

Curras + Baladi: Towards a Levantine Corpus

Karim El Haff
Strasbourg University
France

Mustafa Jarrar
Birzeit University
Palestine

Tymaa Hammouda
Birzeit University
Palestine

Fadi Zaraket
American University of Beirut
Lebanon

Lexical Resources at Birzeit University

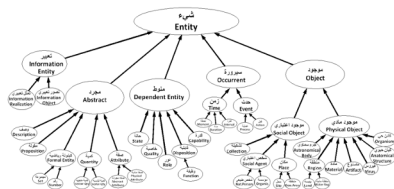
Lexicographic Database



150 lexicons
Largest Arabic
lexicographic
database

WSD 84%

Arabic Ontology/Wordnet



Formal Arabic
Wordnet
with ontologically
clean content

Dialect Corpora



Annotated corpora
each word is annotated
with many morph
features

NLP library



APIs
Linguistic Data, synonyms,
tools, Nested named-
entities, intents, ...

NER 88.4%

Big Linguistic Data Graph

<https://ontology.birzeit.edu>

Curras + Baladi: Towards a Levantine Corpus

Karim El Haff, Mustafa Jarrar, Tymaa Hammouda, Fadi Zaraket

University of Strasbourg, Birzeit University, American University of Beirut
karim.el-haff@etu.unistra.fr, mjarrar@birzeit.edu, 1171779@student.birzeit.edu, fz11@aub.edu.lb

Abstract

This paper presents two-fold contributions: a full revision of the Palestinian morphologically annotated corpus (Curras), and a newly annotated Lebanese corpus (Baladi). Both corpora can be used as a more general Levantine corpus. Baladi consists of around 9.6K morphologically annotated tokens. Each token was manually annotated with several morphological features and using LDC's SAMA lemmas and tags. The inter-annotator evaluation on most features illustrates 87% agreement using the Cohen's Kappa score. Curras was revised by refining all annotations for accuracy, normalization and unification of POS tags, and linking with SAMA lemmas. This revision was also important to ensure that both corpora are compatible and can help the bridge the nuanced linguistic gaps that exist between the two highly mutually intelligible dialects. The papers also presents the helped bridge the nuanced linguistic gaps that exist between the two highly mutually intelligible dialects. Differences and commonalities between both dialects are discussed in the paper. Both corpora are publicly available through a web portal.

Keywords: Arabic morphology, Annotated Corpus, Arabic Dialect, Levantine, Palestinian Arabic, Lebanese Arabic

1. Introduction

Arabic speakers use local dialects in day-to-day communication. Therefore, a question of diglossia arises; we are in a situation where the official language of the state differs from the language spoken in everyday life. We can distinguish several families of dialects: Moroccan, Egyptian, Sudanese, Levantine, Iraqi and Khaliji (Gulf). Arabic dialects tend to diverge from Modern Standard Arabic (MSA) in terms of phonetics, morphology, syntax and vocabulary. This creates a potential lack of inter-comprehension among speakers of different dialect families.

Arabic content was mainly written in MSA. Recently, dialectal content has been increasing massively, especially on social media. MSA is considered among the under-resourced languages by the NLP community (Darwish et al., 2021). Dialectal Arabic (DA) is even less resourced. The resource gap between MSA and the dialects implies a large margin of error when MSA tools are used against dialectal content (Zbib et al., 2012). Thus, it is instrumental to build resources and tools to identify dialects in context and to treat Arabic content based on its unique dialectal identity.

corpus annotations, we have also revised Curras annotations to ensure compatibility with the LDC's SAMA tags and lemmas (Maamouri et al., 2010).

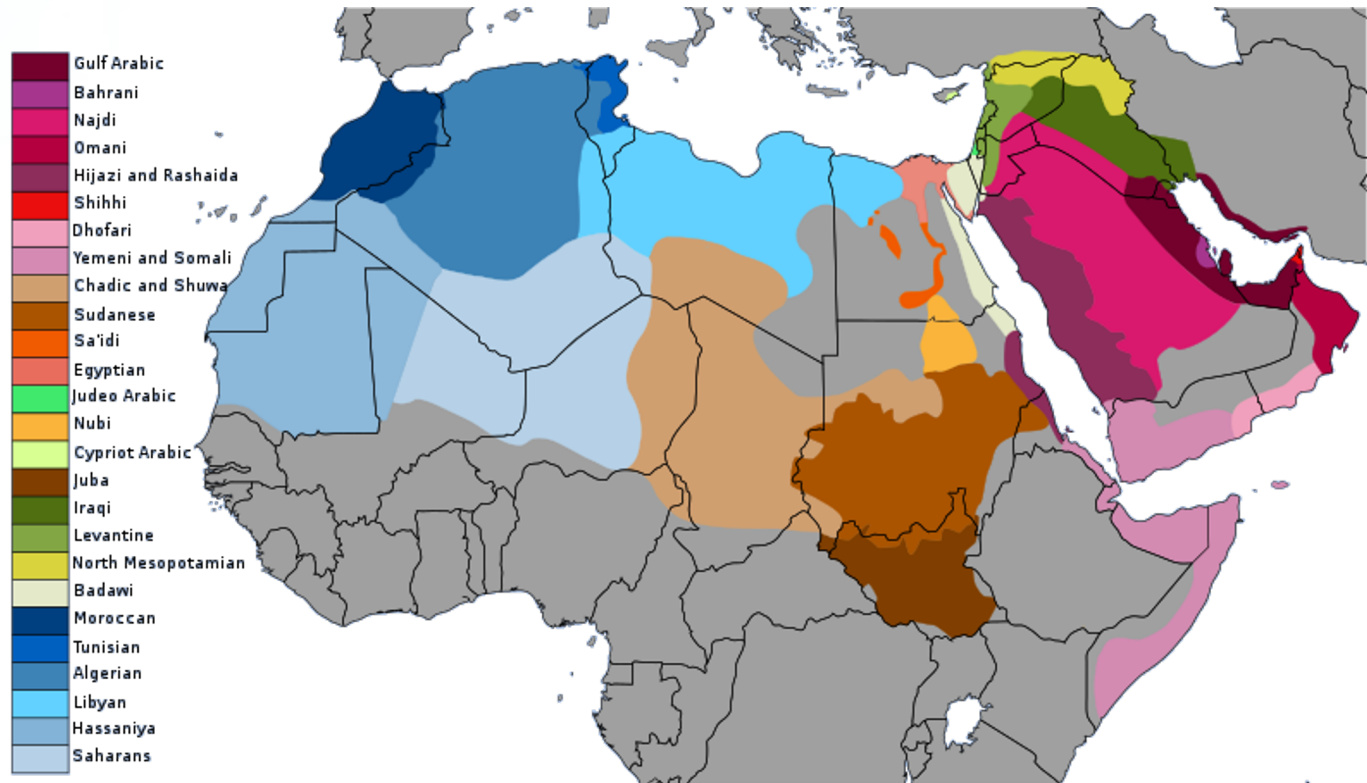
In this paper, we present two-fold contributions:

1. Baladi, a Lebanese morphologically annotated corpus, which consists of 9.6K tokens. Each token was manually annotated with prefixes, suffixes, stem, POS tags, MSA and DA lemmatization, English gloss, in addition to other features such as gender, number, aspect, and person. The corpus was annotated mainly using LDC's SAMA lemmas and tags. The inter-annotator evaluation on most features illustrates 87% agreement using the Cohen's Kappa score (McHugh, 2015).
2. Revision of Curras, by refining all annotations for accuracy, normalization and unification of POS tags, and linking with SAMA lemmas. This revision was also important to ensure that both corpora are compatible and can together form a more general Levantine corpus.

El Haff, K., Jarrar, M., Hammouda, T., Zaraket, F., (2022). [Curras + Baladi: Towards a Levantine Corpus](#). In Proceedings of the International Conference on Language Resources and Evaluation (LREC 2022), Marseille, France.

Arabic is a low resources language

- Classical Arabic
- Modern Standard Arabic (more resources)
- **Arabic Dialects**



Curras, a Palestinian
morphologically
annotated corpus (Revised)



Baladi, a Lebanese
morphologically
annotated corpus



Levantine Corpus

Publicly available:

<https://portal.sina.birzeit.edu/curras>

Baladi Lebanese Corpus

Corpus Collection:

From: Facebook posts, blog posts and traditional poems.

Size: 9.6K tokens, 424 sentences

Annotation Methodology

Tool: a *Smart* Google Sheet (**AnnoSheet**)

Four annotators, over 10 months

Features:

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
الموني	المونة	ال/DET	مون/NOUN	ة/NSUFF_FEM_SG	NOUN	-	-	f	s	مونة 1	مونة 1	provisions

Annotation Guidelines

DA Sentence: بروحي ساكنة ممنون عيونك عيوني وانا عم أخذ المونة رح ظل شوف الدنيا لوما بقيلي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
بروحي	بروحي	ب/PREP	روح/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	s	رُوح 1	رُوح 1	spirit
ساكنة	ساكنة		ساكن/ADJ	ة/NSUFF_FEM_SG	ADJ	-	-	f	s	ساكن 1	ساكن 1	residing
ممنون	ممنون		ممنون/ADJ		ADJ	-	-	m	s	مَمْنُون 1	مَمْنُون 1	grateful
عيونك	عيونك		عيون/NOUN	ك/POSS_PRON_2FS	NOUN	-	-	m	p	عَيْن 1	عَيْن 1	eye
عيوني	عيوني		عيون/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	p	عَيْن 1	عَيْن 1	eye
وانا	وانا	و/CONJ	انا/PRON_1S		PRON	1	-	m	s	أنا 1	أنا 1	I
عم	عم		عم/PROG_PART		PROG_PART	-	i	-	-	عم 0	ظَلَّ 1	continue
أخذ	أخذ	آ/IV1S	أخذ/IV		VERB	1	i	m	s	أَخَذُ 1	أَخَذُ 1	take
الموني	المونة	ال/DET	مون/NOUN	ة/NSUFF_FEM_SG	NOUN	-	-	f	s	مُونَة 1	مُونَة 1	provisions
رح	رح		رح/FUT_PART		FUT_PART	-	-	-	-	رح 0	سَوَّفَ 1	will
ضل	ظل		ظل/IV		VERB	1	i	m	s	ظَلَّ 1	ظَلَّ 1	remain
شوف	شوف		شوف/IV		VERB	1	i	m	s	شَافَ 0	نَظَرَ 1	looking
الدني	الدنيا	ال/DET	الدنيا/NOUN		NOUN	-	-	f	s	دُنْيَا 2	دُنْيَا 2	world
لوما	لوما		لوما/CONJ		CONJ	-	-	-	-	لُوما 0	لُولا 1	except for
بقيل	بقيل		بقيل/IV	ل/IVSUFF SUBJ:1S	VERB	3	i	m	s	بَقِيَ 1	بَقِيَ 1	remain

Annotation Guidelines

DA Sentence: بروحي ساكنة ممنون عيونك عيوني وانا عم آخذ المونة رح ظل شوف الدنيا لوما بقيلي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
بروحي	بروحي	ب/FP	روح/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	s	رُوح 1	رُوح 1	spirit
ساكنة	ساكنة		ساكنة/ADJ	ة/NSUFF_FEM_SG	ADJ	-	-	f	s	ساكن 1	ساكن 1	residing
ممنون	ممنون		ممنون/ADJ		ADJ	-	-	m	s	ممنون 1	ممنون 1	grateful
عيونك	عيونك		عيون/NOUN	ك/POSS_PRON_2S	NOUN	-	-	m	s	عُيون 1	عُيون 1	eyes
عيوني	عيوني		عيون/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	s	عُيون 1	عُيون 1	eyes
وانا	وانا	و/CONJ	انا/PRON_1S		PRON	1	-	m	s	أنا 1	أنا 1	I
عم	عم		عم/PROG_PART		PROG	-	-	m	s	عم 1	عم 1	is doing
آخذ	آخذ	آ/IV1S	أخذ/IV		VERB	-	-	m	s	أخذ 1	أخذ 1	take
الموني	المونة	ال/DET	مون/NOUN	ة/NSUFF_FEM_SG	NOUN	-	-	f	s	مُونَة 1	مُونَة 1	money
رح	رح		رح/FUT_PART		PROG	-	-	m	s	رح 1	رح 1	will do
ضل	ظل		ظل/IV		VERB	-	-	m	s	ظل 1	ظل 1	be present
شوف	شوف		شوف/IV		VERB	1	i	m	s	شاف 0	نظّر 1	looking
الدني	الدنيا	ال/DET	الدنيا/NOUN		NOUN	-	-	f	s	دُنْيَا 2	دُنْيَا 2	world
لوما	لوما		لوما/CONJ		CONJ	-	-	-	-	لُوما 0	لُولا 1	except for
بقيلي	بقيلي		بقي/IV	ي/IVSUFF_SUBJ:1S	VERB	3	i	m	s	بقي 1	بقي 1	remain

CODA (التهجئة الصحيحة)

- The “correct”/unified spelling of the token.
- Followed Curras CODA
- Examples:

Token	CODA
بألكن <i>balkn</i>	بقلكن <i>bqlkn</i>
بقلكون <i>bqlkwn</i>	بقلكن <i>bqlkn</i>
طريء <i>try</i>	طريق <i>tryq</i>
هايدي <i>hāydy</i>	هيدي <i>hydy</i>
عيونن <i>ywnn</i>	عيونن <i>ywnun</i>

Annotation Guidelines

DA Sentence: بروحي ساكنة ممتون عيونك عيوني وانا عم أخذ المونة رح ظل شوف الدنيا لوما بقيلبي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
بروحي	بروحي	ب/PREP	روح/NOUN	ROSS PRON_1S	NOUN	-	-	m	s	روح 1	روح 1	spirit
ساكنة	ساكنة		ساكن/ADJ	/NSUFF_		-	-	f	s	ساكن 1	ساكن 1	residing
ممتون	ممتون		ممتون/ADJ									
عيونك	عيونك		عيون/NOUN									
عيوني	عيوني		عيون/NOUN									
وانا	وانا	و/CONJ	انا/PRON_1S									
عم	عم		عم/PROG_PA									
أخذ	أخذ	آ/IV1S	أخذ/IV									
الموني	المونة	ال/DET	مون/NOUN									
رح	رح		رح/FUT_PART									
ظل	ظل		ظل/IV									
شوف	شوف		شوف/IV									
الدنيا	الدنيا	ال/DET	الدنيا/NOUN									
لوما	لوما		لوما/CONJ									
بقيلبي	بقيلبي		بقيلبي/IV									

Prefixes (سوابق)

- Based on SAMA prefixes tagset, which are:

CONJ { و ، ف }	وحي ، فحي	PART_RESTRICT { الا }	والاهي
CONNEX_PART { ف }	فعلشان	PROG_PART { ب }	بيطلع
DEM_PRON { ه }	هالبلد	REL_PRON { ل ، ما }	ماشاء ، لمين
DET { ال }	الغنية	SUB_CONJ { و ، ما ، إن }	وشكلوا ، مانتني ، إنشاء
EMPHATIC_PART { ت ، ل }	لاوريكي ، تاوريكي	VOC_PART { يا }	يا معلم
FUT_PART { ح ، س }	حبريح ، سبريح	IV1S { آ ، ت }	بامرزم ، آخدها ، بتذكر
INTERROG_PART { با ، شو ، * }	مالقاش ، شويك ، اغنيها ، باغنيها	IV2S { ت }	بتدخل
JUS_PART { ل }	ليقول	IV2P { ت }	بتتذكروا
NEG_PART { لا ، ما ، مش }	ولاشي ، مايقتش ، مشلازم	IV3S { ت ، ي }	بتكي ، بيضحك
PREP	لصاحباتها ، بالبلد ، مشان ، عشان ، علحصان	IV3P { ي }	بيجيوا
{ ل ، ب ، م ، ع ، عد ، * ، ف ، ت ، ك ، مند }	فالبلد ، تاحكيلك ، كدكتجي ، منشان	IV1P { ن }	بتطلع

Annotation Guidelines

DA Sentence: بروحي ساكنة مُمون عيونك عيونني وانا عم آخذ المونة رح ظل شوف الدنيا لوما بقيلي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
بروحي	بروحي	ب/PREP	روح/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	s	روح 1	روح 1	spirit
ساكنة	ساكنة		ساكن/ADJ	ة/NSUFF_FEM	ADJ	-	-	f	s	ساكن 1	ساكن 1	residing
مُمون	مُمون		مُمون/ADJ		ADJ							
عيونك	عيونك		عيون/NOUN		NOUN							
عيونني	عيونني		عيون/NOUN		NOUN							
وانا	وانا	و/CONJ	انا/PRON_1S		PRON_1S							
عم	عم		عم/PROG_PART		PROG_PART							
آخذ	آخذ	آ/IV1S	آخذ/IV		IV							
الموني	المونة	ال/DET	مون/NOUN		NOUN							
رح	رح		رح/FUT_PART		FUT_PART							
ضل	ظل		ظل/IV		IV							
شوف	شوف		شوف/IV		IV							
الدني	الدنيا	ال/DET	الدنيا/NOUN		NOUN							
لوما	لوما		لوما/CONJ		CONJ							
بقيلي	بقيلي		بقي/IV		IV							

suffixes (لواحق)

- Based on SAMA suffix tagset, which are:

PREP {ر}	صحلي	NSUFF_FEM_PL {ة}	اللبنانية
NEG_PART {ش، اش، وش}	بدش، ملقاش، معوش	PVSUFF_DO:3P {هْن، +هْن}	هدهن، هددهن
REL_ADV {يش}	قديش	PVSUFF_DO:2P {كُن، *كُن، *كو، *كن}	عرفتكو، عرفتكين، عرفتكُن
POSS_PRON_1S {ني}	بعدي	PVSUFF_DO:2FS {كي}	شفتاكي
POSS_PRON_2P {كُن، *كُن، *كو، *كن}	بالكو، بالكين، بالكُن	PVSUFF_SUBJ:3P {نو}	طعميتوهن
POSS_PRON_3S {ت}	تاعت	PVSUFF_SUBJ:3S {يت}	طابقيتني
POSS_PRON_3P {هْن، +هْن}	تاعتهن، تاعتهن	IVSUFF_SUBJ:3S {يت}	معجبتني
POSS_PRON_2FS {كي}	عنيكي	IVSUFF_DO:2P {كُن، *كُن، *كو، *كن}	بتقلكو، بتقلكين، تقلكُن
PRON_2FS {كي}	عليكي	IVSUFF_SUBJ:2P {كُن، *كُن، *كو، *كن}	بديكو، بديكين، بديكُن
PRON_2P {كُن، *كُن، *كو، *كن}	نيالكو، نيالكين، نيالكُن		

Annotation Guidelines

DA Sentence: بروحي ساكنة ممنون عيونك عيوني وانا عم آخذ المونة رح ظل شوف الدنيا لوما بقيلي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
بروحي	بروحي	ب/PREP	روح/NOUN	ي/PRON_1S	NOUN	-	-	m	s	رُوح	رُوح	spirit
ساكنة	ساكنة		ساكن/ADJ	ة/NSUFF_FEM_SG		-	-	f	s	ساكن	ساكن	residing
ممنون	ممنون		ممنون/ADJ		ADJ					مُتَمَنِّين	مُتَمَنِّين	grateful
عيونك	عيونك		عيون/NOUN	ك/POSS_PR								
عيوني	عيوني		عيون/NOUN	ي/POSS_PR								
وانا	وانا	و/CONJ	انا/PRON_1S									
عم	عم		عم/PROG_PART									
آخذ	آخذ	آ/IV1S	خذ/IV									
الموني	المونة	ال/DET	مون/NOUN	ة/NSUFF_FEM_SG	NOUN	-	-	f	s	مُونَة	مُونَة	provisions
رح	رح		رح/FUT_PART		FUT_PART	-	-	-	-	رَح	سَوَّف	will
ضل	ظل		ظل/IV		VERB	1	i	m	s	ظَلَّ	ظَلَّ	remain
شوف	شوف		شوف/IV		VERB	1	i	m	s	شَاف	نَظَّر	looking
الدني	الدنيا	ال/DET	الدنيا/NOUN		NOUN	-	-	f	s	دُنْيَا	دُنْيَا	world
لوما	لوما		لوما/CONJ		CONJ	-	-	-	-	لُومَا	لُولا	except for
بقيلي	بقيلي		بقي/IV	ي/IVSUFF_SUBJ:1S	VERB	3	i	m	s	بَقِيَ	بَقِيَ	remain

Stem (الساق)

- <Stem/POS> tagging schema used in Curras
- SAMA POS tagset
- Reused Curras stems (as much as possible)

Annotation Guidelines

DA Sentence: بروحي ساكنة ممنون عيونك عيونني وانا عم أخذ المونة رح ظل شوف الدنيا لوما بقيلي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
بروحي	بروحي	ب/PREP	روح/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	s	رُوح	رُوح	spirit
ساكنة	ساكنة		ساكن/ADJ	ة/NSUFF_FEM_SG	ADJ	-	-	f	s	ساكن	ساكن	residing
ممنون	ممنون		ممنون/ADJ		ADJ	-	-	m	s	مَمْنُون	مَمْنُون	grateful
					VERB	-	-	m	p	عَيْن	عَيْن	eye
					VERB	-	-	m	p	عَيْن	عَيْن	eye
					PRON	1	-	m	s	أنا	أنا	I
					VERB	-	i	-	-	عم	ظَلَّ	continue
					VERB	1	i	m	s	أخذُ	أخذُ	take
					NOUN	-	-	f	s	مُونة	مُونة	provisions
					VERB	-	-	-	-	رح	سَوَّف	will
ضل	ظل		ظل/IV		VERB	1	i	m	s	ظَلَّ	ظَلَّ	remain
شوف	شوف		شوف/IV		VERB	1	i	m	s	شَاف	نَظَّر	looking
الدني	الدنيا	ال/DET	الدنيا/NOUN		NOUN	-	-	f	s	دُنْيَا	دُنْيَا	world
لوما	لوما		لوما/CONJ		CONJ	-	-	-	-	لُوما	لُولا	except for
بقيلي	بقيلي		بقي/IV	ي/IVSUFF_SUBJ:1S	VERB	3	i	m	s	بَقِيَ	بَقِيَ	remain

- Same SAMA tagsets: POS, Person, Aspect, Gender, Number

Annotation Guidelines

DA Sentence: بروحي ساكنة ممنون عيونك عيوني وانا عم آخذ المونة رح ظل شوف الدنيا لوما بقيلي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lamma	MSA Lemma	Gloss
بروحي	بروحي	ب/PREP	روح/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	s	روح	رُوح 1	spirit
ساكنة	ساكنة		ساكن/ADJ	ة/NSUFF_FEM_SG	ADJ	-	-	f	1	ساكن 1	ساكن 1	residing
ممنون	ممنون		ممنون/ADJ		ADJ	-	-	-	s	مَمْنُون 1	مَمْنُون 1	grateful
						-	-	m	p	عَيْن 1	عَيْن 1	eye
						-	-	m	p	عَيْن 1	عَيْن 1	eye
					ART	1	-	m	s	أنا 1	أنا 1	I
					ART	-	i	-	-	عم 0	ظَلَّ 1	continue
						1	i	m	s	أَخَذُ 1	أَخَذُ 1	take
						-	-	f	s	مُؤْنَة 1	مُؤْنَة 1	provisions
					ART	-	-	-	-	رح 0	سَوَّف 1	will
ضل	ظل		ظل/IV		VERB	1	i	m	s	ظَلَّ 1	ظَلَّ 1	remain
شوف	شوف		شوف/IV		VERB	1	i	m	s	شَاف 0	نَظَّر 1	looking
الدني	الدنيا	ال/DET	الدنيا/NOUN		NOUN	-	-	f	s	دُنْيَا 2	دُنْيَا 2	world
لوما	لوما		لوما/CONJ		CONJ	-	-	-	-	لُومَا 0	لُولا 1	except for
بقيلي	بقيلي		بقي/IV	ي/IVSUFF_SUBJ:1S	VERB	3	i	m	s	بَقِيَ 1	بَقِيَ 1	remain

MSA Lemma (مدخلة معجمية فصحي)

- SAMA lemmas
- If no SAMA lemma, add one in the same way, e.g., (0-يوغا)

Annotation Guidelines

DA Sentence: بروحي ساكنة ممنون عيونك عيوني وانا عم أخذ المونة رح ظل شوف الدنيا لوما بقيلي عيون

Token	CODA	Prefix(es)	Stem	Suffix(es)	POS	Person	Aspect	Gender	Number	DA Lemma	MSA Lemma	Gloss
بروحي	بروحي	ب/PREP	روح/NOUN	ي/POSS_PRON_1S	NOUN	-	-	m	s	رُوح 1	رُوح 1	spirit
ساكنة	ساكنة		ساكن/ADJ	ة/NSUFF_FEM_SG	ADJ	-	-	f	s	ساكن 1	ساكن 1	residing
ممنون	ممنون		ممنون/ADJ		ADJ	-	-	m	s	مَمْنُون 1	مَمْنُون 1	grateful
						-	-	m	p	عَيْن 1	عَيْن 1	eye
						-	-	m	p	عَيْن 1	عَيْن 1	eye
						1	-	m	s	أنا 1	أنا 1	I
					RT	-	i	-	-	عم 0	ظَلَّ 1	continue
						1	i	m	s	أَخَذُ 1	أَخَذُ 1	take
						-	-	f	s	مُؤْنَة 1	مُؤْنَة 1	provisions
					RT	-	-	-	-	رح 0	سَوْفَ 1	will
ضل	ظل		ظل/IV		VERB	1	i	m	s	ظَلَّ 1	ظَلَّ 1	remain
شوف	شوف		شوف/IV		VERB	1	i	m	s	شَافَ 0	نَظَرَ 1	looking
الدني	الدنيا	ال/DET	الدنيا/NOUN		NOUN	-	-	f	s	دُنْيَا 2	دُنْيَا 2	world
لوما	لوما		لوما/CONJ		CONJ	-	-	-	-	لُومَا 0	لَوْلَا 1	except for
بقيلي	بقيلي		بقي/IV	ي/IVSUFF_SUBJ:1S	VERB	3	i	m	s	بَقِيَ 1	بَقِيَ 1	remain

Dialect Lemma (مدخلة معجمية عامية)

- from MSA, then its MSA and DA lemmas are the same.
- If not, introduce a new DA lemma

Annotation Guidelines

Frequent functional words in Lebanese (right) and Palestinian (left).

DEM_PRON هيّاني - هيّاني، هيّاني، هيّاني هيداك - هيداك هيدي - هاي، هذي هيدا - هاذا هيديك - هديك، هذيك ياني - ايانبي هياها - هيبي	INTERJ أيه - اه، اها تاري، تخمن - أجرمنعنو لاء - لاء، لع
PRON احنا، نحنا - احنا	REL_PRON تاعول، تاعون - تاعون تبعهن - تبعهنم
ADV بركي - بلكي هلق - هلقيت، هلا، هسا، الحين هون، هنا، هنتياني - هون، هان، هانا، هونا هونيك، هنكه، هنكتياني - هناك، غاد عندهن - عندهم لخالهن - لخالهم بقا - عاد عاللس - عالسكيت	INTERROG_ADV وينكن - وينكنم، وينكو منتلي - منوينلي
	INTERROG_PRON مينو - أنو ميني - انبي بشو - بشو، بايش
	NEG_PART معاش - بلاش

- Randomly selected annotated sentences that together consist of 400 tokens
- Each annotator re-annotated 100 tokens that were annotated by another

❖ Cohen's Kappa inter-annotation agreement

Tag	Values	Agreement	Disagreement	Kappa
Stem	178	357	43	0.884
Prefix	41	380	20	0.860
Suffixes	55	358	42	0.738
POS	22	340	60	0.821
Person	3	359	41	0.629
Aspect	4	384	16	0.911
Gender	3	337	63	0.687
Number	4	347	53	0.741
Overall				0.785

❖ F1-Score, precision and recall metrics

Feature	Precision	Recall	F1-Score
Stem	0.9036	0.8935	0.893
Prefixes	0.964	0.95	0.955
Suffixes	0.948	0.895	0.915
POS	0.898	0.85	0.853
Person	0.928	0.898	0.910
Aspect	0.974	0.96	0.967
Gender	0.845	0.843	0.844
Number	0.881	0.868	0.873
Overall	0.918	0.894	0.901

Curras – Palestinian Dialect Corpus

Originally

- 56K tokens
- morphologically annotated corpus



Jarrar, M., Habash, N., Alrimawi, F., Akra, D., & Zalmout, N. (2016). **Curras: An Annotated Corpus for the Palestinian Arabic Dialect**. *Language Resources and Evaluation*, 50(219), 1-31.

Jarrar, M., Habash, N., Akra, D. F., & Zalmout, N. (2014). **Building a corpus for Palestinian Arabic: a preliminary study**. In *Proceedings – Arabic Natural Language Processing Workshop*, at the Conference on Empirical Methods in Natural Language Processing (EMNLP 2014). Association for Computational Linguistics. (pp. 18-27).

Curras Revisions

Almost all annotations are revised

1. Tokenization and POS

- No parsing errors in tokenization and POS
- Buckwalter transliterations
- CODA corresponds to prefixes + stem + suffixes
- every prefix should be in the predefined set of prefixes,
- every suffix should be in the predefined set of suffixes,
- every stem POS should be in the SAMA POS tagset.

Curras Revisions

Almost all annotations are revised

2. Lemmatization

Manually revised MSA and DA lemmas

- MSA lemmas are linked with SAMA lemmas, new lemmas are marked with “_0”
- DA lemmas are linked with MSA lemmas (if different)

Curras Revisions

Almost all annotations are revised

3. Other features

Person, Aspect, Gender, and Number are verified

Generating Unique Solutions

A minimum set of unique morphological solutions is generated and used to annotated the Lebanese corpus

https://portal.sina.birzeit.edu/curras

Available Online



Curras
Corpus for Palestinian Arabic
مدونة اللهجة العامية الفلسطينية

Word Stem MSA Surface Gloss
 Whole Word Substring
 Palestinian Lebanese

[About](#) [Search](#) [Publications](#) [Download](#) [News](#) [Free Ideas](#)

هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	... معروف فس نفس . بتطلع الخميس ليش هالقدرة الازمة بقولك يوم الخميس دايمًا هيك , هذا يوم في الاسبوع طائر...
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	الجار : يا زلمة دوشتنا طرشتنا يلعن ابو هيك سكة
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	...بلم اللي اكلناه . عزمونا الاسبوع الماضي , عشان نردلهم العزيمة , وهيك بك تحطيلهم ترتزسي
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	رجل 2 : هيك مزبلة بدها هيك ختم
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	رجل 2 : هيك مزبلة بدها هيك ختم
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	... معروف فس نفس . بتطلع الخميس ليش هالقدرة الازمة بقولك يوم الخميس دايمًا هيك , هذا يوم في الاسبوع طائر
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	...ن ما باكل الا الحلال وما يشرب وما يصاحب ! الله يسر عجد , معقول في ناس هيك ما بتكفر بعضها ؟ ! عالم فاضية عككرة بال...
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	...صير تدلني امرار ع اطباق ما فيهن لحمه عشان انا ما كون محدود بخياراتي و هيك .
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	...ها شي المكونات , قائلتي انتبه هيدي يمكن فيها كحول و أنت يمكن ما بتاكل هيك شي , و اختارتي وحدة ما فيها كحول
هيك	هيك	هيك 0	هيكنا 1	هيكنا/اداء استفعال + ظرف	اداء استفعال + ظرف	لا ينطبق	لا ينطبق	"is + this way,thus this way,this,thus"	...ي , ان مش ساعة بيع بائع انتك ... بس في شي تغير , شي مش عارف شو هوي , هيك شي اخفتي مثل لما اكل اكلة مش بالبيت وحس...

Summary

Baladi: Lebanese morphologically annotated corpus (9.6K)

Curras (revised): Palestinian morphologically annotated corpus (56K)

= a more **Levantine Corpus** (65.6K tokens)

Differences and communities are discussed in the paper

Thank You

References