

TTS Survey Report

Enhancing Arabic Speech Technology

with the world's most comprehensive Arabic lexicon

Abstract

The recent advances in deep learning and neural network technology have dramatically improved speech technology [5]. Although the quality of Arabic speech technology has been steadily improving, it still lags significantly behind the other major languages, such as Chinese and English. CJKI's **Comprehensive Arabic Full Form Lexicon**, or *ArabLEX* (1.2 billion entries) [1, 2], is the most comprehensive Arabic computational lexicon ever created. This reports analyzes the significant error rates in the leading TTS systems, and shows how *ArabLEX* can address these shortcomings by dramatically enhancing the quality of both TTS and ASR.

1. Orthographical ambiguity

One reason that Arabic speech technology lags behind is that the Arabic script is highly ambiguous. Words are often written as a string of consonants with no indication of vowels. For example, كاتِب can represent as many as *seven* pronunciations: *kāatib*, *kātibun*, *kātibin*, *kātaba*, *kātibi*, *kātiba* and *kātibu*. Many other characteristics of the Arabic script contribute to a high level of orthographic ambiguity, as described in Halpern's paper on Arabic named entities [4].

The morphological complexity of such cliticized forms as وَلِكَاتِبَاتِهِمَا *walikātibātihima*, and the absence of vowel diacritics, makes Arabic TTS especially challenging. That is, determining the morphological composition of such forms, and the correct vowels for such consonants as ت in وَلِكَاتِبَاتِهِمَا often requires morphological, semantic and contextual analysis which tax the capabilities of state-of-the-art speech technology.

2. Improving TTS accuracy

The extreme orthographic ambiguity of Arabic has led to unacceptably high error rates, even by the TTS systems offered by major players such as Google, Apple and Microsoft.

Our institute has conducted a survey to determine the scope of this problem, some of the results of which are reported in the Appendix. Surprisingly, we discovered that it is not unusual for over 50%, and even 80%, of the words in a sentence to be mispronounced, and that there is a trend for cliticized words to be incorrectly pronounced. For example, the cliticized word وَلِلْكَاتِبِينَ, correctly pronounced *walilkatibīna*, is mispronounced as *walilkatibáyna*.

As can be seen from the Appendix, the error rate of Arabic TTS is unacceptably high. Such a high error rate would be unthinkable in the other major world languages. Another issue is **prosody** (stress and intonation) and **vowel neutralization** (e.g. *na* is written as a long vowel in *أنا* but is shortened in actual pronunciation to *na*). This is a complex issue, described in detail in Halpern's paper on Arabic stress [3].

Speech synthesis, speech recognition and prosody in current Arabic speech technology are, on the whole, inaccurate, unnatural and often unpleasant to the ear. The time has come for developers to make serious efforts to make dramatic improvements.

3. Improving ASR accuracy

ArabLEX includes features specifically designed to support automatic speech recognition (ASR). For speech synthesis (TTS), it is only necessary to *generate* one accurate pronunciation. For example, كاتِبُونَ 'writers' in standard Arabic is pronounced *katibūna*, but for ASR it is also necessary to *recognize* the less formal variant pronunciation *katibūn*. Similarly, the standard pronunciation of أَكْتُبُ 'I write' is *'áktubu*, but the final vowel is often omitted and it is pronounced *'áktub*.

The above alternatives are on a *phonemic* level. That is, the phoneme /na/ is being replaced by the phoneme /n/ as a result of vowel omission. There are also variations on the *phonetic* level; that is, certain phonemes have regional allophones. For example, ج in such words as جَمَل *jamal* is pronounced [gɛ̃mɛl] in Egypt, [dʒɛ̃mɛl] in the Gulf region, and [ʒɛ̃mɛl] in the Levant. It is important to note that this does not refer to the local dialects in those regions, but to a *regional varieties* of Modern Standard Arabic (MSA).

Thus *ArabLEX* not only represents ج in the standard IPA [dʒ] for TTS, it also lists the regional [ʒ] and [g] for ASR training. The goal is to enable the recognition of these allophones, but not to generate them.

4. Benefits to speech technology

One of the key components for training speech technology systems is the *pronunciation dictionary*. A major feature of *ArabLEX* is that it can serve as an extremely comprehensive pronunciation dictionary.

ArabLEX not only maps all unvocalized forms (including all cliticized forms) to their vocalized counterparts and to their lemmas, but also provides precise **phonemic transcriptions** (CARS system) [4] and **phonetic transcriptions** (IPA) that includes precise word stress and vowel neutralization for each entry. For example, in the IPA *wēlikĕːˈtibikumɛ(ˈ)*, the stressed syllable is indicated by (ˈ) (U+0C28), while (ː) (U+02D1) indicates that the final *ɛ* is neutralized vowel of optional half length. These features can help developers significantly enhance the quality of Arabic TTS, and can be used in training ASR systems to achieve higher recognition rates.

To summarize, *ArabLEX* can bring the following benefits to speech technology:

- Covers approximately **1.2 billion entries**, including millions of proper nouns.
- Covers *all combinations* of proclitics and enclitics for inflected wordforms (mostly verbs, nouns, adjectives and proper nouns).
- Tens of millions of orthographic variants for all wordforms.
- Provides an exhaustive list of alternative pronunciations of identical unvocalized strings to enable orthographical disambiguation (e.g. six alternatives for كاتباتك).
- Future versions will provide 'importance flags' to help determine the most likely alternative.
- Highly accurate phonemic transcriptions for all wordforms, including precise stress and vowel neutralization
- Phonetic transcriptions (IPA) indicate the correct allophonic variants in context as well as regional variants for ASR.

References

- [1] Halpern, Jack. [Arablex: Comprehensive Arabic Full Form Lexicon](#)
- [2] Halpern, Jack. [Very Large-Scale Lexical Resources to Enhance Chinese and Japanese Machine Translation](#)
- [3] Halpern, Jack. [Word Stress and Vowel Neutralization in Modern Standard Arabic](#)
- [4] Halpern, Jack. [CJKI Arabic Romanization System](#)
- [5] Wajdan Algihab, Imam Muhammad bin Saud Islamic University. [Arabic Speech Recognition with Deep Learning: A Review](#)

APPENDIX: TTS SURVEY RESULTS

Below are the results of a survey conducted by CJKI (our institute) to compare the TTS systems of Google, Apple (iPhone) and Microsoft (Bing), showing high error rates for all three. The **Unvocalized** field is the original Arabic text, the **Vocalized** field indicates the correct pronunciation, and the **CJKI** field shows the correct pronunciation in CARS phonemic transcription [4]. The CARS transcriptions in these **Google**, **iOS** and **Bing** columns indicate the actual pronunciation by the three TTS engines. Mispronunciations are indicated in red, and the error rate is given in the column headers.

Table 1 and 3 are based on text **composed** for this survey, while tables 2 and 4 use a sentence **extracted** from the web. (It is noteworthy that the error rate for the composed text is actually much lower than for the extracted text.) Tables 1 and 3 compare the results on a word-by-word basis, whereas tables 2 and 4 compare them on a sentence-by-sentence basis, showing the context. The fact that the error rate is sometimes over 80% is surprising and unacceptable to users.

Table 1: Mispronounced Words in Composed Text

Unvocalized	Vocalized	Google (13%)	iOS (31%)	Bing (25%)	CJKI (0%)
عدد	عَدَدَ	ɛádadu	ɛádada	ɛádada	ɛáddada
الكاتب	اَلْكَاتِبُ	lkātibu	lkātibi	lkātibu	lkātibu
ما	مَا	mā	mā	mā	mā
قال	قَالَ	qāla	qāla	qāla	qāla
إن	إِنَّ	ʾinna	ʾinna	ʾinna	ʾinna
هؤلاء	هُؤُلَاءِ	həʾulāʾi	həʾulāʾi	həʾulāʾi	həʾulāʾi
الحكام	اَلْحُكَّامُ	lhukkāmi	lhukkāmi	lhukkāmi	lhukkāma
يفعلونه	يَفْعَلُونَهُ	yafɛalūnahu	yafɛalūnahu	yafɛalūnahu	yafɛalūnahu
في	فِي	fī	fī	fī	fī
الخارج	اَلْخَارِجُ	lkhāriji	lkhārija	lkhāriji	lkhāriji

مثل	مِثْلٌ	míthli	míthli	míthli	míthla
الهجمات	أَلْهَجَمَاتِ	lhajamāti	lhajamāti	lhajamāti	lhajamāti
الإلكترونية	أَلْإِلِكْتُرُونِيَّةِ	l'ilikturuníyyati	l'ilikturuníyyati	l'ilikturuníyyati	l'ilikturuníyyati
ومطاردة	وَمُطَارَدَةٌ	wamuṭārādati	wamuṭārídati	wamuṭārídati	wamuṭārādati
المعارضين	أَلْمُعَارِضِينَ	Imueṣariḍīna	Imueṣariḍīna	Imueṣariḍīna	Imueṣariḍīna
اللاجئين	أَللَّاجِئِينَ	llajjīṭna	llajjīṭna	llajjīṭna	llajjīṭna
في	فِي	fī	fī	fī	fī
العواصم	أَلْعَوَاصِمِ	lɛawāṣimi	lɛawāṣimi	lɛawāṣimi	lɛawāṣimi
الغربية	أَلْغَرْبِيَّةِ	lgharbíyyati	lgharbíyyati	lgharbíyyati	lgharbíyyati
وللكاتبين	وَلِلْكَاتِبِينَ	walilkatībīna	walilkātībáyna	walilkātībáyna	walilkatībīna
من	مِنْ	mína	mína	mína	mína
الصحفيين	أَلصَّحَفِيِّينَ	ṣṣaḥafiyīna	ṣṣaḥafiyīna	ṣṣaḥafiyīna	ṣṣaḥafiyīna
العرب	أَلْعَرَبِ	lɛárabi	lɛárabi	lɛárabi	lɛárabi
صرح	صَرَّحَ	ṣárraḥa	ṣáraḥa	ṣáraḥa	ṣárraḥa
بأن	بِأَنَّ	bi'ánna	bi'ánna	bi'ánna	bi'ánna
عليهم	عَلَيْهِمْ	ɛaláyhim	ɛaláyhim	ɛaláyhim	ɛaláyhim
أن	أَنَّ	'an	'an	'an	'an
يكتبوا	يَكْتُبُوا	yaktúbuwu	yaktúbuwu	yaktúbuwu	yaktúbuwu
ما	مَا	mā	mā	mā	mā
تمليه	تُمْلِيهِ	tumalīhi	tamlīhi	tamlīhi	tumlīhi
عليهم	عَلَيْهِمْ	ɛalayhim	ɛalayhim	ɛalayhim	ɛalayhim
ضمايرهم	ضَمَائِرُهُمْ	ḍamā'íruhum	ḍamā'írihim	ḍamā'írihim	ḍamā'íruhum

Table 2: Mispronounced Words in Extracted Text

Unvocalized	Vocalized	Google (80%)	iOS (90%)	Bing (70%)	CJKI (0%)
الاخوات	الْأَخَوَاتُ	'alikhwātu	'al'akhawāti	'al'akhawātu	'al'akhawātu
المتزوجات	الْمُتَزَوِّجَاتُ	Imutazawwijātu	Imutazawwijāti	Imutazawwijāti	Imutazawwijātu
اللاتي	الَّلَاتِي	lti	llati	llāti	llāti
رزقن	رُزِقْنَ	rízzqin	rúzzqin	rúzzqin	ruzíqna
بابناء	بَابْنَاءَ	babinā'un	bibnā'i	bibnā'i	bi'abnā'a
فليكتبن	فَلْيَكْتُبْنَ	falayiktibna	falktíbnā	falktíbnā	falyaktúbna
اسمائهم	أَسْمَائُهُمْ	'ismā'ahum	smā'ihim	smā'ihim	'asmā'ahum
وسبب	وَسَبَبَ	wasábaba	wasábaba	wasábaba	wasábaba
التسميه	الَّتَسْمِيَةِ	lttasammīhu	lttasammīhi	lttasammīhi	ttasmíyati
رجاءا	رَجَاءًا	rajjā'an	rajā'	rajā'	rajā'an

Table 3: Mispronounced Sentences in Composed Text

TTS	Sentence	Error %
Unvocalized	عدد الكاتب ما قال إن هؤلاء الحكام يفعلونه في الخارج مثل الهجمات الإلكترونية ومطاردة المعارضين اللاجئين في العواصم الغربية. وللكاتبين من الصحفيين العرب صرح بأن عليهم أن يكتبوا ما تمليه عليهم ضمائرهم	-
Vocalized	عَدَدَ الْكَاتِبِ مَا قَالَ إِنَّ هَؤُلَاءِ الْحُكَّامَ يَفْعَلُونَهُ فِي الْخَارِجِ مِثْلَ الْهَجَمَاتِ الْإِلِكْتُرُونِيَّةِ وَمُطَارَدَةِ الْمُعَارِضِينَ الْلَّاجِئِينَ فِي الْعَوَاصِمِ الْغَرْبِيَّةِ. وَلِلْكَاتِبِينَ مِنَ الصَّحَفِيِّينَ الْعَرَبِ صَرَحَ بِأَنَّ عَلَيْهِمْ أَنْ يَكْتُبُوا مَا تُمْلِيهِ عَلَيْهِمْ ضَمَائِرُهُمْ	0%
CJKI	εáddada lkātibu ma qāla 'inna ha'ulā'i lhukkāma yafεalūnahu fī lkhārijī mīthla lhajamāti l'ilikturūniyyati wamuṭārādati lmuεarīdīna llajī'īna fī lεawāṣimi lgharbīyyati. walilkaṭibīna mīna ṣṣaḥafīyīna lεārabi ṣārraḥa bi'ánna	0%

	εalayhim 'an yaktúbuwu mā tumlīhi εalayhim ḍamā'iruhum	
Google	εádadu lkātibu mā qāla 'inna ha'ulā'i lhukkāmi yafεalūnahu fī lkhārijī mīthli lhajamāti l'ilikturuníyyati wamuṭārādati lmuεariḍīna llajī'īna fī lεawāšimi lgharbíyyati. walilkātibīna mīna šṣahafīyīna lεárabi šarraḥa bi'ánna εalayhim 'an yaktúbuwu mā tumalīhi εalayhim ḍamā'iruhum	13%
iOS	εádada lkātibi mā qāla 'inna ha'ulā'i lhukkāmi yafεalūnahu fī lkhārija mīthli lhajamāti l'ilikturuníyyati wamuṭārídati lmuεariḍīna llajī'īna fī lεawāšimi lgharbíyyati. walilkātibáyna mīna šṣahafīyīna lεárabi šáraḥa bi'ánna εalayhim 'an yaktúbuwu mā tamlīhi εalayhim ḍamā'irihim	31%
Bing	εádada lkātibu mā qāla 'inna ha'ulā'i lhukkāmi afealūnahu fī lkhārijī mīthli lhajamāti l'ilikturuníyyati wamuṭārídati lmuεariḍīna llajī'īna fī lεawāšimi lgharbíyyati. walilkātibáyna mīna šṣahafīyīna lεárabi šáraḥa bi'ánna εalayhim 'an yaktúbuwu mā tamlīhi εalayhim ḍamā'irihim	25%

Table 4: Mispronounced Sentences in Extracted Text

TTS	Sentence	Error %
Unvocalized	الاخوات المتزوجات اللاتي رزقن بابناء فليكتبن اسمائهم وسبب التسميه رجاء	-
Vocalized	الْأَخَوَاتُ الْمُتَزَوِّجَاتُ اللَّاتِي رُزِقْنَ بِأَبْنَاءَ فَلْيَكْتُبْنَ أَسْمَاءَهُمْ وَسَبَبُ التَّسْمِيَةِ رَجَاءٌ	0%
CJKI	'al'akhawātu lmutazawwijātu llāṭi ruzíqna bi'abnā'a falyaktúbna 'asmā'ahum wasábaba ttasmíyati rajā'an	0%
Google	'alikhwātu lmutazawwijātu ltī rízqin baḡinā'un falayiktíbna 'ismā'ahum wasábaba lttasammīhu rajjā'an	\80%
iOS	'al'akhawāti lmutazawwijāti llaṭi rúzqin bibnā'i falktíbna smā'ihim wasábaba lttasammīhi rajjā'	90%
Bing	'al'akhawātu lmutazawwijāti llāṭi rúzqin bibnā'i falktíbna	70%

