# Dialectal Arabic Telephone Speech Corpus: Principles, Tool design, and Transcription Conventions

**Mohamed Maamouri, Tim Buckwalter, Christopher Cieri**

**Linguistic Data Consortium**

**University of Pennsylvania**

**maamouri@ldc.upenn.edu, timbuck2@ldc.upenn.edu, ccieri@ldc.upenn.edu**

# PRESENTATION OUTLINE

❖ ARABIC LINGUISTIC BACKGROUND

❖ ARABIC DIALECTAL SPEECH: METHODOLOGICAL TRANSCRIPTION PRINCIPLES AND TECHNOLOGICAL GOALS OF THE PROJECT

❖ *AMADAT*: LDC'S ARABIC MULTI-DIALECTAL TRANSCRIPTION TOOL

❖ METALANGUAGE:  RT-04 ARABIC TELEPHONE SPEECH TRANSCRIPTION CONVENTIONS

❖ BRIEF OVERVIEW OF LEVANTINE ARABIC TRANSCRIPTION GUIDELINES

OUR FOCUS WILL BE ON THE ARABIC DIALECTAL TRANSCRIPTION RATIONALE, THE TECHNOLOGICAL GOALS OF THE PROJECT, THE ANNOTATION TOOL STRUCTURE AND THE LEVANTINE CONVERSATIONAL ARABIC TRANSCRIPTION GUIDELINES
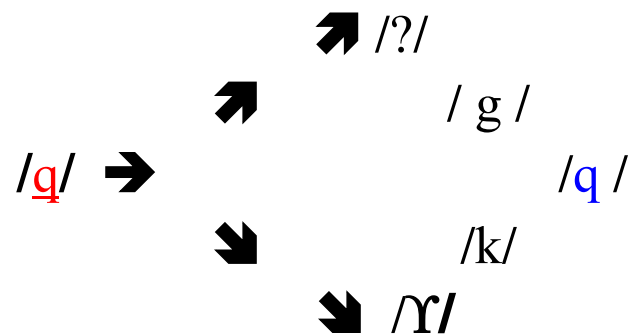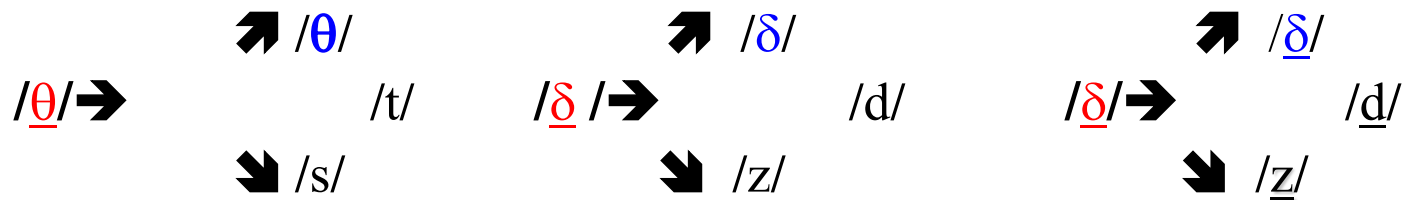
# ARABIC LINGUISTIC BACKGROUND

❖ **"ARABIC LANGUAGE CONTINUUM" WITH ARABIC DIGLOSSIA**
*FUSHA* = Modern Standard Arabic (=MSA) + ARABIC DIALECTS
**+ INTRALINGUAL CODESWITCHING & CODE-MIXING**

❖ **SIGNIFICANT LINGUISTIC DISTANCE BETWEEN MSA & DIALECTS**

❖ **SIGNIFICANT INTER- LINGUISTIC VARIATION AMONG DIALECTS**

❖ **SIGNIFICANT INTRA- LINGUISTIC VARIATION WITHIN DIALECTS**

❖ **IMPORTANT COMMON CORE OF MUTUAL INTELLIGIBILITY**

→ **HIGH LEVEL OF FORM AND STRUCTURE SIMILARITY**
→ **COMMON LEXICAL CORE WITH SIGNIFICANT SEMANTIC DIFFERENTIATION**

# ARABIC LANGUAGE BACKGROUND

❖ **EXISTENCE OF  LIVING MSA WRITING AND READING COMMUNITY**

❖ **INTERNALIZED KNOWLEDGE OF MSA BY EDUCATED AND SEMI-LITERATE NATIVE ARABIC SPEAKERS**

❖ **EXISTENCE OF UNDERLYING MSA COGNATE STRUCTURES**

❖ **USE OF MSA-BASED "ACCOMMODATION FILTERS"**

❖ **DOMINANCE OF MSA-BASED GRAPHEMIC TRADITIONS AND EVIDENCE OF MSA-BASED GRAPHEMIC INTERFERENCE**

❖ **EXISTENCE OF STANDARD MSA-BASED GRAPHEMIC KNOWLEDGE**

→ **PRODUCTIVE BASE FOR CONVERSATIONAL DIALECTAL ARABIC SPEECH-TO-TEXT TRANSCRIPTION SKILLS**

# DIALECTAL ARABIC SOUND CHANGE

## DIALECTAL SOUND CHANGE PATTERNS

$\nearrow$ /θ/ $\qquad$ $\nearrow$ /δ/ $\qquad$ $\nearrow$ /δ̠/

/θ̠/ → /t/ $\qquad$ /δ̠/ → /d/ $\qquad$ /δ̠/ → /d̠/

$\searrow$ /s/ $\qquad$ $\searrow$ /z/ $\qquad$ $\searrow$ /z̠/

---

$\nearrow$ /ʔ/

$\nearrow$ / g /

/q̠/ → $\qquad$ /q /

$\searrow$ /k/

$\searrow$ /ɣ/

# ARABIC DIALECTAL VARIATION

In Egyptian Arabic,MSA /θ/ becomes both /t/ and /s/ while /g / is used to replace / j / and /?/ to replace /q/.  In Sudanese Arabic, MSA /q/ is pronounced /g / and [ ϒ ] while the same <u>phoneme/letter</u> is pronounced /q/, /g/, /?/,and /k/ in Levantine Arabic.

**Example:**  Iraqi.q.h.C.wav

EXISTENCE AND USE OF ARABIC SCRIPT "ARCHIGRAPHEMES"

# LEVANTINE ARABIC EXAMPLE

**Q:** شو ال**ق**صة؟

$w Al**q**Sp?

"What's the story?"

**A/T1:** يا زلمي **ك**لتلك مو**ك**وف مش معت**ك**ل وما في **ك**صّة

yA zlmy **k**ltlk mw**k**wf m$ mEt**k**l wmA fy **k**S~p

**A/T2:** يا زلمي **ق**لتلك مو**ق**وف مش معت**ق**ل وما فيه **ق**صّة

yA zlmy **q**lt**l**k mw**q**wf m$ mEt**q**l wmA fy **q**S~p

"Hey 'dude' I told you arrested not indicted and there
is no story"

❖ **Need to distinguish the transcription approach from the alphabet used.**

- ◆ **Transcription approaches: phonic, orthographic, hybrid**
- ◆ **Alphabets: Arabic, Roman, International Phonetic Alphabet**
- ◆ **One may perform either phonic or orthographic transcription using either Roman or Arabic alphabets**

❖ **Problems with standard approaches**

- ◆ **Alphabets**
  - ▪ IPA is hard to learn
  - ▪ Roman script looks and feels unnatural to Arabic speakers
  - ▪ Few computer systems fully implement Arabic script and bi-directional input.
- ◆ **Transcription Approaches**
  - ▪ MSA lacks conventions for many Levantine forms, does fully not address needs of acoustic modeling
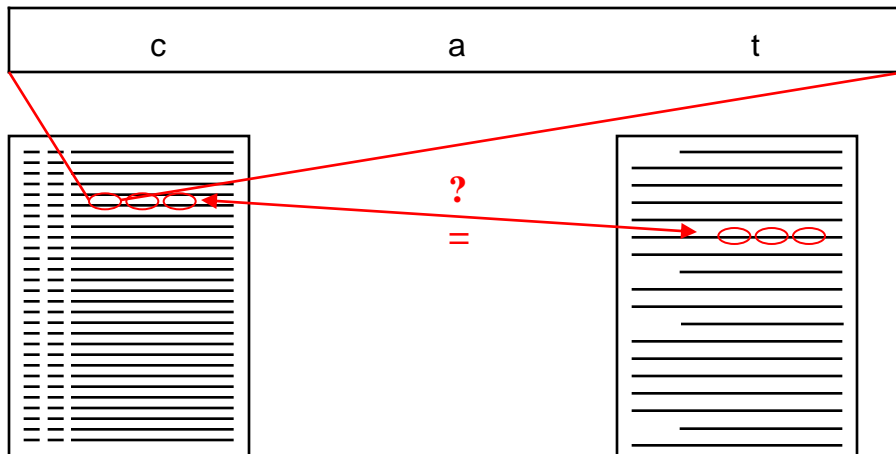  - ▪ purely phonic approach hinders language modeling

# Speech Recognition

❖ **Original Speech**

❖ **Analysis of audio**

| α1 | α2 | α3 | α4 | α5 | α6 | α7 | α8 | A9 |
|----|----|----|----|----|----|----|----|----|
| β1 | β2 | β3 | β4 | β5 | β6 | β7 | β8 | B9 |
| γ1 | γ2 | γ3 | γ4 | γ5 | γ6 | γ7 | γ8 | γ9 |
| δ1 | δ2 | δ3 | δ4 | δ5 | δ6 | δ7 | δ8 | δ9 |

| k | e | k |
|---|---|---|
| p | i | t |
| t | a | p |

❖ **Analysis suggests multiple phonetic interpretations.**

| c | a | t |
|---|---|---|

? =

❖ **Which need to be mapped onto a surface representation**

❖ **Sequences of which are compared against existing text to determine probable accuracy. Off-domain written text often substitutes for rare on-domain transcripts of spoken language.**

# LDC CONVERSATIONAL DIALECTAL ARABIC STT RATIONALE

" **How can we harness the native speaker's knowledge of Arabic orthography conventions and of the MSA linguistic common core to complete a quick, easy, and low-cost Speech-to-Text transcription of <u>Conversational Dialectal Arabic</u> ?"**

## OBJECTIVES OF SPEECH-TO-TEXT TRANSCRIPTION

- ❖ **FRIENDLY TO WRITERS AND READERS: EASY TO LEARN TO WRITE AND READ**
- ❖ **LEXICALLY CONSISTENT: A GIVEN UTTERANCE WILL ALWAYS BE SPELLED THE SAME**
- ❖ **LEXICALLY DISTINCTIVE: DIFFERENT UTTERANCES WILL ALWAYS BE SPELLED DIFFERENTLY**
- ❖ **ACOUSTICALLY CONSISTENT: TRANSCRIPTION/SPELLING PREDICTS PRONUNCIATION**

# CONVERSATIONAL DIALECTAL ARABIC TRANSCRIPTION CHALLENGES

## MSA-BASED/ARABIC ORTHOGRAPHIC SCRIPT-BASED TRANSCRIPTION

## 3 MAJOR CHALLENGES

- ❖ **RARE EVIDENCE OF CONVERSATIONAL DIALECTAL ARABIC TEXT CORPUS** **WITH STABLE MSA-BASED WRITING CONVENTIONS (POETRY, DRAMA, EPISTOLARY, POLITICAL SPEECHES, WEB & INTERNET CHATROOMS)**

- ❖ **DANGER OF INCONSISTENT CONVERSATIONAL DIALECTAL ARABIC MSA-BASED TRANSCRIPTION PRACTICES**

- ❖ **NATIVE LANGUAGE REPRESENTATION: DANGER OF OVER INTERFERENCE OF MSA WRITING CONVENTIONS IN EXISTING CONVERSATIONAL DIALECTAL ARABIC TRANSCRIPTION PRACTICES**

# CONVERSATIONAL DIALECTAL ARABIC STT TRANSCRIPTION OBJECTIVE

**OBJECTIVE:** APPROPRIATE BALANCE BETWEEN THE TWO TENDENCIES BELOW IN ORDER TO AVOID NEGATIVE CONSEQUENCES TO THE SPECIFIC NEEDS OF THE STT SCIENTIFIC RESEARCH COMMUNITY

- **Neither too strict an adherence to the use of MSA-based spelling conventions to reconvert dialectal forms to an unnecessary MSA-representation →** WITH HIGHER RECONSTRUCTION RATE OF 'UNDERLYING' FORMS

- **Nor too cloose an adherence to finer sound /(allo)phonic/ acoustical utterance representation →** LEADING TO AN OUTPUT WITH FINER ACOUSTICAL REPRESENTATION BUT WITH LOWER RATE OF SEMANTIC WORD RECOGNITION

# "*AMADAT*" DESIGN SPECIFICATIONS

- ❖ **ARABIC MULTI-DIALECTAL TRANSCRIPTION AND ANNOTATION TOOL**

- ❖ **TWO TIERS OF TRANSCRIPTION / ANNOTATION**

- ❖ **MODERN STANDARD ARABIC-BASED TRANSCRIPTION (MSAT: '*ORTHOGRAPHIC LEVEL*')**

- ❖ **ARABIC ORTHOGRAPHIC SYSTEM-BASED TRANSLITERATION (AOST: '*SURFACE PHONEMIC LEVEL*' )**

- ❖ **THREE MUTUALLY EXCLUSIVE OPERATION MODES**

# '*AMADAT*' STT TRANSCRIPTION MODES

**MSAT MODE:** QUICK TRANSCRIPTION→'GREEN AREA'

- USE OF NORMAL ARABIC KEYBOARD FOR TRANSCRIPTION
- FIRST PASS WITH MSA-BASED APPLICABLE CONVENTIONS
- METALANGUAGE ANNOTATION (CTS RT-04 ANNOTATION) OBJECTIVE: OPTIMIZED OUTPUT FOR LANGUAGE MODELING

**AOST MODE:** CAREFUL TRANSCRIPTION → 'YELLOW AREA'

- USE OF LATIN KEYBOARD FOR TRANSLITERATION
- USE OF MODIFIED TIM BUCKWALTER CODE WITH SOUND VALUES
- OBJECTIVE: OPTIMIZED OUTPUT FOR ACOUSTIC MODELING

**EDIT MODE:** ANNOTATION CORRECTION → 'RED AREA'

- USE OF LATIN KEYBOARD FOR A TOKEN-BY-TOKEN EDITING
- ACCESS ONLY TO ANNOTATION MANAGEMENT AND QUALITY CONTROL

Annotation File   fsa_10161.txt          Speaker IDs (A/B)

| Begin | End | Track | Transcription in Arabic | Index | Tr |
|---|---|---|---|---|---|
| 147.92 | 149.57 | A | (%أه) بتهم (ضحك) | 98 | (>1 |
| 148.85 | 150.01 | B | بالضبط (أه%) | 99 | (%— |
| 149.80 | 151.01 | A | (إنقطاع) أما (ضحك) | 100 | (Dł |
| 150.43 | 152.40 | B | (أه%) لأ إلا و بأحكيه اللي صحيح أما | 101 | >m |
| 152.04 | 153.42 | A | (إنقطاع)(أه%) صح لأ (%أه) | 102 | (% |
| 152.96 | 154.18 | B | (إنقطاع) معلوم معلوم (إنقطاع) | 103 | (<r |

| Prev | Next | Play | Stop | Bank | Drop | BadSeg |
|---|---|---|---|---|---|---|

Speech Comment

| F1: (breath) | F2: (cough) | F3: (laugh) | F4: (music) | F5: (noise) | F6: (peopletalk) | F7: (sneeze) | F8: (silence) | F9: (pause) |
|---|---|---|---|---|---|---|---|---|

| F10: (%ah) | F11: (%eh) | F12: (%um) | F13: (%ooh) | F14: (%hm) | F15: (noise/) | F16: (overlap) | F17: (overlap/) |
|---|---|---|---|---|---|---|---|

MSA Transcription

(أه%) لأ إلا و بأحكيه اللي صحيح أما

Selected Word                                                          Change Word

Annotation Remark                                                      Insert Word

| Cons Change | Velarized Cons | Voc Variant | Hamzah Drop | Diphthong | -h Deletion | Cons Deletion | -ap Silent | -ap Pronounced | Delete Word |
|---|---|---|---|---|---|---|---|---|---|

# 'MSAT' SPECIFICATIONS AND ISSUES

- ❖ MACHINE-READABLE UNVOCALIZED WRITTEN TEXT DATA
- ❖ NO DIACRITICS IN GENERAL. HOWEVER, USE OF SHADDAH AND INITIAL HAMZA NEED TO BE RE-DISCUSSED BY THE SCIENTIFIC COMMUNITY' USERS
- ❖ FOCUS ON CONSISTENT TRANSCRIPTION OF SAME FORMS
- ❖ FOCUS ON IDENTIFICATION OF SPECIFIC DIALECTAL FORMS (DEFINITIONAL NEEDS TO BE DISCUSSED)
- ❖ ANCHORING OF SOME DIALECTAL FORMS TO MSA-SIMILAR UTTERANCES AND AN 'UNDERLYING' MSA SEMANTIC STRUCTURE
    (DEFINITIONAL NEEDS TO BE DISCUSSED)
- ❖ CAUTIOUS/CONSERVATIVE USE OF RECONSTRUCTED 'UNDERLYING' FORMS: "NO REVERSE MSA ENGINEERING"

# python

File  Mode

Annotation File | fsa_10161.txt | Speaker IDs (A/B) |

| Begin | End | Track | Transcription in Arabic | Index | Tra |
|---|---|---|---|---|---|
| 147.92 | 145.57 | A | (أه%)(صحك) بهم (إصلاح) | 98 | (أ |
| 148.85 | 150.01 | B | بالضبط (أه%) | 99 | (% |
| 149.80 | 151.01 | A | (أ ... أ) (أ ... أ) | 100 | (Pl |

Prev | Next | Play | Stop | Bank | Drop | BadSeg

Speech Comment | F1: (breath) | F2: (cough) | F3: (laugh) | F4: (music) | F5: (noise) | F6: (peopletalk) | F7: (sneeze) | F8: (silence) | F9: (pause)

F10: (%ah) | F11: (%eh) | F12: (%um) | F13: (%ooh) | F14: (%hm) | F15: (noise/) | F16: (overlap) | F17: (overlap/)

Vowelized Arabic Trans

أمَّا صَحِيح اللّي بَحكِيّه و إلّا لأ (%أه)

Linguistic Transliteration

>am~aA SaHiyH Al~ily baHkiy~h wi <il~aA la> (%>h)

Selected Word | >am~aA

Change Word

Insert Word

Annotation Remark

Cons Change | Velarized Cons | Voc Variant | Hamzah Drop | Diphthong | -h Deletion | Cons Deletion | -ap Silent | -ap Pronounced | Delete Word

# 'AOST' SPECIFICATIONS AND ISSUES

❖ **FOCUS ON CLOSE ADHERENCE TO SOUND SPECIFICITIES**

❖ **FOCUS ON FULL FUNCTIONAL VOCALIZATION WITH SUKUN LIMITED TO SYLLABIC DIVISION WHEN NEEDED FOR PRONUNCIATION**

❖ **NO REPRESENTATION OF VOCALIC QUALITY VARIATION BUT LENGTHENING OF UNDERLYING DIPTHONGS**

❖ **INCLUSION OF RELEVANT SOUND FEATURES EXCEPT MORPHOPHONEMIC ASSIMILATION PHENOMENA (EXAMPLE: AL- ), AND EPENTHETIC AND JUNCTURE PHENOMENA**

❖ **USE OF PERSIAN LETTERS FOR CAREFUL TRANSCRIPTION OF UTTERANCES IN WHICH SOUNDS WHICH DO NOT EXIST IN THE ARABIC ORTHOGRAPHY OCCUR**

❖ **WHILE RECORDING AND ANNOTATING DIALECTAL SOUND FEATURES IN AOST, THE LINKED MSAT TOKENS AND QUICK TRANSCRIPTION BASELINE REMAIN UNCHANGED/STABLE**

Annotation File | fsa_10161.txt

Speaker IDs (A/B) |

| Begin | End | Track | Transcription in Arabic | Index | Tra |
|-------|-----|-------|-------------------------|-------|-----|
| 147.92 | 149.57 | A | (صحك) بهم (إصلاح) | 98 | (P |
| 148.85 | 150.01 | B | بالضبط (%أه) | 99 | (% |
| 149.80 | 151.01 | A | بالنار (أ ا.ب) | 100 | (Ph |

| Prev | Next | Play | Stop | Bank | Drop | BadSeg |
|------|------|------|------|------|------|--------|

Speech Comment

| F1: (breath) | F2: (cough) | F3: (laugh) | F4: (music) | F5: (noise) | F6: (peopletalk) | F7: (sneeze) | F8: (silence) | F9: (pause) |

| F10: (%ah) | F11: (%eh) | F12: (%um) | F13: (%ooh) | F14: (%hm) | F15: (noise/) | F16: (overlap) | F17: (overlap/) |

Vowelized Arabic Trans

أمّا صَحِيح اللّي بَحكِيّه و إلّا لأ (%أه)

Linguistic Transliteration

>am~aA SaHiyH Al~ily baHkiy~h wi <il~aA la> (%>h)

Selected Word |

Change Word

Insert Word

Annotation Remark

| Cons Change | Velarized Cons | Voc Variant | Hamzah Drop | Diphthong | -h Deletion | Cons Deletion | -ap Silent | -ap Pronounced | Delete Word |

# RT-04 CONVERSATIONAL ARABIC TRANSCRIPTION CONVENTIONS

## DISFLUENT SPEECH

- **FILLED PAUSES AND HESITATION SOUNDS**
- **PARTIAL WORDS AND RESTARTS**
- **CONTRACTED WORDS**
- **MISPRONOUNCED WORDS**
- **HARD-TO-UNDERSTAND SECTIONS**
- **BACKGROUND NOISES**
- **SPEAKER-PRODUCED NOISES**

## LINGUISTIC MARKUP

- **LINGUISTIC CHANGE FEATURES**
- **SOCIO-LINGUISTIC VARIATION FEATURES**
- **FOREIGN WORDS**

## MSA-based orthography

*"whenever possible, follow the spelling conventions and word segmentation of MSA."* <span style="color:red">*Like this:*</span>

قلت لك    `/?ultil:ak/`

مضبوط    `/mazbu:T/`

مثل    `/mitl/`

مثلا    `/masalan/`

*"whenever possible, follow the spelling conventions and word segmentation of MSA."* <span style="color:red">*Avoid this:*</span>

ألتلك   /?ultil:ak/

مزبوط   /mazbu:T/

متل   /mitl/

مسلا   /masalan/

## Exceptions

*"Note, however, the following exceptions…"*

**1** **list of high-frequency colloquial words**

**2** **conjugation paradigms of colloquial verbs**

**3** *nunation (-an -in -un)* **is transcribed if heard**

## Exception 1

*High-Frequency Colloquial Words (c. 120)*

| | | | | |
|---|---|---|---|---|
| علشان | زلمه | بعدين | إيد | إحنا |
| عم | زي | بكرة | أيش | اللي |
| فا | شو | بلكي | إيمتى | إمبيرح |
| فيه | شوي | بينات | أيوه | إنتي |
| فيش | شوية | جوا | برا | إنتوا |
| فين | عشان | دغري | بس | أنو |

# Exception 2

## *Colloquial Verbs Conjugation Paradigm*

| | | | | | | |
|---|---|---|---|---|---|---|
| هو | بيشوف | ما بيشوفش | بيجي | ما بيجيش | بيقرى | ما بيقراش |
| هي | بتشوف | ما بتشوفش | بتيجي | ما بتيجيش | بتقرى | ما بتقراش |
| هم | بيشوفوا | ما بيشوفوش | بيجوا | ما بيجوش | بيقروا | ما بيقروش |
| إنت | بتشوف | ما بتشوفش | بتيجي | ما بتيجيش | بتقرى | ما بتقراش |
| إنتي | بتشوفي | ما بتشوفيش | بتيجي | ما بتيجيش | بتقري | ما بتقريش |
| إنتوا | بتشوفوا | ما بتشوفوش | بتيجوا | ما بتيجوش | بتقروا | ما بتقروش |
| أنا | بشوف | ما بشوفش | بجي | ما بجيش | بقرى | ما بقراش |
| إحنا | بنشوف | ما بنشوفش | بنيجي | ما بنيجيش | بنقرى | ما بنقراش |

## Exception 2

*Colloquial Verbs Conjugation Paradigm*

| هو | شاف | ما شافش | إجى | ما جاش | قرى | ما قراش |
|---|---|---|---|---|---|---|
| هي | شافت | ما شافتش | إجت | ما جتش | قرت | ما قرتش |
| هم | شافوا | ما شافوش | إجوا | ما جوش | قروا | ما قروش |
| إنت | شفت | ما شفتش | جيت | ما جيتش | قريت | ما قريتش |
| إنتي | شفتي | ما شفتيش | جيتي | ما جيتيش | قريتي | ما قريتيش |
| إنتوا | شفتوا | ما شفتوش | جيتوا | ما جيتوش | قريتوا | ما قريتوش |
| أنا | شفت | ما شفتش | جيت | ما جيتش | قريت | ما قريتش |
| إحنا | شفنا | ما شفناش | جينا | ما جيناش | قرينا | ما قريناش |

## Exception 3

*Nunation (tanween) should reflect actual pronunciation*

مرحباً /marHaban/

مرحبا /marHaba/

أهلاً وسهلاً /?ahlan wa-sahlan/

أهلا وسهلا /?ahla wa-sahla/

**Issues:** *Choose the variant with the highest frequency of usage*

| | | | |
|---|---|---|---|
| 45 | أنا باحكي | 455 | أنا بحكي |
| 1,420 | أيوا | 2,530 | أيوه |
| 2,540 | برضه | 3,180 | برضو |

## Issues:

*Transcribe hamza when it is pronounced*

ممتاز يا ابو محمد    /mumta:z y-abu muHam:ad/

أهلين أبو طارق    /?ahle:n ?abu Ta:riq/

أنا والاولاد    /?ana wa-liwla:d/

الأب والأولاد    /?il-?ab wa-l-?awla:d/

# CONCLUSION:  Collection Update

**[September 19, 2004]**

- ❖ **13604 Recruits** (Domestic, International) / **11450 active callers**
- ❖ **2184 calls** completed
- ❖ **1662 are available as of today.**
- ❖ **1400 of them have more than 8 minutes** speech.
- ❖ **Male-Female ratio among the 2184 calls where the genders of both speakers are available :   M M  710   /   F F 300 /     M F 354       / F M 398**
    **Male to female ratio is: 1086 to 676 = 61.6% to 38.4%**

  **[ Note that when calls involve speakers with no gender information, those  calls are excluded from the calculations above].**
- ❖ **2305 speakers were used for the 2184 calls**. 1251 speakers only appeared in 1 call; 381 appeared in 2 calls; 488 appeared in 3 calls.

  **[1  times 1251; 2  times  381;  3  times  488;  4  times  117;  5  times 41]**

## 2 hrs EVALUATION SET/2 hrs DEVELOPMENT SET
## 68 hours + 32 hours TRAINING SET

**For more information, go to:**

http://www.ldc.upenn.edu/Projects/EARS/Arabic/Guidelines_Levantine_MSA.htm