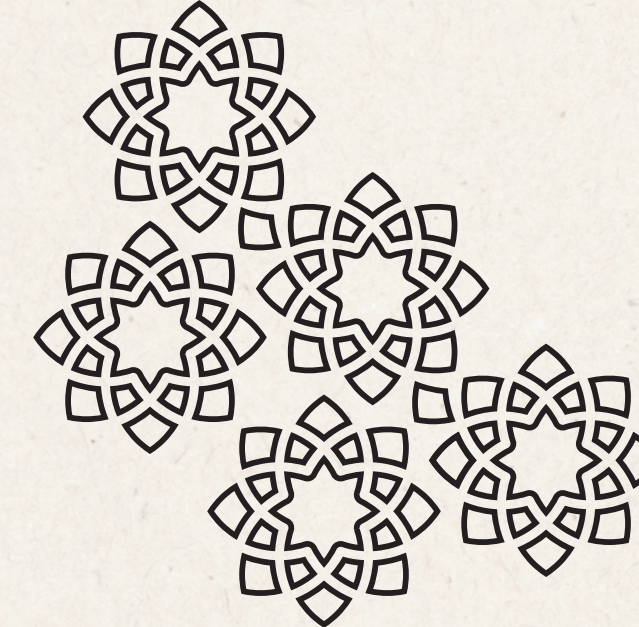




Institut für Maschinelle Sprachverarbeitung |
Summer Semester 24



LINGUISTIC ROADMAP OF DARIJA

Helping Darija Linguists Grow

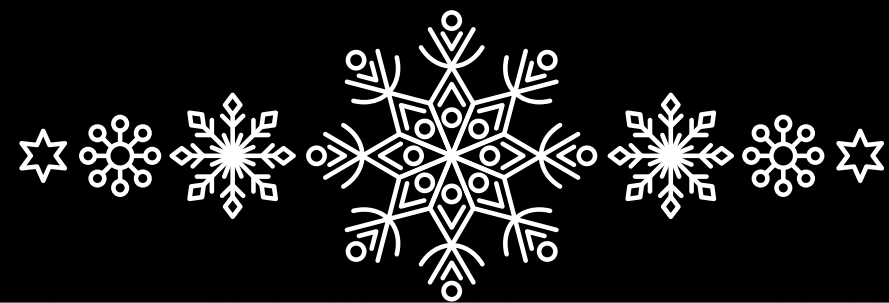
COURSE:
Text Technology

STUDENTS:
Yassir EL ATTAR & Shawn Chua

PROFESSOR:
Kerstin Jung



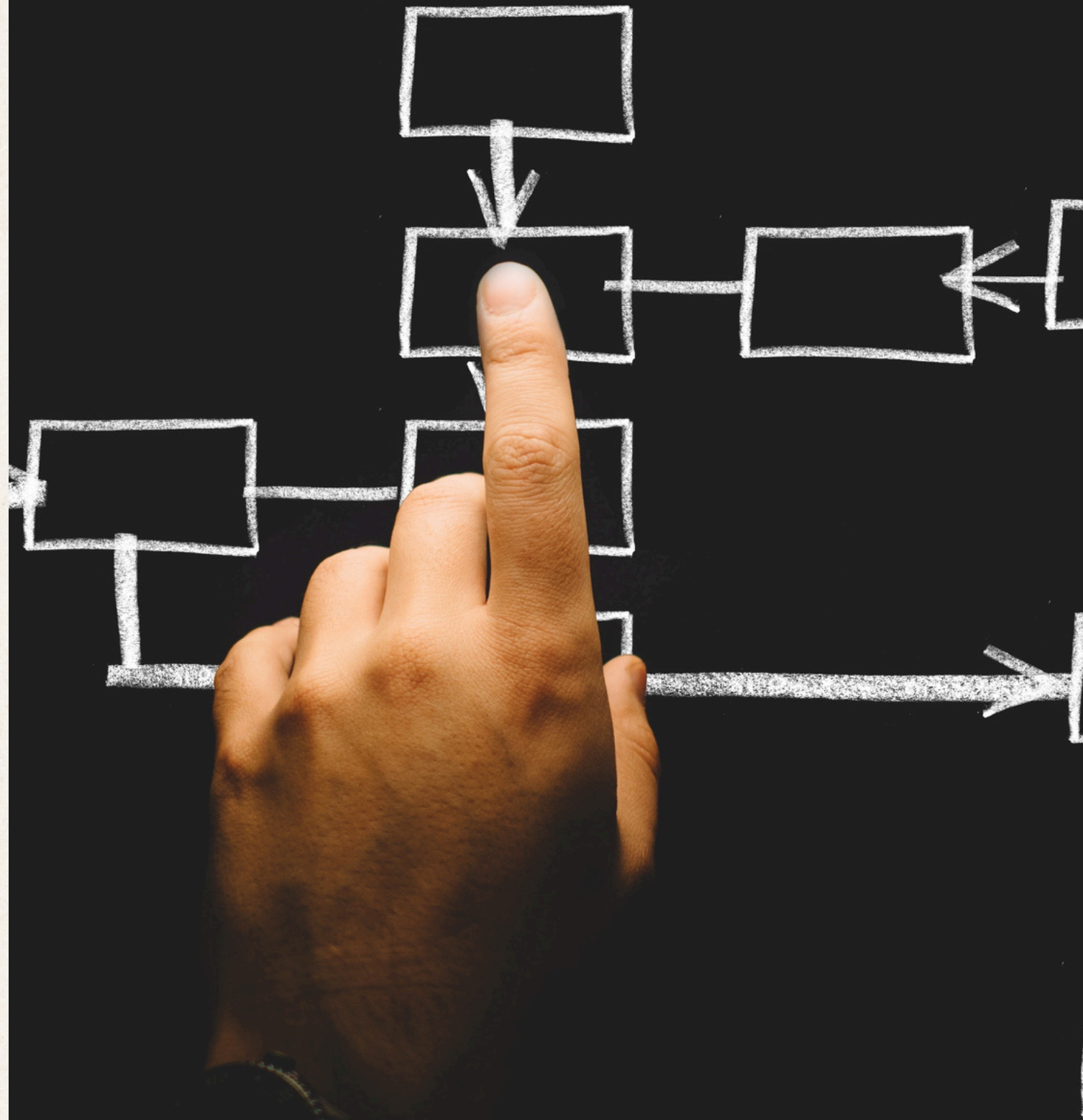
Agenda



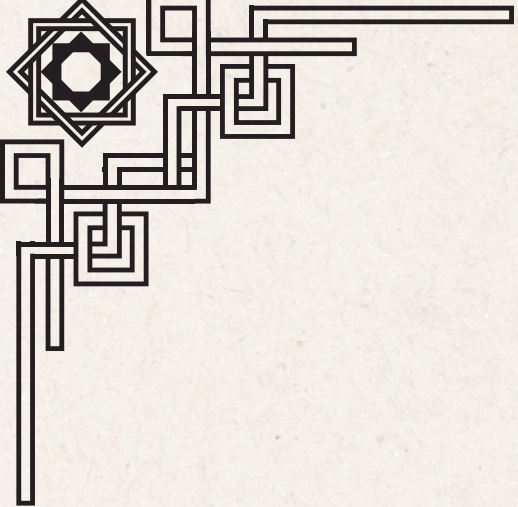
03	Overview
05	Objectives
06	Timeline
07	Collect
09	Prepare
11	Access
12	Extensions

Overview

Where did this
idea come
from?



- 01 Discussion of our native languages and their variations
- 02 Darija is a new and developing language in Morocco
- 03 Many linguistics aspects of Darija are still uncovered

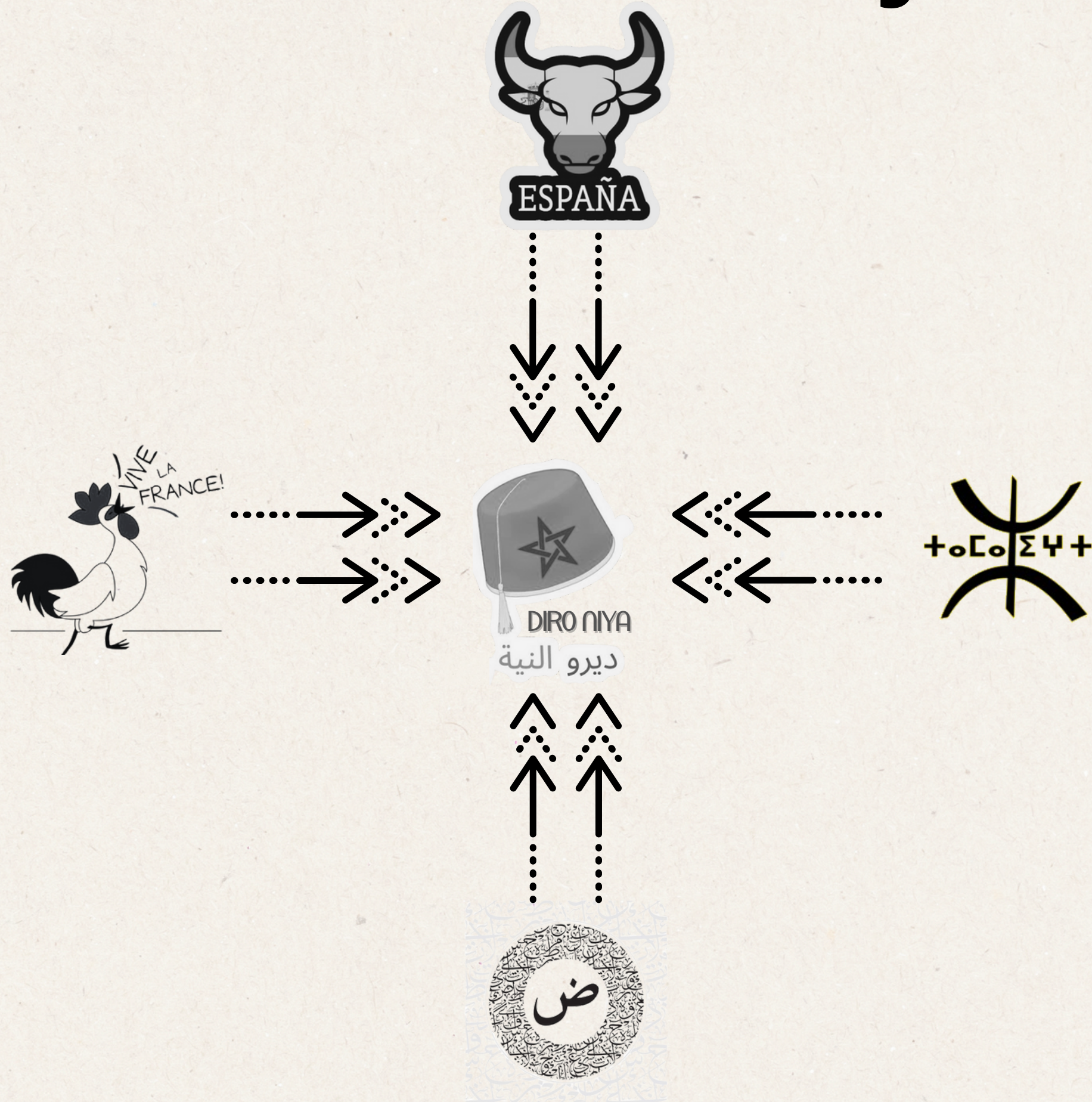


Why Darija?

04/13

What resulted?

Arabic	Latin	Darija_Ar	Darija
ق	-	ق	9/ q
ح	-	ح	7
ع	-	ع	3
س	-	ص	S
ص	-	ص	D
ط	-	ط	T
-	g	ڭ	g
-	p	پ	p
-	v	ف	v



Objectives and Goals

By the end of our project, we will create a platform to:



Goal # 1

Be an interactive database for linguists interested in Darija



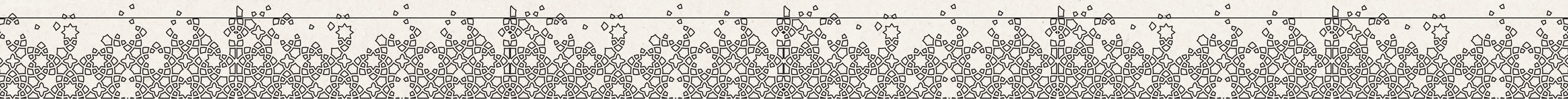
Goal # 2

Improve the current Dataset and solve the common problems



Goal # 3

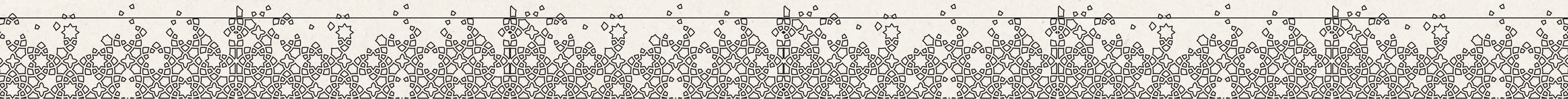
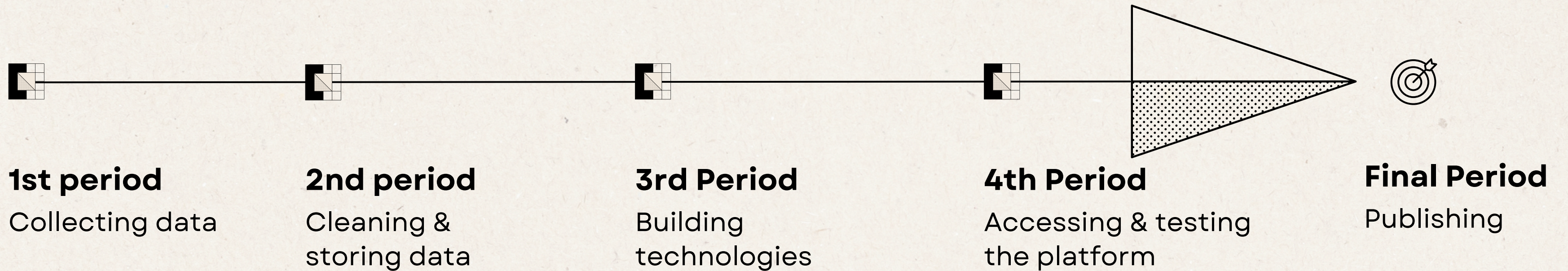
Help linguists contribute and access the database

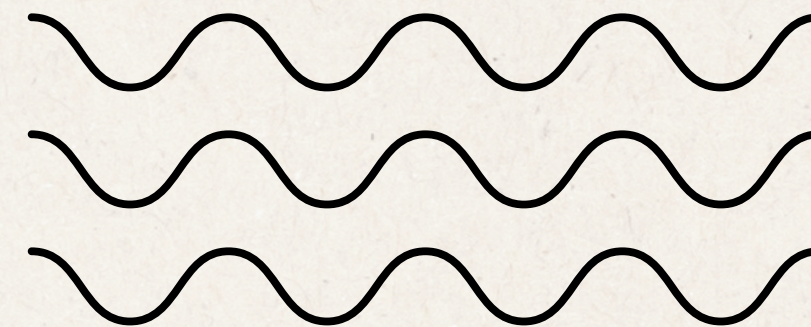
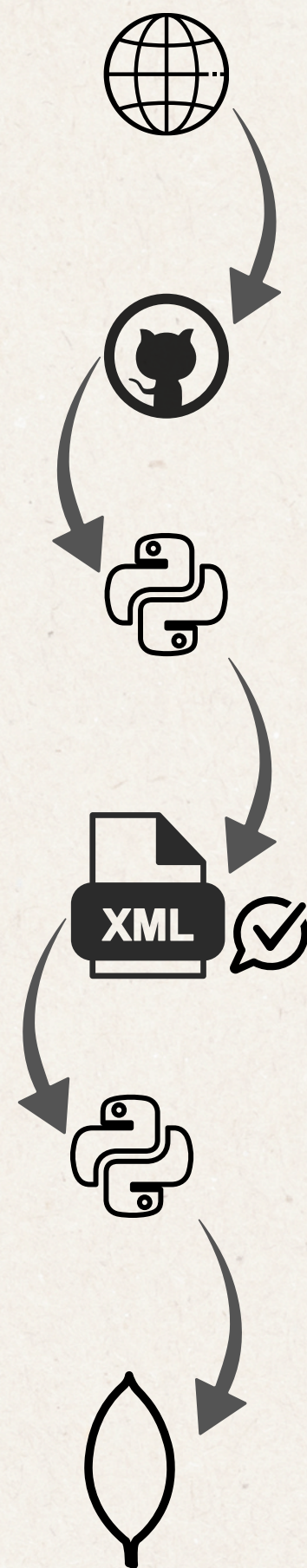
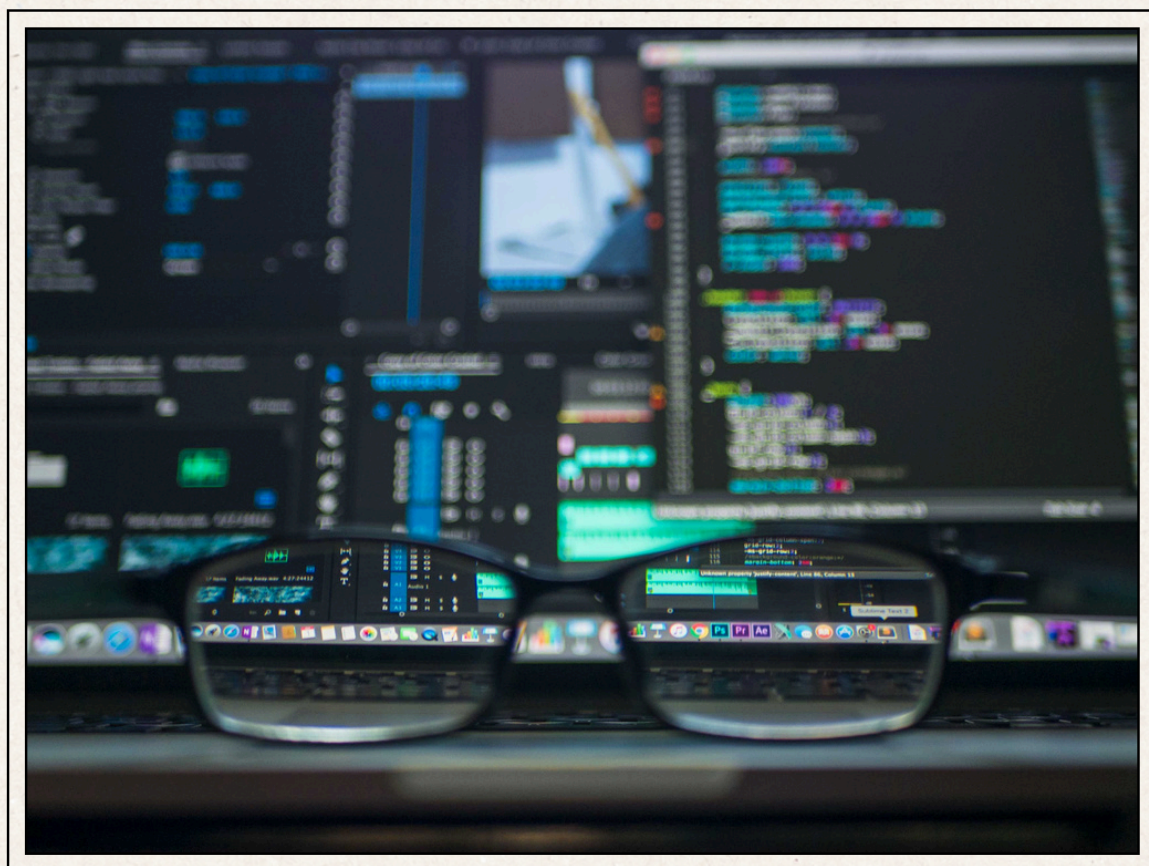


Timeline

06/13

Project timeline





Collect

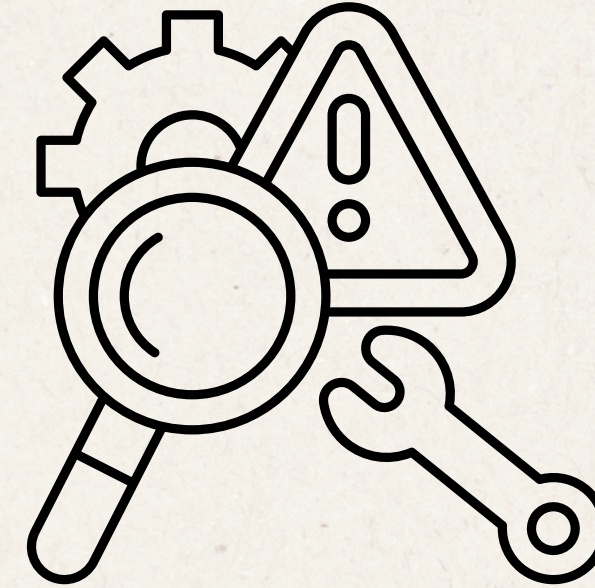
- From online search to •
- Github Darija Open Dataset (DODa) •
- Python code to convert to XML •
- Validation of XMLs with XSDs •
- Python to convert from XML to MongoDB •



Problems encountered

08/13

The original dataset still had some issues to deal with



Issue #1

Not all entries in the dataset have the same number of spelling forms

Issue #2

Missing values for entries (English/and Darija_Ar)

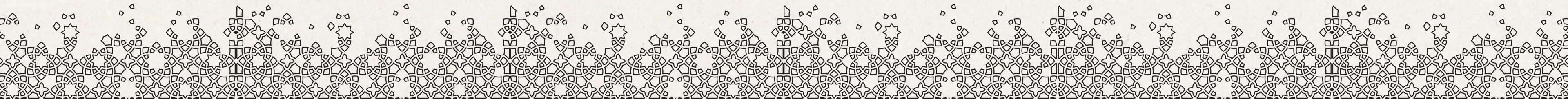
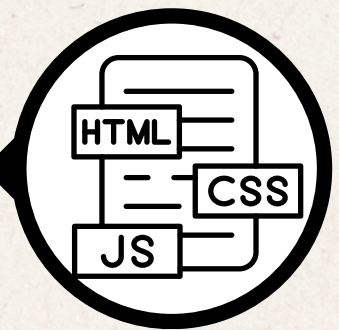
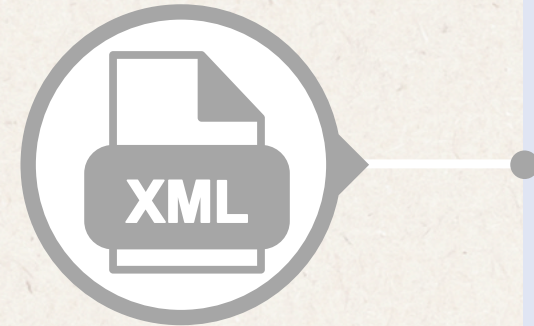
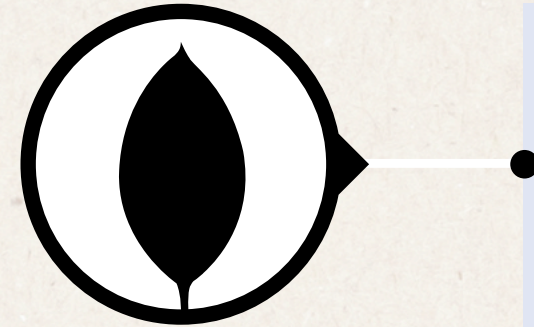
issue #3

Some records/entries appeared in multiple collections/tables

1	n1	n2	n3	n4	darija_ar	eng
2	mousi9a	musiqa	mosi9a		موسيقا	music
3	rasm				رسم	painting
4	kamanja				كامانجا	violin
5	3oud				عود	oud
6	chTi7	chTih			شطيح	dance
7	ghna	ghona			غنا	singing
8	chi3r				شعر	poetry
9	fenn	fann	fnn		فَنّ	art
10	msr7	masra7	masrah	mesre7	مسرح	theater

Prepare

- MongoDB and Nodejs interaction using Expressjs Framework
- Creating necessary queries of NoSQL
- XSDs to validate XML files for each collection
- Converting XML files to JSON using JS and add them to database
- HTML page to display & access the data



Problems encountered

10/13

we encountered several issues when preparing the platform, such as:

Issue #1

Creating XSDs to validate XML files according to the norms of the data structure

Issue #2

Connecting MongoDB to NodeJS using ExpressJS Server

Issue #3

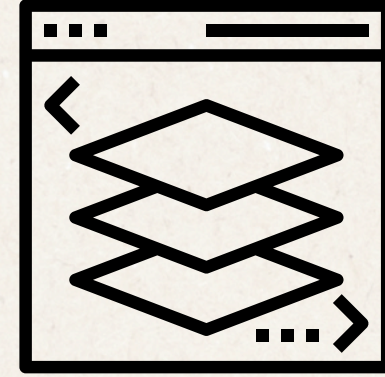
Duplicate entries and records in database cause incorrect results

Issue #4

Problems with using JS to validate XML files with XSDs

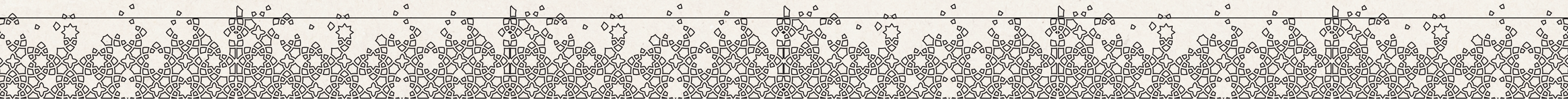
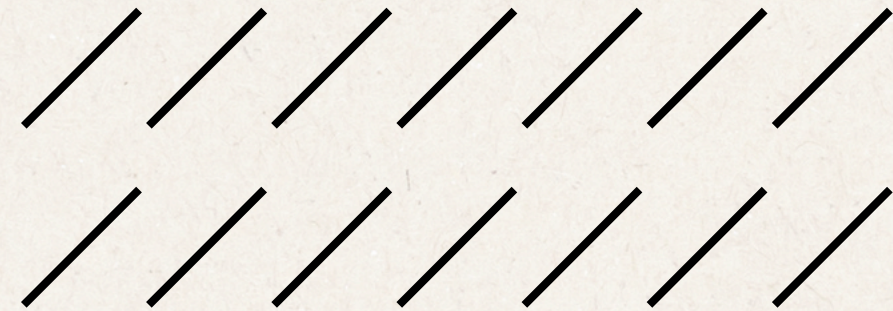
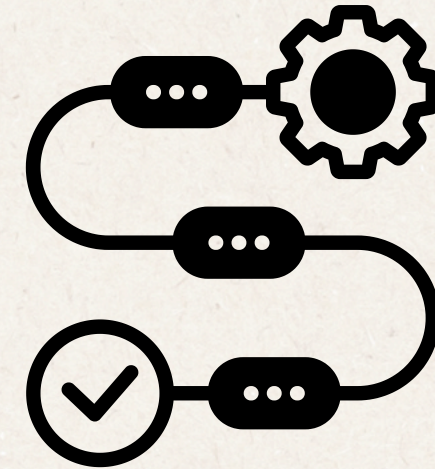
The screenshot displays a web interface titled "English to Darija Search". At the top, there is a search input field containing the word "black" and a "Search" button. Below the search bar, the results are presented under the heading "Search Results:". The results are organized into two rows of three columns each. Each column represents a different transliteration or spelling form of the word "black".

Spelling form n1	Spelling form n2	Spelling form n3	Spelling form n1	Spelling form n2	Spelling form n3	Spelling form n1	Spelling form n2	Spelling form n3
k7l	k7al	k7el	3zzi	3azzi	3ezzi	black	black	black
كحل			عزي	عززي	عزّي	English : black	English : black	English : black
English : black						Arabic Darija : عزي	Arabic Darija : عززي	Arabic Darija : عزّي

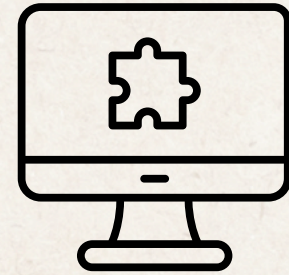


Access

- HTML Page to access data
- Contribution using XML files
- Validation of XML files
- XML to MongoDB
- Statistics of the database



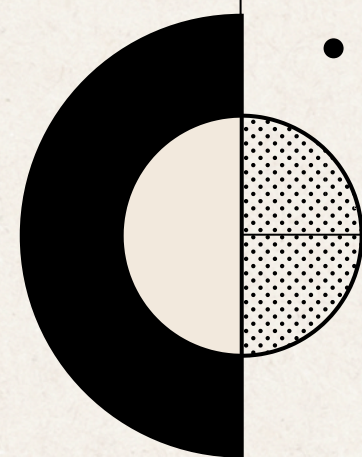
Extensions



12/13

ADDITIONAL TECHNOLOGIES AND TECHNIQUES USED

<ul style="list-style-type: none">• XSD additional options	<i>Validate XML element to select from options</i>
<ul style="list-style-type: none">• Javascript	<i>Javascript manipulation of JSON</i>
<ul style="list-style-type: none">• NodeJS & ExpressJS	<i>Easy interactions between MongoDB (noSQL) & FrontEnd Page (HTML)</i>
<ul style="list-style-type: none">• NoSQL	<i>More detailed noSQL queries for specific tasks on MongoDB</i>



Thank you



CONTACT US

E-mail {yassir.el.attar,shawn.chua}@ims.uni-stuttgart.de

Program Computational Linguistics M.Sc.

Department IMS - Faculty 5

References

DATABASE & CONTENT

Darija Open Dataset. [Online] Available at: <https://darija-open-dataset.github.io/>

GitHub - darija-open-dataset/dataset. [Online] Available at: <https://github.com/darija-open-dataset/dataset/> .

MongoDB. [Online] Available at: <https://www.mongodb.com> .

IMAGES & ICONS

Icons - Available at: <https://www.canva.com>

Spain Symbol - Available at: <https://www.pinterest.com/pin/spain-spanish-flag-souvenir-bull-head-sticker-1089378597347561283/>

France Symbol - Available at: <https://qph.cf2.quoracdn.net/main-qimg-967efe8c5aae711f6c6d7cddf5bace2c-lq>

Tamazight Symbol - Available at: https://tunisie-amazigh.blogspot.com/2013/12/blog-post_26.html

Arabic Symbol - Available at: <https://www.nawa3em.com/UserFiles/002-daddd.jpg>

Darija 'Diro Niya' Symbol - Available at:
<https://www.spreadshirt.fr/shop/design/diro+niya+morocco+qatar+world+cup+avocay+autocollant-D63adb9d93664df7481c8091c?sellable=qrjDBE5oEMHqQp5wykXn-1459-215>

