

# A Lexical Database for Modern Standard Arabic Interoperable with a Finite State Morphological Transducer

Mohammed Attia, Pavel Pecina, Antonio Toral, Lamia Tounsi, Josef van Genabith

National Centre for Language Technology (NCLT),  
School of Computing, Dublin City University

Funded by:

Enterprise Ireland, the Irish Research Council for Science  
Engineering and Technology (IRCSET), and  
the EU projects PANACEA and META-NET

# Contribution

- We develop a finite state morphological transducer for Modern Standard Arabic
  1. Open source, distributed under the GPLv3 license
  2. Large scale, more than 30,000 lemmas
  3. Corpus based, truly representative of Modern Standard Arabic and not Classical Arabic.
  4. Compatible with Foma, an open-source fst compiler

# Short Tutorial

(1) Download Foma

<http://foma.sourceforge.net>

(2) Download AraComLex

<http://aracomlex.sourceforge.net>

(3) Build the transducer: README

# The transducer online

- You can test the transducer online:  
<http://www.cngl.ie/aracomlex>

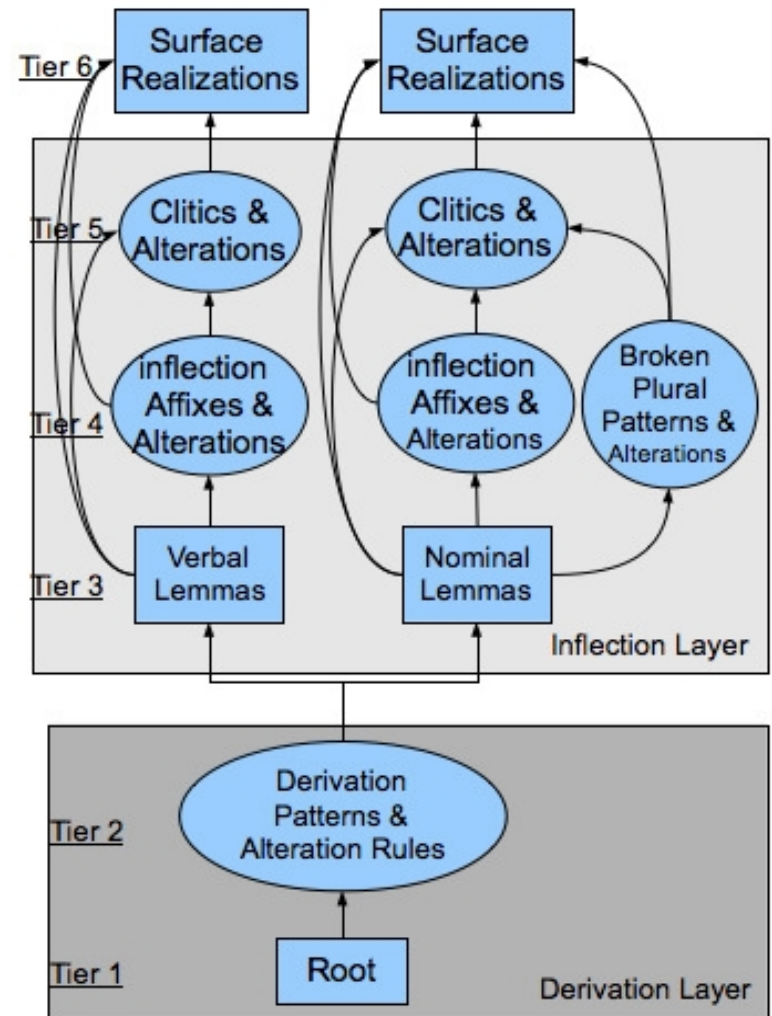
# Introduction

- Modern Standard Arabic vs. Classical Arabic
- Current State of Arabic Lexicography
  - Lexicons are not corpus-based
  - Buckwalter Arabic Morphological Analyser
- Importance of Lexical Resources

# Introduction

- Arabic Morphotactics

Root	درس drs			
Pattern	$R_1aR_2aR_3a$	$R_1aR_2R_2aR_3a$	$R_1\bar{a}R_2iR_3$	$muR_1aR_2R_2iR_3$
POS	V	V	N	N
lemma	d a r a s a 'study'	d a r r a s a 'teach'	d ā r i s 'student'	m u d a r r i s 'teacher'



# Aim

- Constructing a lexical database for Modern Standard Arabic
- Building a finite-state morphological transducer

# Methodology

- Using a medium-scale manually created lexicon of 10,799 lemmas, with detailed info for:
  - Nouns (human/nonhuman, POS, Continuation Classes)
  - Verbs (transitive/intransitive, allow passive, allow imperative)
- Using statistics from a 1 billion word corpus
  - 90% from the LDC's Arabic Gigaword
  - 10% collected from the Al-Jazeera website
- Using a pre-annotation tool: MADA+TOKAN



# Methodology

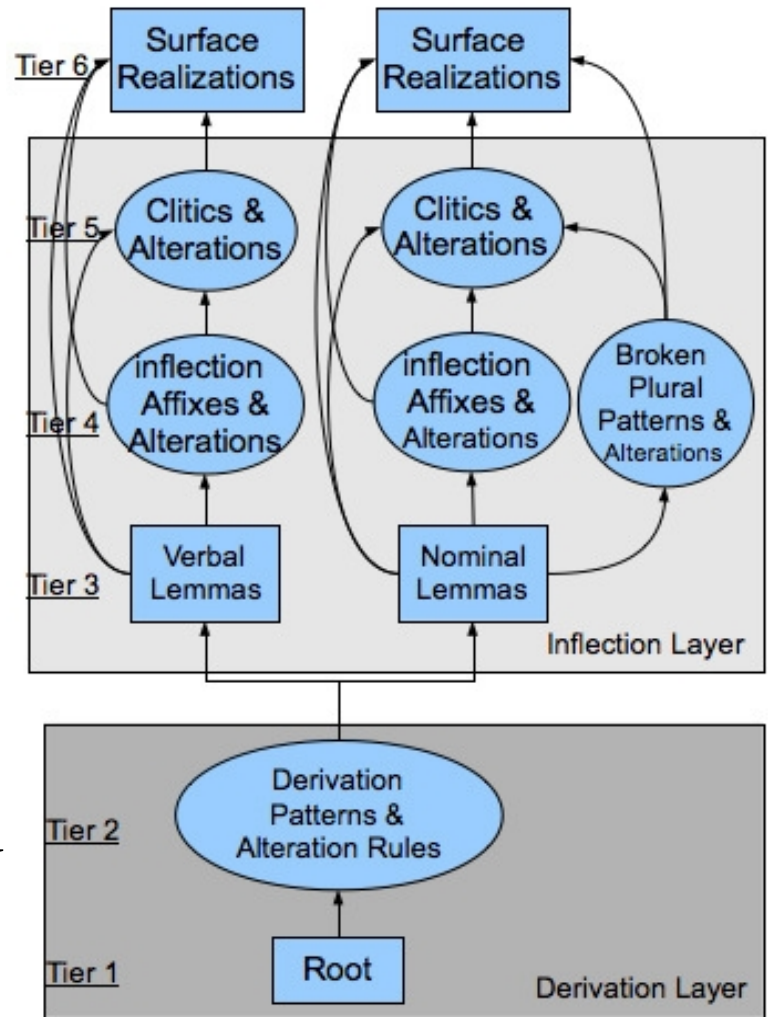
- Using Finite State Technology (XFST)
  - Bidirectional: Suitable for analysis and generation
  - handles concatenative and non-concatenative morphotactics
  - Speed and efficiency in dealing with millions of paths
  - Handles separated dependencies.
  - Handles phonological and orthographic changes through alteration rules.

# Methodology

- Design Approach:  
Three approaches
  - Root-based Morphology  
Xerox Arabic FTM
  - Stem-based morphology  
Buckwalter
  - Lemma-based morphology

\$kr      \$akar      PV    thank;give thanks

\$kr      \$okur      IV    thank;give thanks



# Methodology

## Our Approach: Lemma-based morphology

### (1) LEXICON Verbs

^ss^شَكَرَ[‘thank’]^se^	Transitive;
^ss^فَرِحَ[‘be-happy’]^se^@D.V.P@	Intransitive;
^ss^أَمَرَ[‘order’]^se^@D.M.I@	Transitive;
^ss^قَالَ[‘say’]^se^	Intransitive;

### (2) LEXICON Nouns

+masc+human^ss^مُعَلِّمٌ[‘teacher’]^se^	FemMascduFemduFemplMascpl;
+masc+human^ss^طَالِبٌ[‘student’]^se^	FemMascduFemduFempl;
+masc+nonhuman^ss^كِتَابٌ[‘book’]^se^	Mascdu;
+fem+nonhuman^ss^كُرْسَاتَةٌ[‘notebook’]^se^	DualFempl;

	Masculine Singular	Feminine Singular	Masculine Dual	Feminine Dual	Masculine Plural	Feminine Plural	Continuation Class
1	مُعَلِّمٌ <i>mu'allim</i> 'teacher'	مُعَلِّمَةٌ <i>mu'allimat</i>	مُعَلِّمَانِ <i>mu'allimān</i>	مُعَلِّمَاتَانِ <i>mu'allimātān</i>	مُعَلِّمُونَ <i>mu'allimuwn</i>	مُعَلِّمَاتٌ <i>mu'allimāt</i>	Fem-Mascdu-Femdu-Mascpl-Fempl
2	طَالِبٌ <i>tālib</i> 'student'	طَالِبَةٌ <i>tālibat</i>	طَالِبَانِ <i>tālibān</i>	طَالِبَاتَانِ <i>tālibātān</i>	-	طَالِبَاتٌ <i>tālibāt</i>	Fem-Mascdu-Femdu-Fempl
3	تَحْضِيرِيٌّ <i>taḥḍiryiy</i> 'preparatory'	تَحْضِيرِيَّةٌ <i>taḥḍiryiyat</i>	تَحْضِيرِيَّانِ <i>taḥḍiryiyān</i>	تَحْضِيرِيَّاتَانِ <i>taḥḍiryiyātān</i>	-	-	Fem-Mascdu-Femdu
4	-	بَقْرَةٌ <i>baqarat</i> 'cow'	-	بَقْرَتَانِ <i>baqaratān</i>	-	بَقْرَاتٌ <i>baqarāt</i>	Femdu-Fempl
5	تَنَازُلٌ <i>tanāzul</i> 'concession'	-	-	-	-	تَنَازُلَاتٌ <i>tanāzulāt</i>	Fempl
6	-	ضَعِيَّةٌ <i>ḍaḥiyyat</i> 'victim'	-	ضَعِيَّاتَانِ <i>ḍaḥiyyatān</i>	-	-	Femdu
7	مَحْضٌ <i>maḥḍ</i> 'mere'	مَحْضَةٌ <i>maḥḍat</i>	-	-	-	-	Fem
8	إِمْتِحَانٌ <i>imtiḥān</i> 'exam'	-	إِمْتِحَانَانِ <i>imtiḥānān</i>	-	-	إِمْتِحَانَاتٌ <i>imtiḥānāt</i>	Mascdu-Femdu
9	طَيَّارٌ <i>ṭayyār</i> 'pilot'	-	-	-	.tayyAruwn	-	Mascdu-Mascpl
10	كِتَابٌ <i>kitāb</i> 'book'	-	كِتَابَانِ <i>kitābān</i>	-	-	-	Mascdu
11	دِيمُقْرَاطِيٌّ <i>diy-muqrāṭiy</i> 'democrat'	-	-	-	دِيمُقْرَاطِيُّونَ <i>diy-muqrāṭiyuwn</i>	-	Mascpl
12	خُرُوجٌ <i>ḥuruwġ</i> 'exiting'	-	-	-	-	-	NoNum
13	مَبَاحِثٌ <i>mabāḥit</i> 'investigators'	-	-	-	-	-	Irreg_pl

# Methodology

## Alteration Rules:

Alteration Rules are used for handling discrepancies between surface forms and underlying representation or lemmas. We have 130 replace rules.

$a \rightarrow b \parallel L \_ R$

# Results to Date

- Start-off with a seed lexicon
  - Four Lexical Databases, manually constructed
    - 5,925 nominal lemmas
    - 1,529 verb lemmas
    - 490 patterns (456 for nominals and 34 for verbs)
    - lemma-root look up database

# Results to Date

- Automatically Extending the Lexical Database: Lexical Enrichment
  - Data-driven filtering technique
    - 40,648 lemmas (in Buckwalter or SAMA 3.1)
    - Statistics from three web search engines
    - Statistics from the corpus annotated by MADA
    - 29,627 lemmas (left after filtering)

# Results to Date

## Automatically Extending the Lexical Database: Feature Enrichment

- Machine Learning
- Multilayer Peceptron classification algorithm
- Training Data: 4,816 nominals and 1,448 verbs
- Classes for nominals: continuation classes (or inflection paths), the semantico-grammatical feature of humanness, and POS (noun or adjective)
- Classes for verbs: transitivity, allowing the passive voice, and allowing the imperative mood
- We feed these datasets with frequency statistics from the corpus and build a vector grid.



# Results to Date

- Extending the Lexical Database
  - Feature enrichment using Machine Learning

Exp.	Classes	Features	P	R	F
	<b>Nominals</b>				
1	Continuation Classes: 13 classes	number, gender, case, clitics	0.62	0.65	0.63
2	Human: yes, no, unspecified		0.86	0.87	0.86
3	POS: noun, adjective		0.85	0.86	0.85
	<b>Verbs</b>				
4	Transitivity: transitive, intransitive	number, gender, person, aspect, mood, voice, clitics	0.85	0.85	0.84
5	Allow passive: yes, no		0.72	0.72	0.72
6	Allow imperative: yes, no		0.63	0.65	0.64

# Results to Date

- Extending the Lexical Database
  - With Machine Learning we add:
    - 18,000 new lemmas:
      - 12,974 nominals
      - 5,034 verbs

# Results to Date

- AraComLex Lexicon Writing Application

**AraComLex: Arabic Computer Lexicon**

[Nominal Lemmas](#) || [Verb Lemmas](#) || [Templates](#) || [FST Morphological Analyser](#)

### Nominal Lemmas

Returned records: 3 [View](#) [Add New](#)

---

form\_id: 1, arabicUnpointed: خادم, arabicPointed: : خادم, gloss\_bw: (web)\_server\_(computer)  
 lemma\_bw: xAdim\_1, partOfSpeech\_pw: noun, Repeated records: 0, hasARoot: xdm, template\_auto: "@A@i@", template\_regex: ".A.i".

partOfSpeech_modif: <input type="text" value="noun"/>	lemma_modif: <input type="text" value="xAdim_1"/>	gloss_modif: <input type="text" value="(web)_server_(computer)"/>	lemma_morph: <input type="text" value="+masc"/>	partOfSpeech_ma: <input type="text" value="Noun"/>	continuationClass: <input type="text" value="FemMascdFemduMascdFempl"/>	human: <input type="text" value="yes"/>
lemma_extra: <input type="text" value="unspec"/>	irreg_plural: <input type="text" value="خادم#خدمة"/>	irregp_morph: <input type="text" value="unspec"/>	matched: <input type="text" value="1"/>	deleted: <input type="text" value="0"/>	reviewed: <input type="text" value="0"/>	<a href="#">Add Copy</a> <a href="#">Remove</a>

**Statistics:**  
 lemma\_freq: 33174, masc\_sg: 31986, masc\_dl: 231, masc\_pl: 0, fem\_sg: 0, fem\_dl: 0, fem\_pl: 957, prc0: 2059, prc1: 808, prc2: 351, prc3: 0, enc0: 428

---

form\_id: 1, arabicUnpointed: خادم, arabicPointed: : خادم, gloss\_bw: servant;attendant  
 lemma\_bw: xAdim\_1, partOfSpeech\_pw: noun, Repeated records: 0, hasARoot: xdm, template\_auto: "@A@i@", template\_regex: ".A.i".

partOfSpeech_modif: <input type="text" value="noun"/>	lemma_modif: <input type="text" value="xAdim_1"/>	gloss_modif: <input type="text" value="servant;attendant"/>	lemma_morph: <input type="text" value="+masc"/>	partOfSpeech_ma: <input type="text" value="Noun"/>	continuationClass: <input type="text" value="FemMascdFemduMascdFempl"/>	human: <input type="text" value="yes"/>
lemma_extra: <input type="text" value="unspec"/>	irreg_plural: <input type="text" value="خادم#خدمة"/>	irregp_morph: <input type="text" value="unspec"/>	matched: <input type="text" value="1"/>	deleted: <input type="text" value="0"/>	reviewed: <input type="text" value="0"/>	<a href="#">Add Copy</a> <a href="#">Remove</a>

**Statistics:**  
 lemma\_freq: 3006, masc\_sg: 0, masc\_dl: 0, masc\_pl: 86, fem\_sg: 2776, fem\_dl: 144, fem\_pl: 0, prc0: 1477, prc1: 130, prc2: 120, prc3: 0, enc0: 387

Find:     Highlight all  Match case

# Results to Date

- FST Morphology Coverage and RPW Results
  - a test corpus of 800,000 words, divided as
    - 400,000 for Semi-Literary text
    - 400,000 for General News texts.

Morphology	No. of Lemmas	General News		Semi-Literary	
		Coverage	Rate per word	Coverage	Rate per word
AraComLex 1.0	10,799	79.68%	1.67	69.37%	1.62
AraComLex 2.0	28,807	86.89%	2.10	85.14%	2.09
AraComLex 2.1	30,587	87.13%	2.09	85.73%	2.08
SAMA	40,648	88.13%	5.32	86.95%	5.30

# Future Work

- Going beyond SAMA
- Including Named Entities and MWEs
- Building a spell checker

# Conclusion

- Open-source finite state transducer for Modern Standard Arabic (AraComLex) distributed under the GPLv3 license.
- We successfully use machine learning to predict morpho-syntactic features for newly acquired words.
- Comparing our morphological transducer to SAMA, we find that we achieve comparable coverage and lower rate of analyses per word.