

Camel Morph MSA: A Large-Scale Open-Source Morphological Analyzer for Modern Standard Arabic

Christian Khairallah, Salam Khalifa, Reham Marzouk, Mayar Nassar and Nizar Habash

Computational Approaches to Modeling Language Lab, New York University Abu Dhabi, UAE
{christian.khairallah,nizar.habash}@nyu.edu

Motivations

- Morphological analyzers add value to neural models in several NLP tasks:
 - Grammatical Error Correction (Alhafni et al., 2023)
 - Morphological Disambiguation (Inoue et al., 2022)
 - Machine Translation (Oudah et al., 2019)
- Make building analysis/generation tools easier and more consistent
- Make expansion of lexicons more efficient

Contributions

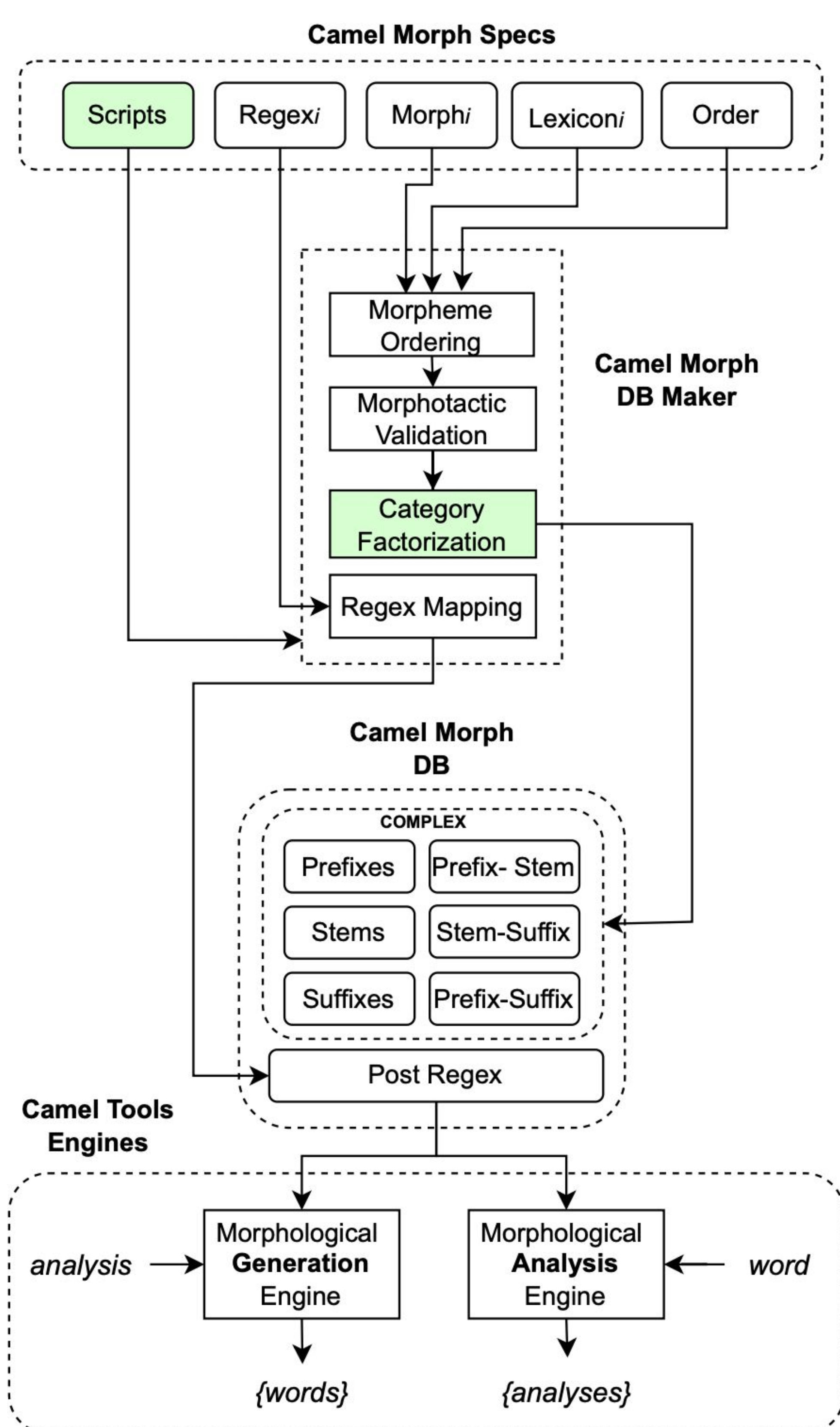
- Released largest open-source Arabic morphological analyzer/generator
- Developing an extendable large-scale implementation using **Camel Morph**
- Benchmarking our models against a publicly available analyzer
- Our data and code are publicly available



morph.camel-lab.com

Camel Morph Framework (Habash et al., 2022)

- **Compile morphological specifications into optimized DBs**
 - **From**: linguistic representations and rules
 - **To**: lists of complex prefixes and suffixes, stems, and their compatibilities



Camel Morph Specifications & Database

- **Example of Morphology Specification and Validation**
 - Morph order defines the full space of all morphemes that can co-occur by their class.
 - Each form (allomorph) sets some truth conditions to be true.
 - For a word to be valid, the required truth conditions of every form (allomorph) in it must be already set by some other allomorph.

	Morph Order										<div>sufaraA⁺⁺+a</div> <div><div>safiy^r+aAt⁺⁺+him</div><div><div>sufaraA⁺+ŷ⁺+i+him</div><div>Al⁺+safiy^r+aAt⁺⁺</div></div></div>				O1	O2	O2	O3																																																																																																																																																																																																																																																																																																																																																																																																					
	DBPrefix		DBStem				DBSuffix																																																																																																																																																																																																																																																																																																																																																																																																																
O1	[Conj]	[Prep]	[NomStem]				[NomBuff]				[NomSuff.IG]																																																																																																																																																																																																																																																																																																																																																																																																												
O2	[Conj]	[Prep]	[NomStem]				[NomBuff]				[NomSuff.CG] [Pronoun]																																																																																																																																																																																																																																																																																																																																																																																																												
O3	[Conj]	[Prep]	[Determiner]		[NomStem]				[NomBuff]				[NomSuff.DG]																																																																																																																																																																																																																																																																																																																																																																																																										

	Count	Verb Example	Count	Noun Example	Count	Particle Example
Morphology Specifications	Lemma	1 ramay 'throw' رمى	1 safiyr 'ambassador' سفير	1 calay 'on, upon' على		
	Proclitics	30 wa, fa, Aa, ... و, ف, أ, ...	19 Al, wa, fa, li, bi, ... ال, و, ف, ل, ب, ...	9 wa, fa, fa, Aa, ... و, ف, أ, ...		
	Prefixes	13 ya, ta, ... ي, ت, ...	N/A	N/A		
	Pre-Buffers	1 a ا	N/A	N/A		
	Stems	3 زم (ماضي), زم (مضارع), زم (امر) ram (perfect), r.m (imperfect/command)	2 سفير, سفيرا سفيرا (base), sufaraA (broken plural)	4 على, على, علام, calay, çalay, çalaAma, ...		
	Post-Buffers	11 ay, ø, ay, iy, ... ي, ء, ي, ي, ...	5 w, ' , y, ø, ... و, ء, ي, ø, ...	0		
	Suffixes	100 lu, ta, naA, at, ... ل, ت, ن, أ, ت, ...	79 a, i, i, ø, ahu, aAili, ... ا, ي, ي, ø, ا, هـ, ا, ا, ل, ي, ...	13 hu, hi, hum, haA, ... هـ, هـ, هـ, هـ, ...		
	Enclitics	18 hu, hi, hum, him ... هـ, هـ, هـ, هـ, ...	19 hu, hi, hum, him ... هـ, هـ, هـ, هـ, ...	0		
	Conditions	25 #ay (defective), ...	25 FP (+At.suffix), ...	6 encø (presence of object clitic)		
	Order Seqs.	69 verbal orders	21 nominal orders	2 particle orders		
DB	Complex					
	Prefixes	3,129 wa, fa, Aa, aaya, ... و, ف, أ, و, ...	199 wa, waAl, wali, ili, ... و, و, ال, و, ل, ل, ...	14 wa, fa, >a, >awa, ... و, ف, أ, و, ...		
	Stems	10 ramay, ramay, ... رمى, رمى, ...	4 safiyr, sufaraW ... سفير, سفارو, ...	4 calay çalay, çalaAma, ... على, على, علام, ...		
	Suffixes	972 ta, tahu, hi, ... ت, ت, هـ, هـ, ...	243 aAtu, aAthim, aki, ... ا, ا, ت, ا, ت, هـ, ا, ك, ...	14 hu, hi, hum, haA ... هـ, هـ, هـ, هـ, ...		
Compatibility Combinations	2,687 Stem-Suffix, Prefix-Stem, Prefix Suffix	394 Stem-Suffix, Prefix-Stem, Prefix Suffix	14 Stem-Suffix, Prefix-Stem, Prefix Suffix			
Word Forms (unique analyses)	42,588 Word: Åaramay.tahu أرمنية 'did you throw it?' Features: verb, perfective, active, 2nd person, masculine, singular, interrogative, 2nd person singular	12,942 Word: walisafiyaAthim وسفيراتهم 'and for their ambassadors' Features: noun, feminine, plural, genitive, construct, 3rd person ps possessive, 3rd person singular	224 Word: façalay.hA فاضلها '3rd person fs pron, conjunction' Features: prep, 3rd person fs pron, conjunction			

Evaluation

- Compare to Calima MSA (Taji et al. 2018)~SAMA (Graff et al., 2009)
- Coverage evaluation of Penn Arabic Treebank (Maamouri et al., 2004): Recall 95.9% of analyses; 90% of mismatches due to gold errors.

		Camel Morph MSA Specs		Camel Morph MSA DB		SAMA/CALIMA DB			
(a)	Lemmas (Stems)	105,102	(140,612)	105,102	(154,573)	42,218	(71,466)	Lemmas (Stems)	(c)
	<i>Verbs</i>	9,333	(38,156)	9,333	(47,540)	9,279	(26,343)	<i>Verbs</i>	
	<i>Nominals</i>	33,267	(39,837)	33,267	(44,414)	32,701	(44,742)	<i>Nominals</i>	
	<i>Others</i>	230	(347)	230	(347)	238	(381)	<i>Others</i>	
	<i>Proper Nouns - Annex</i>	62,272	(62,272)	62,272	(62,272)				
(b)	Prefix Morphs (Allom.)	60	(65)	14,726		4,640		DBPrefix Sequences	(d)
	Suffix Morphs (Allom.)	205	(406)	12,724		1,191		DBSuffix Sequences	
	Stem Buffers	111		15,044		979		Compatibility Entries	
	Unique Condition Terms	88		535,186,314	(242,824,398)	70,250,488		Unique Diacritized Forms (no Annex)	(e)
	Morph Order Sequences	122		1,447,312,125	(630,731,386)	227,471,211		Unique Analyses (no Annex)	
					8,070.764	(3,214.695)	1,514.577		Unique Analyses w/o Clitics (no Annex)

	Camel Morph MSA		SAMA/CALIMA	
	MSA	CA	MSA	CA
Run Time (sec)	12,293	4,667	4,231	1,960
Type OOV	67.9%	34.7%	75.1%	43.2%
Token OOV	2.3%	1.5%	3.5%	2.5%
Analyses/Type	18.9	21.2	13.7	15.2
Analyses/Token	38.6	45.7	18.9	20.3