

LREC 2008

Arabic Dialect Processing

Mona Diab Nizar Habash
Center for Computational Learning Systems
Columbia University

Tutorial Contents

- Introduction
- Description of MSA Phenomena
- Description of Dialectal Phenomena
- Sample Applications
- Resources and References

Introduction

- Arabic is a Semitic language
- Forms of Arabic
 - Classical Arabic (CA)
 - Classical Historical texts
 - Liturgical texts
 - Modern Standard Arabic (MSA)
 - News media & formal speeches and settings
 - Only written standard
 - Dialectal Arabic (DA)
 - Predominantly spoken vernaculars
 - No written standards
- Dialect vs. Language
 - Linguistics vs. Politics

3

Introduction

- ~300M people worldwide speak Arabic
- Arabic is the/an official language of 23 countries
- No native speakers of CA nor MSA
- In the Arabic speaking world, MSA and CA are the only Arabic taught in schools

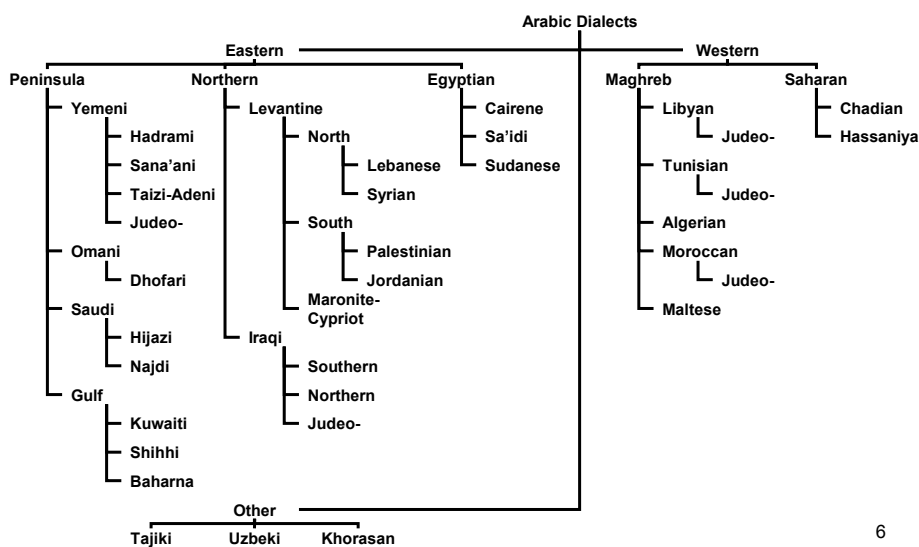
4

Introduction

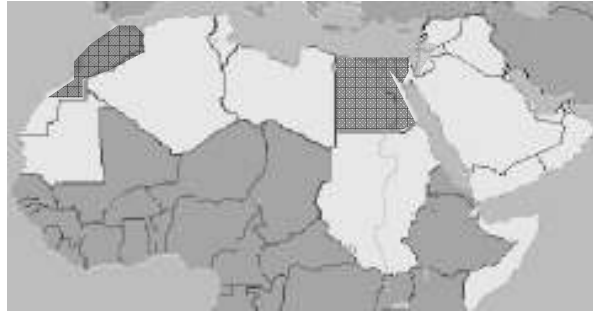
- Arabic Diglossia
 - Diglossia is where two forms of the language exist side by side
 - MSA is the formal public language
 - Perceived as "language of the mind"
 - Dialectal Arabic is the informal private language
 - Perceived as "language of the heart"
- General Arab perception: dialects are a deteriorated form of Classical Arabic
- Continuum of dialects

5

Geographical Continuum



6



lam jaftari nizār ṭawilatan	ζadīdatan	لو يشتري نزار طاولة جديدة
didn't buy	Nizar table	new
nizār maṣṭarāf ṭarabēza gidīda	●	نزار ماشراف طريزة جديدة
nizār maṣṭarāf ṭawile	●	نزار ماشراف طاولة جديدة
nizar maṣṭarāf mida	●	نزار ماشراف ميده جديدة
Nizar	not-bought-not	table new

7

Social Continuum

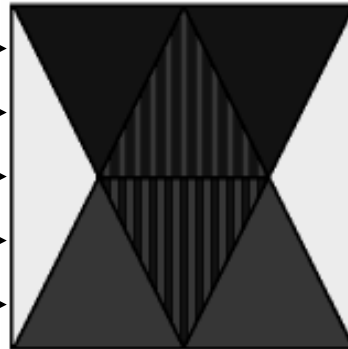
- Factors affecting dialect
 - Lifestyle
 - Bedouin, urban, rural
 - Education & Social Class
 - Religion
 - Muslim, Christian, Jewish, Druze, etc.
 - Gender

8

Social Continuum

- Badawi's levels

- Traditional Arabic →
- Modern Arabic →
- Educated Colloquial →
- Literate Colloquial →
- Illiterate Colloquial →



Classical Dialect Foreign

- Polyglossia

Why Study Arabic Dialects?

- **Almost no** native speakers of Arabic sustain continuous spontaneous production of MSA
- Ubiquity of Dialect
 - Dialects are the primary form of Arabic used in all unscripted spoken genres: conversational, talk shows, interviews, etc.
 - Dialects are increasingly in use in new written media (newsgroups, weblogs, etc.)
 - Dialects have a direct impact on MSA phonology, syntax, semantics and pragmatics
 - Dialects lexically permeate MSA speech and text
- Substantial Dialect-MSA differences impede direct application of MSA NLP tools

10

Why Study Arabic Dialects?

- Degrees of linguistic distance

	Syntax	Morphology	Lexicon	Phonology
MSA-Dialect	++	+++	++++	++++
Inter-Dialect	+	+++	++++	++++
Intra-Dialect	0	0	+	+

- Lack of standards for the dialects
- Lack of written resources

11

A Note on Romanization

- Phonological Transcription

- IPA

- Transliteration

- Strict (one-to-one)

- Buckwalter Encoding



- Loose

- Many spelling variants

- Qadafi, kadaphi, kaddafy, etc.

- This tutorial's examples are in

- Arabic script

- Transcription (IPA)

- Transliteration (Buckwalter) slAm

ل	*	ذ	ز	ع
م	x	ر	ا	آ
ن	z	ز	>	أ
ه	s	س	س	ؤ
و	S	ش	<	إ
ي	S	ص	ص	ئ
ي	D	ض	ا	أ
ف	T	ط	ب	ب
N	Z	ظ	ة	ة
K	E	ع	ت	ت
a	g	غ	ث	ث
u	—	—	ج	ج
i	f	ف	ح	ح
~	q	ق	خ	خ
o	k	ك	د	د

12

Tutorial Contents

- Introduction
- Description of MSA Phenomena
 - Orthography
 - Morphology
 - Syntax
- Description of Dialectal Phenomena
- Sample Applications
- Resources and References

13

Arabic Script

الخطُ العَرَبِي

- An alphabet
- Written right-to-left
- Letter have allographic variants
- Optional diacritics
- Common ligatures
- Used to write many languages in addition to Arabic: Persian, Kurdish, Urdu, Pashto, etc.

14

Arabic Script

Alphabet

- letter forms

ع ط ص س ر د ح ب ا
ء ي و ه ن م ل ف

- letter marks

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

15

Arabic Script

Alphabet

- letters (form+mark)

- Distinctive

ب ت ث س ش
/θ/ /t/ /b/
/ʃ/ /s/

- Non-distinctive

أ إ آ إ ء

/ʔ/
glottal stop aka hamza

16

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
بُ /ban/	بَ /ba/
بُنْ /bun/	بُ /bu/
بِ /bin/	بِ /bi/

17

Arabic Script

Diacritics

- No-vowel marker (*sukun*)

مَكْتَبٌ /maktab/ *office*

- Double consonant marker (*shadda*)

كَتَّبَ /kattab/ *to dictate*

- Combinable

بُبُّ بَّبُّ بَبُّ
/bbu/ /bbin/ /bban/

No Vowel
بْ /b/

Double Consonant
بَّبُّ /bb/

18

Arabic Script

Putting it together

Simple combination

Arab /ʕarab/ عَرَب ← عَرَب = عَرَب

West /ʕarb/ غَرَب ← غَرَب = غَرَب

Ligatures

Peace /salām/ س ل ا م ← س ل ا م سلام

19

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

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

20

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



21

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 /ħubb/ love حَبَّ /ħabb/ seed
 - Dual use of ا, و, ي as consonant and long vowel
 - ا (/ʔ/, /ā/) و (/w/, /ū/) ي (/j/, /ī/)

22

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 بهاءة بهاءة (ء أأؤئ) /baha'/ + 3MascSing (his glory)

23

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

24

Tutorial Contents

- Introduction
- Description of MSA Phenomena
 - Orthography
 - Morphology
 - Syntax
- Description of Dialectal Phenomena
- Sample Applications
- Resources and References

25

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

26

Derivational Morphology

- Templatic Morphology

- Root

ك ت ب
k t b

- Pattern



- Lexeme

مكتوب
maktūb
written

كاتب
kātib
writer

Lexeme.Meaning =

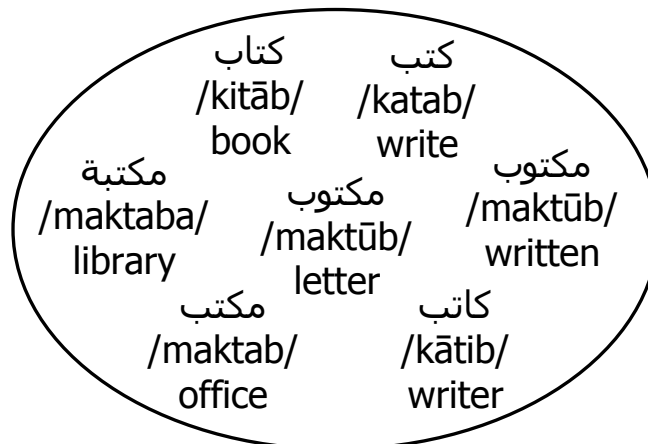
(Root.Meaning+Pattern.Meaning)*Idiosyncrasy.Random

27

Derivational Morphology

Root Meaning

- ك ت ب KTB = notion of "writing"



28

Root Polysemy

LHM-1 لحم

"meat"

لحم /laħm/

Meat

لحام /laħħām/

Butcher



LHM-2 لحم

"battle"

ملحمة /malħama/

Fierce battle

Massacre

Epic



LHM-3 لحم

"soldering"

لحم /laħam/

Weld, solder,

stick, cling



29

Derivational Morphology

Pattern Meaning

- Verb Pattern Meaning is hard to define systematically

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

30

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
 - Definiteness: definite, indefinite
 - Case: nominative, accusative, genitive
 - Possessive clitic

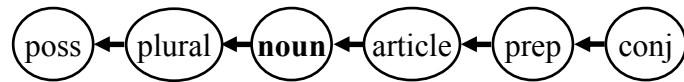
31

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

32

Inflectional Morphology Nouns

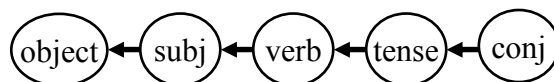


<p>و كيبوتنا /wakabiyūtinā/ و + ك + بيوت + نا wa+ka+biyūt+nā and+like+houses+our <i>And like our houses</i></p>	<p>وللمكتبات /walilmaktabāt/ و + ل + ال + مكتبة + ات wa+li+al+maktaba+āt and+for+the+library+plural <i>And for the libraries</i></p>
---	--

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

33

Inflectional Morphology Verbs



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

- Morphotactics
- Subject conjugation (suffix or circumfix)

34

Inflectional Morphology

- Perfect verb subject conjugation (*suffixes only*)

	Singular	Dual	Plural
1	كُتِبْتُ katabtu	كُتِبْنَا katabnā	
2	كُتِبْتَ katabta	كُتِبْتُمَا katabtumā	كُتِبْتُمْ katabtum
3	كُتِبَ kataba	كُتِبَا katabā	كُتِبُوا katabtū

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

	Singular	Dual	Plural
1	أَكْتُبُ aktubu	نَكْتُبُ naktubu	
2	تَكْتُبُ taktubu	تَكْتُبَانِ taktubān	تَكْتُبُونَ taktubūn
3	يَكْتُبُ yaktubu	يَكْتُبَانِ yaktubān	يَكْتُبُونَ yaktubūn

35

Feminine form and other verb moods not shown

Morphological Ambiguity

- Derivational ambiguity
 - قاعدة: basis/principle/rule, military base, Qa'ida/Qaeda/Qaida
- Inflectional ambiguity
 - تَكْتُبُ: you write, she writes
 - Segmentation ambiguity
 - وَجَدَ: he found; وَجَدَ: and+grandfather
- Spelling ambiguity
 - Optional diacritics
 - كَاتِبُ: /kātib/ writer, /kātab/ to correspond
 - Suboptimal spelling
 - Hamza dropping: اِ، اُ → ا
 - Undotted ta-marbuta: ه → ة
 - Undotted final ya: ي → ى

36

Tutorial Contents

- Introduction
- Description of MSA Phenomena
 - Orthography
 - Morphology
 - Syntax
- Description of Dialectal Phenomena
- Sample Applications
- Resources and References

39

Morphology and Syntax

- Rich morphology crosses into syntax
 - Pro-drop / Subject conjugation
 - Verb subcategorization and object clitics
 - Verb_{transitive}+subject+object
 - Verb_{intransitive}+subject
 - Verb_{passive}+subject
- 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.

40

Morphology and Syntax

- Morphological interactions with syntax (continued)
 - Case
 - MSA is case marked: 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 (*minimal segmentation*) is necessary

41

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

42

Sentence Structure

- Verbal sentences
 - Verb agreement with gender only
 - كتب الولد\الاولاد wrote_{3MascSing} the-boy/the-boys
 - كتبت البنت\البنتات wrote_{3FemSing} the-girl/the-girls
 - Pronominal subjects are conjugated
 - كتبتُ wrote-you_{MascSing}
 - كتبتُم wrote-you_{MascPl}
 - كتبوا wrote-they_{MascPl}
 - Passive verbs
 - Same structure: Verb_{passive} Subject_{underlyingObject}
 - Agreement with surface subject

43

Sentence Structure

- Verbal sentences
 - Common structural ambiguity
 - *Third masculine/feminine singular are structurally ambiguous*
 - Verb_{3MascSing} Noun_{Masc}
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

44

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*

45

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_{3MascPl}-they poems]] *The boys wrote the poems*
 - Full agreement in this order (SVO)
 - الاشعار كتبها الاولاد
[the-poems [wrote_{3MascSing}-it the boys]] *The poems, the boys wrote*

46

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_{fem} the-ambitious_{fem}
 - الكاتبات الطموحات the-writer_{femPl} the-ambitious_{femPl}

47

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

48

Phrase Structure

- Morphological *definiteness* interacts with syntactic structure

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

49

Tutorial Contents

- Introduction
- Description of MSA Phenomena
- Description of Dialectal Phenomena
 - Orthography
 - Lexicon
 - Morphology
 - Syntax
 - Code switching
- Sample Applications
- Resources and References

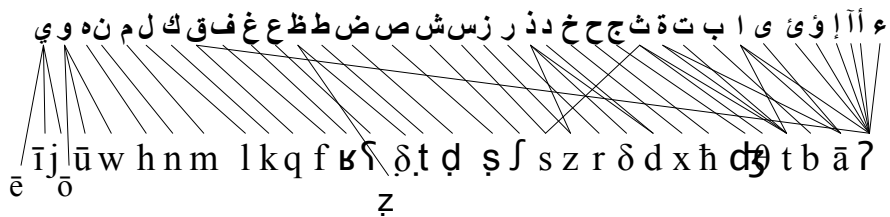
50

Phonological Variations

MSA



LEV



- phoneme quality differences

51

Phonological Variations

- Major variants

MSA		Dialects
ق	/q/	/q/, /k/, /ʔ/, /g/, /dʒ/
ث	/θ/	/θ/, /t/, /s/
ذ	/ð/	/ð/, /d/, /z/
ج	/dʒ/	/dʒ/, /g/

- Some of many limited variants
 - /l/ → /n/ MSA: /burtuqāl/ → LEV: /burtʔān/ 'orange'
 - /ʁ/ → /ħ/ MSA: /kaʁk/ → EGY: /kaħk/ 'cookie'
 - Emphasis add/delete: MSA: /fustān/ → LEV: /fuʂʔān/ 'dress'

52

Script Choices

- Arabic script:
 - + continuity with MSA
 - + masks the vocalic and some consonantal difference across dialects
 - ambiguity
- Latin script
 - + precision
 - lose connections among dialects (within dialects)
 - politically loaded
- Other scripts
 - Hebrew and Syriac
 - Different religious/ethnic preferences

53

Arabic Script Orthographic Variants

	IRQ	LEV	EGY	TUN	MOR
/dʒ/	ج	ج	چ	ج	ج
/g/	گ	چ	ج	ق	ك
/tʃ/	چ	تش	تش	تش	تش
/p/	پ	پ	پ	پ	پ
/v/	ف	ف	ف	پ	پ

- Historical variants: MSA (ق, ف) = MOR (ف, ب)
- Modern proposals: LEV /ʔ/ ق[□], /ē/ ی, /ō/ و (Habash 1999)

54

Syrian Arabic in Arabic Script

رح إحكي عنا نحن السوريين ..المعروفين بمأكولاتنا الشهية
واللذيذة والمميزة...مو بس هيك كل الخير فيها..دسمة وتقيلة
وعين الله ما بينقصها شي من المكسرات و...و...واللي لا
يمكن ترحمنا إذا ما رحمنا حالنا ..فبتلاقينا منهجم عالآكل يا
قاتل يا مقتول حتى التلت اللي لازم نتركه للنفس بديق بعيننا
و منعبيه أكل

55

<http://www.soriagate.net/showthread.php?t=32678>

Latin Script

- Several proposals to the Arabic Language Academy in the 1940s
- Said Akl Experiment (1961) →
- Web Arabic (Arabish, Franco-arabe)
 - No standard, but common conventions

عربي	IPA	Latin	عربي	IPA	Latin
أأعوى	/ʔ/	ʔ 2 Ø	ث	/θ/	th
ة	/a/,/t/	a t	ط	/t/	t T 6
ح	h	H h 7	ع	/ʕ/	ʕ 3 Ø
خ	/x/	kh 7' x 8	غ	/ɣ/	g gh 3'
ذ	/ð/	th	ق	/q/	q
ش	/ʃ/	sh ch	ي	/y/ /ay/ /ī/ /ē/	y, i, e, ai, ei, ...

Akl 1961

ʕ caʕoof	F fo
B be	V ve
P po	Q qoof
T to	L laam
T tahh	M miim
J jiin	N nnun
X xo	H ho
K ko	W waaw
D daal	A a
D daal	A a
R ro	I i
Z zayn	E e
Z zahh	E e
S siin	O o
S saad	U u (ou)
C ciin	U u
Y yayn	Y yu
G gayn	
G zo (zoo)	

56

Egyptian Arabic in Latin Script

nadeity bsho2 nadeit
olteely ta3ala geit
laha3atbek 3alli fat
wala 7atta haloom 3aleiky
adeeni rge3telek
adeeni bein edeiky
kefaya dmoo3 ba2a
mush 3aref ashoo3 3eneiky

57

http://www.jencomics.com/artist_a/amr_diab_lyrics/adeeni_rege3tilik_lyrics.html

The Case of Maltese

- An Arabic dialect that is considered a separate language
- Standardized Latin-based orthography

Kulhadd hu intitolat għal dawn il-jeddijiet u l-libertajiet imxandra f'din l-Istqarrija, bla ebda għażla, bħal ta' razza, lewn, sess, ilsien, reliġjon, opinjoni politika jew kull opinjoni oħra, oriġini nazzjonali jew soċjali, proprjetà, twelid jew kull qagħda oħra. Mhux biss, iżda l-ebda għażla m'għandha ssir fuq bażi tal-qagħda politika, ġuridika jew internazzjonali tal-pajjiż jew territorju li minnu tiġi l-persuna kemm jekk ikun indipendenti, kemm jekk ikun fdat lil xi pajjiż ieħor, m'għandux gvem tiegħu jew għandu xi limiti oħra fis-sovraniġġ tiegħu.

58

<http://www.language-museum.com/m/maltese.htm>

Hebrew Script

- Example from Tunisian Judeo-Arabic

"The Ballad of Hannah and her Seven Sons "

קעת חנה וזכריה
א אסמעה קולי אנא חנה ואנערד מא ג'רא לי
לי סבע בנן באל כרם ועז ובאל ד'אלי. וכאן
ביהום ולד זג'יר ונהו יע'וי כאל הלאלי ווקעו פי יד כאפר
מא יכאף מן רב אל עאלי. יעלח לנא לנבכי טול אל
איאם וליאלי

59

<http://www.uwm.edu/~corre/jatexts/>

Lack of Orthographic Standards

- Orthographic inconsistency
- Egyptian /mabin?ulhalakf/

- | | |
|---------------------|----------------|
| - mA binqwlhA lak\$ | ما بنقولها لكش |
| - mAbin&ulhalak\$ | مابنؤلهاالكش |
| - mA bin}ulhAlak\$ | ما بنئلهاالكش |
| - mA binqulhA lak\$ | ما بنقلها لكش |
| - ... | |

60

Spelling Inconsistency I

في البدايا خلق الله (السَّمَا) والأرض. والأرض
كانت خَرَبَانِي وفاضيي وعلى وُشْنُ الغمق عتيمي وروح
الله يرفرق على وُشْنُ المويي. وقال الله خَلِّي يصير ضَوء
وصار(ضوء) وشاف الله (الضَوء) أنو شي ظريف وفرَّق
الله بين الضَّوء والعتيمي. وسَمَّى الله الضَّوء نهار
والعتيمي سَمَّاها ليل وكان(مسا) وكان صباح يوم واحد.
وقال الله خَلِّي يصير جَوِّ في وسط المويي ويصير
فَاصِل بين المُوَيِّي ومُوَيِّي. وعمل الله الجَوِّ وفرَّق بين
المُوَيِّي اللَّيِّ تحت الجَوِّ والمُوَيِّي فوقَ الجَوِّ وهيك صار.
وسَمَّى الله الجَوِّ(سما) وكان(سما) وكان صباح يوم ثاني.

61

<http://www.language-museum.com/a/arabic-north-levantine-spoken.php>

Spelling Inconsistency II

- ya alain lesh el 2aza
ti7keh 3anneh kaza w kaza
iza bidallak ti7keh hek
2areeban ra7 troo7 3al 3aza

chi3rik 3emilleh na2zeh
li2anneh manneh mi2zeh
bass law baddik yeha 7arb
fikeh il layleh ra7 3azzeh

62

<http://www.onelebanon.com/forum/archive/index.php/t-8236.html>

Tutorial Contents

- Introduction
- Description of MSA Phenomena
- Description of Dialectal Phenomena
 - Orthography
 - Lexicon
 - Morphology
 - Syntax
 - Code switching
- Sample Applications
- Resources and References

63

Lexical Variation

- Arabic Dialects vary widely lexically

English	Table	Cat	Of	I_want	There_is	There_isn't
MSA	Tāwila طاولة	qīTTa قطة	idafa Ø	'uridu اريد	yūjadu يوجد	lā yujadu لا يوجد
Moroccan	mida ميدة	qeTTa قطة	dyāl ديال	byīt بغيت	kāyn كاين	mā kāynš ما كاينش
Egyptian	Tarabēza طربيزة	'oTTa قطة	bitāç بتاع	çāwez عاوز	fī في	mafīš مفيش
Syrian	Tāwle طاولة	bisse بسة	tabaç تبع	biddi بدي	fī في	mā fi ما في
Iraqi	mēz ميز	bazzūna بزونة	māl مال	'arīd اريد	aku اكو	māku ما

- Arabic orthography allows consolidating some variations

64

Lexical Variation

- خلف EGY:reproduce - GLF: give condolences
- مكوى EGY:press iron - GLF:buttocks
- براد EGY:kettle - LEV:fridge
- مرا EGY:prostitute - LEV:woman
- ماشي EGY/LEV:okay - MOR:not
- بسط EGY/LEV:make happy - IRQ:beat up
- العافية EGY/LEV:health - MOR:hell fire
- بلش LEV:start - SUD:end

65

Foreign Borrowings

- أوكي >wky okay
- مرسي mrsy merci
- بندورة bndwrp pomodoro (italian)
- بيرا byrA birra (italian)
- فرمت frmt format
- تلفون tlfwn telephone
- تلفن talfan to phone

66

Tutorial Contents

- Introduction
- Description of MSA Phenomena
- Description of Dialectal Phenomena
 - Orthography
 - Lexicon
 - Morphology
 - Syntax
 - Code switching
- Sample Applications
- Resources and References

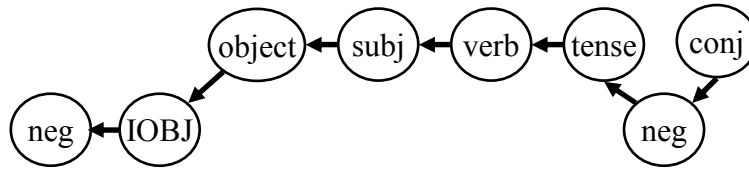
67

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 (2nd, 3rd person)
 - Loss of morphological moods
 - Subjunctive/jussive form dominates in some dialects
 - Indicative form dominates in others
- Other aspects increase in complexity

68

Morphological Variation Verb Morphology



MSA
ولم تكتبوها له
/walam taktubū+hā la+hu/
/wa+lām taktubū+hā la+hu/
and+not_past write_you+it for+him

EGY
وماكنتبوتها لوش
/wimakatabtuhalūʃ/
/wi+ma+katab+tu+ha+lū+ʃ/
and+not+wrote+you+it+for_him+not

And you didn't write it for him

69

Morphological Variation

	<i>Perfect</i>		<i>Imperfect</i>		
	Past	Subjunctive	Present habitual	Present progressive	Future
MSA	كتب /kataba/	يكتب /jaktuba/	يكتب /jaktubu/		سيكتب /sajaktubu/
LEV	كتب /katab/	يكتب /jiktob/	بيكتب /bjoktob/	عم بيكتب /ʕam bjoktob/	حيكتب /ħajiktob/
EGY	كتب /katab/	يكتب /jiktib/	بيكتب /bjiktib/		هيكتب /hajiktib/
IRQ	كتب /kitab/	يكتب /jiktib/	ديكتب /dajiktib/		رح يكتب /raħ jiktib/
MOR	كتب /kteb/	يكتب /jekteb/	كيكتب /kjektib/		غيكتب /ʕajektib/

70

Morphological Variation

Verb conjugation

	Perfect			Imperfect		
	1S	2S♂	2S♀	1S	1P	2S♀
MSA	كُتِبْتُ /katabtu/	كُتِبْتَ /katabta/	كُتِبْتِ /katabti/	اُكْتُبُ /aktubu/	نُكْتُبُ /naktubu/	تُكْتُبِينَ /taktubīna/ تُكْتُبِي /taktubī/
LEV	كُتِبْتُ /katabt/		كُتِبْتِي /katabti/	اُكْتُبْ /aktob/	نُكْتُبْ /noktob/	تُكْتُبِي /toktobi/
IRQ	كُتِبْتُ /kitabit/		كُتِبْتِي /kitabti/	اُكْتُبْ /aktib/	نُكْتُبْ /niktib/	تُكْتُبِينَ /tikitbīn/
MOR	كُتِبْتُ /ktebt/	كُتِبْتِي /ktebti/		نُكْتُبْ /nekteb/	نُكْتُبُوا /nektebu/	تُكْتُبِي /tektebi/

71

Tutorial Contents

- Introduction
- Description of MSA Phenomena
- Description of Dialectal Phenomena
 - Orthography
 - Lexicon
 - Morphology
 - Syntax
 - Code switching
- Sample Applications
- Resources and References

72

Idafa Construction

- Genitive/Possessive Construction
- Both MSA and dialects
 - Noun1 Noun2
 - ملك الاردن
king Jordan
the king of Jordan / Jordan's king
- Ta-marbuta allomorphs

	Idafa	No Idafa	Waqf
MSA	+at		+a
EGY	+it	+a	

- Dialects have *an additional* common construct
 - Noun1 <exponent> Noun2
 - LEV: الملك تبع الاردن *the-king belonging-to Jordan*
 - <exponent> differs widely among dialects

73

Demonstrative Articles

- Forms

	Proclitic	Word	
		Proximal	Distal
MSA	-	هذا, هذه, هؤلاء	ذلك, تلك, اولئك
EGY	-	ده, دي, دول	
LEV	+هـ	هدا, هادي, هدول	هداك, هديك, هدوك

- Word Order (Example: *this man*)

	Pre-nominal	Post-nominal
MSA	هذا الرجل	X
EGY	X	الرجل ده
LEV	هدا الرجال	الرجال هدا

74

Negation of Declarative Verbal Sentences

	Pre	Circum	Post
MSA	لا, لم, لن, ما lA, l _m , l _n , mA	X	X
EGY	مش m\$	ما ... ش mA ... \$	X
LEV	ما, مش mA, m\$	ما ... ش mA ... \$	ش \$

75

Sentence Word Order

- Verbal sentences
 - The boys wrote the poems
 - MSA
 - Verb Subject Object (Partial agreement)
كتب الاولاد الاشعار
wrote_{masc} the-boys the-poems
 - Subject Verb Object (Full agreement)
الاولاد كتبوا الاشعار
the-boys wrote_{mascPl} the-poems
 - LEV, EGY
 - Subject Verb Object
الاولاد كتبوا الاشعار
The-boys wrote_{mascPl} the-poems
 - Less present: Verb Subject Object
كتبوا الاولاد الاشعار
wrote_{mascPl} the-boys the-poems
 - Full agreement in both orders

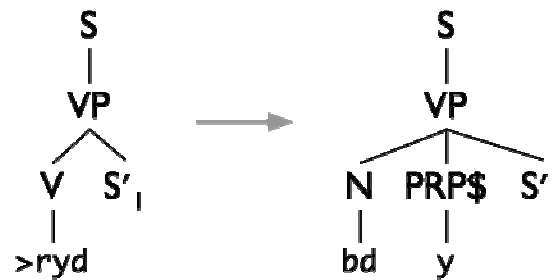
	V-S <i>explicit subject</i>	V(S) <i>pro dropped subject</i>	S-V <i>explicit subject</i>
MSA	35%	30%	35%
LEV	10%	60%	30%

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

76

Lexico-syntactic Variation

- 'want' (Levantine)



77

Tutorial Contents

- Introduction
- Description of MSA Phenomena
- Description of Dialectal Phenomena
 - Orthography
 - Lexicon
 - Morphology
 - Syntax
 - Code switching
- Sample Applications
- Resources and References

78

Code Switching

MSA
LEV

MSA and Dialect mixing in speech

- phonology, morphology and syntax

لا أنا ما بعتمد لأنه عملية اللي عم ببعارضوا اليوم تمديد للرئيس لحد هم اللي طالبوا بالتمديد للرئيس الهراوي وبالتالي موضوع منه موضوع ميدني على الأرض أنا بحترم أنه يكون في نظرة ديمقراطية للأمر وأنه يكون في احترام للعبة الديمقراطية وأن يكون في ممارسة ديمقراطية وبعتمد 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 بتتمير جهود الوطنية الشاملة كي يظل في مصالحة وطنية كي يظل في توافق ما بين المسلم والمسيحي في لبنان يحضن أبناء هذا البلد ما يترك المسار يروح باتجاه الخطأ نعم إنما خطاب القسم كان موضوع مبادئ طرحت هو ملتزم فيها اللي مشبوها معه وأمنوا 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 بنفهم 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 قناعتني في هذا الموضوع.

79

Aljazeera Transcript http://www.aljazeera.net/programs/op_direction/articles/2004/7/7-23-1.htm

Code Switching with English

- Iraqi Arabic Example
 - ya ret 3inde hech sichena tit7arrak wa77ad-ha , 7atta ma at3ab min asawwe zala6a yomiyya :D
 - 3ainee Zainab, tara hathee technology jideeda, they just started selling it !! Lets ask if anybody knows where do they sell them ! :

80

<http://www.aliraqi.org/forums/archive/index.php/t-16137.html>

Dialectal Impact on MSA

- Loss of case endings and nunation in read MSA
/fī bajt dʒadīd/
instead of /fī bajtin dʒadīdin/
'in a new house'
- A shift toward SVO rather than VSO in written MSA

81

Dialectal Impact on MSA

- Structure borrowing
- Example: monies and properties of the company

- اموال الشركة وممتلكاتها
• /ʔamwālu ʃʃarikati wamumtalakātuhā/
• *monies the-company and-properties-its*

- اموال وممتلكات الشركة
• /ʔamwālu wamumtalakātu ʃʃarikati/
• *monies and-properties the-company*



82

Dialectal Impact on MSA

- Code switching in written MSA
- Dialectal lexical and structural uses
 - Example Newswire Alnahar newspaper (ATB3 v.2)

فأخذ على خاطر الأخوان ومن حقهم ان يزععلوا

f>x* EIY xATr AlAxwAn wmn hqhm An yzElw

*then-was-taken upon self the-brothers and-from right-their
to be-angry*

'they were upset, and they had the right to be angry'

83

Tutorial Contents

- Introduction
- Description of MSA Phenomena
- Description of Dialectal Phenomena
- Sample Applications
 - Automatic speech recognition
 - Dictionary creation
 - Morphological analysis
 - Part-of-speech tagging
 - Syntactic parsing
 - Machine translation
- Resources and References

84

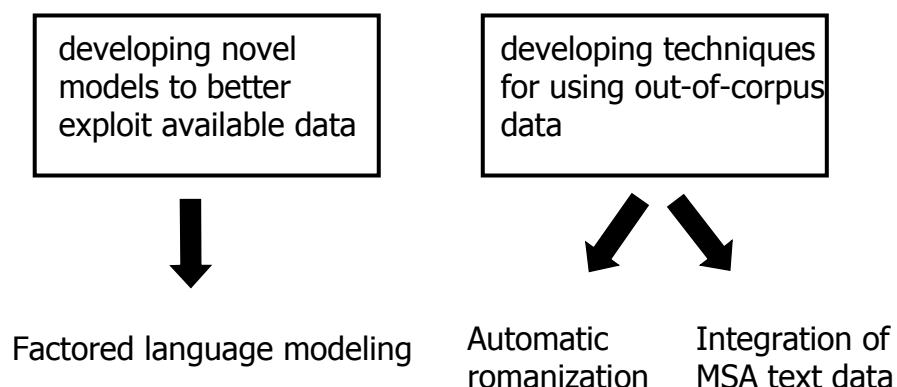
Arabic ASR: State of the Art

- BBN TIDESOnTap: 15.3% WER
- BBN CallHome system: 55.8% WER
- JHU WS 2002: 53.8% WER
- WER on conversational speech noticeably higher than for other languages
(eg. 30% WER for English CallHome)

85

JHU WS02 Approach

improvements to Arabic ASR through



86

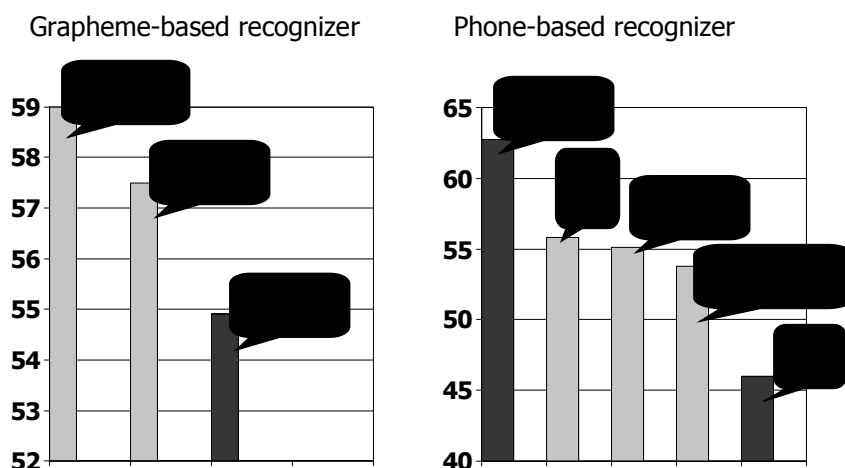
Slide courtesy of (Kirchhoff et al.2002)

Approach Details

- Factored Language Models
 - complex morphological structure leads to large number of possible word forms
 - break up word into separate components
 - build statistical n-gram models over individual morphological components rather than complete word forms
- Automatic Vowelization/Diacritization
 - try to predict vowelization automatically from data and use result for recognizer training
- Integrate data from MSA written sources

87

JHU WS02 Results (WER)



88

Slide courtesy of (Kirchhoff et al.2002)

Dialect-MSA Dictionary

- Problem: total lack of Dialect-MSA resources
 - No Dialect-MSA parallel text
 - No paper dictionaries for Dialect-MSA
- Dialect-MSA dictionary is required for many NLP applications exploiting MSA resources
 - e.g., to translate dialect sentences to MSA before parsing them with an MSA parser

89

Levantine-MSA Dictionary

- **The Automatic-Bridge dictionary (AB)**
 - English as a bridge language between MSA and LA
- **The Egyptian-Cognate dictionary (EC)**
 - Levantine-Egyptian cognate words in Columbia University Egyptian-MSA lexicon (2,500 lexeme pairs)
- **The Human-Checked dictionary (HC)**
 - Human cleanup of the union of AB and EC
 - Using lexemes speeded up the process of dictionary cleaning
 - reducing the number of entries to check
 - minimizing word ambiguity decisions
 - Morphological analysis and generation are required to map from inflected LA to inflected MSA
- **The Simple-Modification dictionary (SM)**
 - Minimal modification to LA inflected forms to look more MSA-like
 - Form modification: (أغنيا >gnyA 'rich pl.') is mapped to (أغنياء >gnyA')
 - Morphology modification: (أشرب b\$rb 'I drink') is mapped to (أشرب >\$rb)
 - Full translation: (كمان kmAn 'also') is mapped to (ايضا AyDAF)

90

(Maamouri et al. 2006)

Dialectal Morphological Analysis

- **MAGEAD** (Habash and Rambow 2006)
 - Morphological Analysis and GEneration for Arabic and its Dialects
- **Levels of Morphological Representation**
 - Lexeme Level
 - Aizdaha_r₁ PER:3 GEN:f NUM:sg ASPECT:perf
 - Morpheme Level
 - [zhr,1tV2V3,iaa] +at
 - Surface Level
 - Phonology: /izdaharat/
 - Orthography: Aizdaharat (إِزْدَهَرَات)

91

The Lexeme

- Lexeme is an abstraction of all inflectional variants of a word
 - ... كتابان الكتابين كتبهم للكتب كُتِبَ كِتَابٌ |مُتَابِ|
- For us, lexeme is formally a triple
 - Root or NTWS
 - Morphological behavior class (MBC)
 - {بيت بيوت} 'house' vs. {بيت ابیات} 'verse'
 - Meaning index
 - |قاعدة قواعد1| : 'rule'
 - |قاعدة قواعد2| : 'military base'

92

Morphological Behavior Class

- MBC::Verb-I-au (*katab/yaktub*)

cnj=wa	→	wa+
tense=fut	→	sa+
per=1 + num=sg	→	'+
per=1 + num=pl	→	n+
mood=indic	→	+u
mood=sub	→	+a
aspect=imper	→	V12V3
aspect=perf	→	1V2V3
voice=act	→	a-u
voice=pass	→	u-a
obj=3FS	→	hA
obj=1P	→	nA

...

93

Morphological Behavior Class

- MBC::Verb-I-au (*katab/yaktub*)

cnj=wa	→	wa+
tense=fut	→	sa+
per=1 + num=sg	→	'+
per=1 + num=pl	→	n+
mood=indic	→	+u
mood=sub	→	+a
aspect=imper	→	V12V3
aspect=perf	→	1V2V3
voice=act	→	a-u
voice=pass	→	u-a
obj=3FS	→	hA
obj=1P	→	nA

...

وَسَنَكْتُبُهَا

wasanaktubuhA

We will write it

94

MSA EGY

Morphological Behavior Class

- MBC::Verb-I-au (*katab/yaktub*)

cnj=wa	→	wa+ wi+	
tense=fut	→	sa+ Ha+	وَسَنَكْتُبُهَا
per=1 + num=sg	→	'+	
per=1 + num=pl	→	n+ n+	wasanaktubuhA
mood=indic	→	+u +0	wiHaniktibhA
mood=sub	→	+a	
aspect=imper	→	V12V3 V12V3	
aspect=perf	→	1V2V3	وَحَنَكْتُبُهَا
voice=act	→	a-u i-i	
voice=pass	→	u-a	
obj=3FS	→	hA hA	We will write it
obj=1P	→	nA	

...

⁹⁵
MSA EGY

Morphological Behavior Class

- MBC::Verb-I-au (*katab/yaktub*)

cnj=wa	→	wa+ wi+	→ [CONJ:wa]
tense=fut	→	sa+ Ha+	→ [PART:FUT]
per=1 + num=sg	→	'+	
per=1 + num=pl	→	n+ n+	→ [SUBJ_PRE_1P]
mood=indic	→	+u +0	→ [SUBJ_SUF_Ind]
mood=sub	→	+a	
aspect=imper	→	V12V3 V12V3	→ [PAT:I-IMP]
aspect=perf	→	1V2V3	
voice=act	→	a-u i-i	→ [VOC:Iau-ACT]
voice=pass	→	u-a	
obj=3FS	→	hA hA	→ [OBJ:3FS]
obj=1P	→	nA	

...

⁹⁶
MSA EGY

Morphological Behavior Class

- **MBC::Verb-I-au** (*katab/yaktub*)
 - cnj=wa → [CONJ:wa]
 - tense=fut → [PART:FUT]

 - per=1 + num=pl → [SUBJ_PRE_1P]
 - mood=indic → [SUBJ_SUF_Ind]

 - aspect=imper → [PAT:I-IMP]

 - voice=act → [VOC:Iau-ACT]

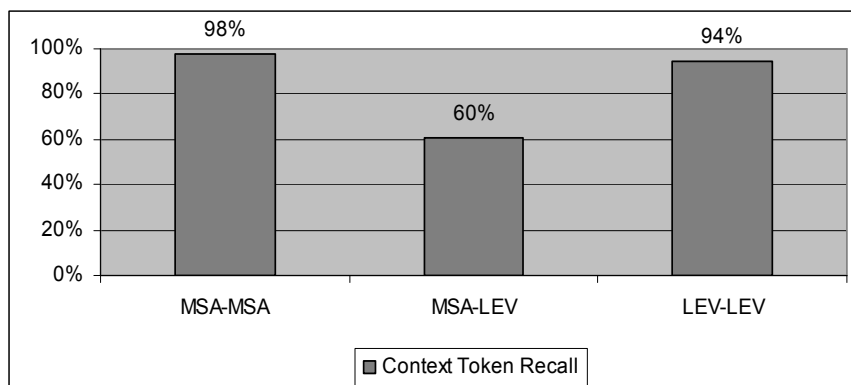
 - obj=3FS → [OBJ:3FS]

...

97

Levantine Evaluation

- Results on Levantine Treebank



98

Arabic Dialect POS Tagging

- Duh and Kirchhoff 2005; Duh and Kirchhoff 2006
 - Egyptian Arabic and Levantine Arabic
 - Minimal supervision
 - dialectal text
 - and MSA morphological analyzer
 - Cross-dialect sharing techniques
- Rambow et al. 2005
 - Levantine Arabic
 - LEV-MSA transduction using LEV-MSA lexicon
 - MSA POS Tagging
 - Projection of MSA tags onto LEV

99

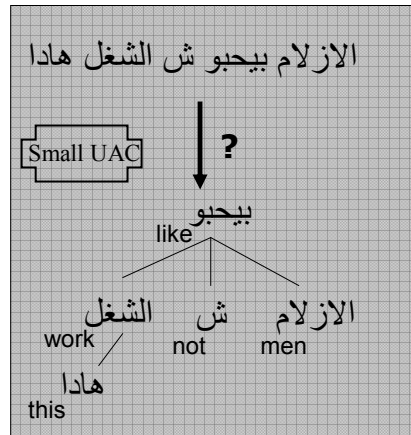
Arabic Dialect Parsing

- Possible Approaches
 - Annotate corpora ("Brill Approach")
 - Too expensive
 - Leverage existing MSA resources
 - Difference MSA/dialect not enormous
 - Linguistic studies of dialects exist
 - Too many dialects: even with dialects annotated, still need leveraging for other dialects

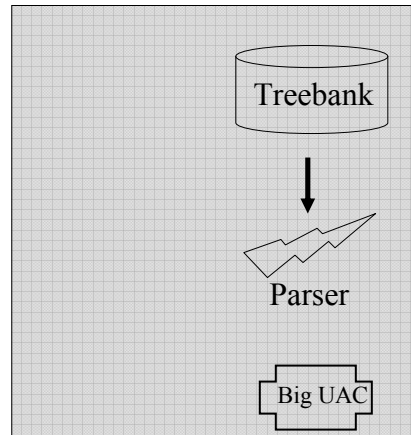
100

Parsing Arabic Dialects: The Problem

- Dialect -



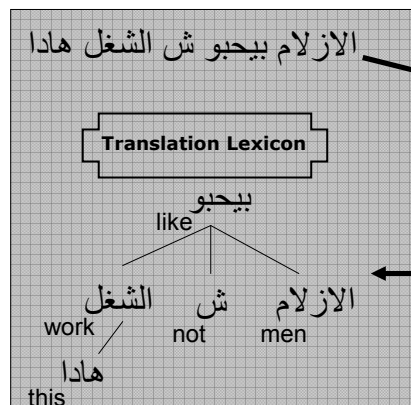
- MSA -



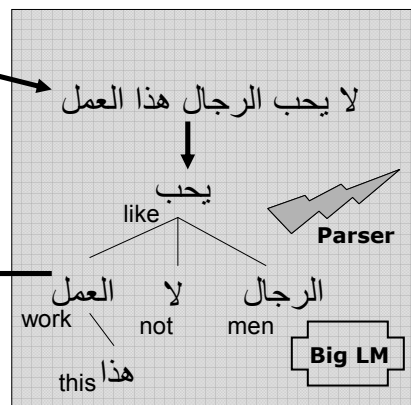
101

Sentence Transduction Approach

- Dialect -



- MSA -

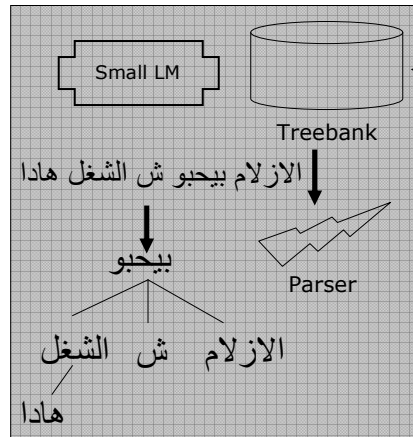


102

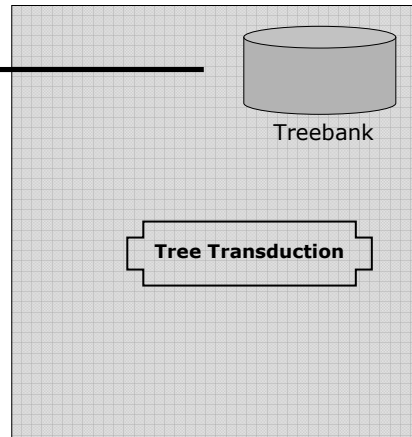
(Rambow et al. 2005; Chiang et al. 2006)

MSA Treebank Transduction

- Dialect -



- MSA -

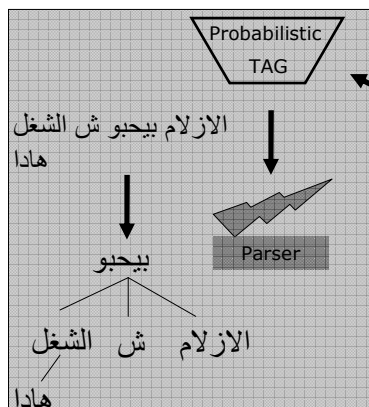


103

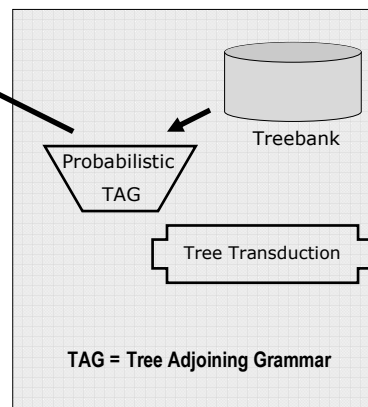
(Rambow et al. 2005; Chiang et al. 2006)

Grammar Transduction

- Dialect -



- MSA -



104

(Rambow et al. 2005; Chiang et al. 2006)

Dialect Parsing Results

Absolute/Relative F-1 improvement

	No Tags	Gold Tags
Sentence Transduction	4.2/9.0%	3.8/9.5%
Treebank Transduction	3.5/7.5%	1.9/4.8%
Grammar Transduction	6.7/14.4%	6.9/17.3%

Dialect-MSA dictionary was the biggest contributor to improved parsing accuracy: more than a 10% reduction on F1 labeled constituent error

105

(Rambow et al. 2005; Chiang et al. 2006)

Arabic Dialect Machine Translation

- **Problems**
 - Limited resources
 - Non-standard Orthography
 - Morphological complexity
- **Solutions**
 - Rule-based segmentation (Riesa et al. 2006)
 - Minimally supervised segmentation (Riesa and Yarowsky 2006)
 - Spelling normalization (Riesa et al. 2006)
 - Leveraging MSA resources (Riesa et al. 2006, Zollman et al. 2006, Rambow et al. 2005)
 - Dialect-MSA lexicons (Rambow et al. 2005, Chiang et al. 2006, Maamouri et al. 2006)

106

Arabic Dialect Machine Translation

- **TransTac: DARPA Program on Translation System for Tactical Use**
 - Iraqi \leftrightarrow English speech-to-speech MT
 - Phraselator: <http://www.phraselator.com/>
- **MT as a component**
 - JHU Workshop on Parsing Arabic dialect (Rambow et al. 2005, Chiang et al. 2006)

107

Dialect 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)

108

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

109

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)

110

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](#)

111

Resources

Morphology

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

Dialect Resources

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

112

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)
- [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)

113

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 \(COLING/ACL 1998\)](#)

114

References

Books

- Bateson, Mary Catherine. Arabic Language Handbook. Georgetown University Press. 2003.
- Brustad, Kristen E. The Syntax of Spoken Arabic: A Comparative Study of Moroccan, Egyptian, Syrian, and Kuwaiti Dialects. Georgetown University Press. 2000.
- Holes, Clive. Modern Arabic: Structures, Functions, and Varieties. Georgetown University Press. 2004.

Conference Papers and Journal Articles

- Alexandrescu, A. and K. Kirchhoff. 2006. Factored Neural Language Models. HLT.
- 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.

115

References

- Diab, M., K. Hacioglu and D. Jurafsky. 2004. Automatic Tagging of Arabic Text: From raw text to Base Phrase Chunks. HLT-NAACL.
- Diab, M., K. Hacioglu and D. Jurafsky. "Automated Methods for Processing Arabic Text: From Tokenization to Base Phrase Chunking." Book Chapter. In Arabic Computational Morphology: Knowledge-based and Empirical Methods. Editors A. van den Bosch and A. Soudi.
- Duh, K. and K. Kirchhoff. 2005. POS Tagging of Dialectal Arabic: A Minimally Supervised Approach. ACL Workshop on Semitic Languages.
- Duh, K. and K. Kirchhoff. 2006. Lexicon Acquisition for Dialectal Arabic Using Transductive Learning. EMNLP.
- Elming, J. and N. Habash. 2007. Combination of Statistical Word Alignments Based on Multiple Preprocessing Schemes. NAACL.
- Fischer, W. 2001. A Grammar of Classical Arabic. Yale Language Series. Yale University Press. Translated by Jonathan Rodgers.
- Habash, N. Issues in Palestinian Arabic Spelling Standardization. NACAL 27, 1999. Baltimore, MD.
- 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.
- Habash, N. and O. Rambow. 2006. MAGEAD: A Morphological Analyzer and Generator for the Arabic Dialects. ACL.
- Habash, N. and O. Rambow. 2007. Arabic Diacritization through Full Morphological Tagging. NAACL.
- 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.

116

References

- 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., A. Soudi and T. Buckwalter. 2006. "On Arabic Transliteration." Book Chapter. In Arabic Computational Morphology: Knowledge-based and Empirical Methods. Editors A. van den Bosch and A. Soudi.
- Habash, N. Arabic and its Dialects. Multilingual Magazine. #81, July/August 2006.
- 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.
- Habash, N., B. Dorr and C. Monz. 2006. Challenges in Building an Arabic Generation-heavy Machine Translation System and Extending it with Statistical Components. AMTA.
- Hwa, R., C. Nichols and K. Sima'an. 2006. Corpus Variations for Translation Lexicon Induction. AMTA.
- Ittycheriah, A. and S. Roukos. 2005. A maximum entropy word aligner for Arabic-English machine translation. HLT and EMNLP.
- 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 Workshop. ICASSP.
- Kirchhoff, K. and D. Vergyri. 2004. Cross-dialectal acoustic data sharing for Arabic speech recognition. ICASSP.
- Lee, Y., K. Papineni, S. Roukos, O. Emam and H. Hassan. 2003. Language Model Based Arabic Word Segmentation. ACL.

References

- 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.
- Maamouri, M., T. Buckwalter, and C. Cieri. Dialectal Arabic Telephone Speech Corpus: Principles, Tool Design, and Transcription Conventions. NEMLAR 2004.
- Maamouri, M., D. Graff, H. Jin, C. Cieri, and T. Buckwalter. Dialectal Arabic Orthography-based Transcription and CTS Levantine Arabic Collection. EARS RT-04 Workshop.
- 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.
- Riesa, J. and D. Yarowsky. Minimally Supervised Morphological Segmentation with Applications to Machine Translation. AMTA06.
- Riesa, J., B. Mohit, K. Knight and D. Marcu. Building an English-Iraqi Arabic Machine Translation System for Spoken Utterances with Limited Resources. Interspeech 2006.
- Rogati, M., S. McCarley, and Y. Yang. 2003. Unsupervised Learning of Arabic Stemming Using a Parallel Corpus. ACL.
- Sadat, Fatiha and Nizar Habash. 2006. Combination of Preprocessing Schemes for Statistical MT. ACL.
- Smith N., D. Smith, and R. Tromble. 2005. Context-Based Morphological Disambiguation with Random Fields. HLT-EMNLP.
- Smrz, O. and P. Zemánek. 2002. Sherds from an Arabic Treebanking Mosaic. Prague Bulletin of Mathematical Linguistics, (78).
- Snider, N. and M. Diab. 2006. Automatic Discovery of Verb Classes in Modern Standard Arabic. ACL.

References

- Snider, N. and M. Diab. 2006. Unsupervised Induction of Arabic Verb Classes. NAACL.
- Soudi, A., V. Cavalli-Sforza, and A. Jamari. 2001. A Computational Lexeme-Based Treatment of Arabic Morphology. ACL workshop on Arabic Natural Language Processing.
- Vergyri, D., K. Kirchhoff, K. Duh and A. Stolcke. 2004. Morphology-based language modeling for Arabic speech recognition. ICSLP.
- Vergyri, D. and K. Kirchhoff. 2004. Automatic Diacritization of Arabic for Acoustic Modeling in Speech Recognition. COLING Workshop on Arabic-script Based Languages.
- Xu J. 2002. UN Parallel Text (Arabic-English), LDC Catalog No.: LDC2002E15.
- Yang, M. and K. Kirchhoff. Phrase-Based Backoff Models for Machine Translation of Highly Inflected Languages. EACL.
- Žabokrtský, Z. and O. Smrž. 2003. Arabic Syntactic Trees: from Constituency to Dependency. EACL.
- Zitouni, I., J. Sorensen, and R. Sarikaya. 2006. Maximum entropy based restoration of Arabic diacritics. ICC3 and ACL.
- Zitouni, I., J. Olive, D. Iskra, K. Choukri, O. Emam, O. Gedge, M. Maragoudakis, H. Tropf, A. Moreno, A. Rodriguez, B. Heuft and R. Siemund. 2002. OrienTel: Speech-Based Interactive Communication Applications for the Mediterranean and the Middle East. ICSLP.
- Zollmann, A., A. Venugopal and S. Vogel. 2006. Bridging the Inflection Morphology Gap for Arabic Statistical Machine Translation. NAACL.

119

References

Conference/Institution/Program 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
ICASSP	= International Conference on Acoustics, Speech and Signal 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
NACAL	= North America Conference on Afro-asiatic Languages
EARS	= DARPA Program (Efficient, Affordable, Reusable Speech-to-Text)

120