

# A Robust Morphological Analysis System for the Moroccan Dialect: Design and Evaluation

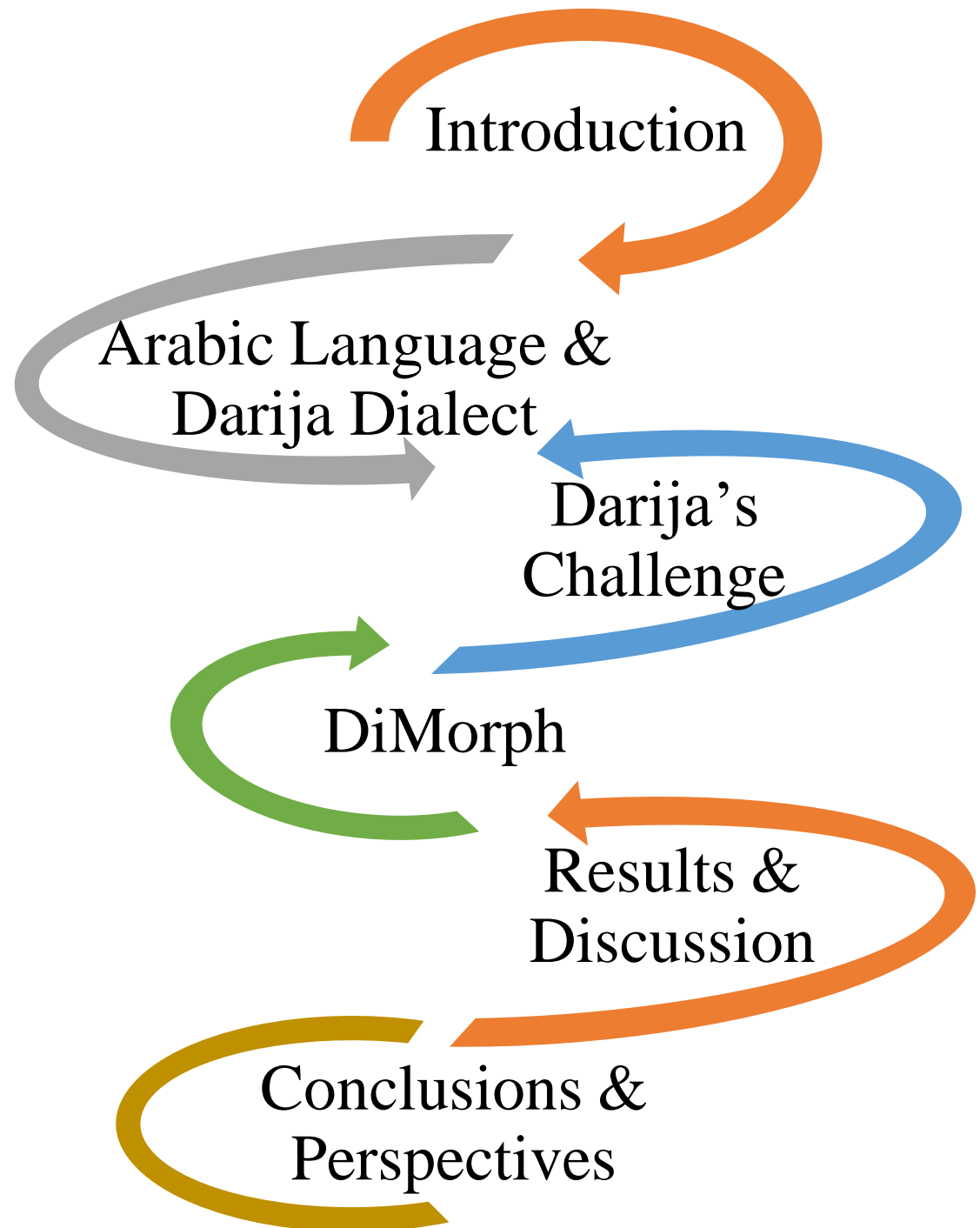


Nadia Khlif  
Azzeddine Mazroui  
Ouafae Nahli

[nadia.khlif@ump.ac.ma](mailto:nadia.khlif@ump.ac.ma)  
[azze.Mazroui@gmail.com](mailto:azze.Mazroui@gmail.com)  
[ouafae.nahli@ilc.cnr.it](mailto:ouafae.nahli@ilc.cnr.it)



## PLAN



- Arabic NLP faces significant challenges due to diglossia with Modern Standard Arabic (MSA) used in writing and dialects in spoken communication.
- With social media, dialects now appear widely in written form, but their lack of standardization and limited resources hinder NLP development.
- To address this gap, our work within the *CWALM* project focuses on developing tools for Contemporary Written Arabic.
- As part of this effort, we have developed **DiMorph**, a morphological analyzer specifically designed to dialectal Arabic.

# Arabic Language & Darija dialect

## Classical Arabic

- The Language of the Quran
- Adherence to Strict Grammatical Rules
- Religious Significance and Unique Lexicon

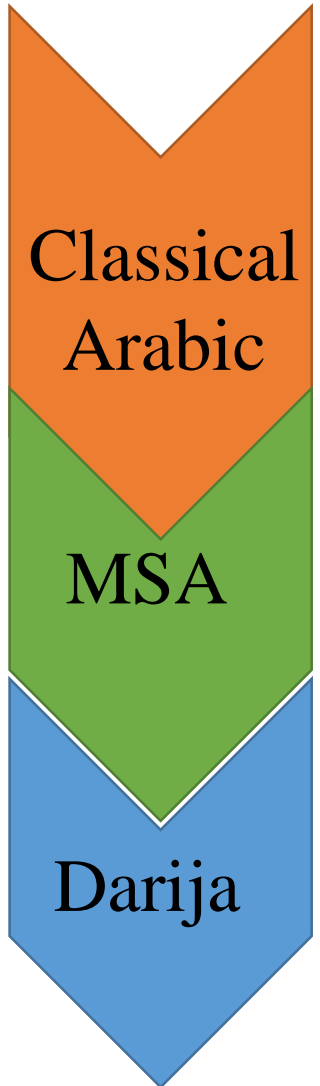
## Modern Standard Arabic (MSA)

- Foundation in Classical Arabic
- Linguistic Authorities and Standardization
- Uniformity in Written Communication

## Arabic Dialects

- Localization and Informality
- Organic Evolution within Communities
- Lack of the formal recognition and standardized grammar rules

# Arabic Language & Darija dialect



/qiTar/ قطار

Camel walking in sequence



/sayya:rat/ سيارة

Camel walking in simultaneous

Word Time Shifts



/qiTar/ قطار

Train

/tora:n/ تران

Train

/sayya:rat/ سيارة

Car

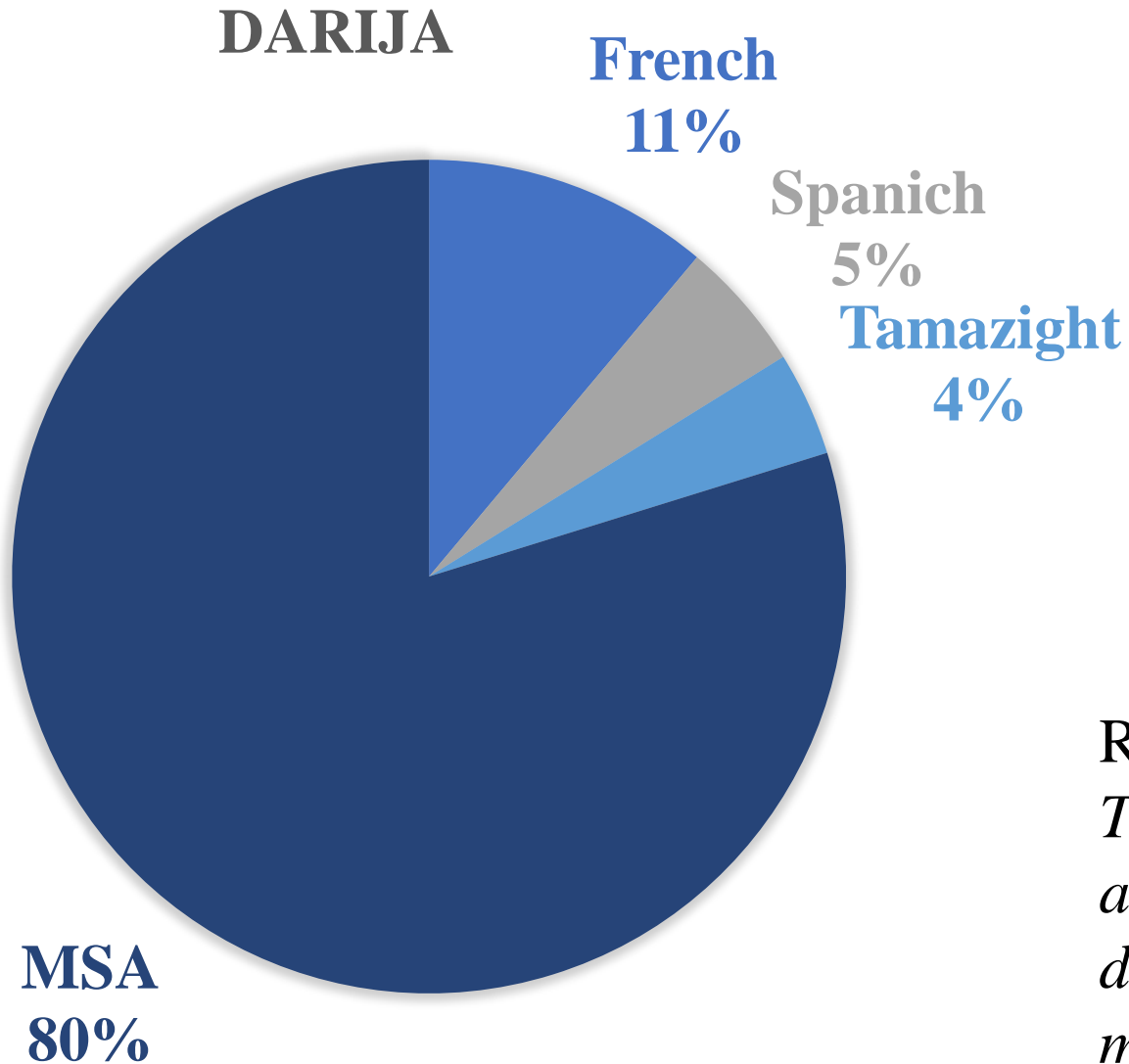


/Tu:mu:bi:l/ طوموبيل /sayya:rat/ سيارة

Car

Lexical borrowing

# Arabic Language & Darija dialect

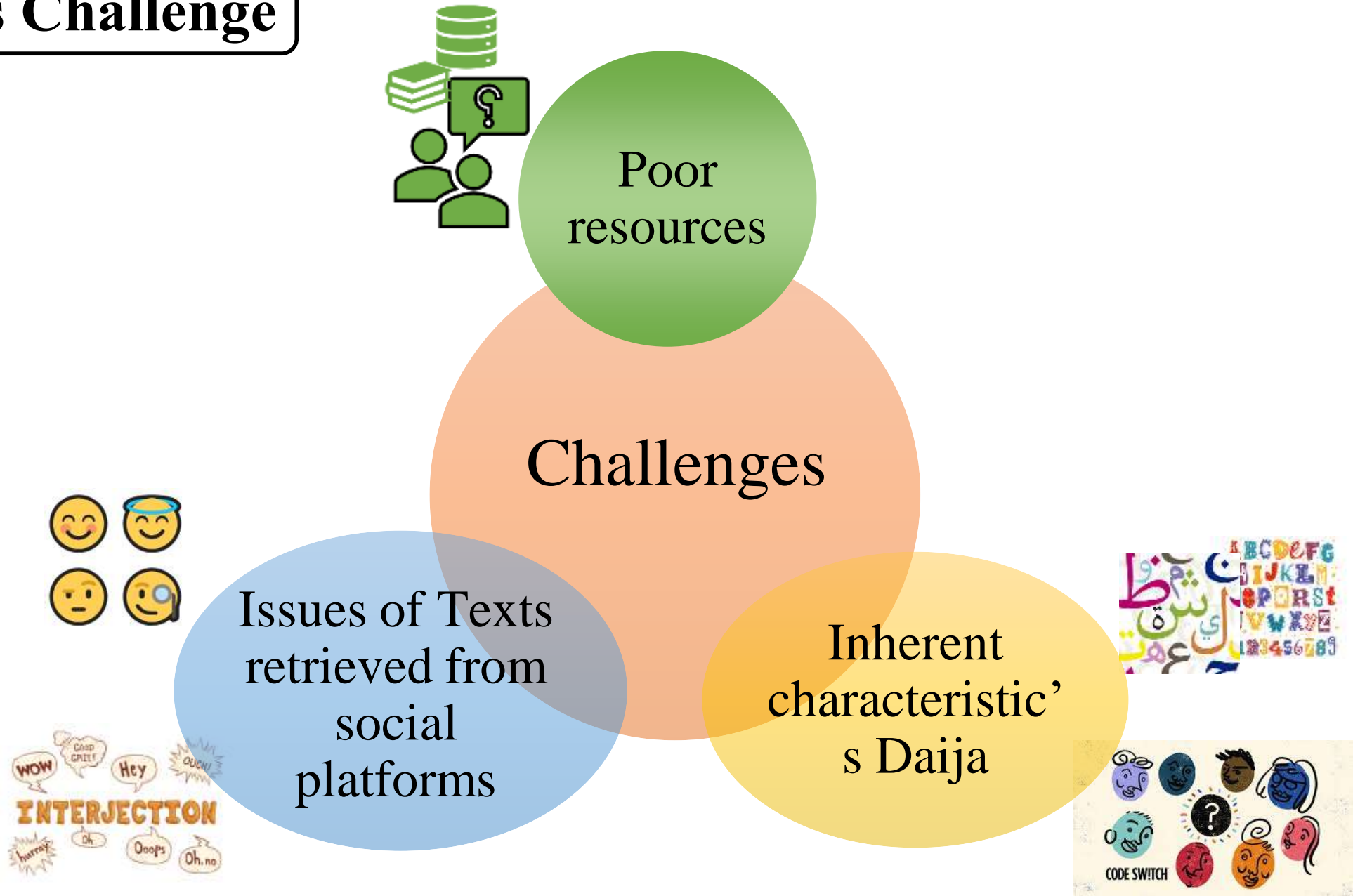


Reference article:

*Tachicart, Ridouane, Karim Bouzoubaa, and Hamid Jaafar. 2016. "Lexical differences and similarities between moroccan dialect and Arabic".*



# Darija's Challenge



# Darija's Challenge

## Poor resources

- The lack of available linguistic resources and annotated corpora significantly hinders progress.

## Issues of Texts retrieved from social platforms

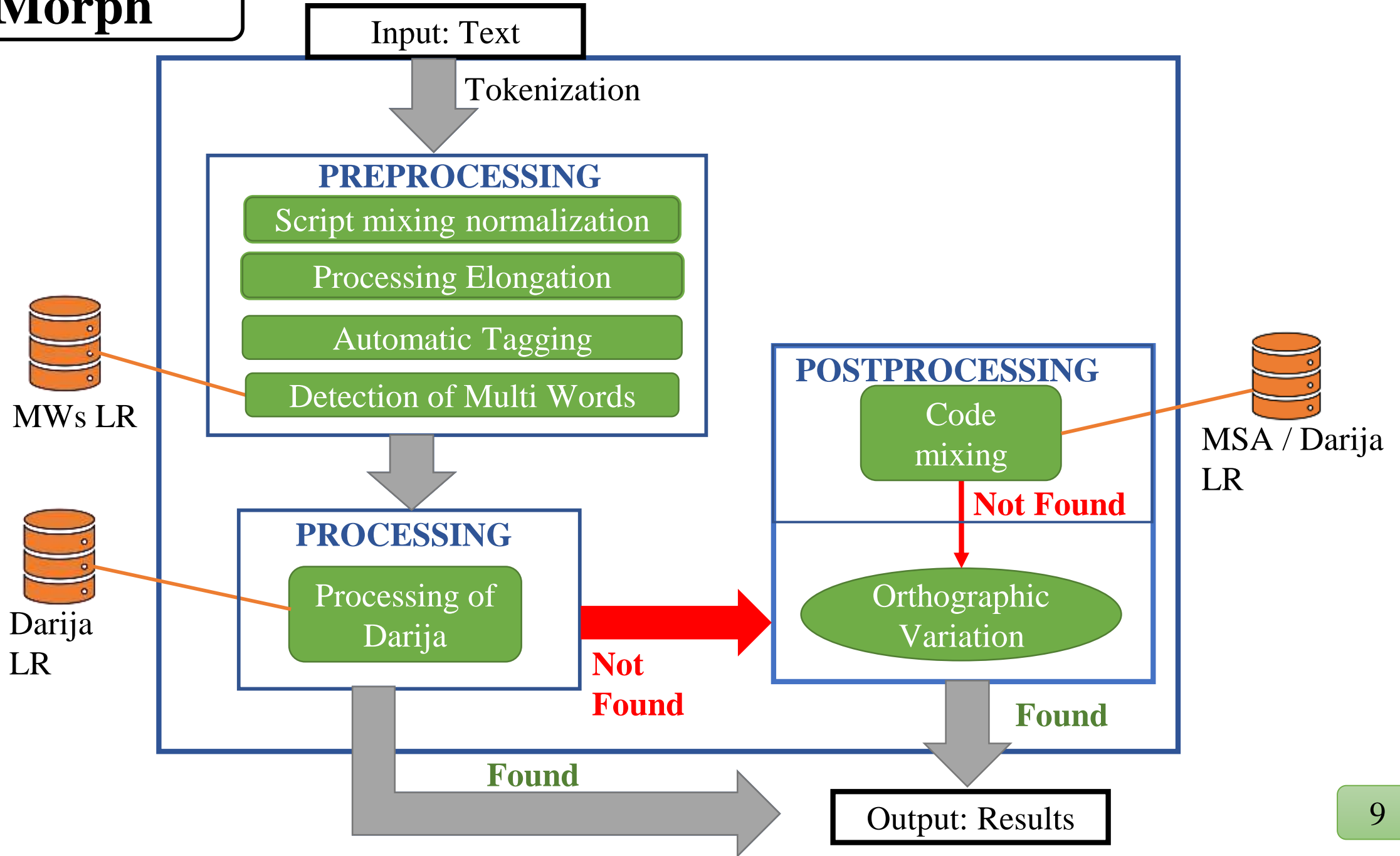
- Elongation: مبرووووك 'congratulation'
- Interjection: اوف 'ohf'
- Emoticon.

## Inherent characteristic's Daija

- Lack of Orthographic and Grammatical Rules:  
راس, رءس, رأس
- Early Technology Limitations: ع <- 3, ح <- 7
- Code/Script-mixing: غاطو



# DiMorph



# DiMorph



## Script mixing normalization

### Numeral character processing:

For example:

- ها 7 فت => فتها /ftaHħa:/ `he opened it'

### Literal character processing:

For example:

- طو g ا => گاطو /ga:Tu:/ `cake'



## Elongation processing

- بازااا bazza:a:f => بازا 'much'



## Automatic tagging

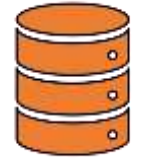
Implementation of an automatic tagging system to identify:

- Punctuation/number
- Emoticons
- Interjections
- Word Foreign

Where the foreign number or word is preceded by a dialectal prefix (e.g., ب/b/ - ف/f/ - و/w/ - ال/Al/).

For example:

- لPaola /l=Paola/ : ل/PREP+WORD\_FOREIGN `for Paola’.
- ب5000 /b=5000/ : ب/PREP+NUMBER `with 5000’.



قولي اسي أحمد اش طرى وجرا في قضية الفقيه بن صالح  
Tell me Mr Ahmed what happened concerning the Sidi Kacem case

Detection and  
segmentation

الفقيه بن صالح  
Fquih Ben Salah

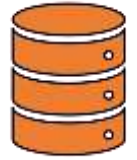
Compound proper noun

اش طرى وجرا  
what happened

Multi-word expression



## Multi-Word detection



MWs LR

	<b>Multi-Word expressions</b>	<b>Compound proper nouns</b>
Number of Units	153	511

## Multi-Word detection

الفقيه بن صالح  
Fquih Ben Salah

Compound proper noun

Analyzing

اش طرى وجرا  
what happened

Multi word expression

```
TOKEN:   الفقيه بن صالح  Alfqyh bn SALH
Arabic-msd: =الفقيه/B-GPE:CITY=/NSUFF+
Buckwalter-msd: =lafoqiyh/B-GPE:CITY=/NSUFF+
Arabic-msd: =بن/I-GPE:CITY=/NSUFF+
Buckwalter-msd: =ban/PROP_N I-GPE:CITY=/NSUFF+
Arabic-msd: =صالح/I-GPE:CITY=/NSUFF+
Buckwalter-msd: =SalaH/I-GPE:CITY=/NSUFF+
Vocalised PROP_N:   lafoqiyh ban SalaH
                  الفقيه بن صالح
Glosses:   Fquih Ben Salah+
```

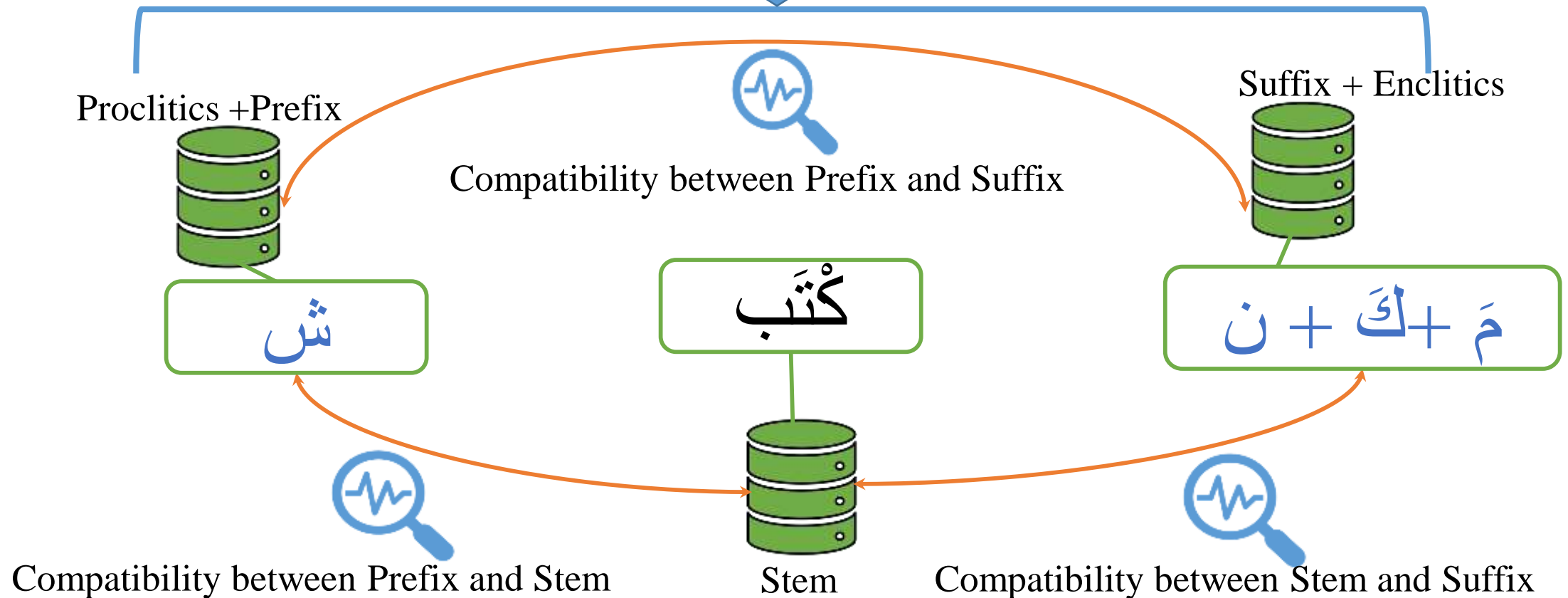
```
TOKEN:   اش طرا وجرا  A$ TrA wjrA
Vocalised Token:   أَشْ طَرَا وَجَرَا
Arabic-msd:   أَشْ طَرَا وَجَرَا/MULTIWORD+
Buckwalter-msd: >a$ ToraA wajorA/MULTIWORD+
Glosses:   what happen?
```

# DiMorph

Processing

مَكْنَزُ كُتَبِشْ

~~DET + Verb~~







## Linguistic Resources & Processing

DictStem	NOUNS	VERBS	ADJECTIFS	ADVERBS	PROPER NOUNS	PRONOMS	FUNCTION WORDS
<b>Darija</b>	8760	11647	1674	214	931	92	153
<b>Foreign</b>	675	122	33	8	100	-	3

DictPrefix	DictSuffix	Compatibility Tables		
proclitics + prefixes	suffixes + enclitics	Prefix=Stem	Stem=Suffix	Prefix=Suffix
394	641	770	780	1067



**MSA Token Processing in Darija:** Detecting Code-Mixing Tokens in DiMorph through Analysis of Clitics in Darija and Stem in MSA.

For example:

کيستفزاها /kayastafazzha/ “He is provoking her”

### Orthographic Variation

- Standardize written forms according to orthographic rules, ensuring that spelling variations are unified into a single, consistent form.

Type	Variation	Example
Hamza Normalization	unified form → ء, أ, إ, ئ, و	سأل → سأل, سأل
Phonological Shifts	ث → ت	ثلاثة → ثلاثة
	ذ → د	ذهب → ذهب
	ظ → ض	ظلام → ضلام
Word-Final Changes	و → ه	عنده → عندو
	ة → ا	ركبة → ركبا

## Results & discussion

### Evaluation

- **In-Vocabulary Rate (INV Rate):** Measures the percentage of tokens successfully analyzed by the system.
- **Out-of-Vocabulary Rate (OOV Rate):** Measures the percentage of tokens the system could not analyze.

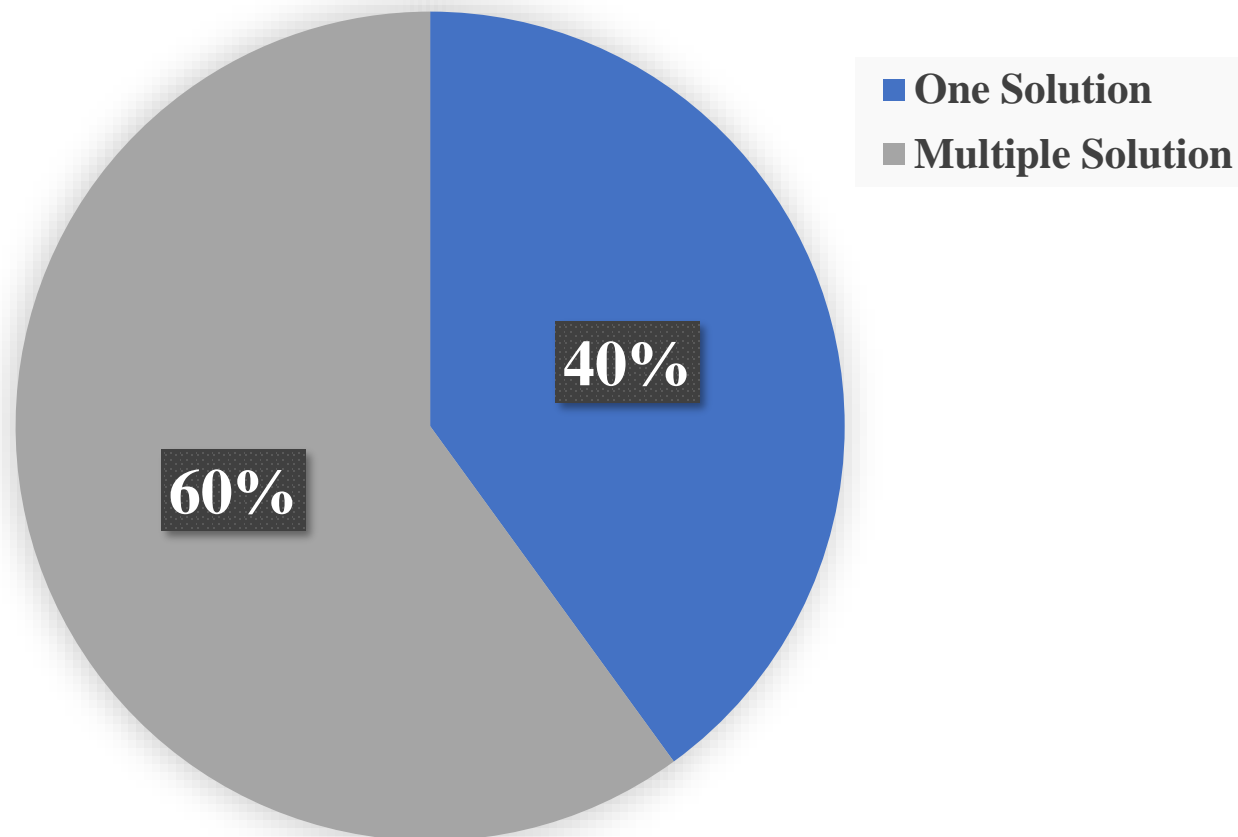
	<b>Total Tokens</b>	<b>INV rate</b>	<b>OOV rate</b>
DiMorph with Preprocessing and Postprocessing	11 085	<b>97.84%</b>	2.16%

# Results & discussion

## Evaluation

Further statistical analysis reveals that:

### Ambiguity Rate



- 89.71% of the cases corresponds to instances of homography.
- 9.31% of the cases corresponds to instances of polysemy.
- 0.98% of the ambiguity is due to the absence of vocalization.

# Conclusion

## DiMorph's Analysis Capacity

### Average Analyses per Token:

- DiMorph provides, on average, 1.45 possible analyses for each token.
- This metric highlights the system's current capability to generate multiple solutions per token, indicating flexibility in interpretation.

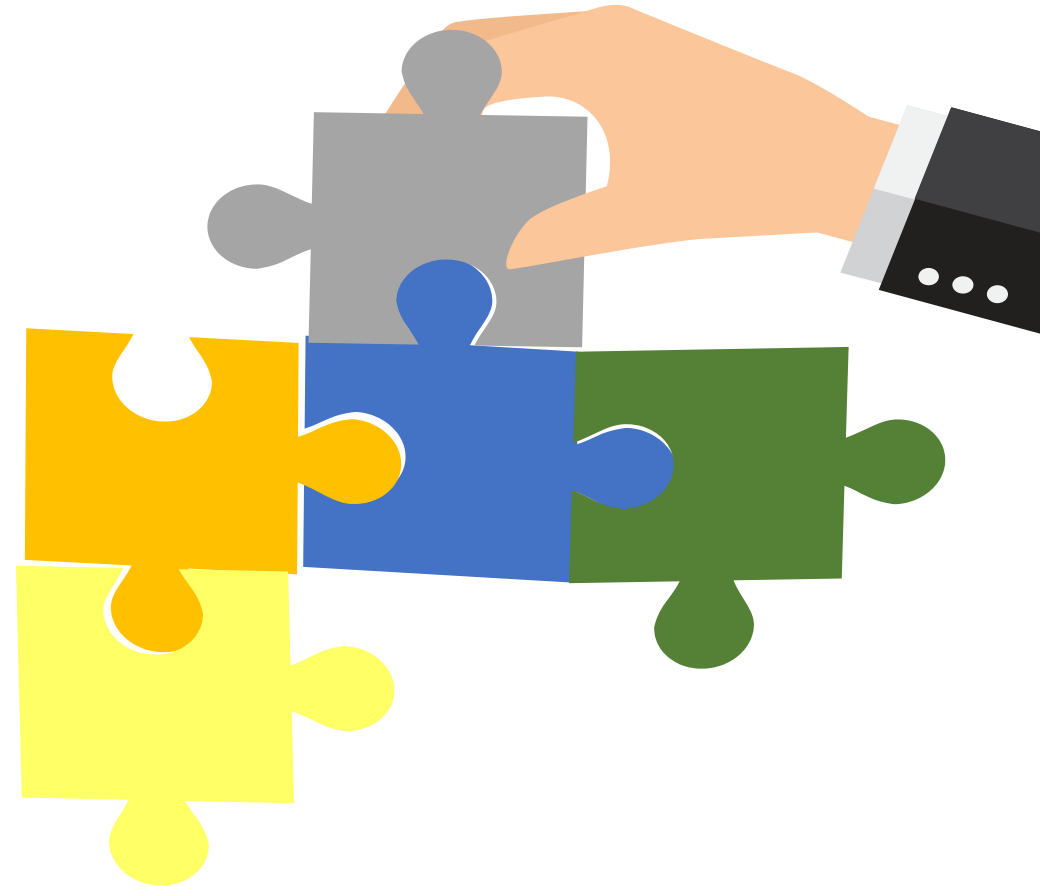
### Challenges:

- Out-of-Context Issue: While DiMorph generates multiple analyses, it often lacks the ability to determine the correct analysis in context, leading to potential ambiguity.

# Perspectives

● Enrich the Moroccan DiMorph linguistic resources.

● Apply deep learning models to provide context-aware solutions and accurately annotate the Moroccan corpus.





*Thank you for  
your attention*

Nadia khlif  
Azzeddine Mazroui  
Ouafae Nahli

[nadia.khlif@ump.ac.ma](mailto:nadia.khlif@ump.ac.ma)  
[azze.Mazroui@gmail.com](mailto:azze.Mazroui@gmail.com)  
[ouafae.nahli@ilc.cnr.it](mailto:ouafae.nahli@ilc.cnr.it)

