# Challenges and Advances in Constructing Arabic Dialect Corpora and Linguistic Tools for the Moroccan Dialect
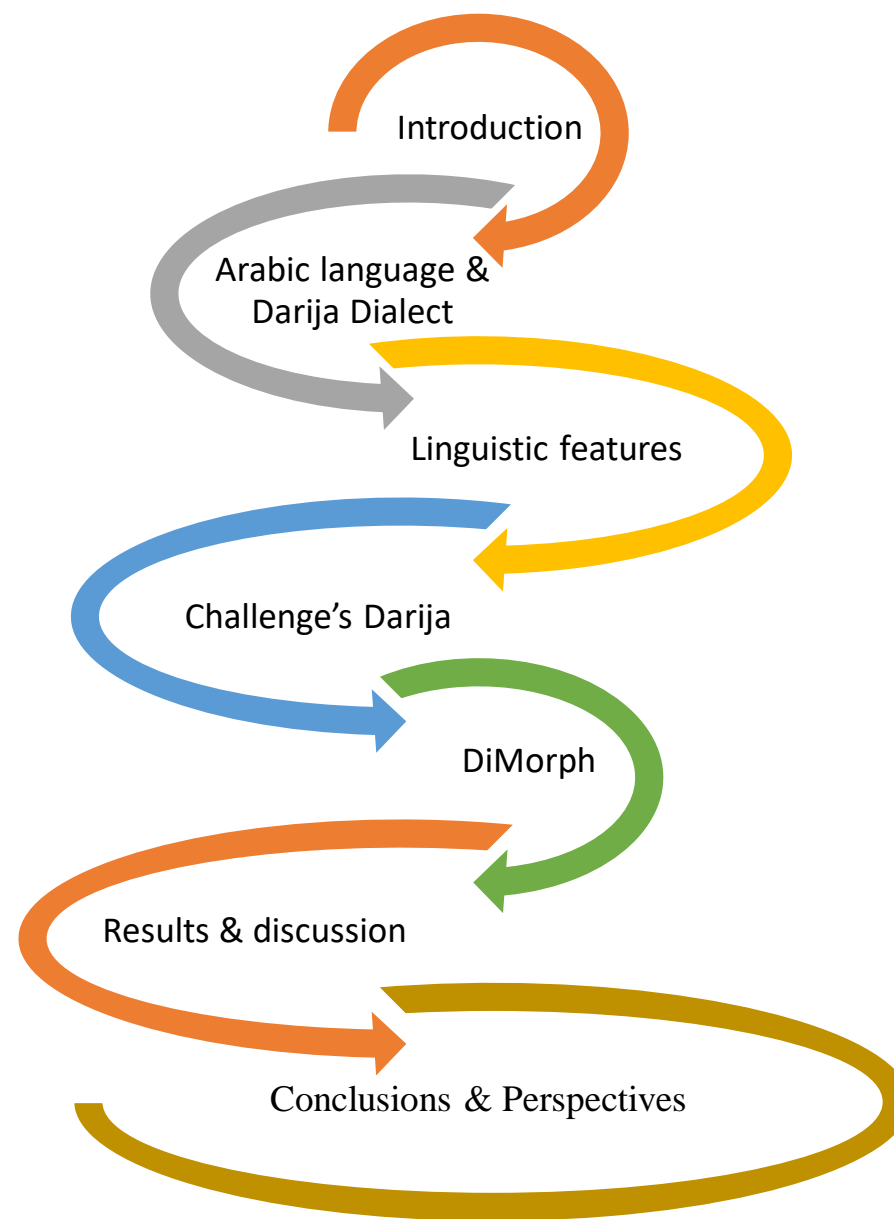
Nadia khlif
nadia.khlif@ilc.cnr.it

AIUCD2024

07/05/2024

PLAN

Introduction

Arabic language & Darija Dialect

Linguistic features

Challenge's Darija

DiMorph

Results & discussion

Conclusions & Perspectives

A LEXICAL CORPUS-BASED MODEL OF CONTEMPORARY WRITTEN ARABIC

**Data Collection**

**Step1:**

Collection of a Representative Corpus: establish a genuinely representative corpus of Colloquial Arabic Varieties.

**Tools: DiMorph**

**Step2:**

Tools Adaptation: enhance linguistic annotations within our corpora through the adaptation of the morphological Analyzer Aramorph for processing written dialectal words.

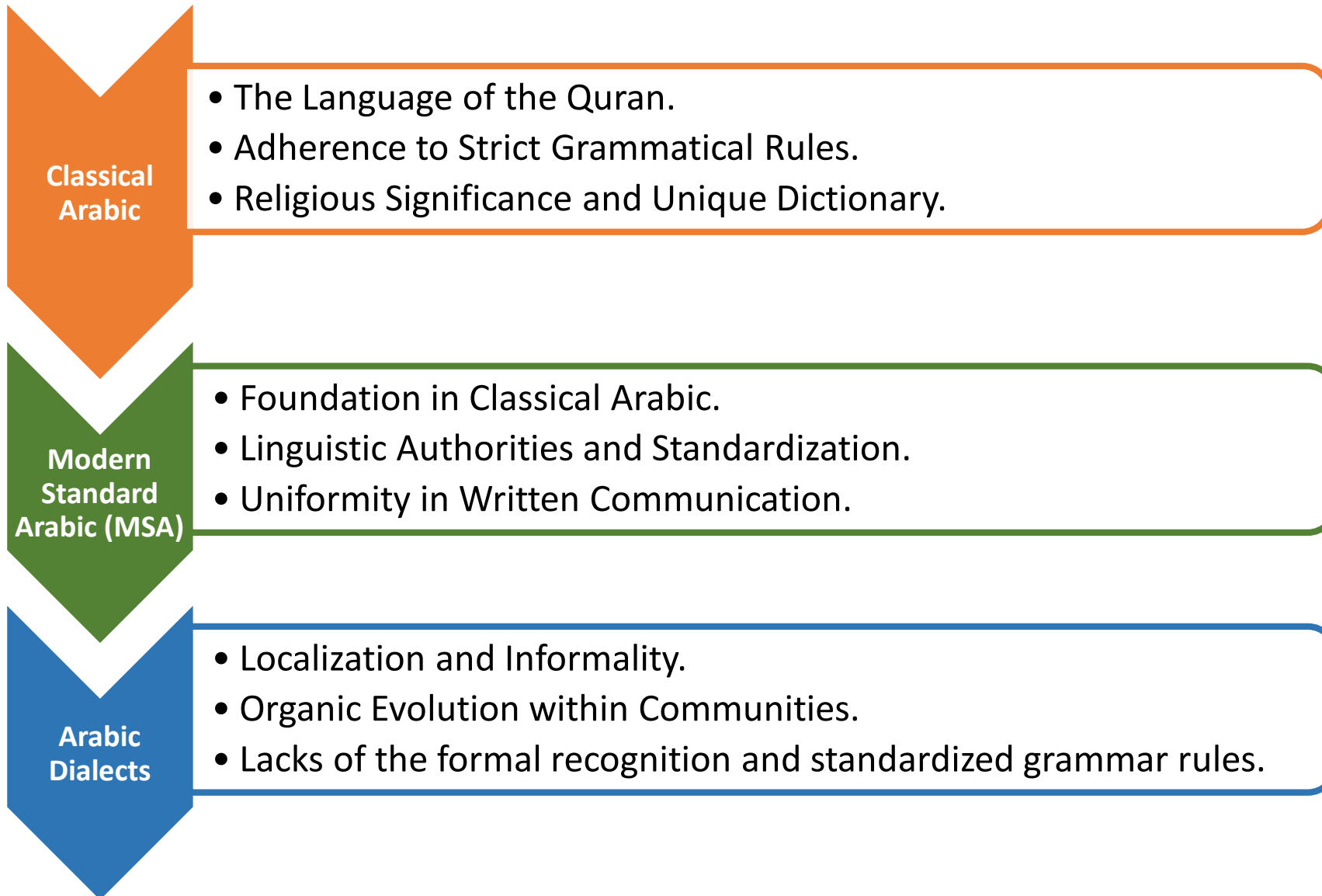**Manual disambiguation of a Subcorpus**

**Step3:**

Corpus Annotation: manual disambiguation of a subcorpus collected and analyzed to adapt methodologies in deep learning for the Automatic Annotation of the Entire corpus.
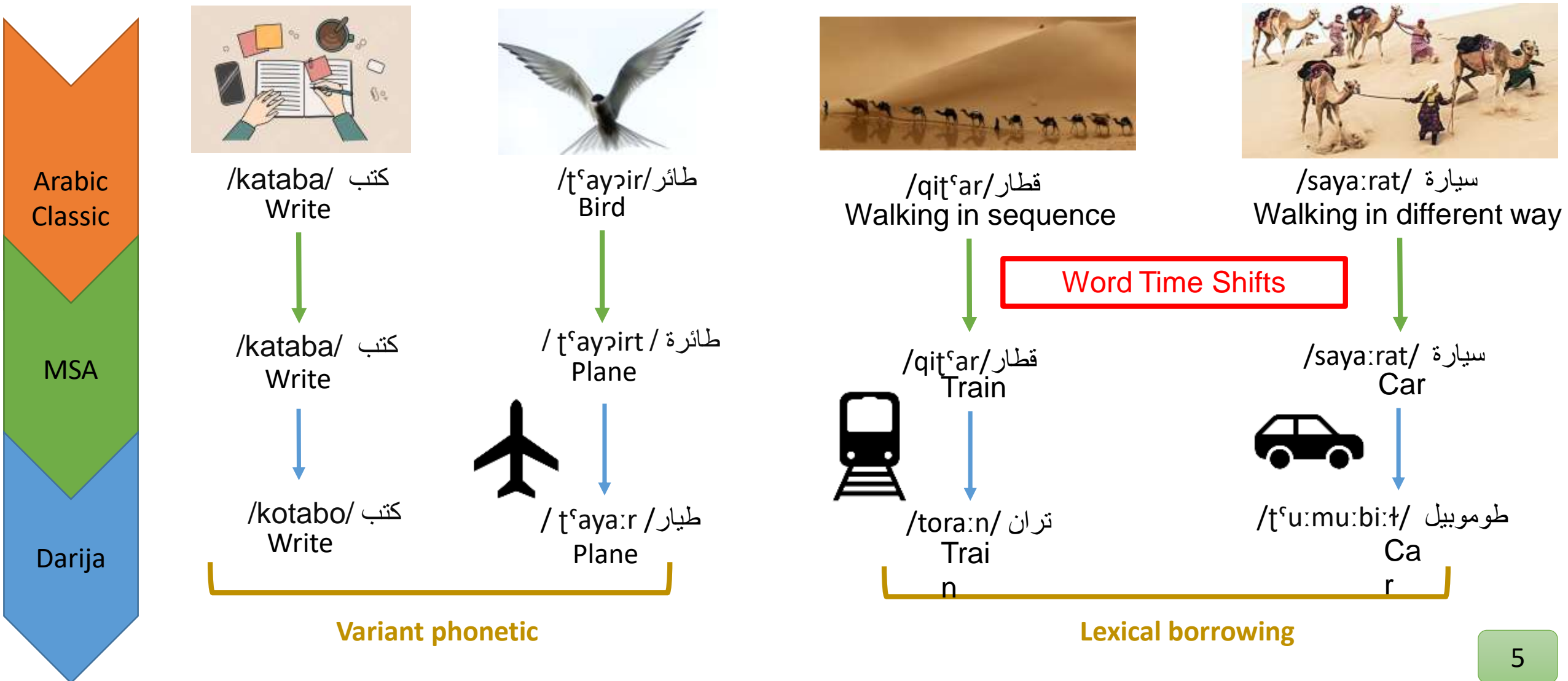
**Lexical Model**

**Step4:**

Developing a Lexical Model: Bridging Corpus Data and Existing Lexical Sources.
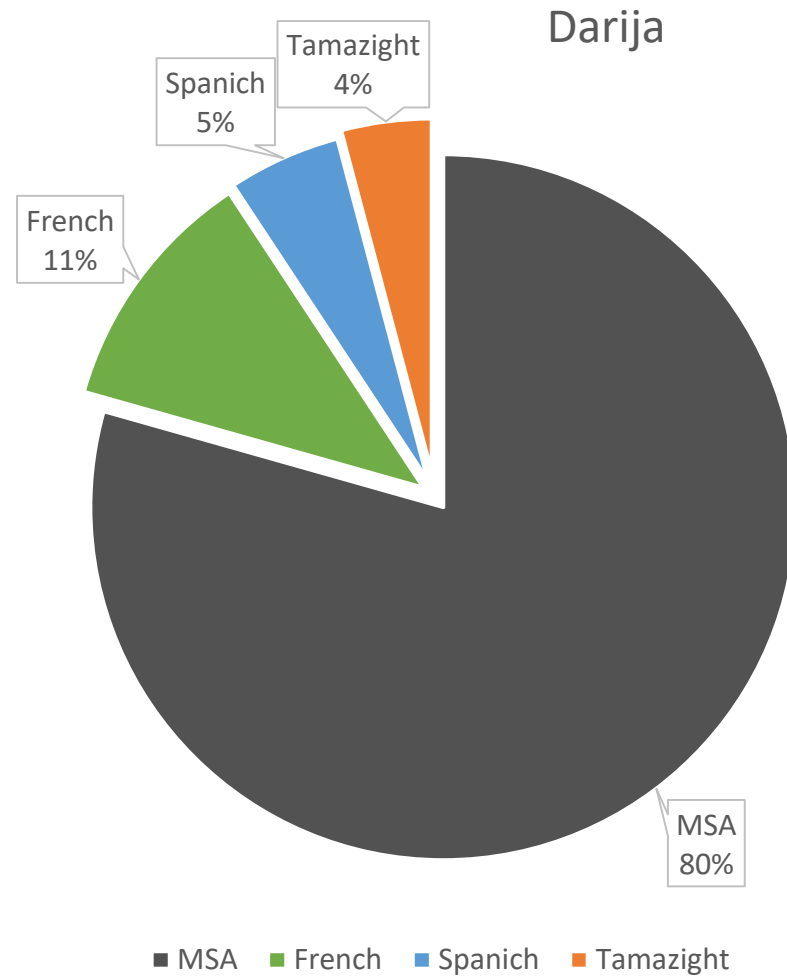
# Arabic Language & Darija dialect

**Classical Arabic**
- The Language of the Quran.
- Adherence to Strict Grammatical Rules.
- Religious Significance and Unique Dictionary.

**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.
- Lacks of the formal recognition and standardized grammar rules.

4

# Arabic Language & Darija dialect

**Arabic Classic**

كتب /kataba/
Write

طائر/ṭʕayʔir/
Bird

قطار/qiṭʕar/
Walking in sequence

سيارة /saya:rat/
Walking in different way

**MSA**

كتب /kataba/
Write

طائرة / ṭʕayʔirt /
Plane

Word Time Shifts

قطار/qiṭʕar/
Train

سيارة /saya:rat/
Car

**Darija**

كتب /kotabo/
Write

طيار / ṭʕaya:r /
Plane

تران /tora:n/
Train

طوموبيل /ṭʕu:mu:bi:ł/
Car

**Variant phonetic**

**Lexical borrowing**

5

## Arabic Language & Darija dialect

Darija



Spanich 5%

Tamazight 4%

French 11%

MSA 80%

■ MSA  ■ French  ■ Spanich  ■ Tamazight

Article de reference:
Tachicart, Ridouane, Karim Bouzoubaa, and Hamid Jaafar. 2016. "Lexical differences and similarities between moroccan dialect and Arabic".
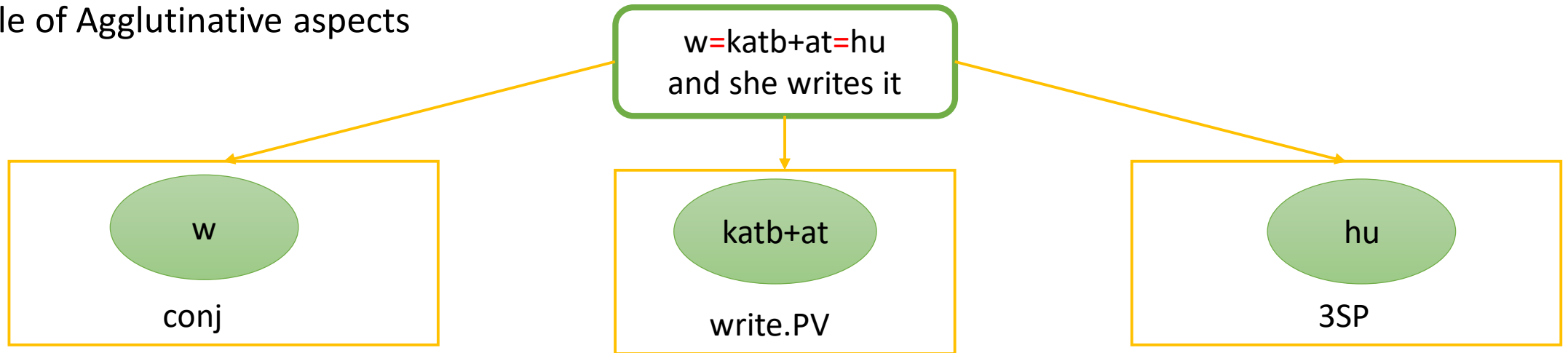
6

- The morpho-syntactic layer combines the inflected form with clitics (prepositions, conjunctions, definite articles, etc.) to shape a rich and complex surface form.

- The inflectional layer is the one where the lemma combines with inflectional affixes to give inflectional forms.

- The derivation layer is the deepest one. At this level, the root combines with the vowels, according to determined patterns, to produce a verbal or a nominal lemma.
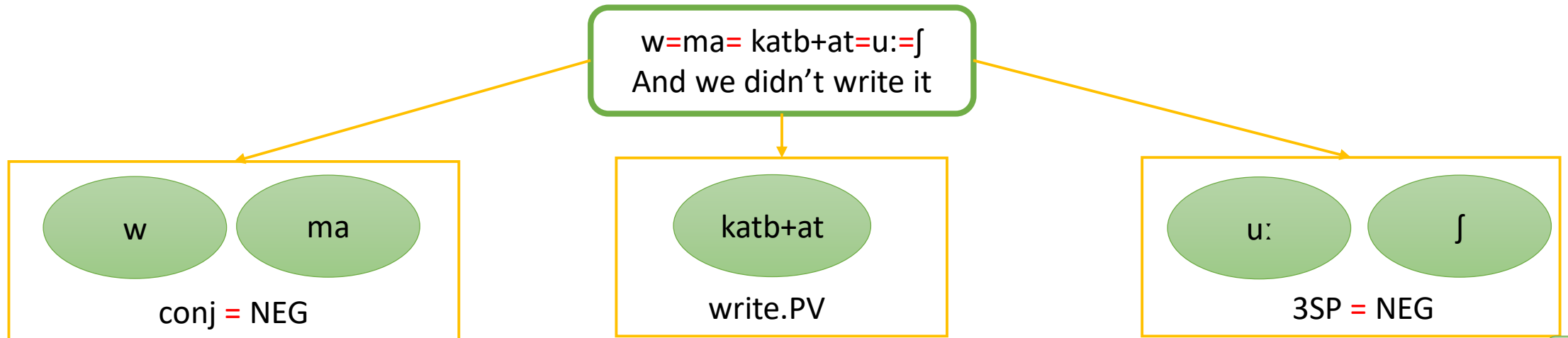


Surface realisation

Alteration Rules
Clitics
Verbal Inflected Form
Nominal Inflected Form

Alteration Rules
Inflection Affixes
Verbal Lemma
Nominal Lemma

Alteration Rules
Pattern
Root
Vocalic

Surface realisation

Alteration Rules

Clitics

Verbal Inflected Form

Nominal Inflected Form

Alteration Rules

Inflection Affixes

Verbal Lemma

Nominal Lemma

Alteration Rules

Pattern

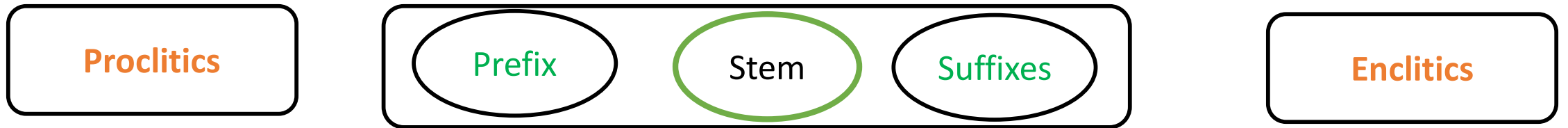Root

Vocalic

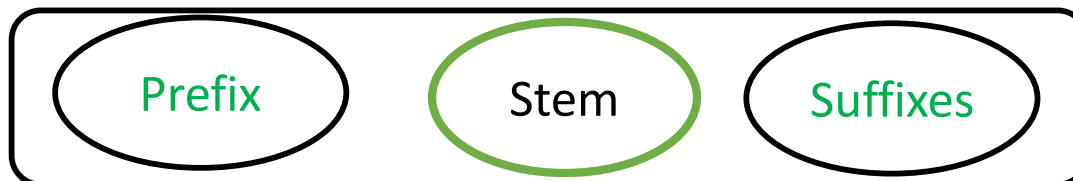# Linguistic features

Example of Agglutinative aspects



Example of Grammatical Aspects



8

In the two examples, the inflected form is surrounded by clitics and the morphological structure is:

**Proclitics**

Prefix   Stem   Suffixes

**Enclitics**

By removing clitics, the remaining word form is a minimally autonomous inflected form whose structure consists of:

Prefix   Stem   Suffixes
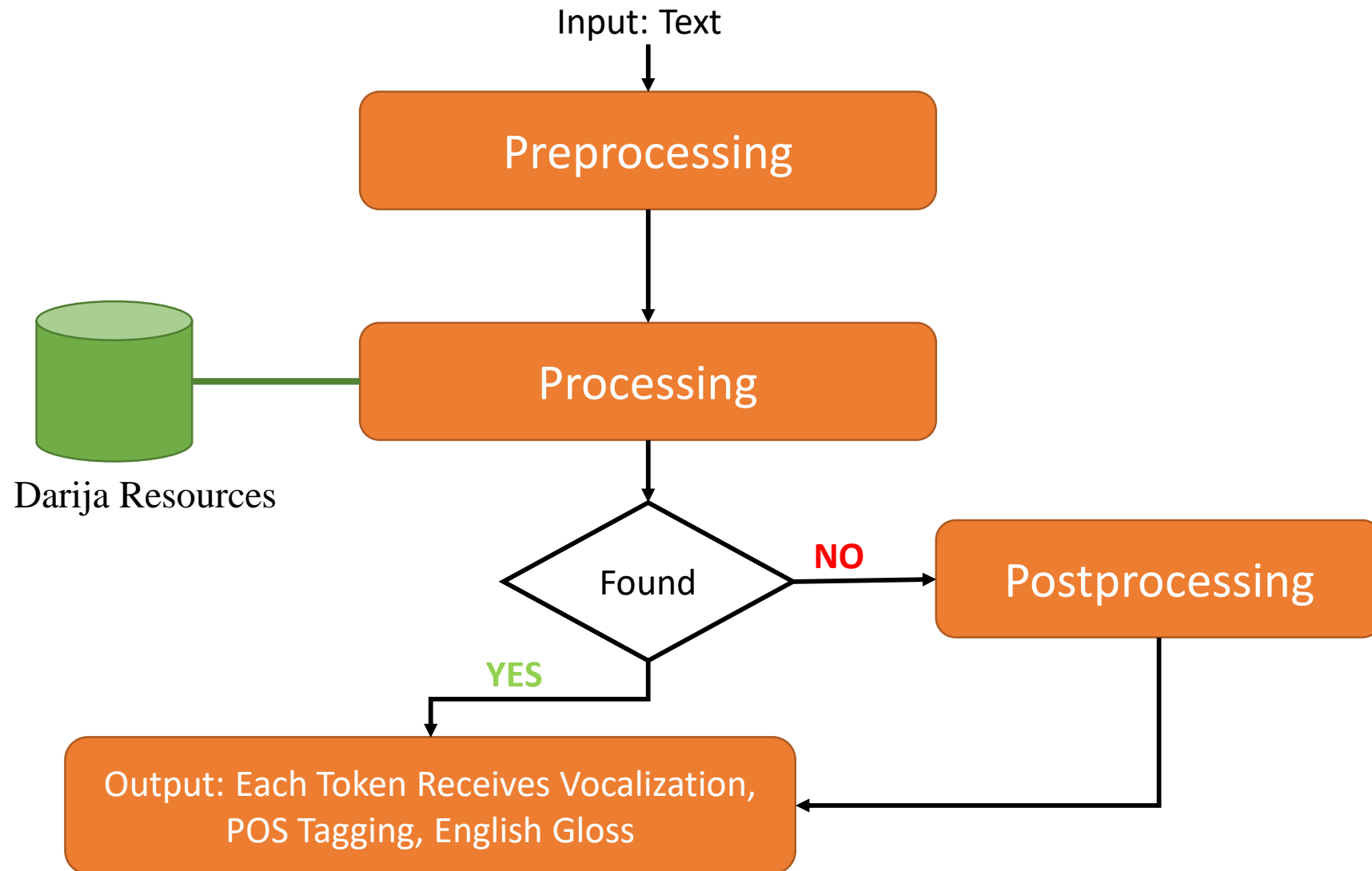
9

Challenge's Darija
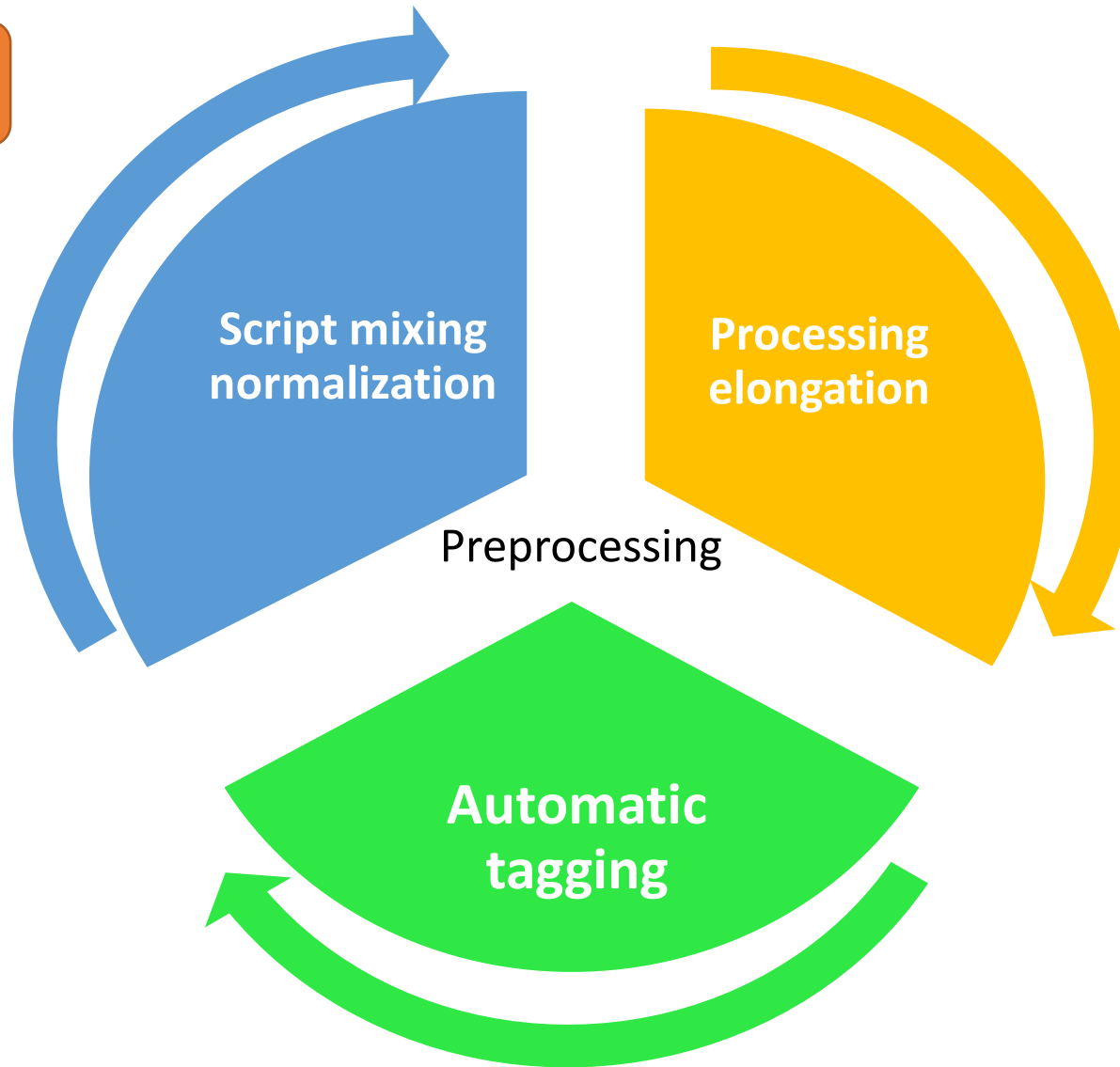
Poor resources

Challenges

Issues of Texts retrieved from social platforms

Inherent characteristic 's Daija

## Challenge's Darija

Examples :

| | Phenomenon | Examples | English |
|---|---|---|---|
| Script mixing | Is a mixing letter between Arabic and Latin | gال<br>قال<br>/qa:l/ | 'say' |
| Code switching | Is a switching between Moroccan dialect and French words 'les valises'. | كنخوي لي فاليز<br>/kanaxuːiː liː faːliːz/ | 'I unload the luggage' |
| Code mixing | Is a French word 'La valise' with MSA enclitic 'ال' and feminine suffix 'ة'. | الفاليزة<br>/aː=lfaːliːz-t/ | 'The luggage' |

11

DiMorph

Input: Text

Preprocessing

Processing

Darija Resources

Found

NO

Postprocessing

YES

Output: Each Token Receives Vocalization, POS Tagging, English Gloss

12

DiMorph

Preprocessing

Script mixing normalization

Processing elongation

Preprocessing

Automatic tagging

13

Processing elongation

| Elongated word | Normalized word | English translation |
|:---:|:---:|:---:|
| مبرووووك<br>mabru:u:u:u:k | مبروك<br>mabru:k | `congratulation' |
| بزااف<br>bazza:a:f | بزاف<br>bazza:f | `much' |

14

Automatic tagging

Implementation of an automatic tagging system to identify:
- Punctuation/ number.
- Emoticons.
- Interjections.

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'.

Script mixing normalization

Automated Substitution for One-to-One Correspondences: When each number or Latin letter can be directly replaced by an Arabic letter, we implement automated substitution. For example:

- فت7ها <= فتحها   /ftaHħa:/  `he opened it'

## DiMorph

Semi-Automated Substitution for One-to-Many Correspondences: The specific case involving the Latin letter "g" or the Persian letter "گ", both representing the phoneme /g/, isn't automatically handled due to the varied correspondences with Arabic phonemes such as "ق" /q/ or "ج" /j/. Instead, these cases are managed at the lexical level, considering the multiple possible phonetic mappings of the previous letter /g/.For example:

| origin | Term in Darija | Cognate Arabic | Script in Arabic and Latin alphabet | Script in Arabic and Persian alphabet | English translation |
|--------|----------------|----------------|-------------------------------------|---------------------------------------|---------------------|
| Arabic | /ga:l/ | قال/qa:/ | gال | گال | 'say' |
| Arabic | /glas/ | جلس/ʒlas/ | gلس | گلس | 'sit' |
| Foreign | /ga:tˤu:/ | | gاطو | گاطو | 'cake' |
| Foreign | /ga:ʕ/ | | gاع | گاع | 'never' |

17

## DiMorph

**Linguistic Resources**

| Minimal word | | | **Prefix** | **Stem** | **Suffix** | | |
|---|---|---|---|---|---|---|---|
| Maximal word | Proclitic1 | Proclitic2 | Prefix | Stem | Suffix | Enclitic1 | Enclitic2 |
| | **DictPrefix** | | | **DictStem** | | **DictSuffix** | |

DiMorph

Linguistic Resources & Processing

|  | NOUNS | VERBS | ADJECTIFS | ADVERBS | PRONOMS | FUNCTION WORDS |
|---|---|---|---|---|---|---|
| Darija | 5169 | 3132 | 1128 | 146 | 28 | 156 |
| Foreign | 295 | 28 | 16 | 6 | - | 4 |

| DictPrefix | DictSuffix | Compatibility Tables | | |
|---|---|---|---|---|
| proclitics + prefixes | suffixes + proclitics | AB | BC | AC |
| 297 | 570 | 770 | 780 | 1067 |

Input: Text

Tokenization

**PREPROCESSING**

Script mixing normalization

Processing Elongation

Automatic Tagging

**ANALYZING**

Processing of Darija

Darija LR

Not Found

Found

Output:
Vocalization,
POS Tagging,
English Gloss

**POSTPROCESSING**

Code switching

MSA LR

Not Found

Code mixing

MSA/Darija LR

Not Found

Orthographic Variation

Not Found

Found

Manual correction

Found

Postprocessing

MSA Detection in DiMorph: Identifying Code-Switching Tokens Not Analyzed as MSA.
For example:

- سنستورد **/sa=na-stawrid-u/  "We will import".**

MSA-Darija Identification: Detecting Code-Mixing Tokens in DiMorph through Analysis of Clitics in Darija and Stem in MSA.
For example:

- غنستوردو **/ʁa=nstawrd-uː/  "We will import".**

Code switching

MSA Resources

Code mixing

MSA : DictStem +
Darija: -DictPrefix
        - DictSuffix
        - Compatibility Tables

Postprocessing

Orthographic variation

**Phonological simplification:** This refers to the process whereby a language alters its sound system, resulting in the reduction or elimination of certain phonetic features.

| MSA | *Darija* | script 1 | script 2 | English translation |
|---|---|---|---|---|
| [kaθar-a] | [ktar] | كْثَر [kθar] | كْتَر [ktar] | abound |
| [ðˤalaːm] | [dˤlaːm] | ظْلَام [ðˤlaːm] | ضْلَام [dˤlaːm] | darkness |
| [ʕaðaːb] | [ʕdaːb] | عْذَاب [ʕðaːb] | عْذَاب [ʕdaːb] | tribulation |

**Postprocessing**

📎 Orthographic variation

**Strict spelling conventions:** include the representation of the glottal stop, known as the *hamza*.

| precedent vowel | lengthening letter | MSA script 1 | phonetical script 2 | English translation |
|---|---|---|---|---|
| [a] | ا | رَأْس [raʕs] | رَاس [raːs] | 'head' |
| [u] | و | مُؤْمِن [muʕmin] | مُومَن [muːman] | 'believer' |
| [i] | ي | دَافِئ [daːfiʕ] | دَافِي [daːfiː] | 'warm' |

**DiMorph**

Postprocessing

Manual correction

**The spelling errors:** were identified as 'Not Found' during analysis and subsequently they were corrected by the annotators.
For example:

- تدحك **/tadħak/** is corrected in تضحك **/tadˤħak/** `she laughs'.
- عندومشاكيل **/ʕanduːmaʃaːkiːl/** is corrected in عندو مشاكيل **/ ʕanduː maʃaːkiːl /** `he has a problem'.

Size of the corpus is: 91.592 tokens

## Results discussion

Size of the corpus is: 91.592 tokens
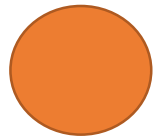
# Results discussion

Size of the corpus is: 91.592 tokens



DiMorph

1%

99%

Found
Not Found

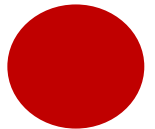# Conclusion

The creation of a morphological analyzer customized to dialectal characteristics.

The construction of linguistic resources specifically for dialects.
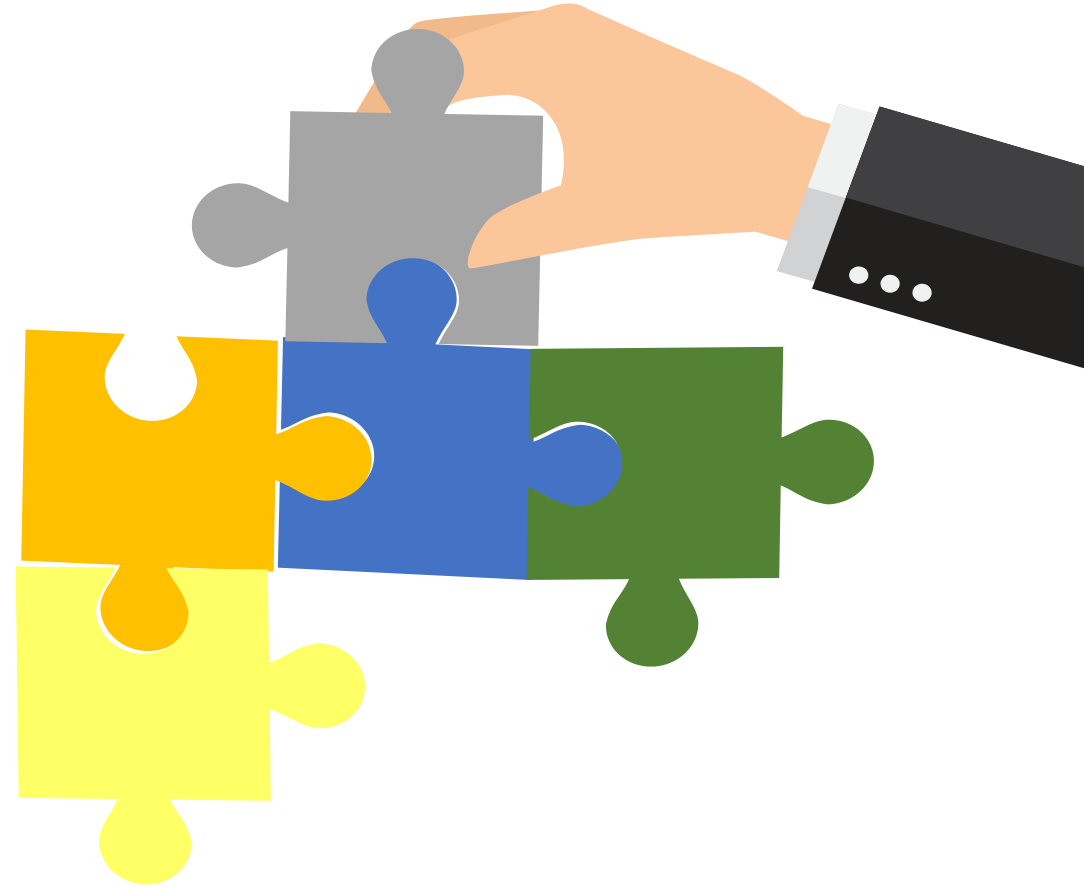
The construction of corpus annotated and vocalized.

Enrich the Moroccan DiMorph linguistic resources.

Apply deep learning strategies to annotate the Moroccan corpus.

# *Thank you for your attention*

Nadia khlif          nadia.khlif@ilc.cnr.it
Giulia Benotto       giulia.benotto@ilc.cnr.it
Ouafae Nahli         ouafae.nahli@ilc.cnr.it