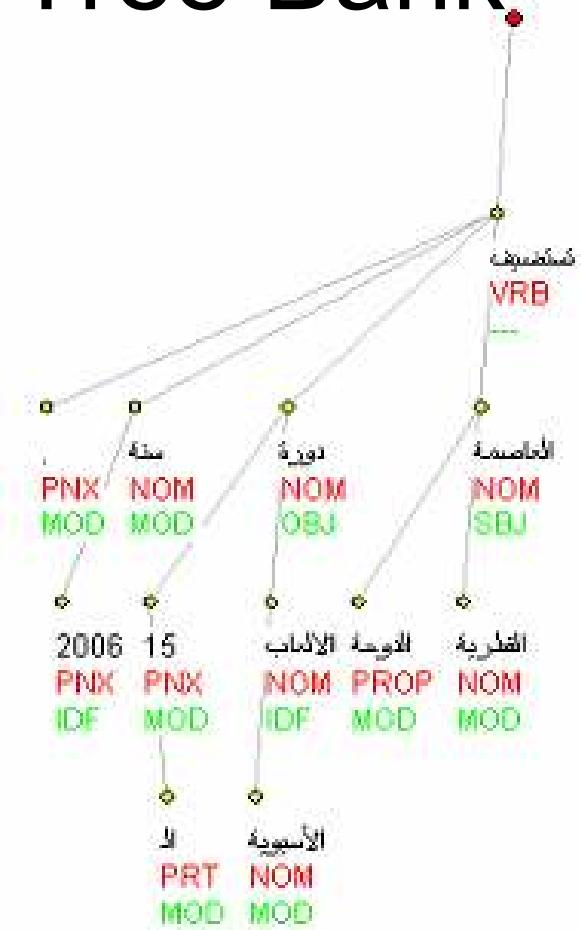
# Computational Resources

- Prague Dependency Treebank
- Partial overlap with PATB and Arabic Gigaword
  - Agence France Presse, AlHayat and Xinhua
- Morphological analysis
  - Extends on PATB
- Dependency representation



# CATiB:Columbia Arabic Tree Bank

- **CATiB Representation**
  - *Lite Dependency Syntax*
  - *Tokenization*
    - CONJ PART BASE PRON
  - *Part of Speech Tag set*
    - Six tags: VRB, VRB-pass, NOM, PROP, PRT, PNX
  - *Syntactic Annotation*
    - Eight dependency relations: SBJ, OBJ, TPC, PRD, IDF, TMZ, MOD, FLAT
- **Practical Considerations**
  - Less information to annotate
  - Dependencies easier to annotate than phrase structure
  - Terminology and representation close to traditional Arabic grammar



# Computational Resources

- Applications using Arabic treebanks
  - Statistical parsing
    - Bikel's parser (Bikel 2003)
      - Same engine used with English, Chinese and Arabic
    - Nivre's MALT parser (Nivre et al. 2006)
  - Base-phrase Chunking
    - (Diab et al, 2004; Diab et al. 2007)
  - POS tagging and morphological disambiguation
    - (Diab et al, 2004; Diab et al. 2007; Habash and Rambow, 2005a; Smith et al., 2005; Roth et al. 2008)
    - Other non-treebank-based POS tagging efforts: (Khoja, 2001)
- Formalism conversion
  - Constituency to dependency (Žabokrtský and Smrž 2003; Habash et al. 2007; Tounsi et al., 2009)
  - Tree-adjoining grammar extraction (Habash and Rambow 2004)
- Automatic diacritization
  - Zitouni et al. (2006); Habash&Rambow (2007); Shaalan et al (2008) among others
  - Diacritization for MT (Diab et al. 2007)

# **Other Tutorial Slides**

- **Columbia's Arabic Dialect Modeling Group (CADIM)**
  - <http://www1.ccls.columbia.edu/~cadim/>
    - Presentations

**MEADR 2009**

**Cairo, Egypt**  
**April 21, 2009**

# Introduction to Arabic Natural Language Processing

**Nizar Habash**

Columbia University

Center for Computational Learning Systems

*habash@ccls.columbia.edu*

