

College of
Computing
UM6P

Building efficient Natural Language Processing systems for low-resource languages

Abdellah EL MEKKI

College of Computing, Mohammed VI Polytechnic University

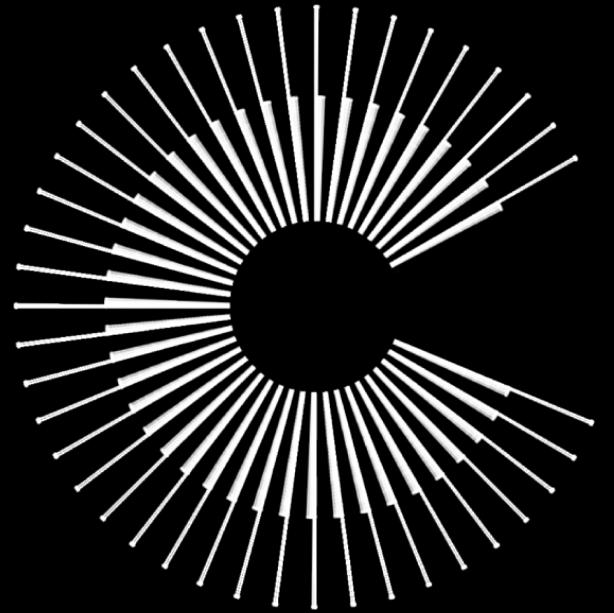
Who am I?

Abdellah EL MEKKI

[2018 - 2023] PhD. Mohammed VI
Polytechnic University, Morocco

- Presented my work at several international conferences (NAACL, COLING, CAAI).
- Ranked Top-3 in various NLP challenges.

[2022 - 2023] Visiting Researcher
University of British Columbia, Canada



College of
Computing
UM6P



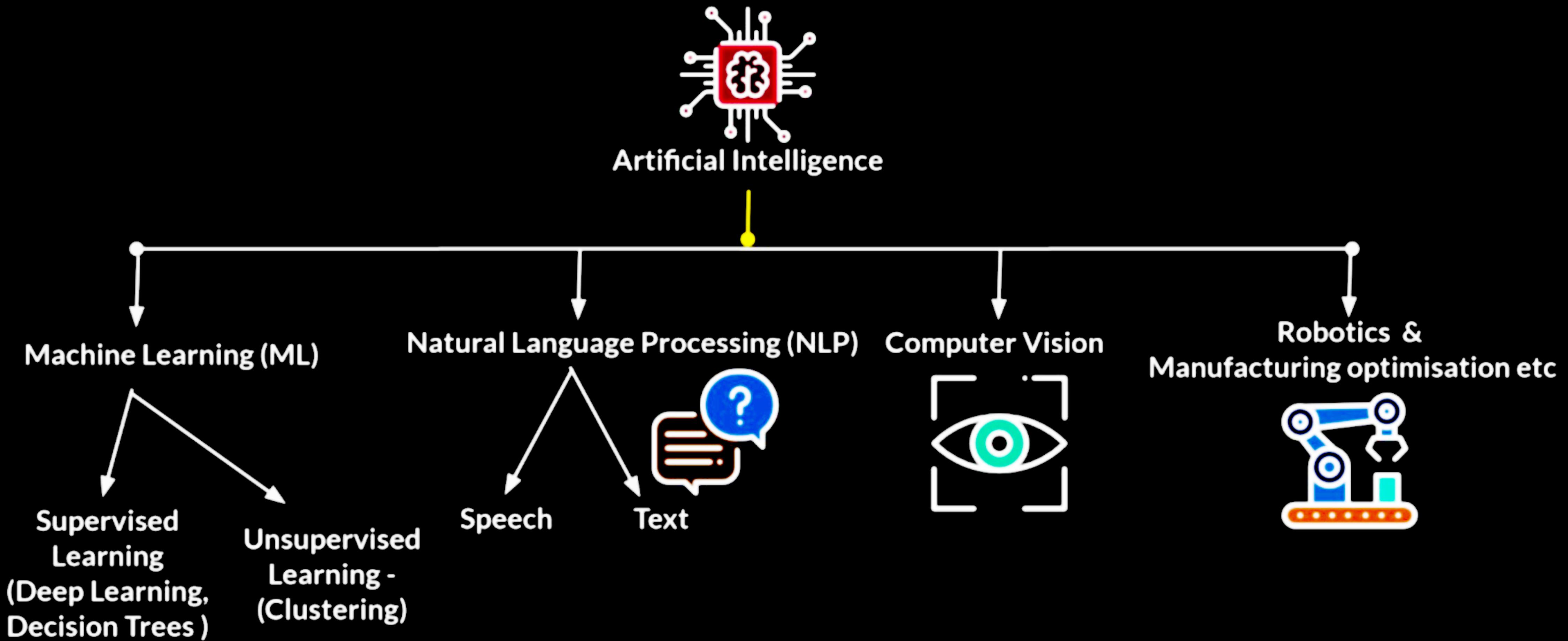
Who am I?



Topics

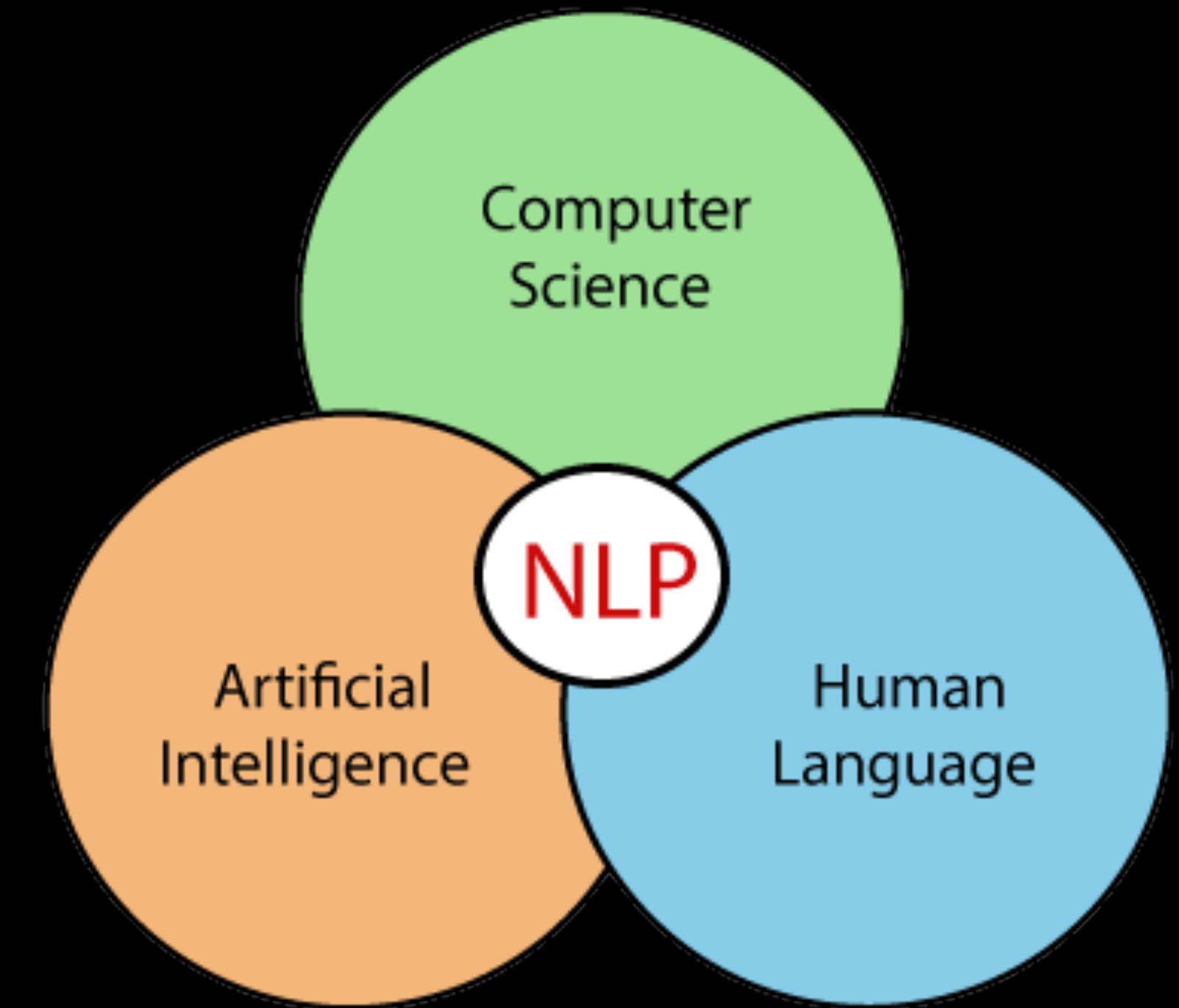
- NLP for low-resource languages.
- Multilingual NLP.
- Transfer Learning in NLP.
- Speech Recognition.
- And now ... Large Language Models.

NLP and AI



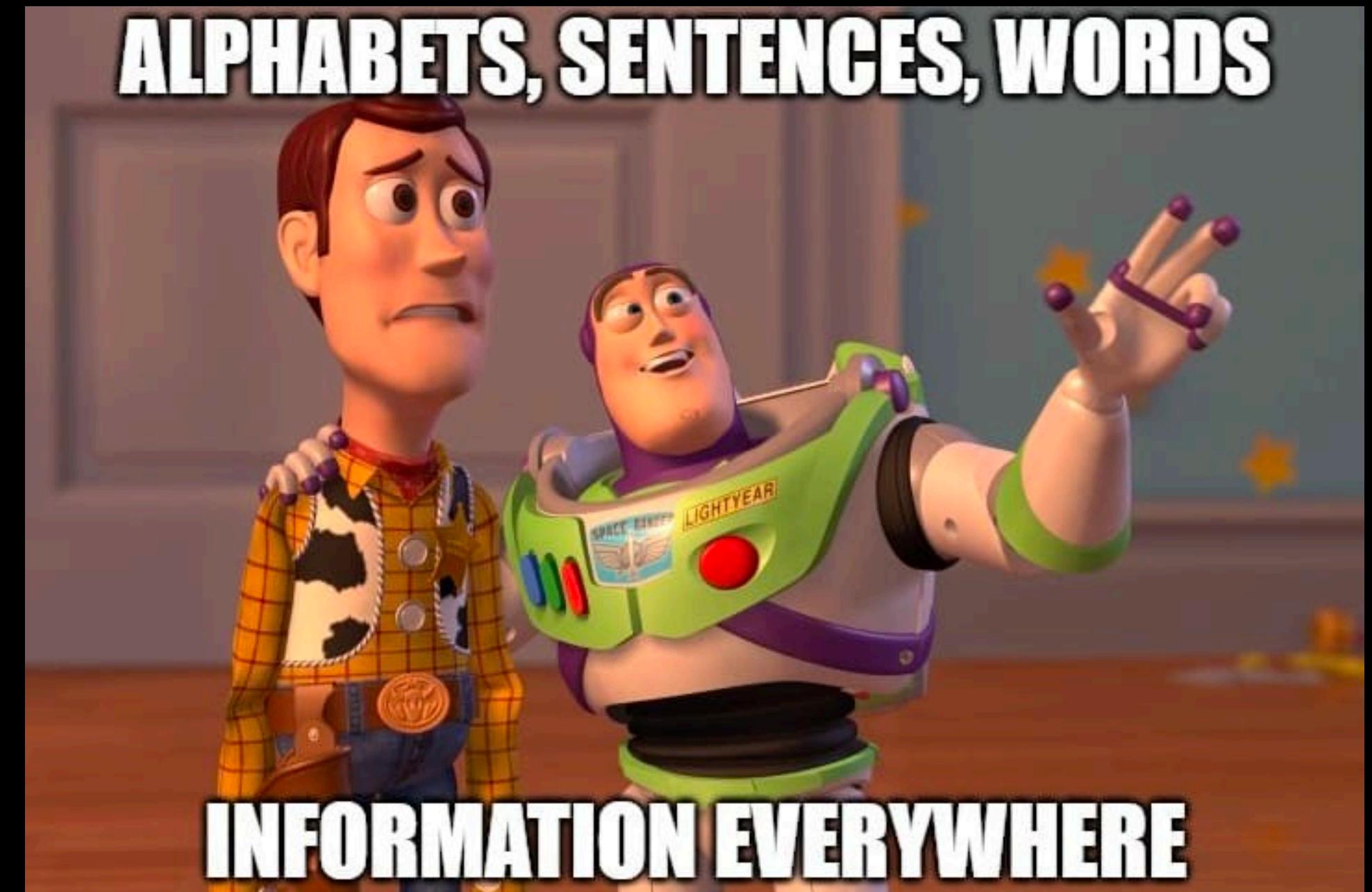
What is NLP?

★**Wiki: Natural language processing**
(NLP) is a field of computer science, artificial intelligence, and computational linguistics concerned with the interactions between computers and human **(natural)** languages.



What is NLP?

- Identify the structure and meaning of words, sentences, texts and conversations.
- Deep understanding of broad language.
- NLP is all around us.



Machine Translation

gute nacht

All Images Videos News More

About 15,600,000 results (0.41 seconds)

German – detected English

gute nacht Good night

Open in Google Translate • Feedback

← Tweet

Soul @aelgatri ...

L'UM6P lance son data center abritant le plus puissant supercalculateur d'Afrique [ledesk.ma/encontinu/lum6...](#) via [@LeDesk_ma](#)

Translated from French by Google

UM6P launches its data center housing the most powerful supercomputer in Africa [ledesk.ma/encontinu/lum6...](#) via [@LeDesk_ma](#)

ledesk.ma
L'UM6P lance son data center abritant le plus puissant su...
L'Université Mohammed VI Polytechnique (UM6P) de Benguerir a procédé, vendredi, à l'inauguration de son ...

10:37 AM · Feb 20, 2021

Question answering

- What does “divergent” mean?
- What year was Abraham Lincoln born?
- How many states were in the United States that year?
- How much Chinese silk was exported to England at the end of the 18th century?
- What do scientists think about the ethics of human cloning?



Google x | microphone | search

All News Images Shopping Maps More Settings Tools

About 258,000,000 results (0.63 seconds)

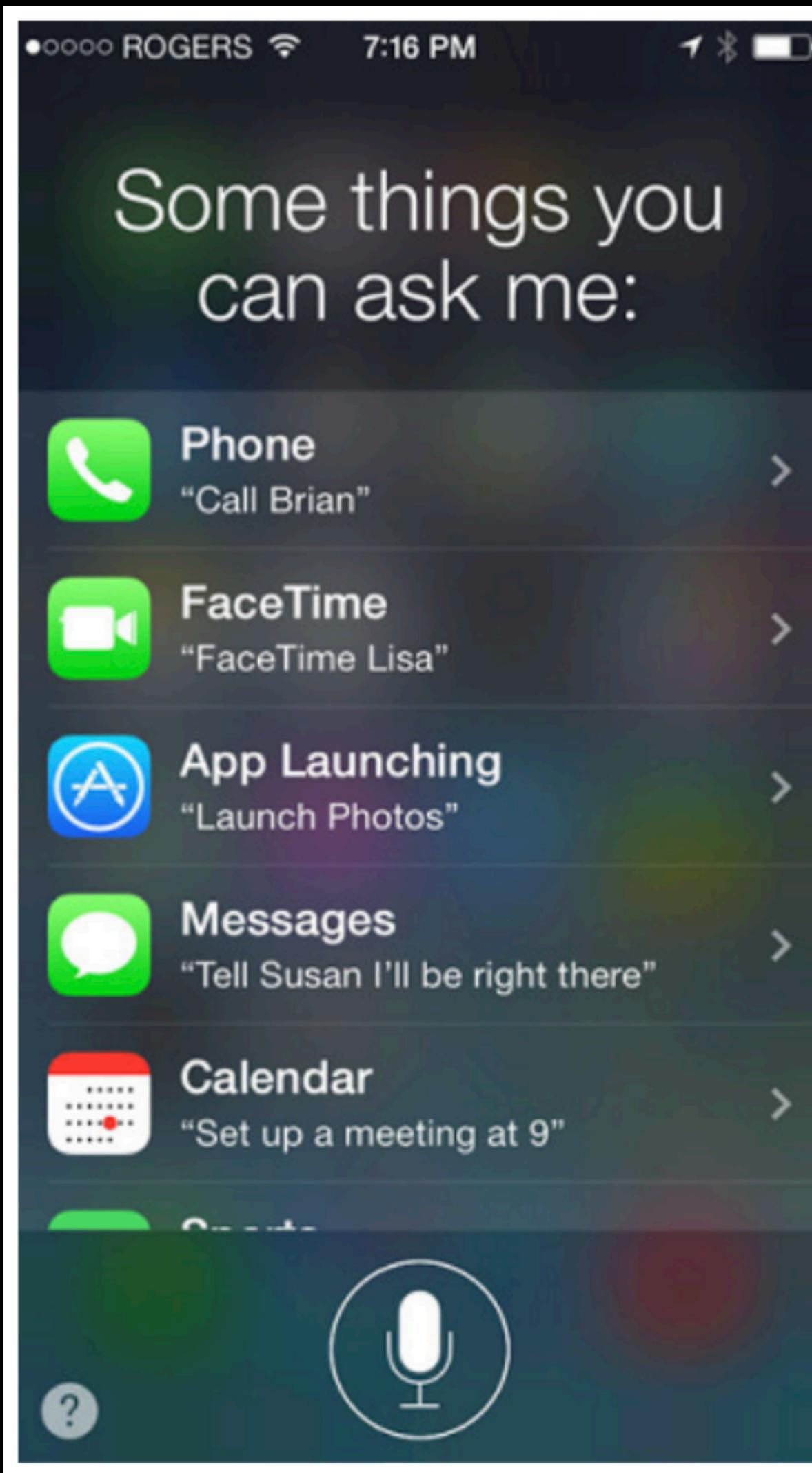
[en.wikipedia.org › wiki › List_of_U.S._states_by_date... ▾](#)

[List of U.S. states by date of admission to the Union - Wikipedia](#)

... states ratified the 1787 Constitution, then the order in which the others were admitted to the Union. A state of the United States is one of the 50 constituent entities that shares its sovereignty with the federal government. Americans are citizens of both the federal republic and of the state in which ...

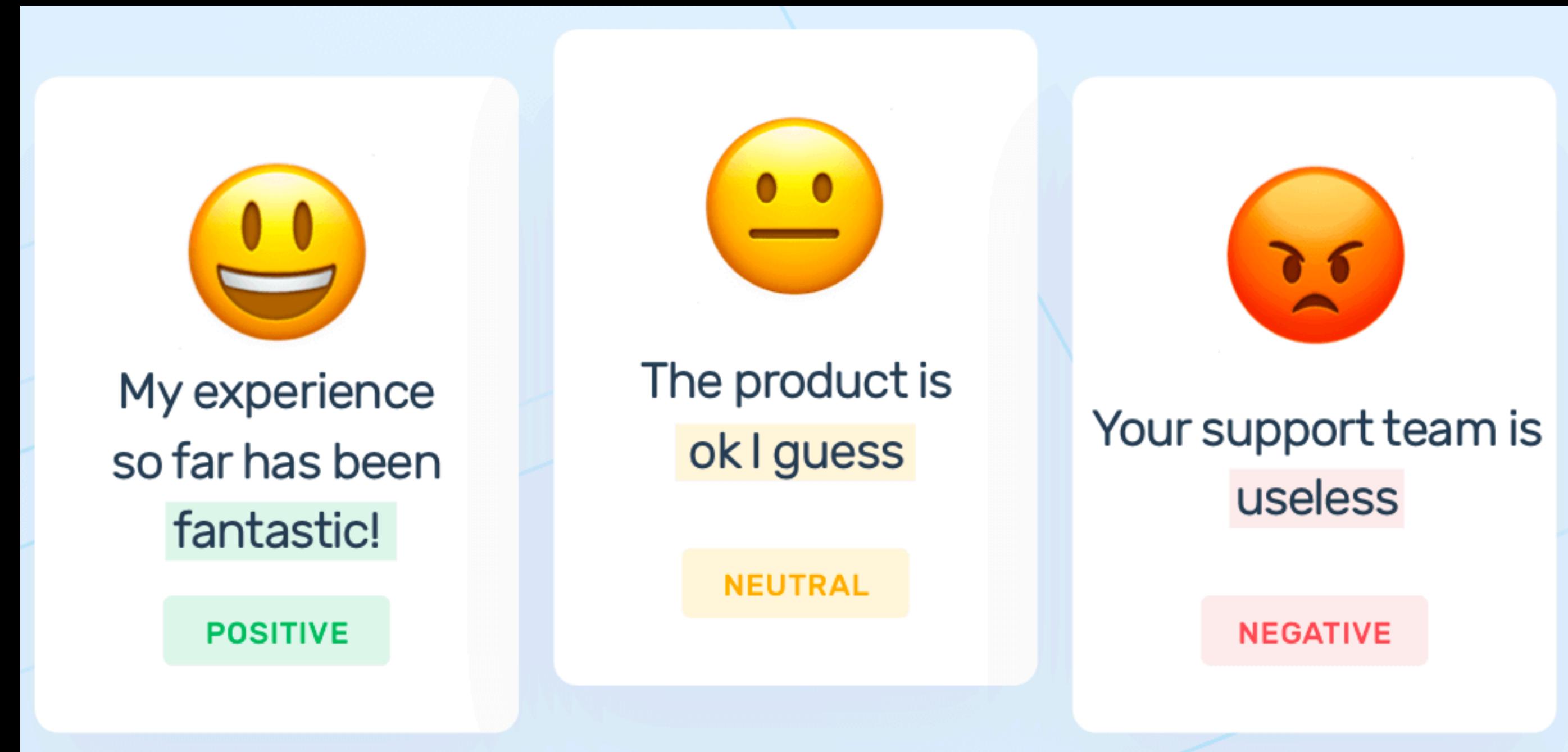
[List of U.S. states · Articles of Confederation ... · See also · Notes](#)

Smart voice assistants



- **Siri contains:**
 - Speech recognition
 - Language analysis
 - Dialog processing
 - Text to speech

Sentiment/opinion analysis



How to build an NLP system?

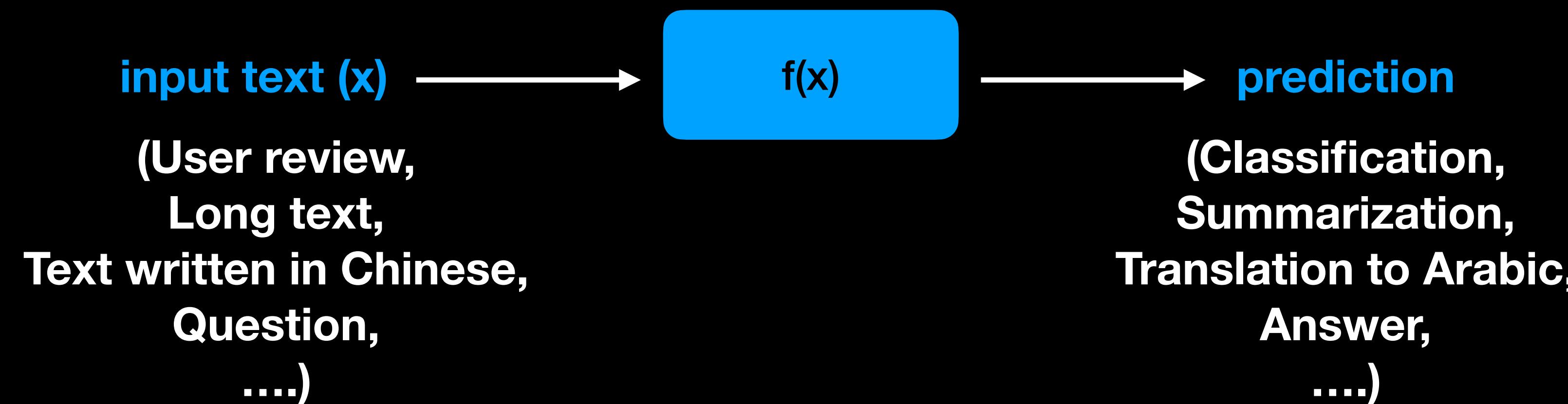
input text (x)

(User review,
Long text,
Text written in Chinese,
Question,
....)

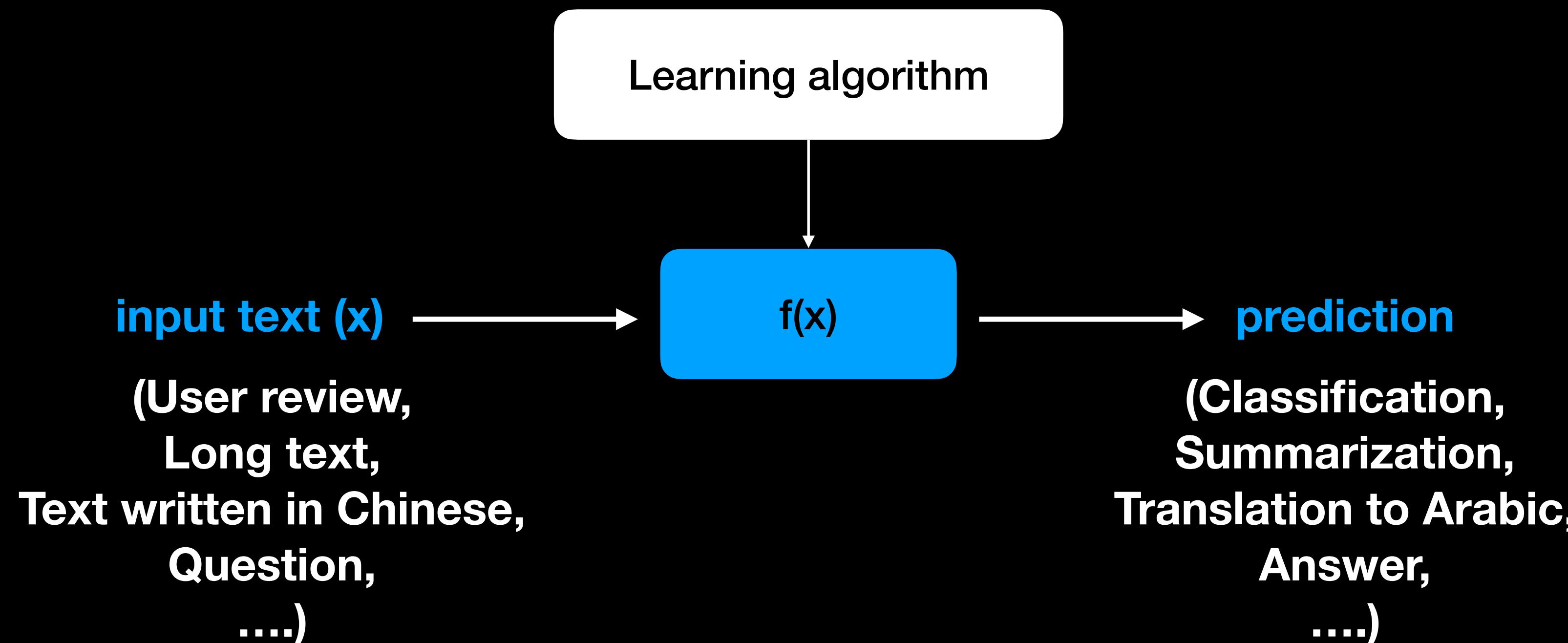
prediction

(Classification,
Summarization,
Translation to Arabic,
Answer,
....)

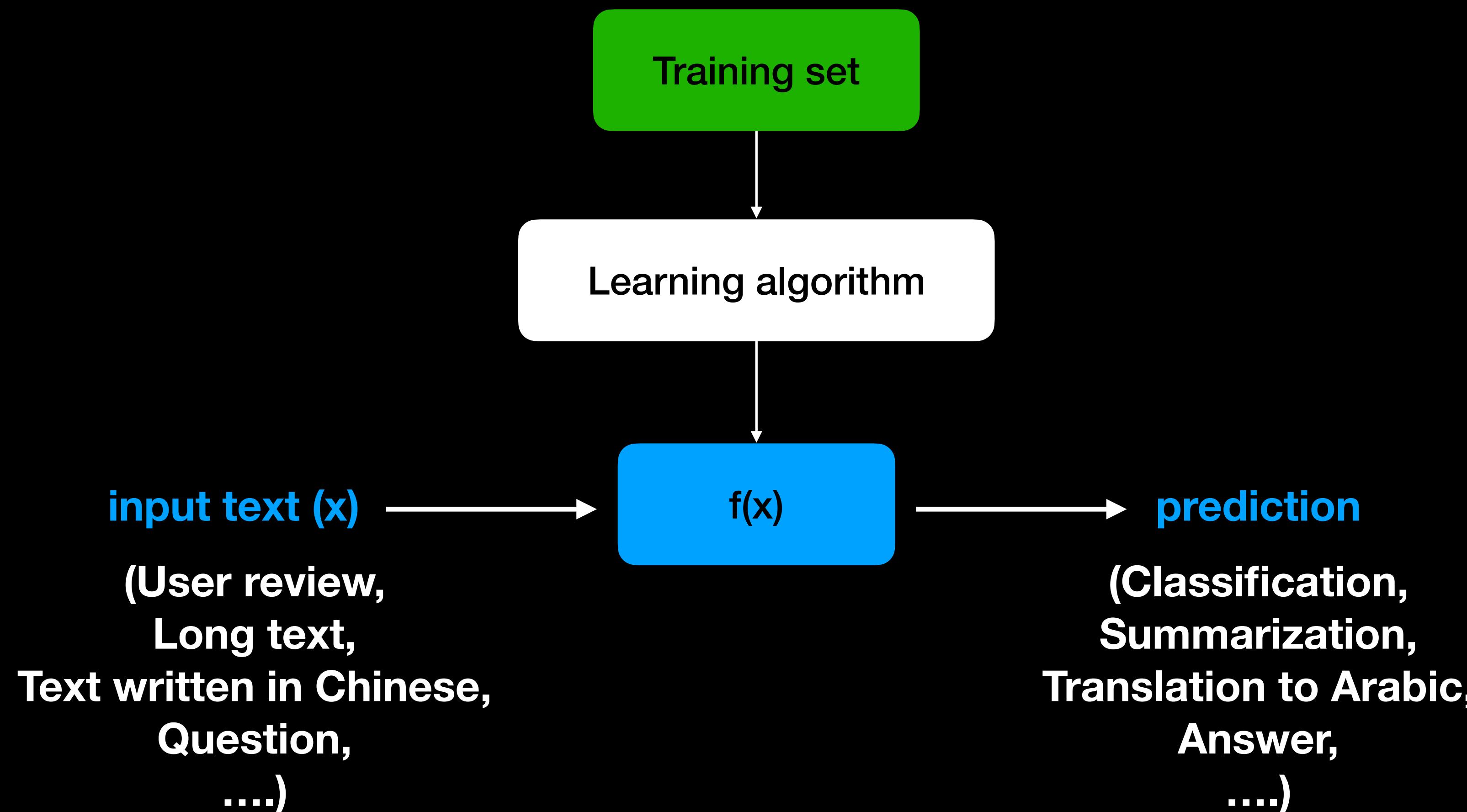
How to build an NLP system?



How to build an NLP system?



How to build an NLP system?



History of NLP

- **Hand-crafted Systems:** Knowledge Engineering [1950s-]
 - Rules written by hand; adjusted by error analysis.
 - Require experts who understand both the systems and domain.
 - Iterative guess-test-tweak-repeat cycle.
- **Automatic Trainable (Machine Learning)** Systems with engineered features [1985s-2012]
 - The tasks are modeled in a statistical way.
 - More robust techniques based on rich annotations.
 - Perform better than rules (Parsing 90% vs. 75% accuracy).
- **Automatic Trainable Neural architectures** with no/limited engineered features [2012-]

Components of NLP

- Natural Language Understanding
 - ◆ Mapping the given input in the natural language into a useful representation.
 - ◆ Different level of analysis required: morphological, syntactic, semantic.
- Natural Language Generation
 - ◆ Producing output in the natural language from some internal representation.
 - ◆ Different level of synthesis required:
 - Deep planning (what to say)
 - Syntactic generation

What is NLP?

- $NL \in \{\text{English, Arabic, French, Hindi, ..., Mandarin}\}$
- Automation of NLs:
 - ◆ Analysis ($NL \rightarrow \text{Representation}$)
 - ◆ Generation ($\text{Representation} \rightarrow NL$)
 - ◆ Acquisition of Representation from knowledge and data.

Why NLP is HARD?

- Richness: there are many ways to express the same meaning, and immeasurably many meaning to express.
- Linguistic diversity across languages, dialects, genres styles.

1. Ambiguity

2. Sparsity

3. Variation

4. Expressivity

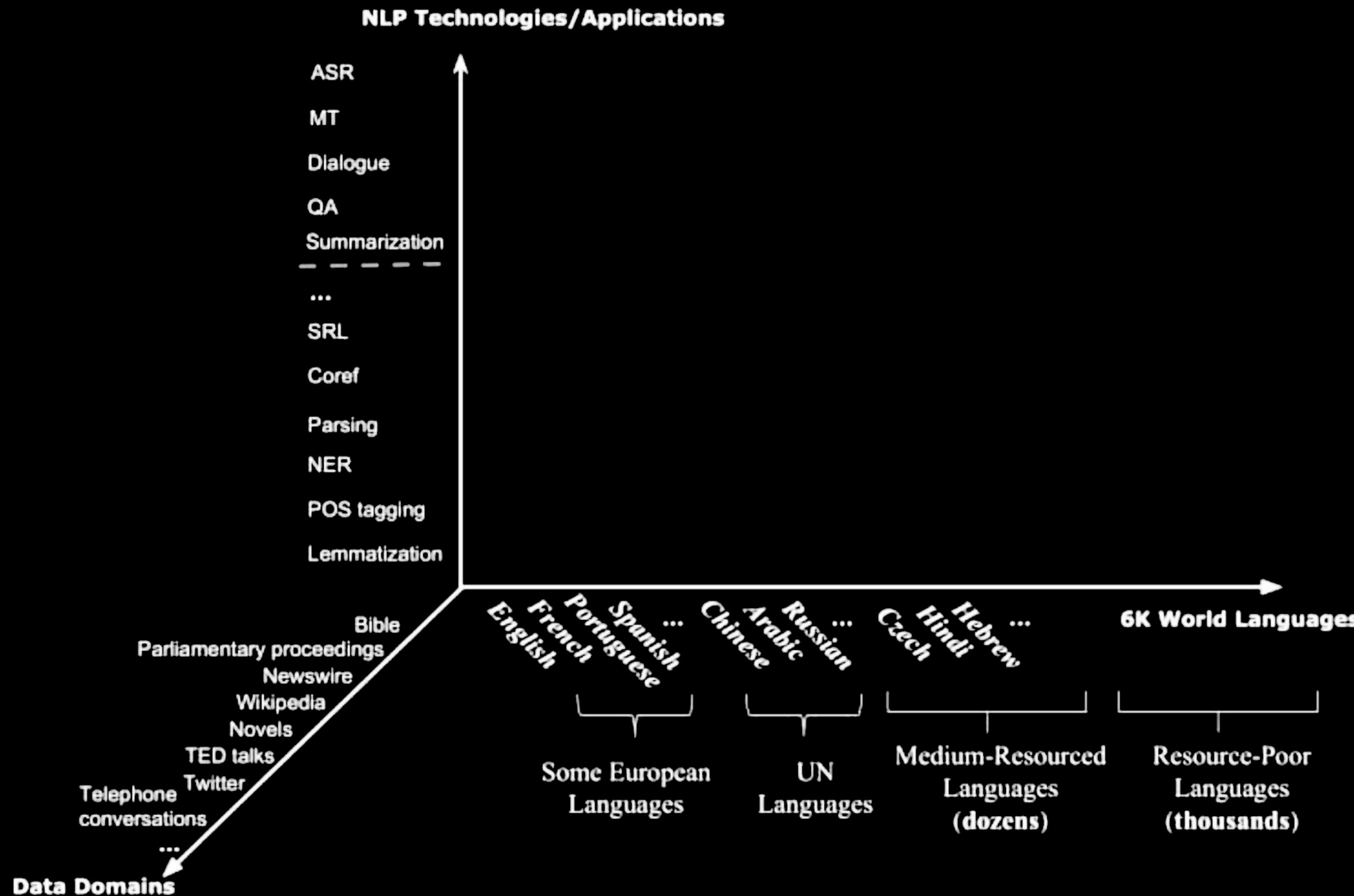
Why NLP is HARD?

Ambiguity

- Ambiguity at multiple levels:
 - Word senses: **bank** (finance or river?)
 - Part of speech: **chair** (noun or verb?)
 - Syntactic structure: **I can see a man with a telescope**
 - Multiple: **I saw her duck**



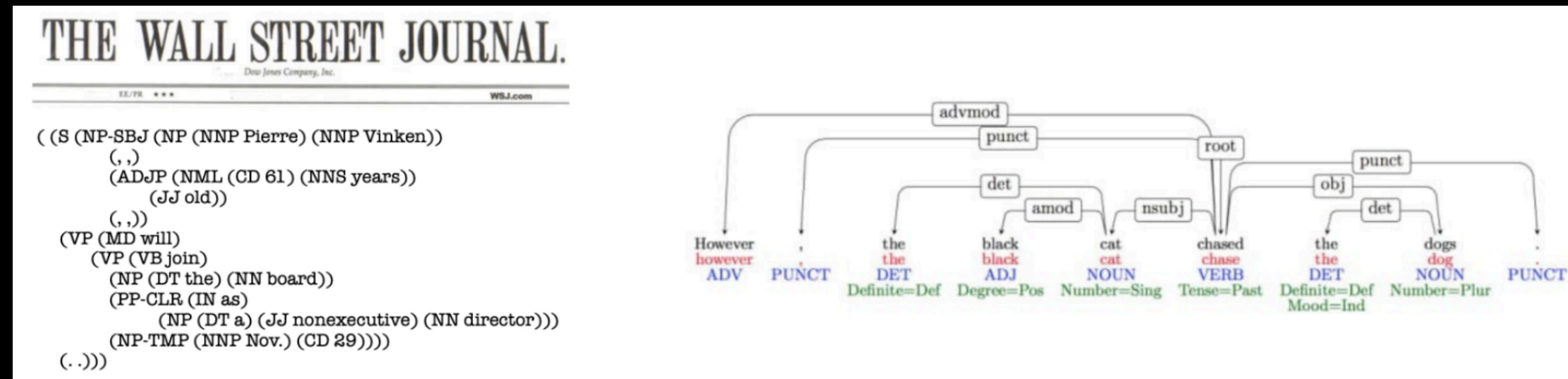
Why NLP is HARD?



Why NLP HARD

Variation

- Suppose we train a part of speech tagger or a parser on the Wall Street Journal...



- What will happen if we try to use this tagger/parser on social media?

Trinity (@christinedarvin)

Hayyy, naminiss ko na yung chicken wings sa UN 😔

Translate Tweet

4:57 AM · Aug 17, 2020 · Twitter for Android

Domo (@djxdomo)

I ain't never met a gatekeeper in my life because I'm finna do whatever tf I'm finna do.

4:22 PM · Aug 31, 2020 · Twitter for iPhone

Why NLP is HARD?

Expressivity

- Not only can one form have different meanings (ambiguity) but the same meaning can be expressed with different forms:

How Google Translation Works



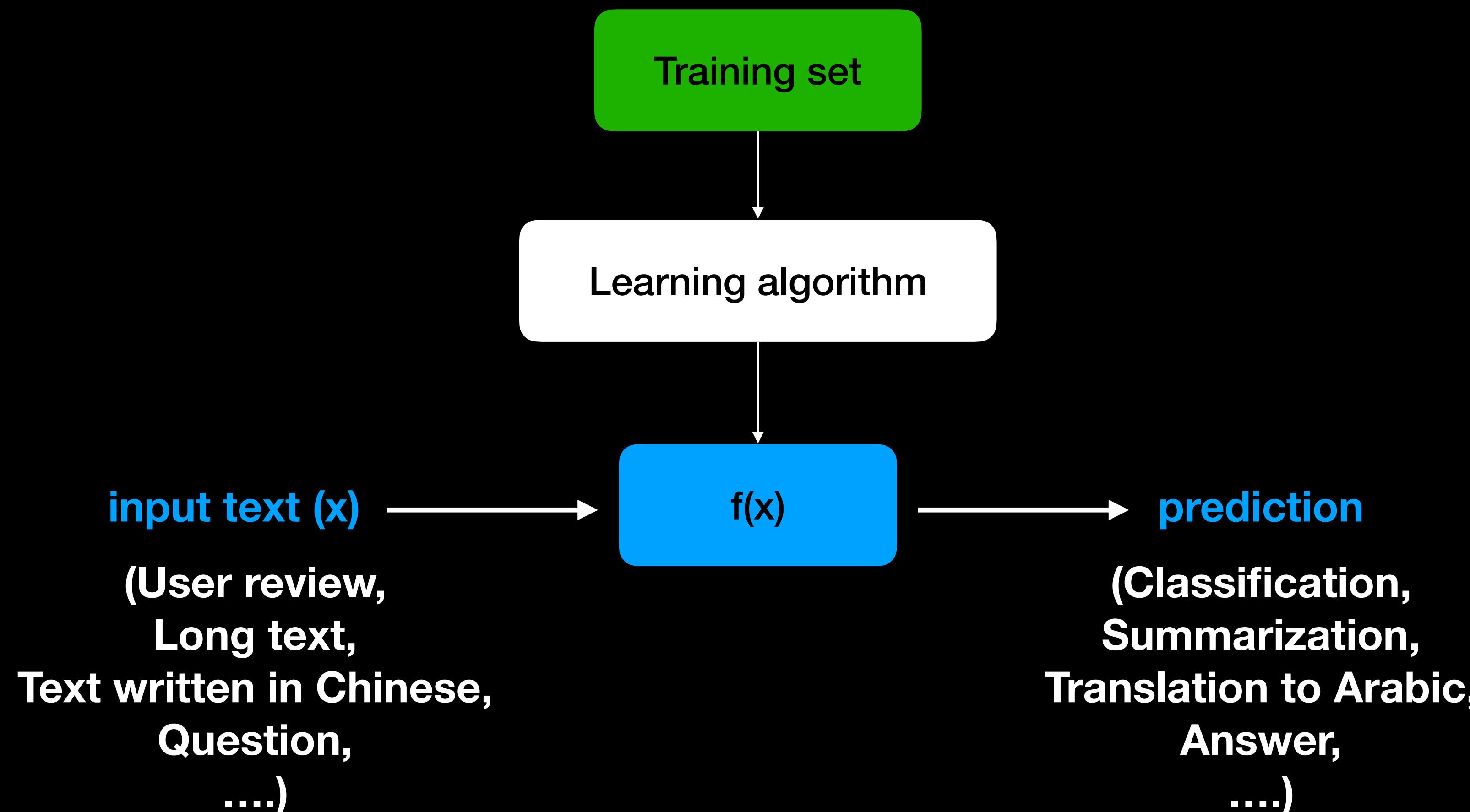
Google Translate

She gave the book to Deni. vs. She gave Deni the book.

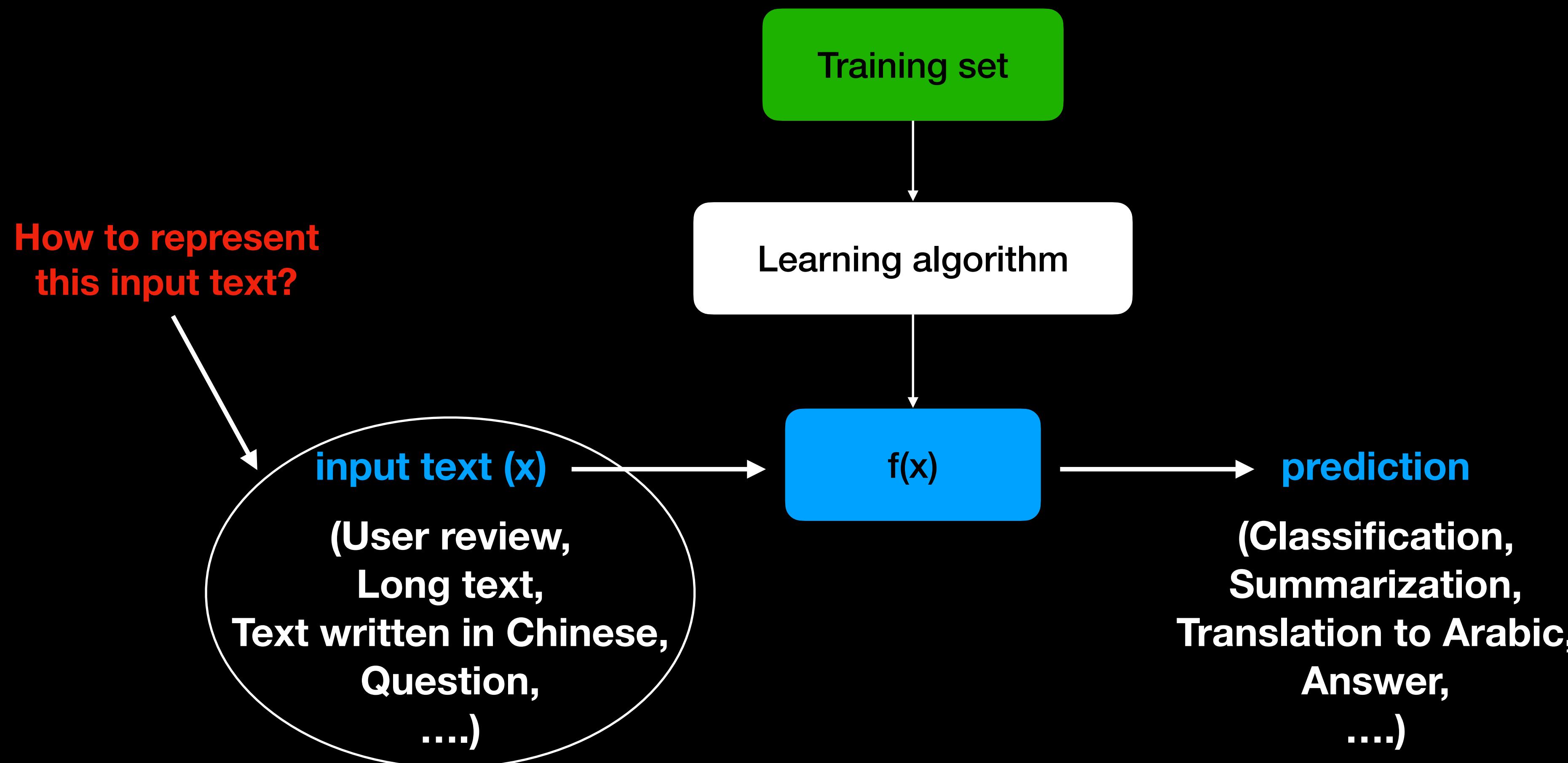
Some kids popped by. vs. A few children visited.

Is that window still open? vs. Please close the window.

How to build an NLP system?

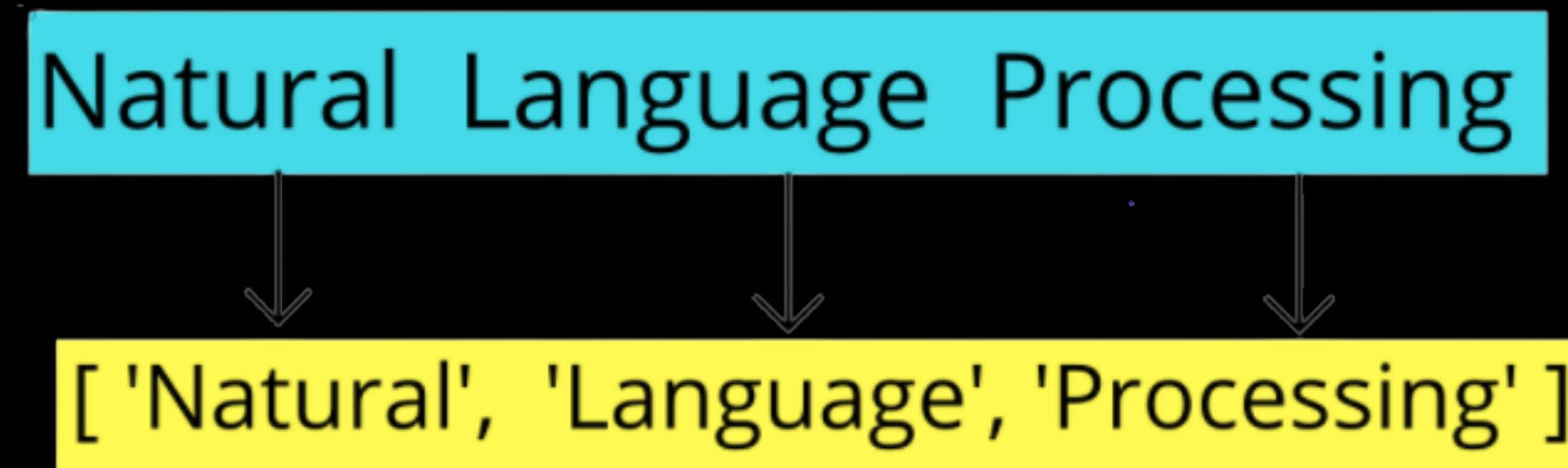


How to build an NLP system?



Tokenization

- Break a complex sentence into words.
- Understand the importance of each of the words with respect to the sentence.
- Produce a structural description of an input sentence.



One-hot encoding

- We've now converted our text to tokens. How can we convert a single token a feature vector?
 - Use **one-hot encoding**.
 - Length of feature vector = #distinct tokens (so can be pretty large, e.g. 10-100K).
 - Map token i to a feature vector x with entry i equal to 1 and other entries 0.
 - E.g. **movie** might map to feature vector [0, 1, 0, 0, 0, 0].

good	movie	not	a	did	like
0	1	0	0	0	0

Bag-of-Words

- How can we map a **sequence** of tokens (sentence, paragraph, document) to a feature vector.
 - Count number of occurrences of each token in our text -> later we'll then try to identify marker words like **excellent** or **disappointed**.
 - Length of feature vector = #distinct tokens
 - Entry i of vector is set equal to number of times token i appears in text. E.g.

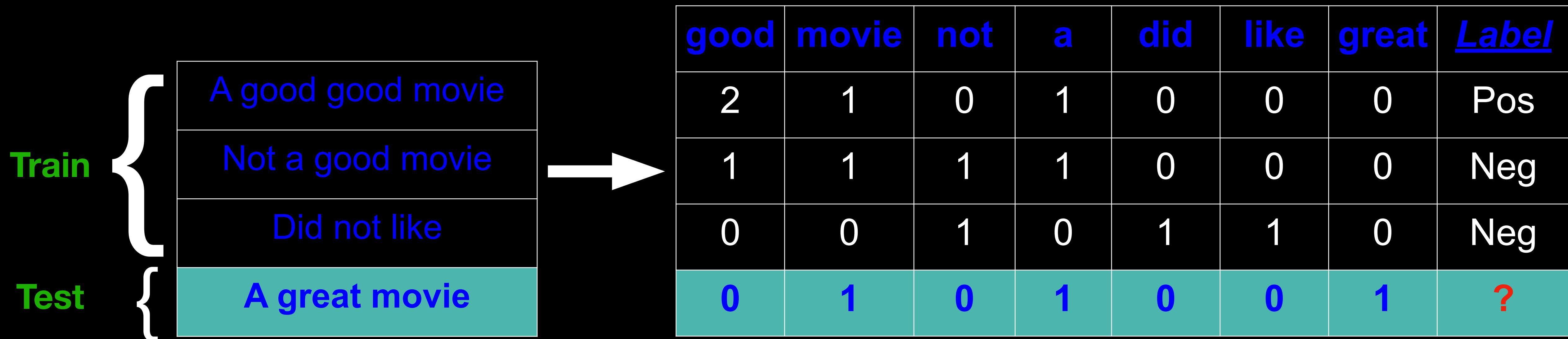
A good good movie
Not a good movie
Did not like



good	movie	not	a	did	like
2	1	0	1	0	0
1	1	1	1	0	0
0	0	1	0	1	1

- Information about the word order is removed, that's why it's called a "bag of words".

Vector-based similarity

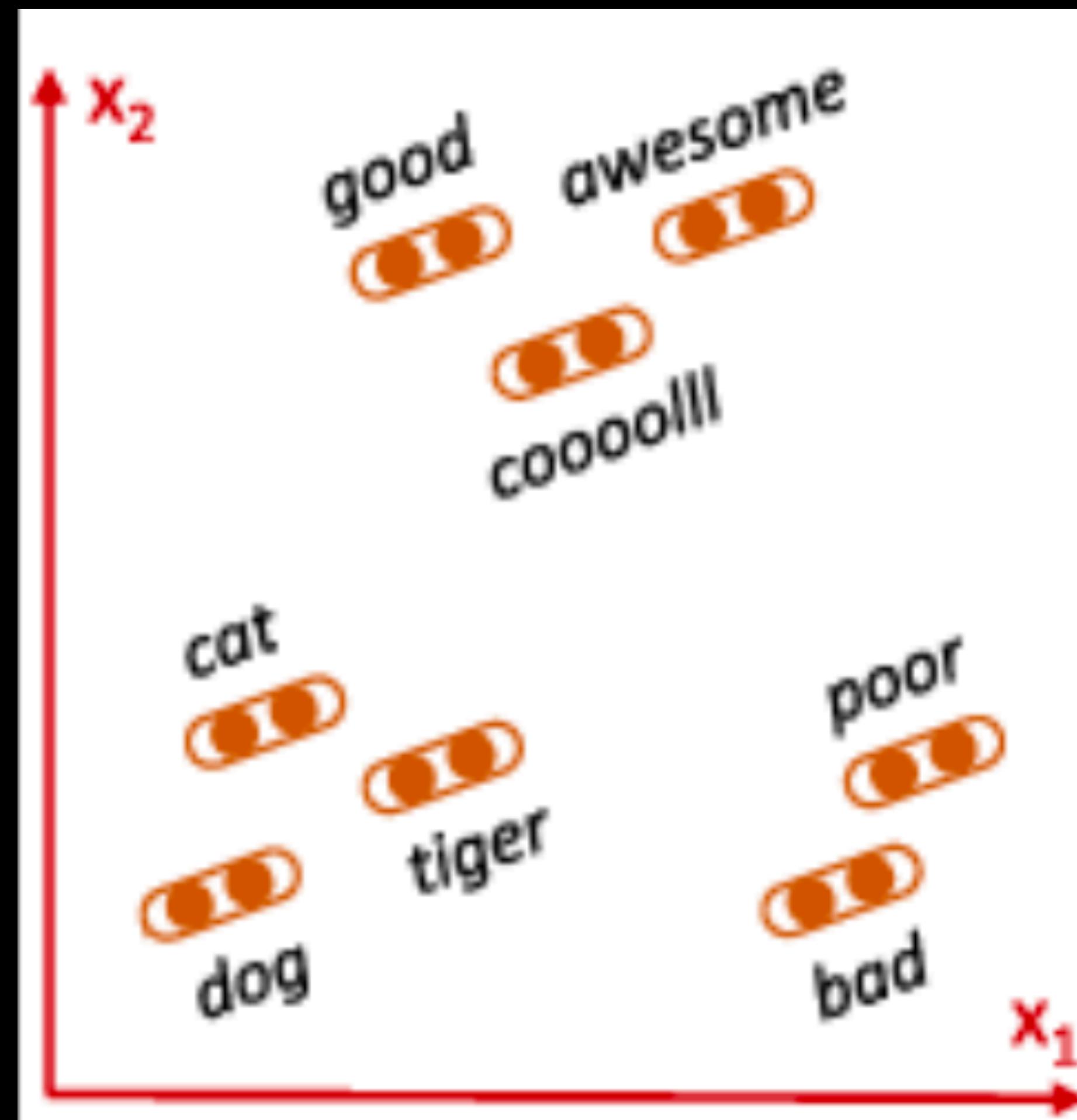


- What would an ML model predict for the last sentence?

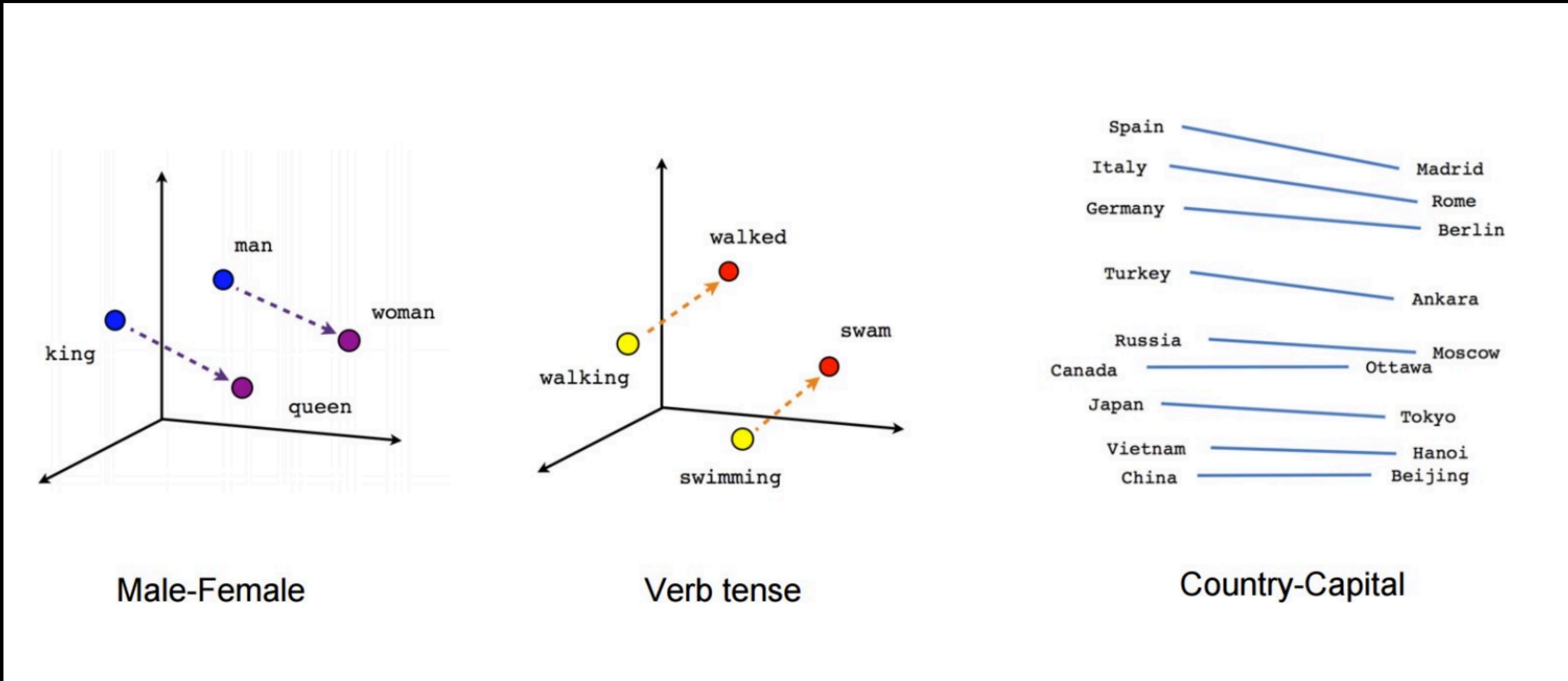
Cons of bag-of-words

- Large vocabulary \Leftrightarrow Large dimensions.
- All word vectors are orthogonal (no similarity).

Word Embeddings



Word Embeddings

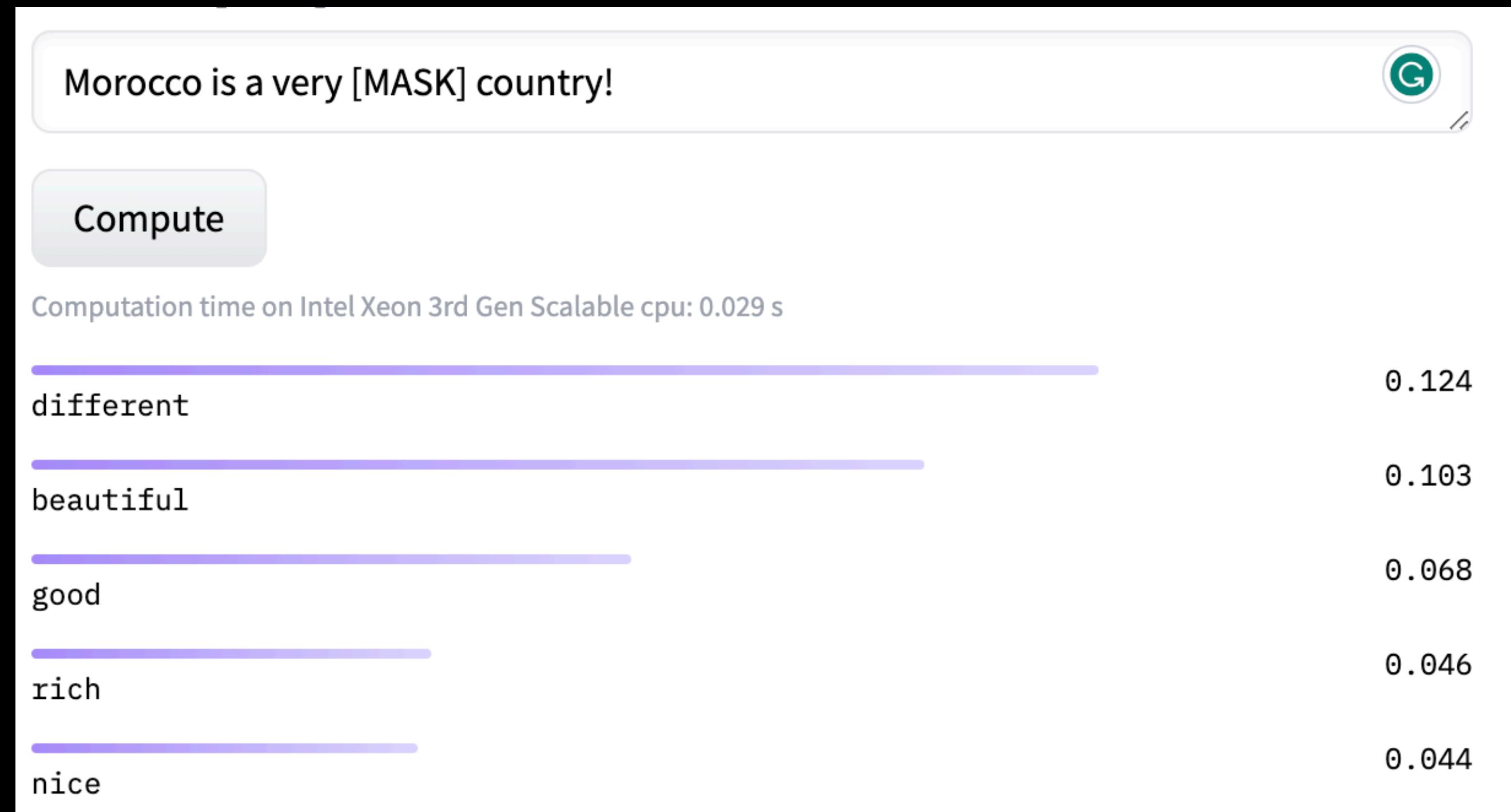


Word Embedding - Fill the Blank

Morocco is a very ... country!

Word Embedding - Fill the Blank

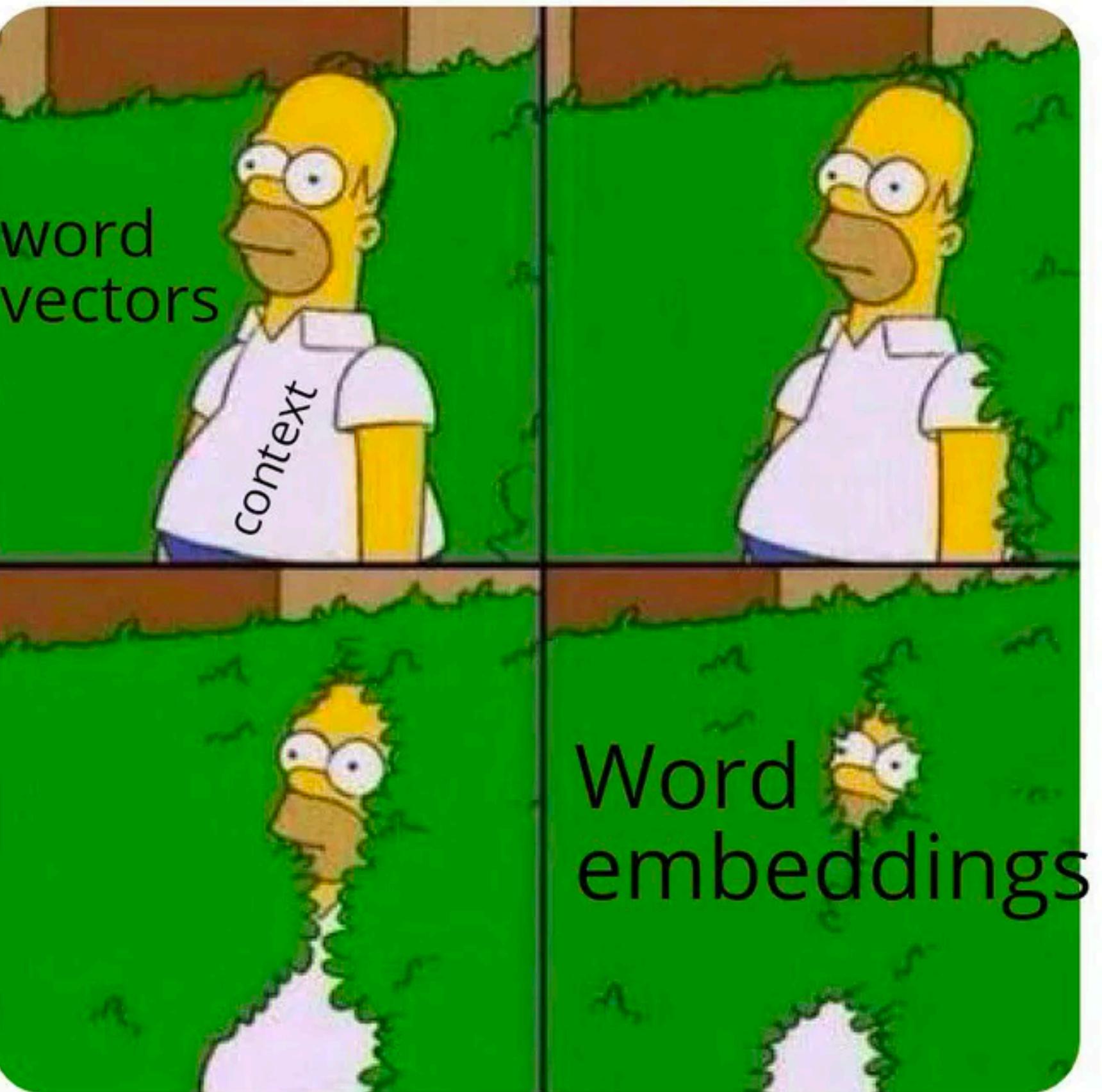
Morocco is a very ... country!



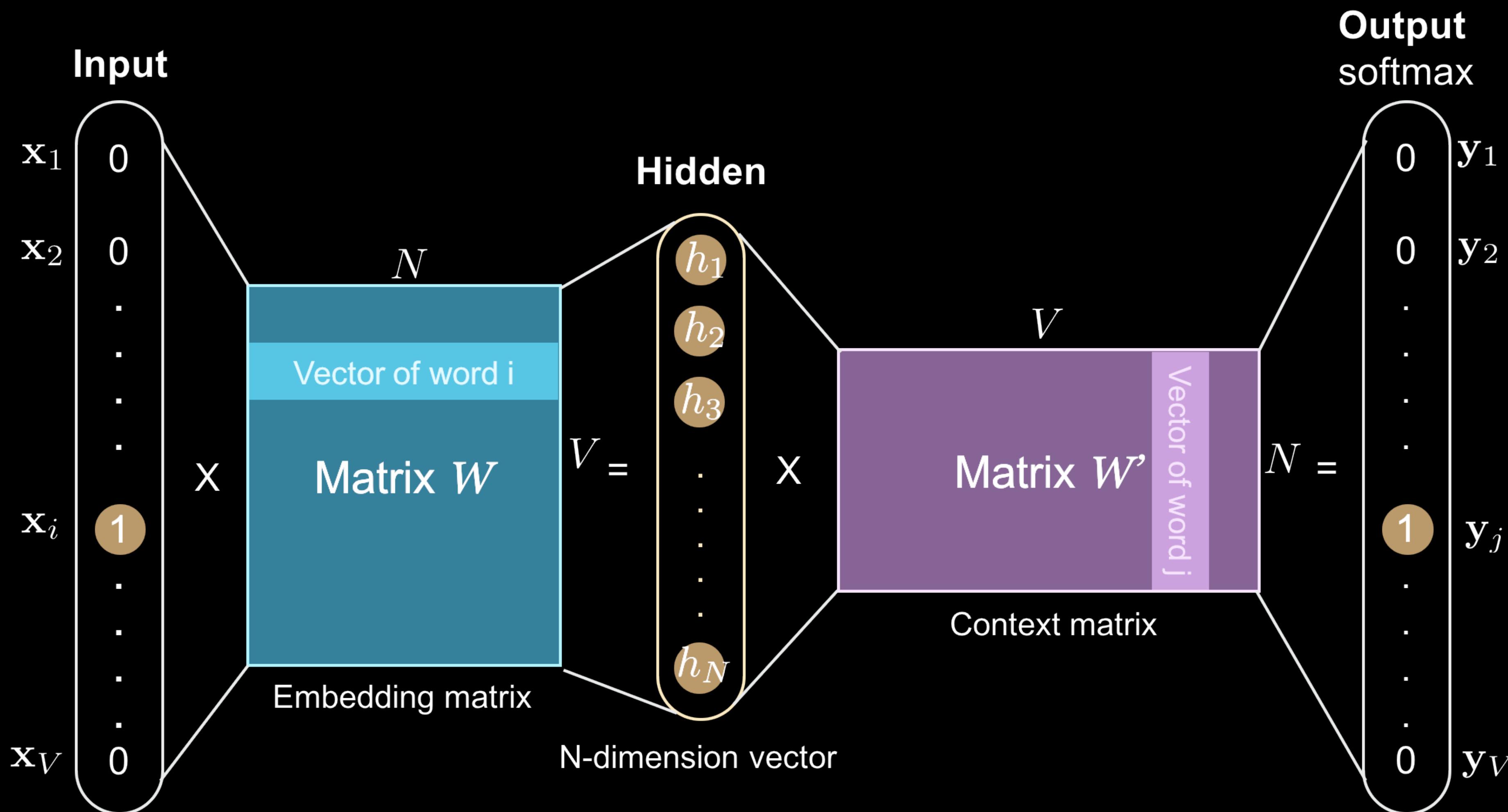
Word Embeddings

- Word2vec
 - CBOW
 - Skip-Gram
- Glove
- FastText

Word embeddings
in a nutshell



Word2vec (Skip-Gram)



How to use Word Embeddings?

This movie is good.

I did not liked this movie.

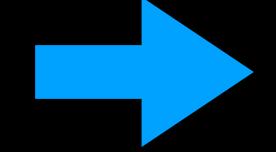
This movie is great.

How to use Word Embeddings?

This movie is good.

[This, movie, is, good.]

I did not liked this movie.



[I, did, not, liked, this, movie.]

This movie is great.

[This, movie, is, great.]

How to use Word Embeddings?

This movie is good.

I did not liked this movie.

This movie is great.

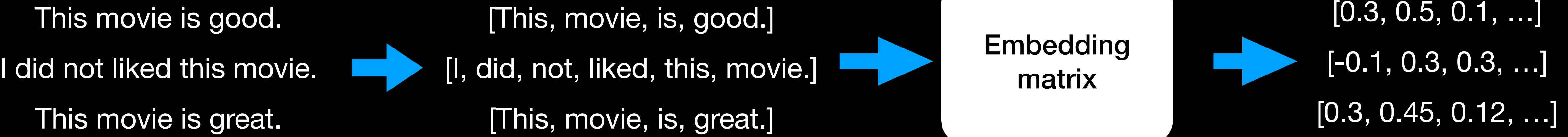
[This, movie, is, good.]

[I, did, not, liked, this, movie.]

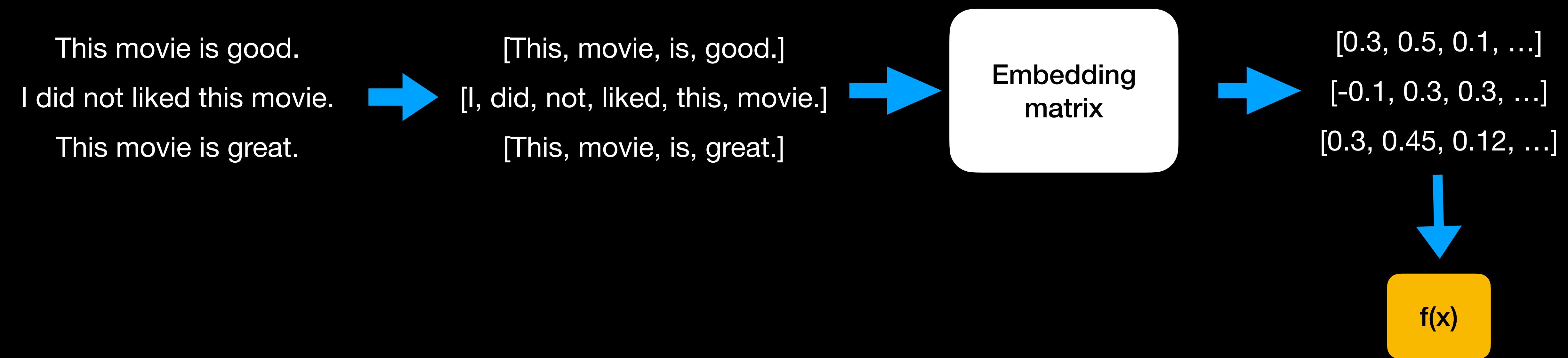
[This, movie, is, great.]

Embedding
matrix

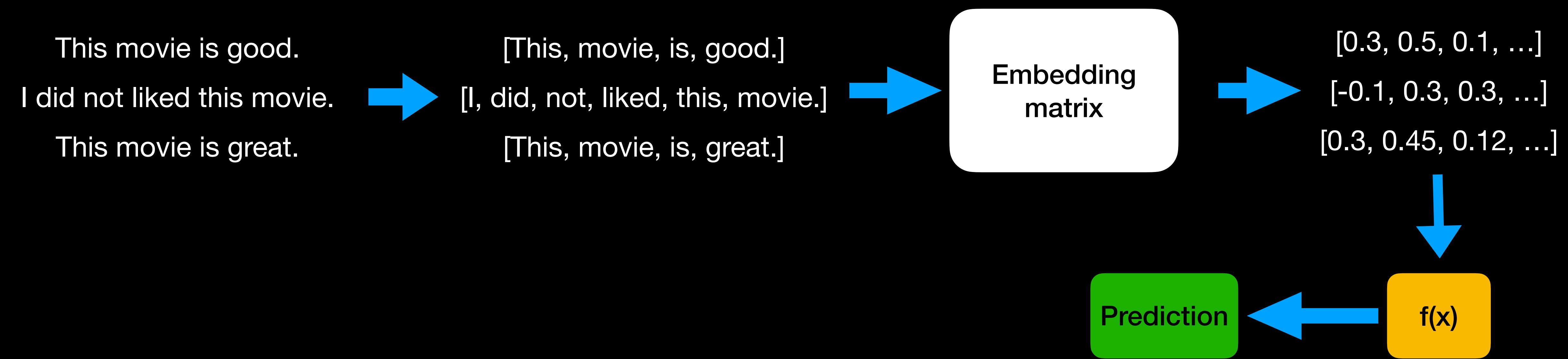
How to use Word Embeddings?



How to use Word Embeddings?



How to use Word Embeddings?



Transformers and Attention

We went to the river **bank**.



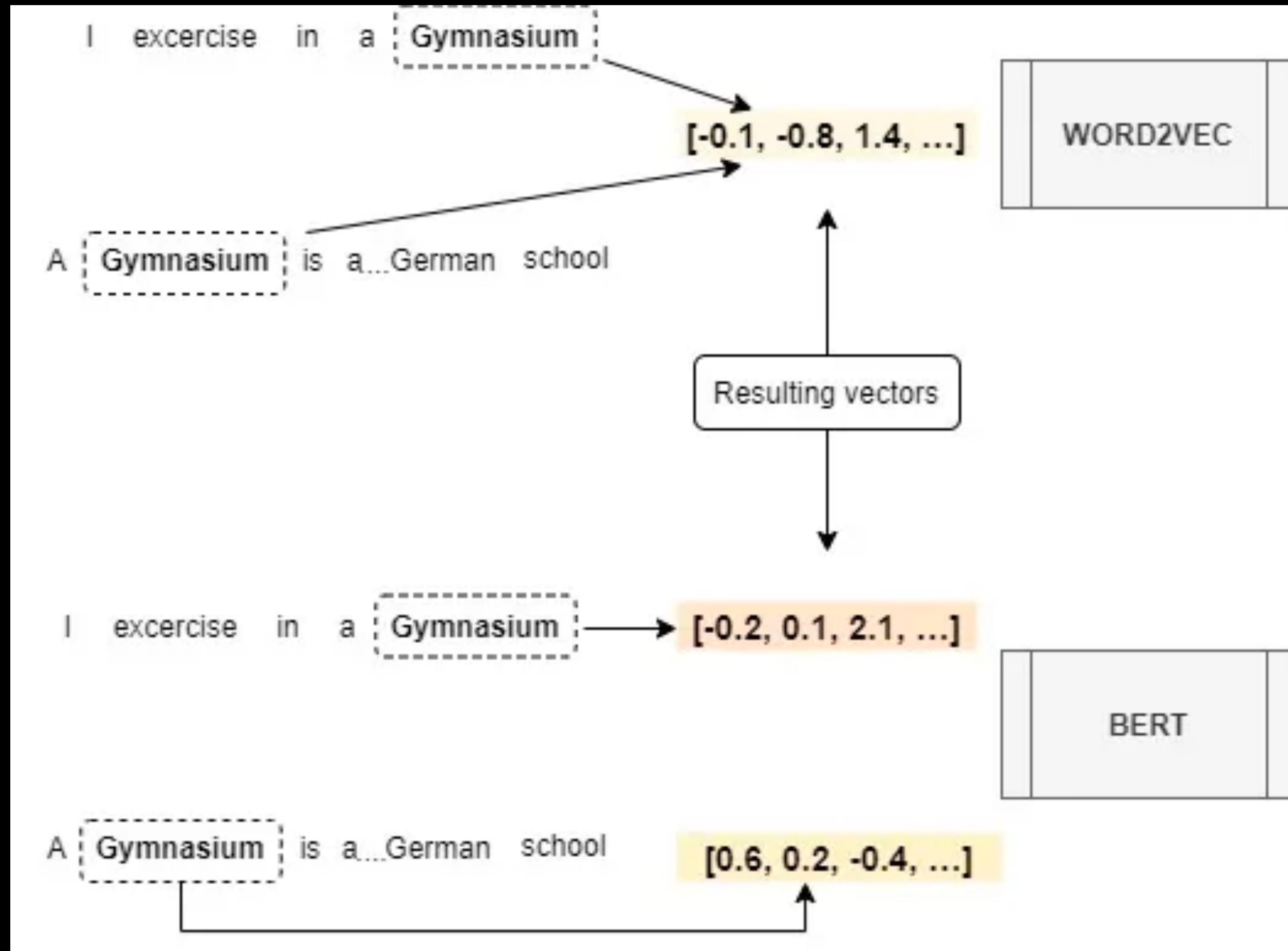
$[-0.1, 0.3, 1.5, \dots]$

Word2vec

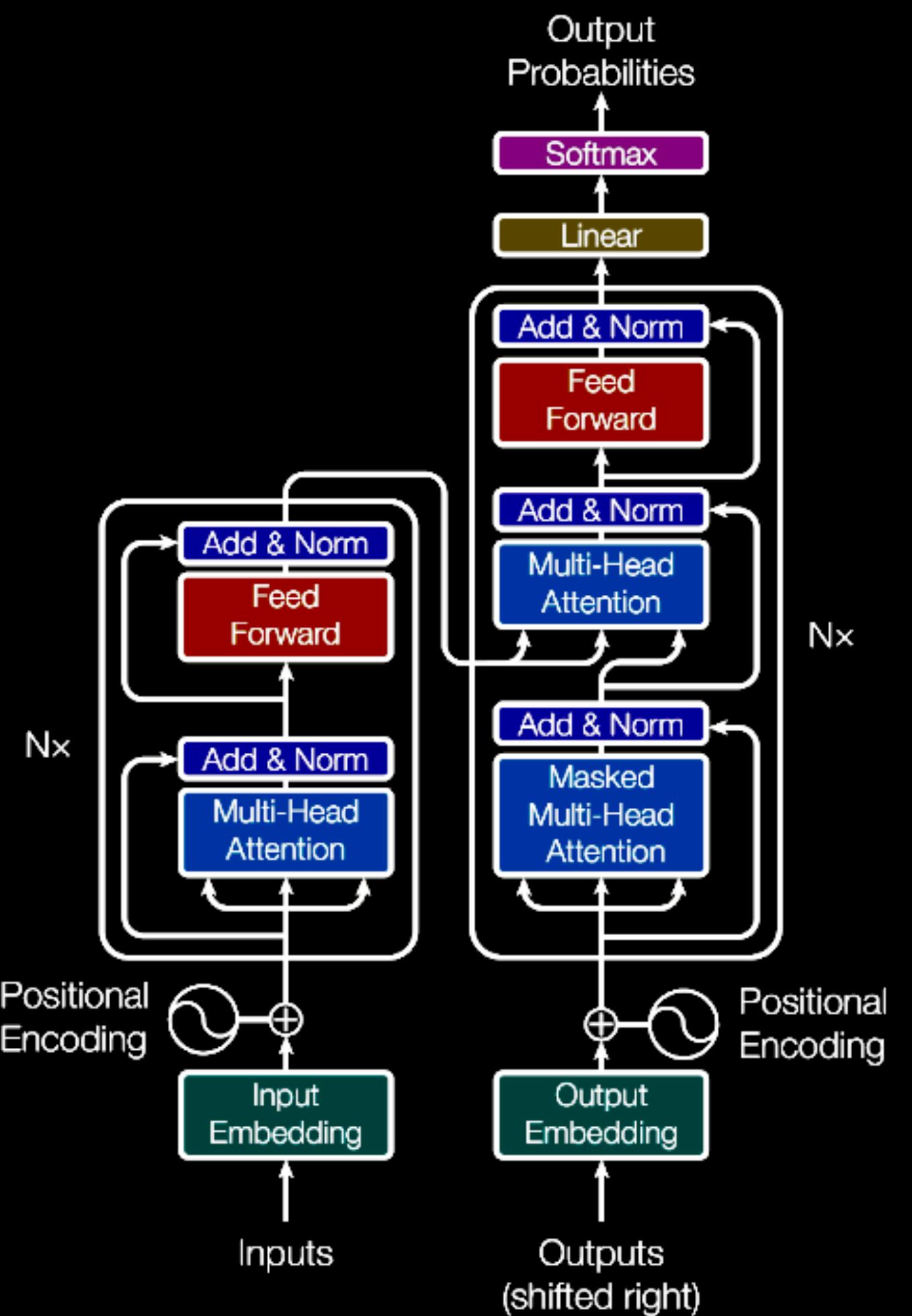


I need to go to the **bank** to make a deposit.

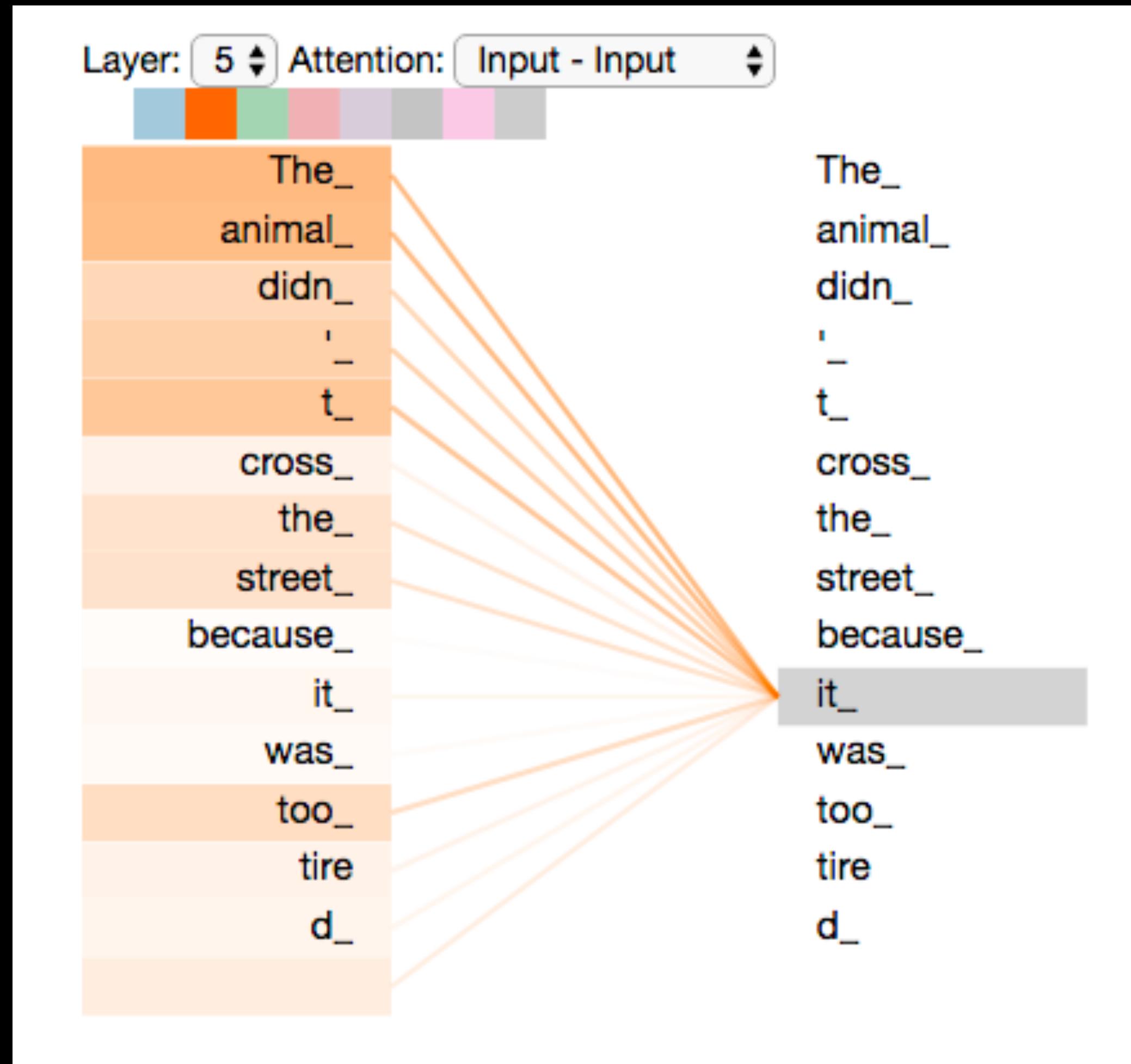
Transformers and Attention



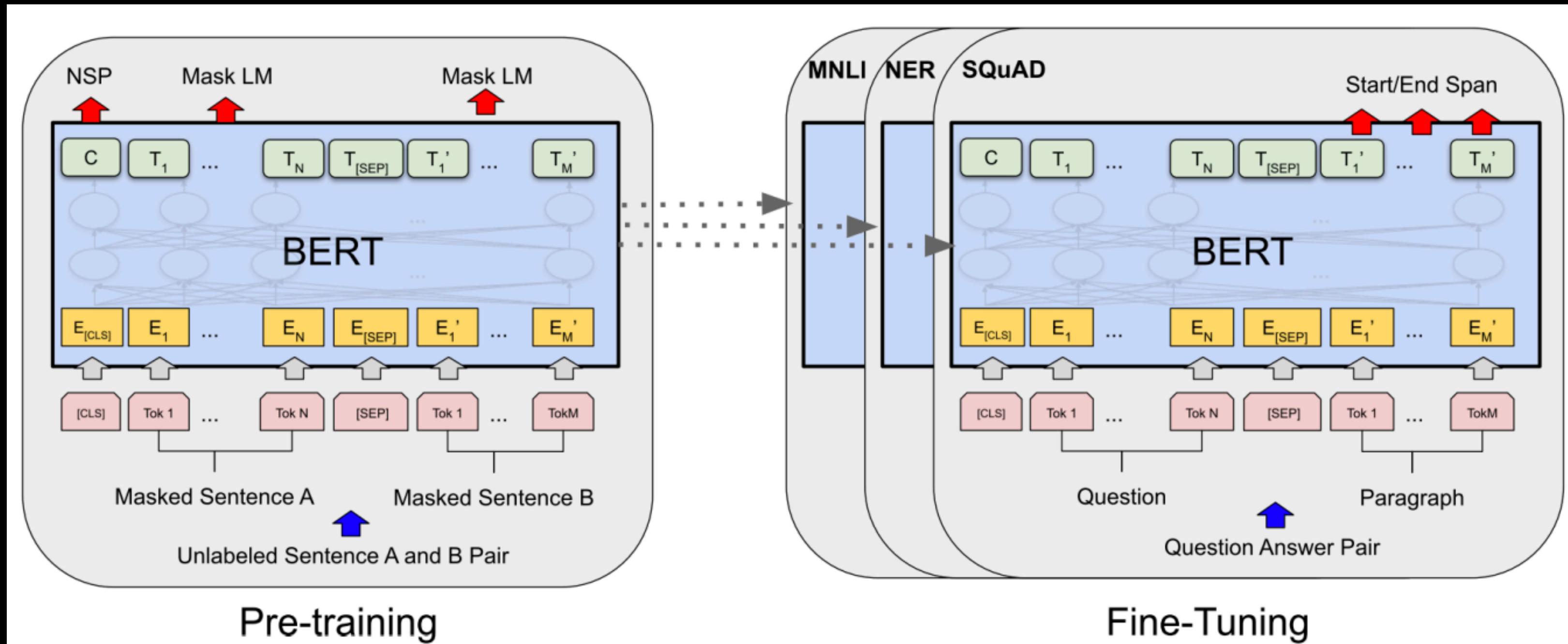
Transformers and Attention



Transformers and Attention



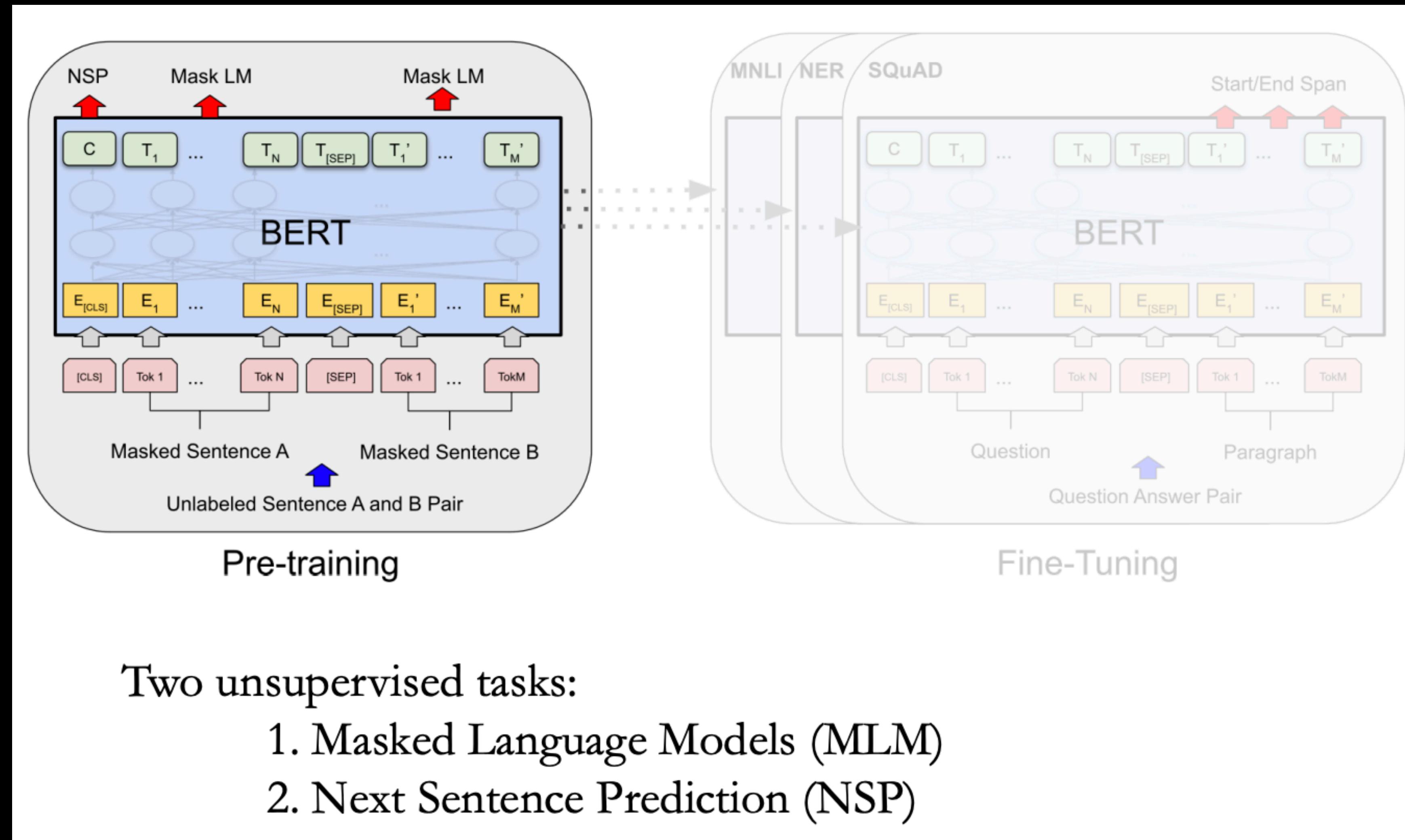
BERT: Framework



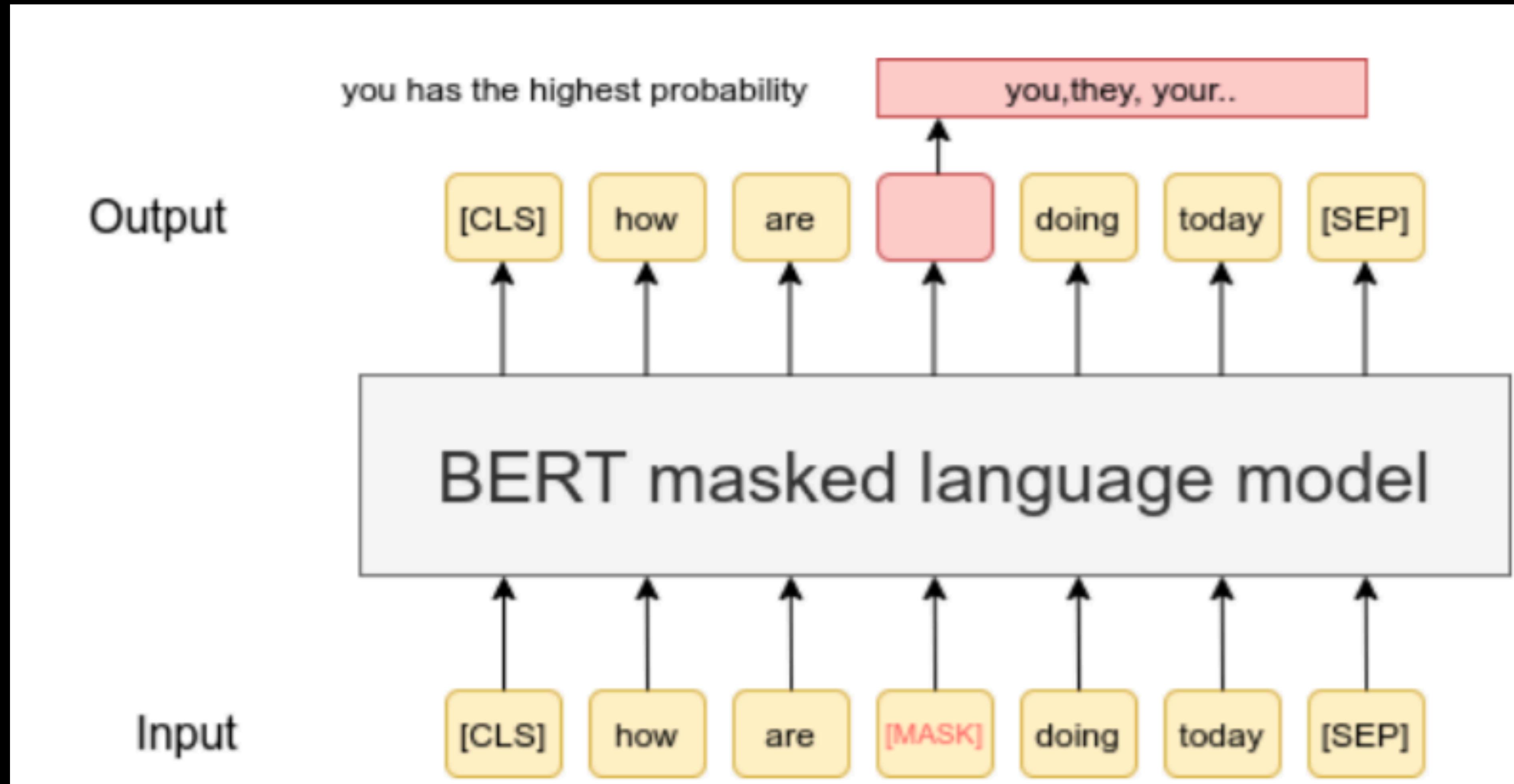
Pre-training : trained on unlabeled data over different pre-training tasks.

Fine-Tuning : fine-tuned parameters using labeled data from the downstream tasks.

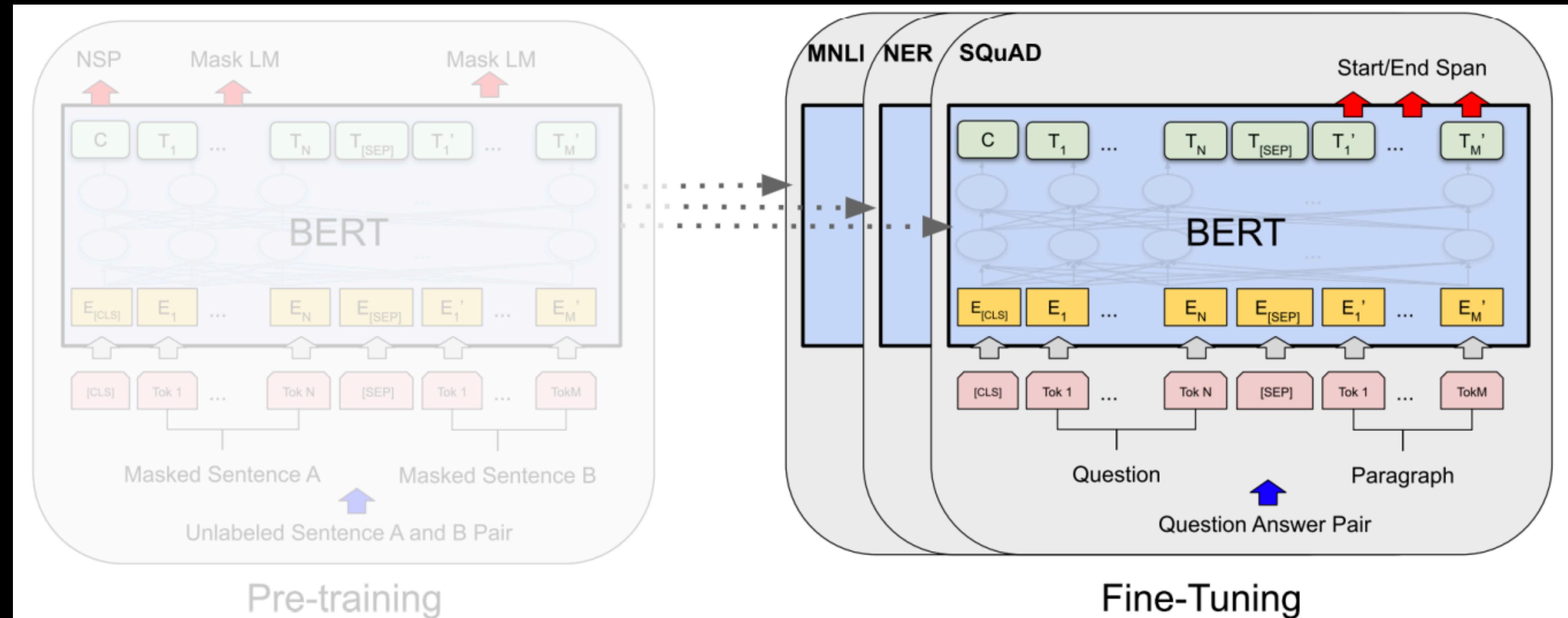
BERT: Pre-Training



BERT: Masked Language Modeling (MLM)



BERT: Fine-Tuning



Fine-Tuning : fine-tuned parameters using labeled data from the downstream tasks.

GTP-3: Training

Text: Second Law of Robotics: A robot must obey the orders given it by human beings



Generated training examples

Example #	Input (features)	Correct output (labels)
1	Second law of robotics :	a
2	Second law of robotics : a	robot
3	Second law of robotics : a robot	must
...		

GPT-3: Base Model

What is the capital of Morocco? (Hint: The word Morocco is borrowed from Arabic.)



Stop by the Towers for evening ice cream service from 6:00 PM to 6:30 P.M. and enjoy the historic elegance of the old century room which features wood moldings and a stone fireplace. We will close at 8:45 PM due to the last train from Skyline Drive at 9:15.

After dinner, we will head out to see the movie. Click [here](#) for a trailer...

GPT-3: ChatGPT

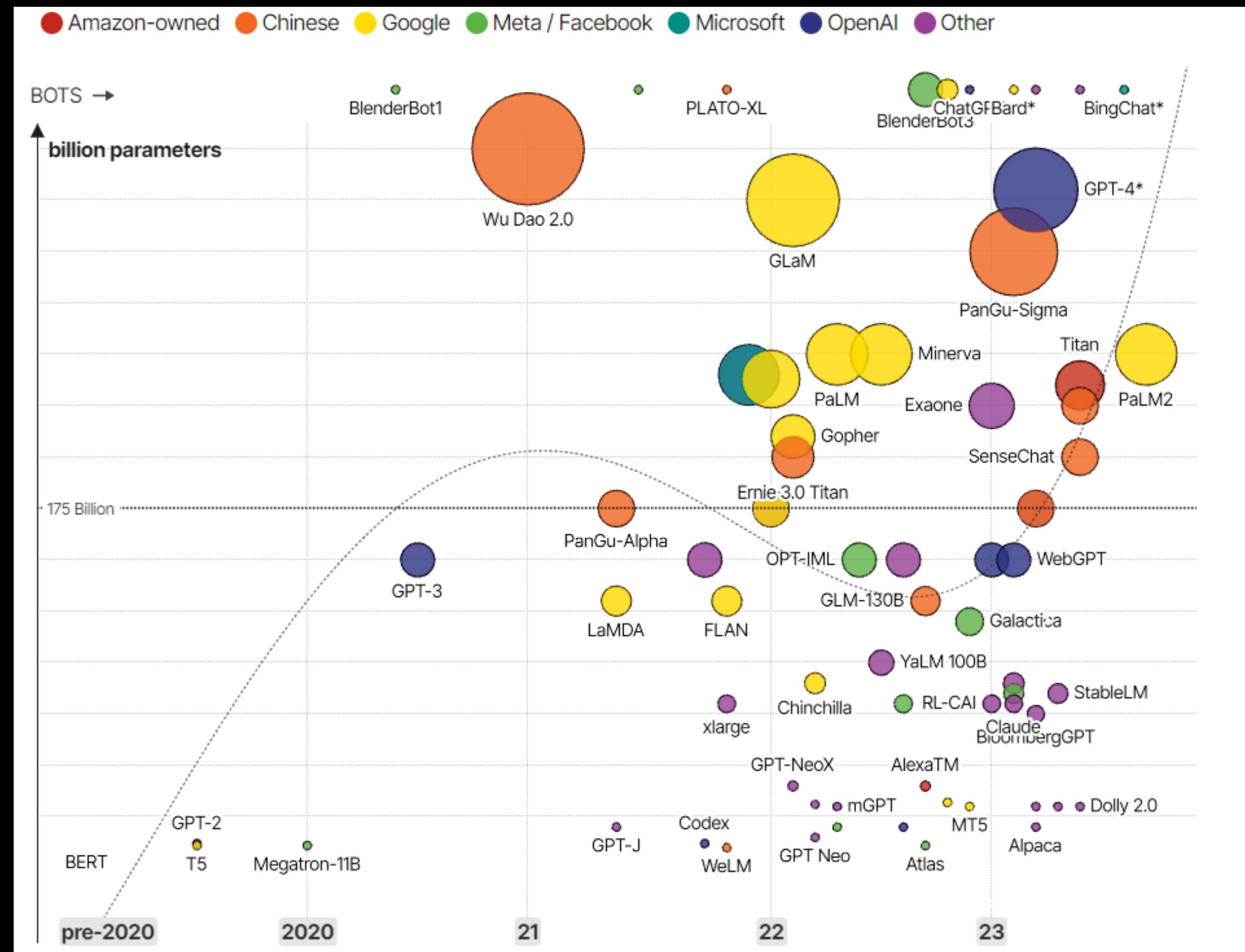
What is the capital of Morocco?



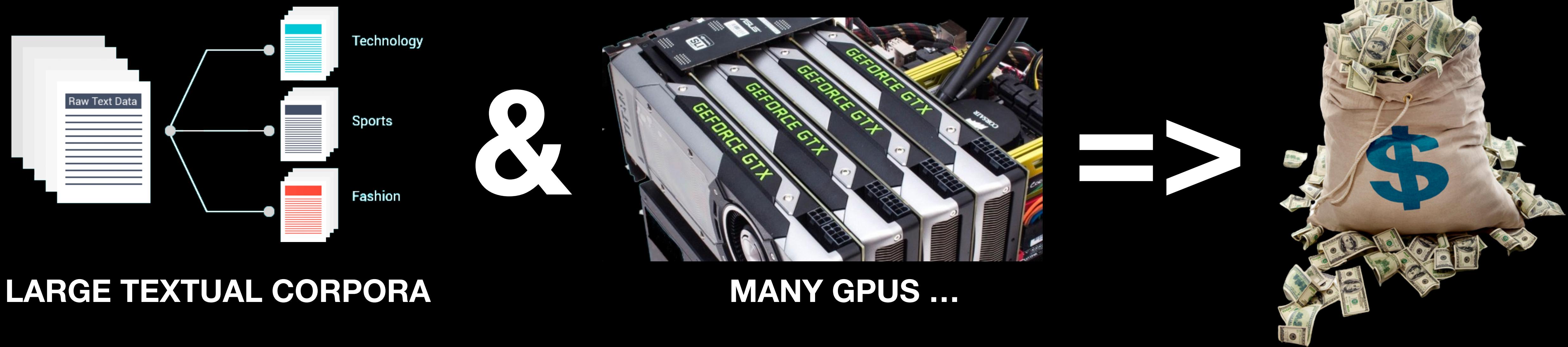
Rabat is the capital of Morocco.

How can we leverage NLP
achievements for
underrepresented languages?

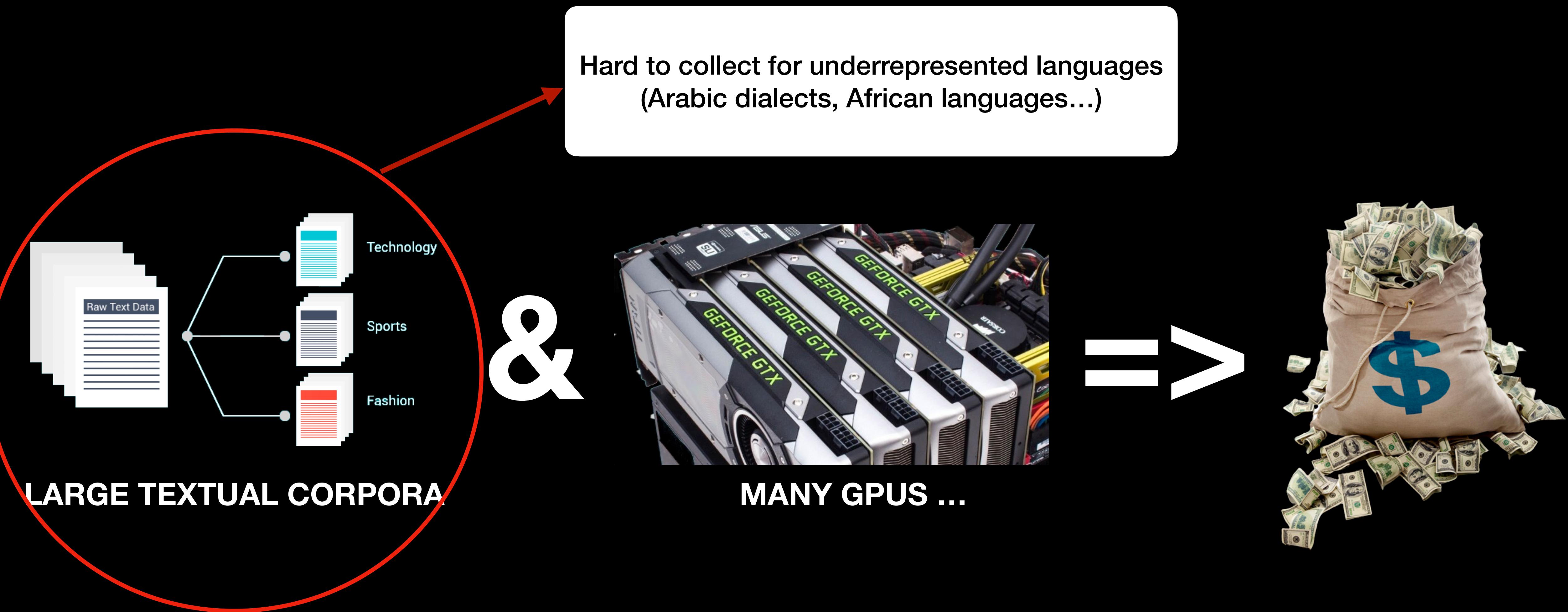
NLP \leftrightarrow Large Language Models



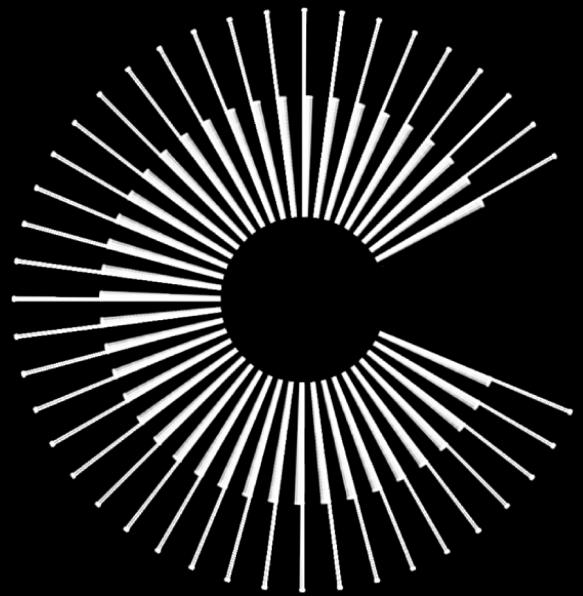
How to build a similar language model?



How to build a similar language model?



Our team: UM6P-CC



College of
Computing
UM6P

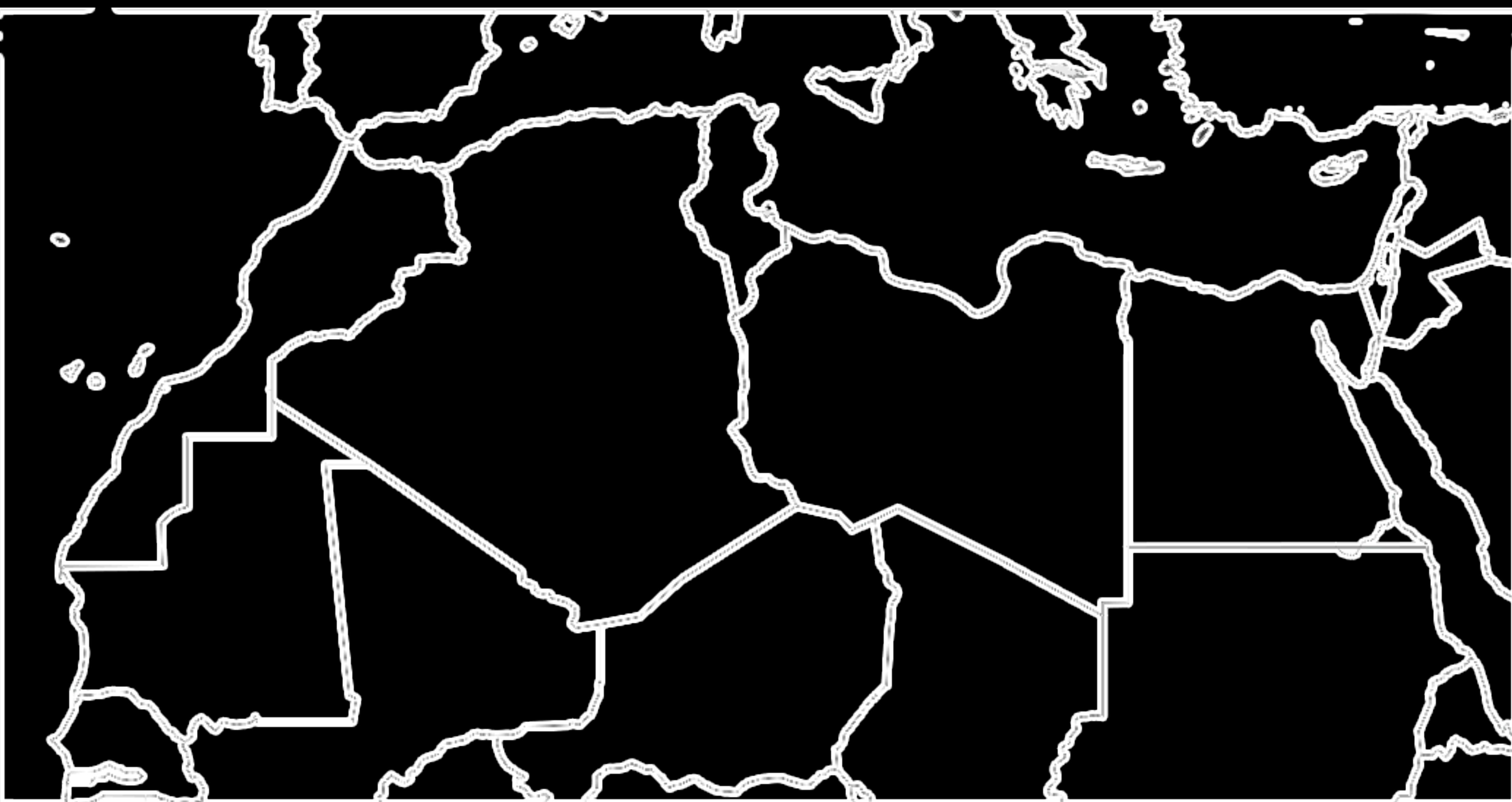


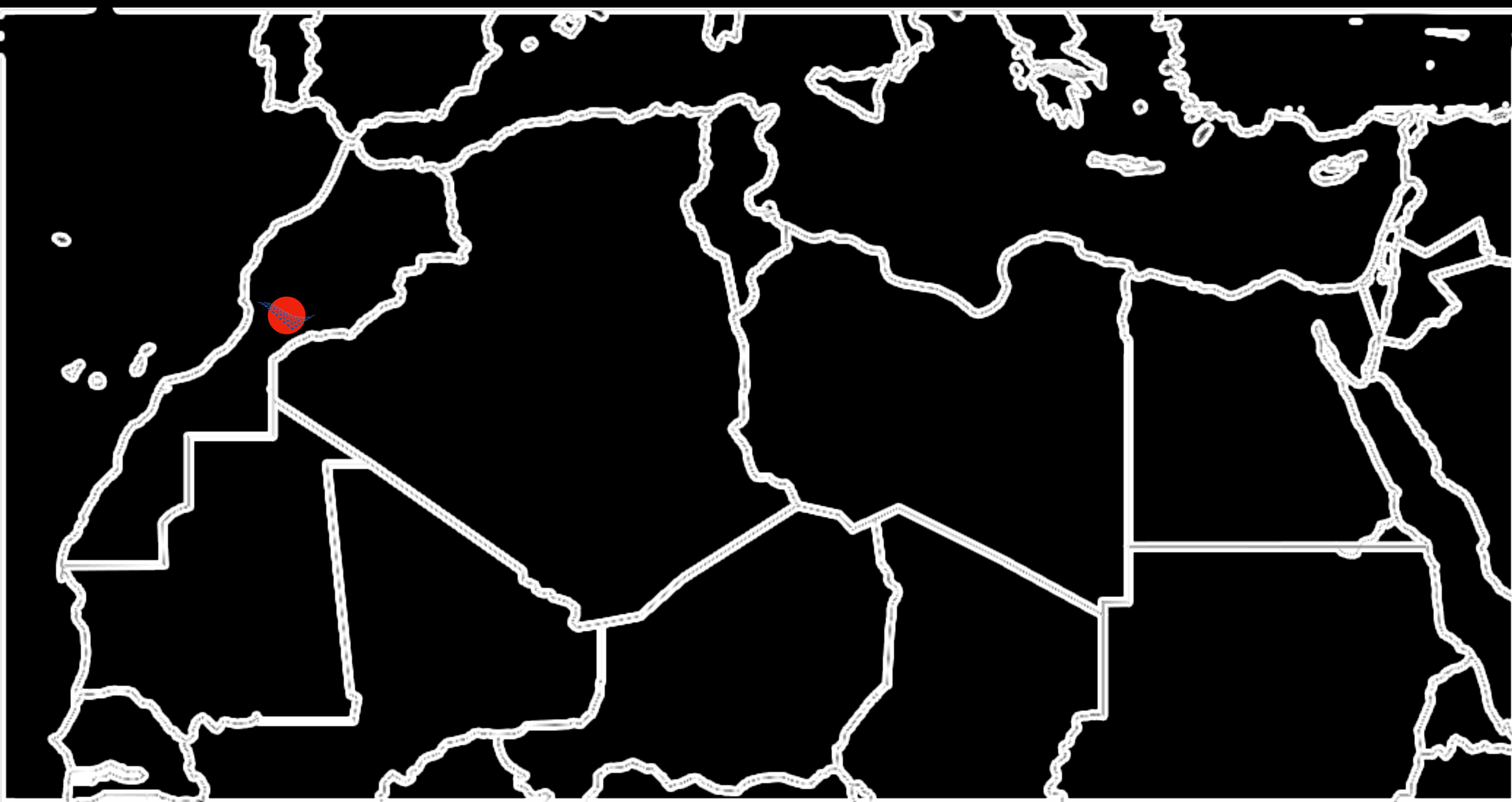
- Dr. Abdelkader EL MAHDAOUY
- Post-Doctoral Researcher In Natural Language Processing / Machine Learning
- University Mohammed VI Polytechnic, Morocco
- Dr. Ismail BERRADA
- UM6-CC professor and researcher in applied ML
- University Mohammed VI Polytechnic, Morocco
- Abdellah EL MEKKI
- PhD Student in Natural Language Processing / Machine Learning
- University Mohammed VI Polytechnic, Morocco

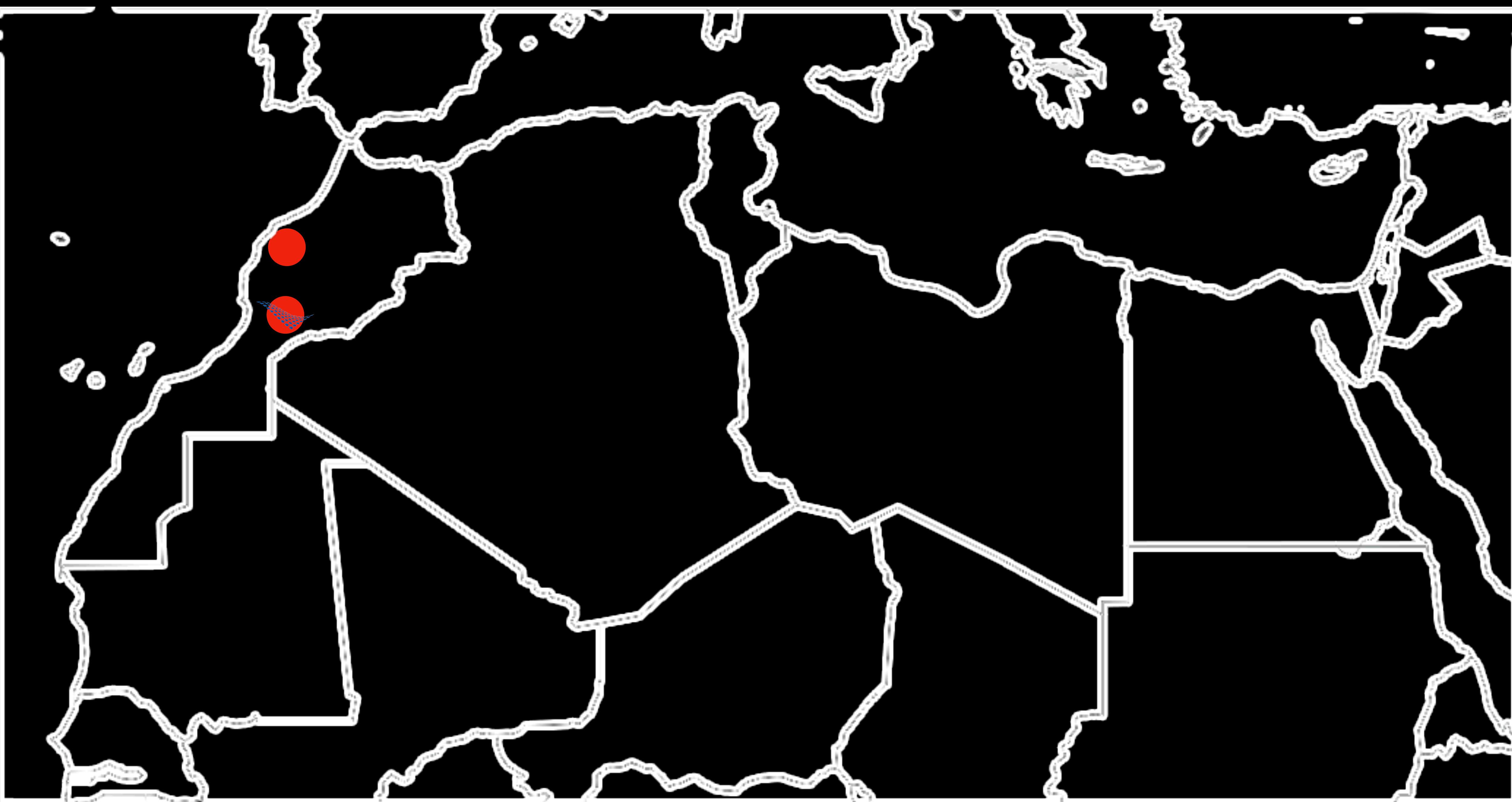
Modern Standard Arabic (MSA) NLP

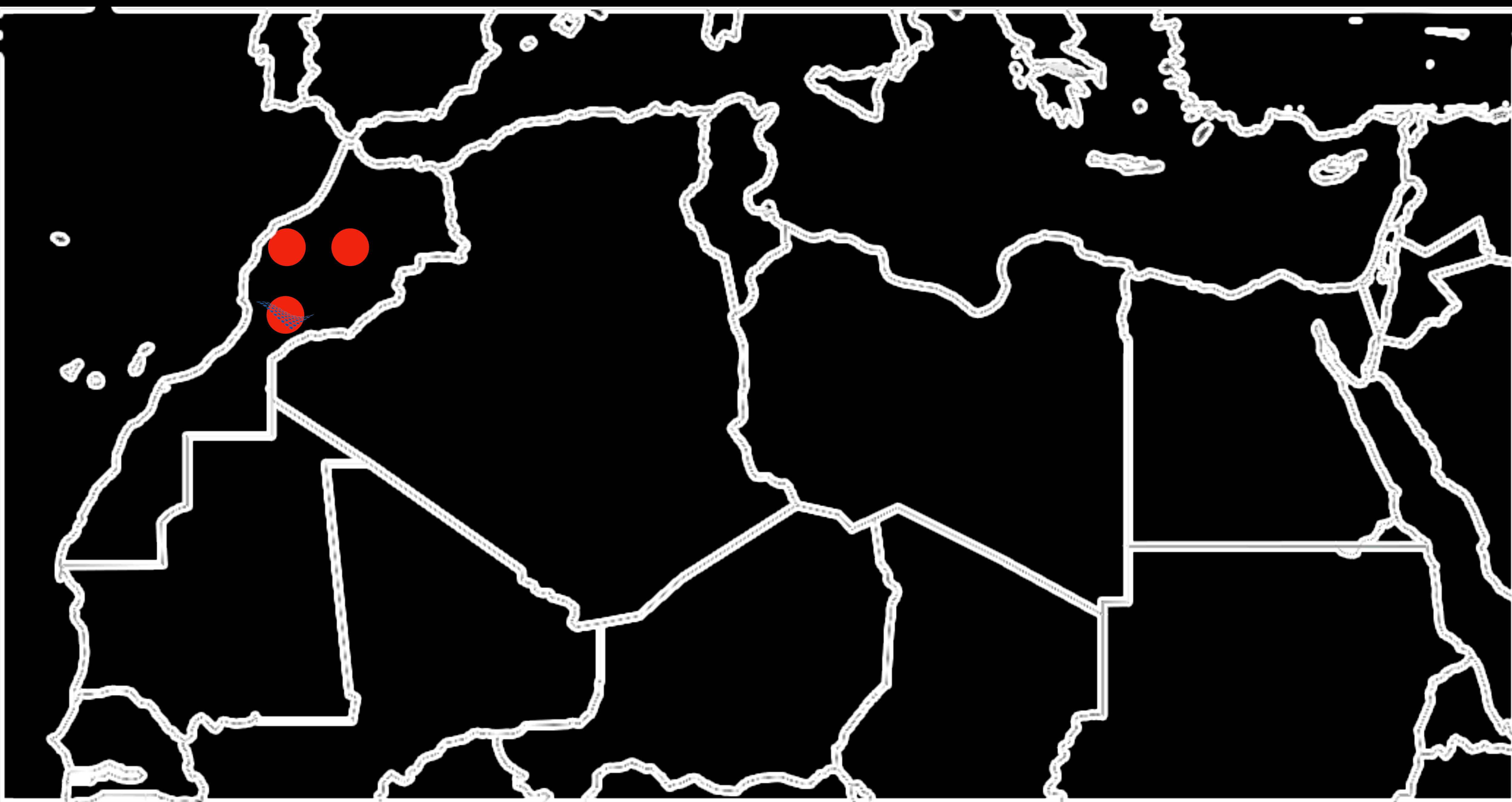
- Good Machine Translation results.
- Good conversation systems.
- Availability of corpora to train Arabic NLP systems on.

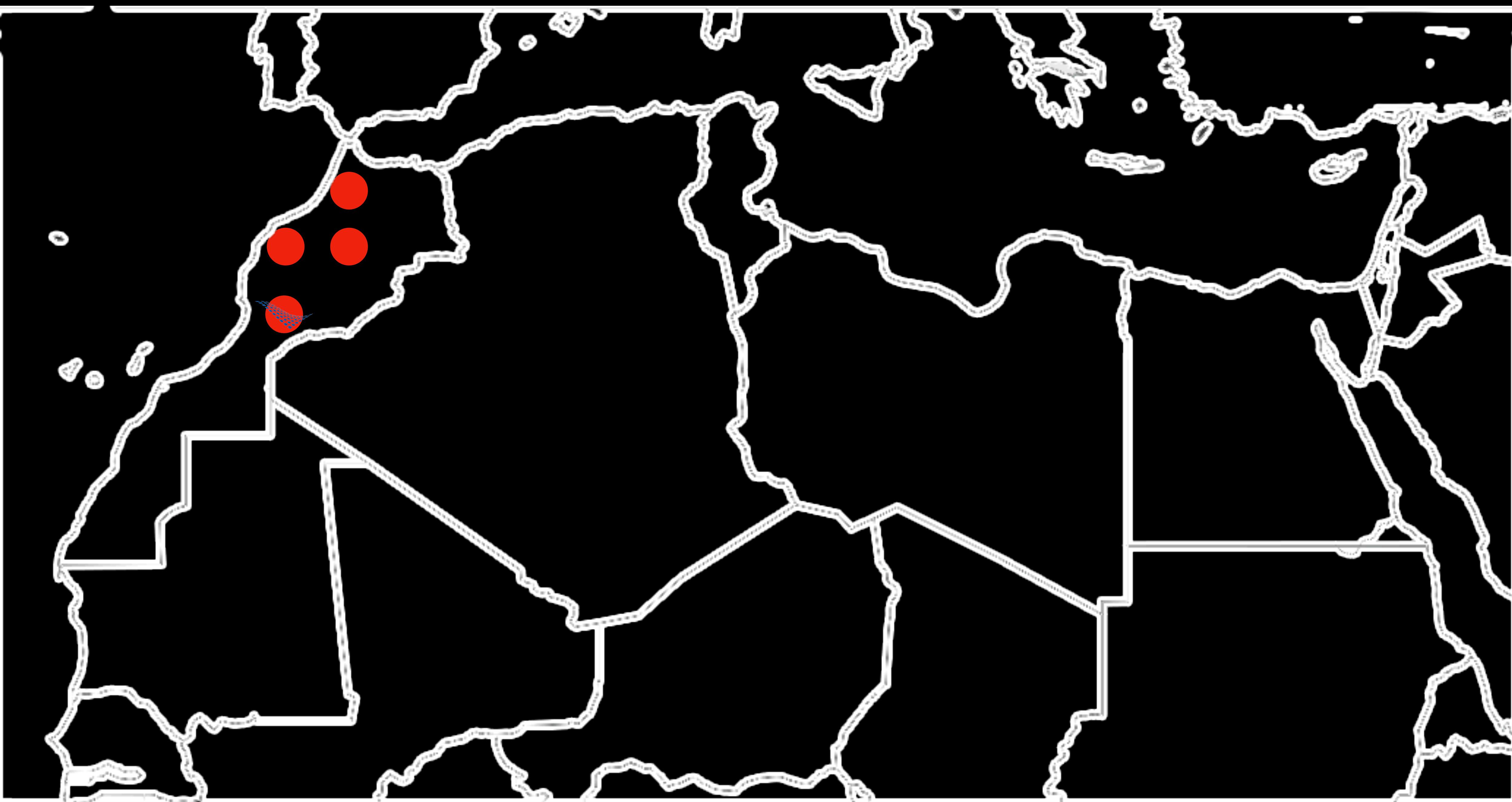


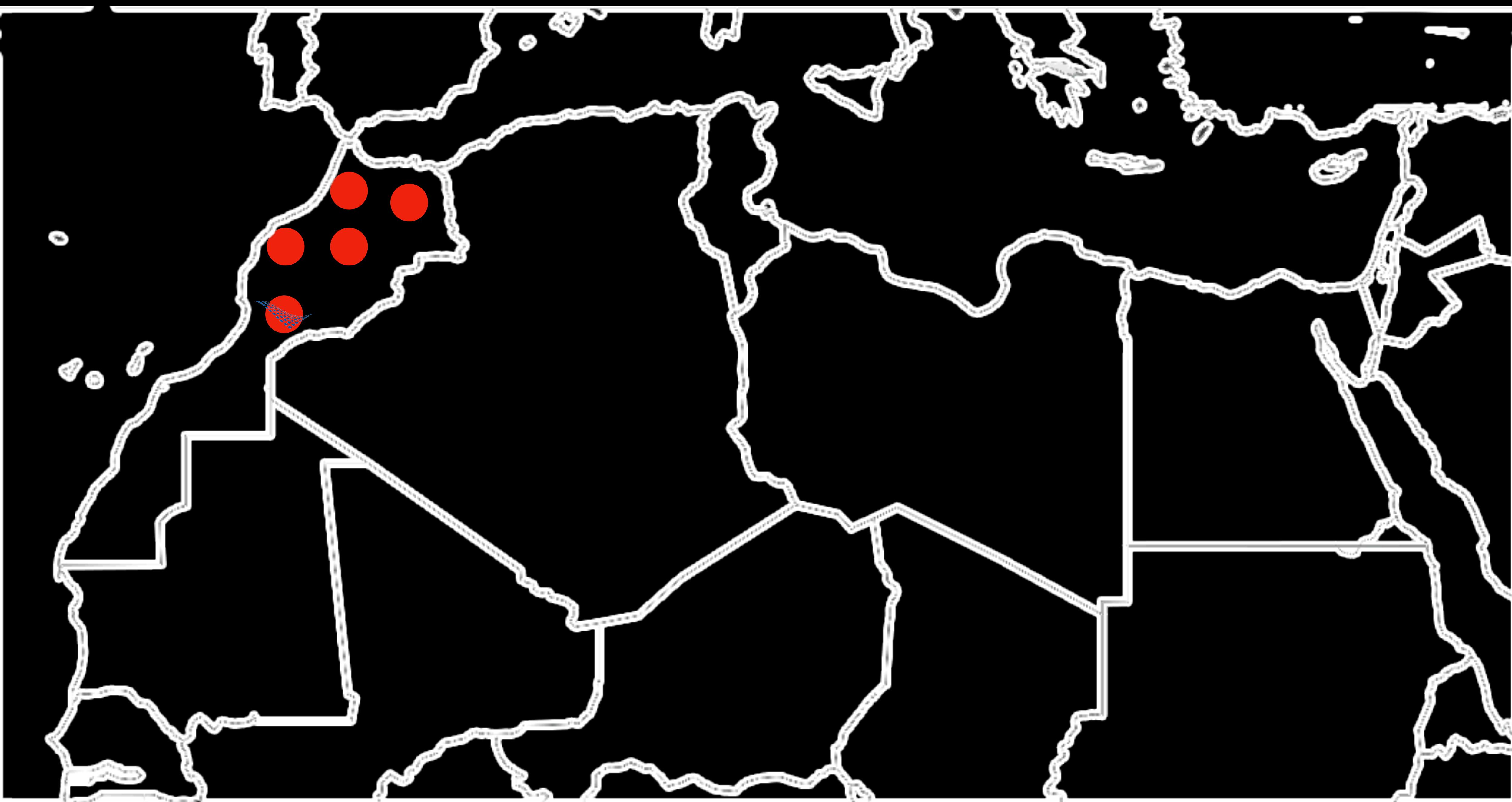


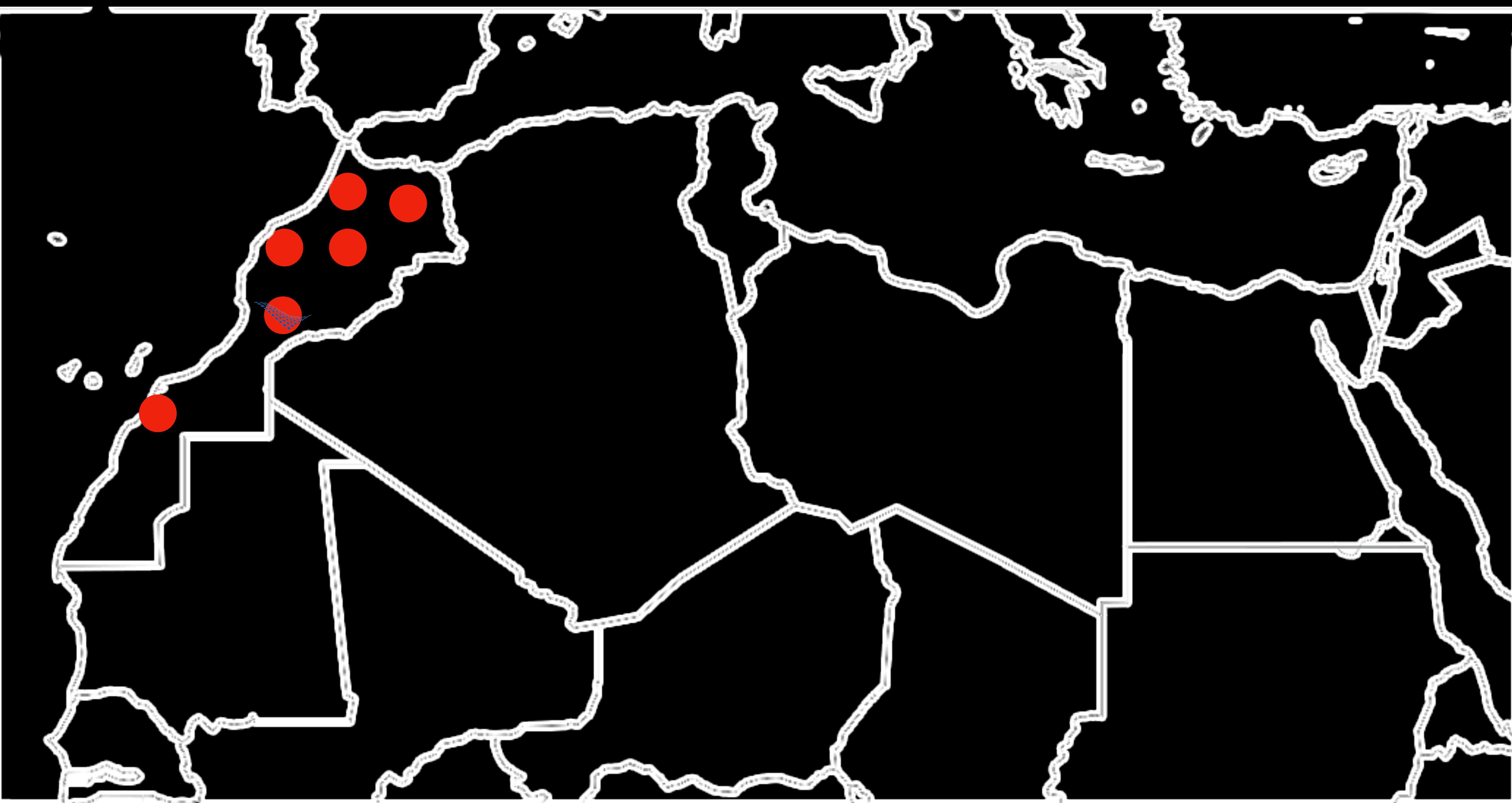


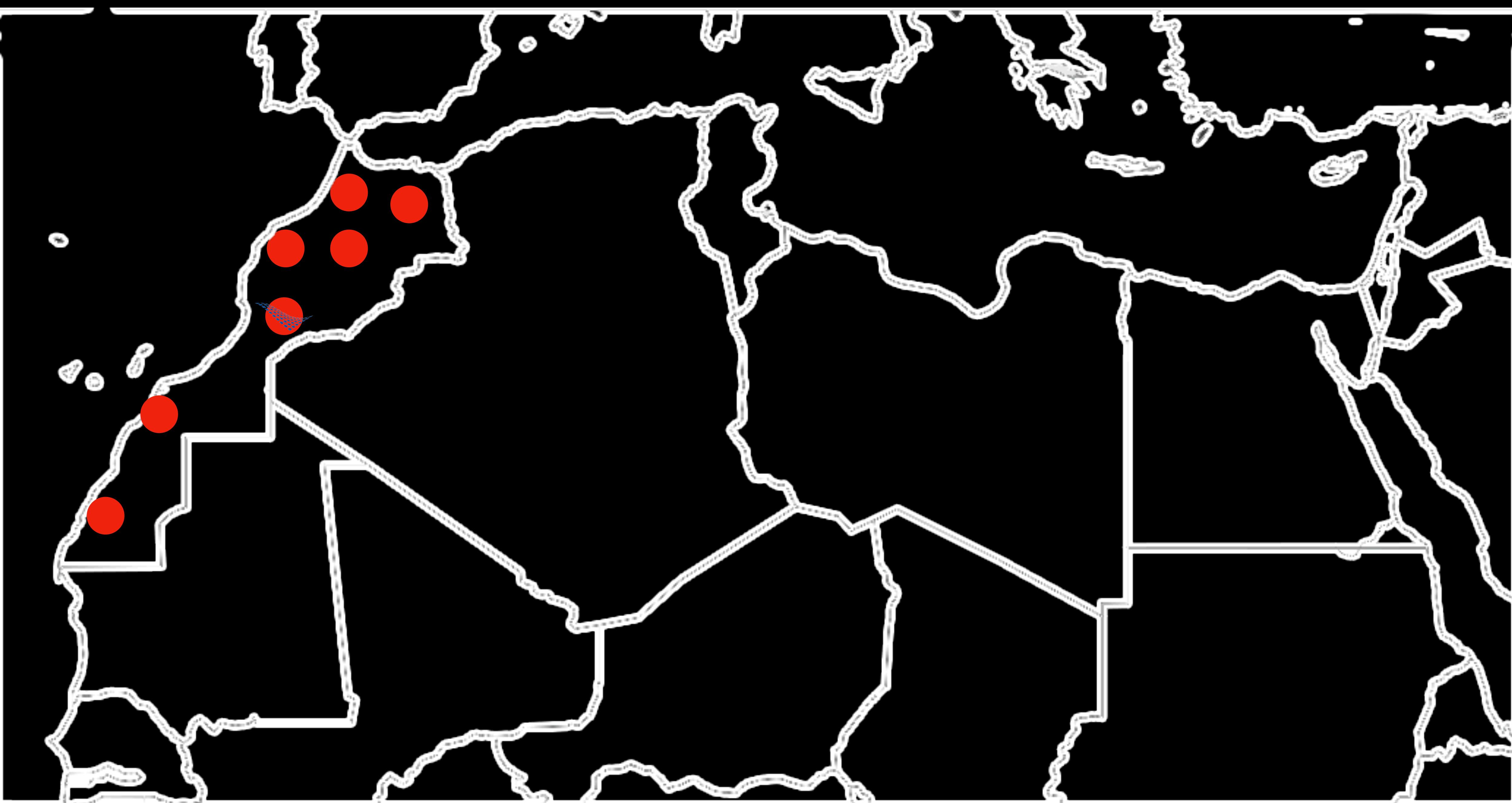


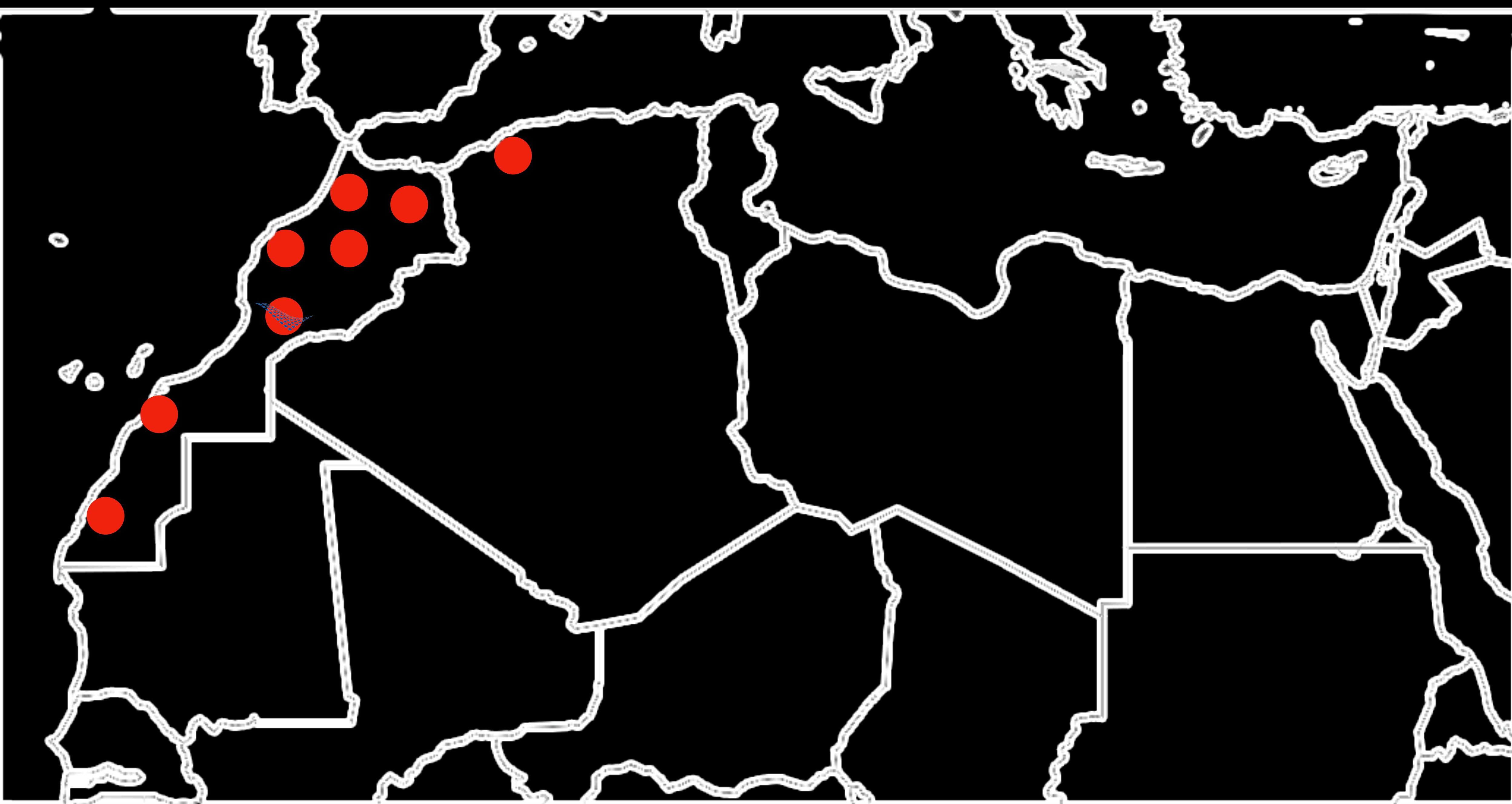


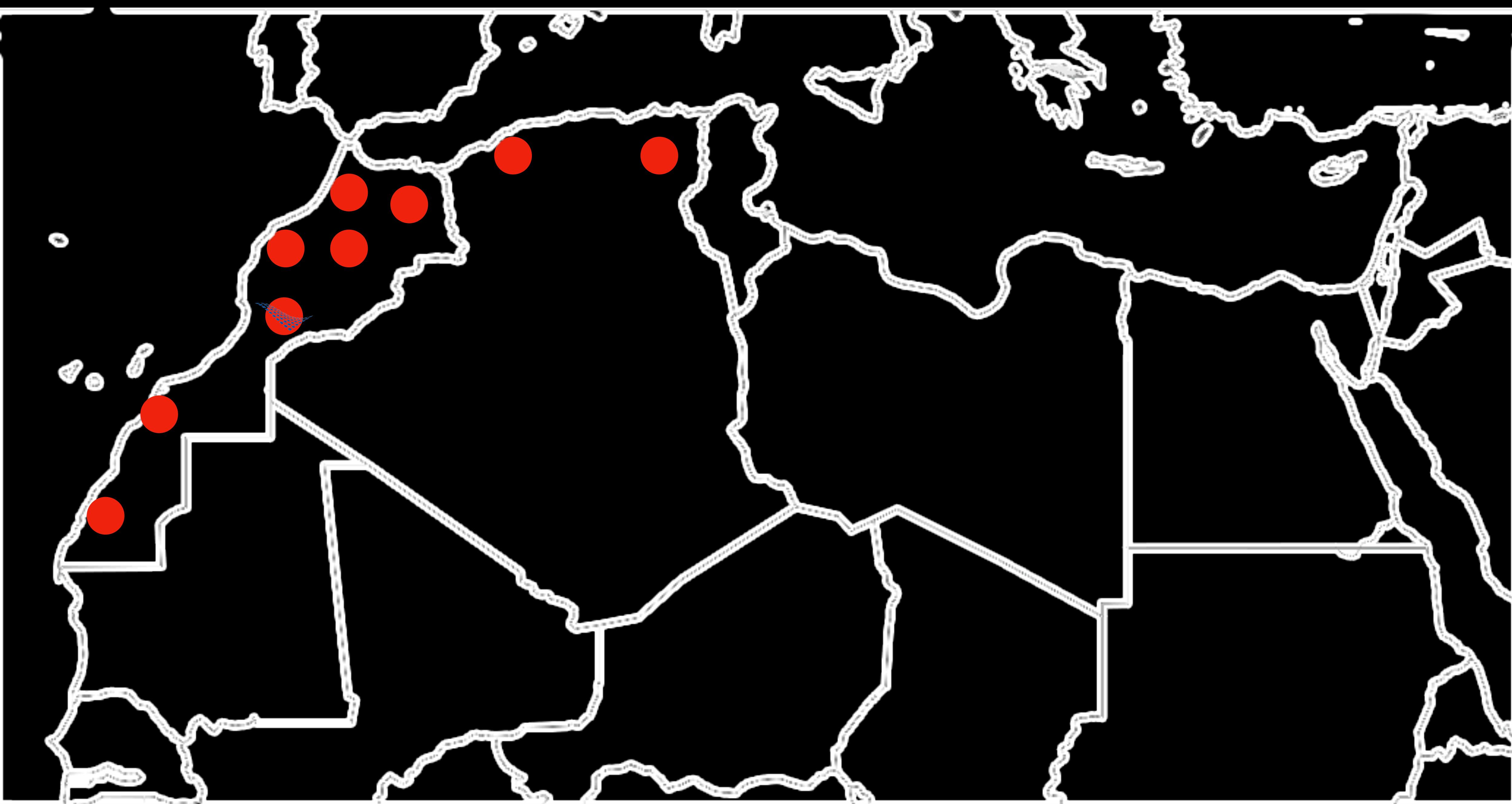


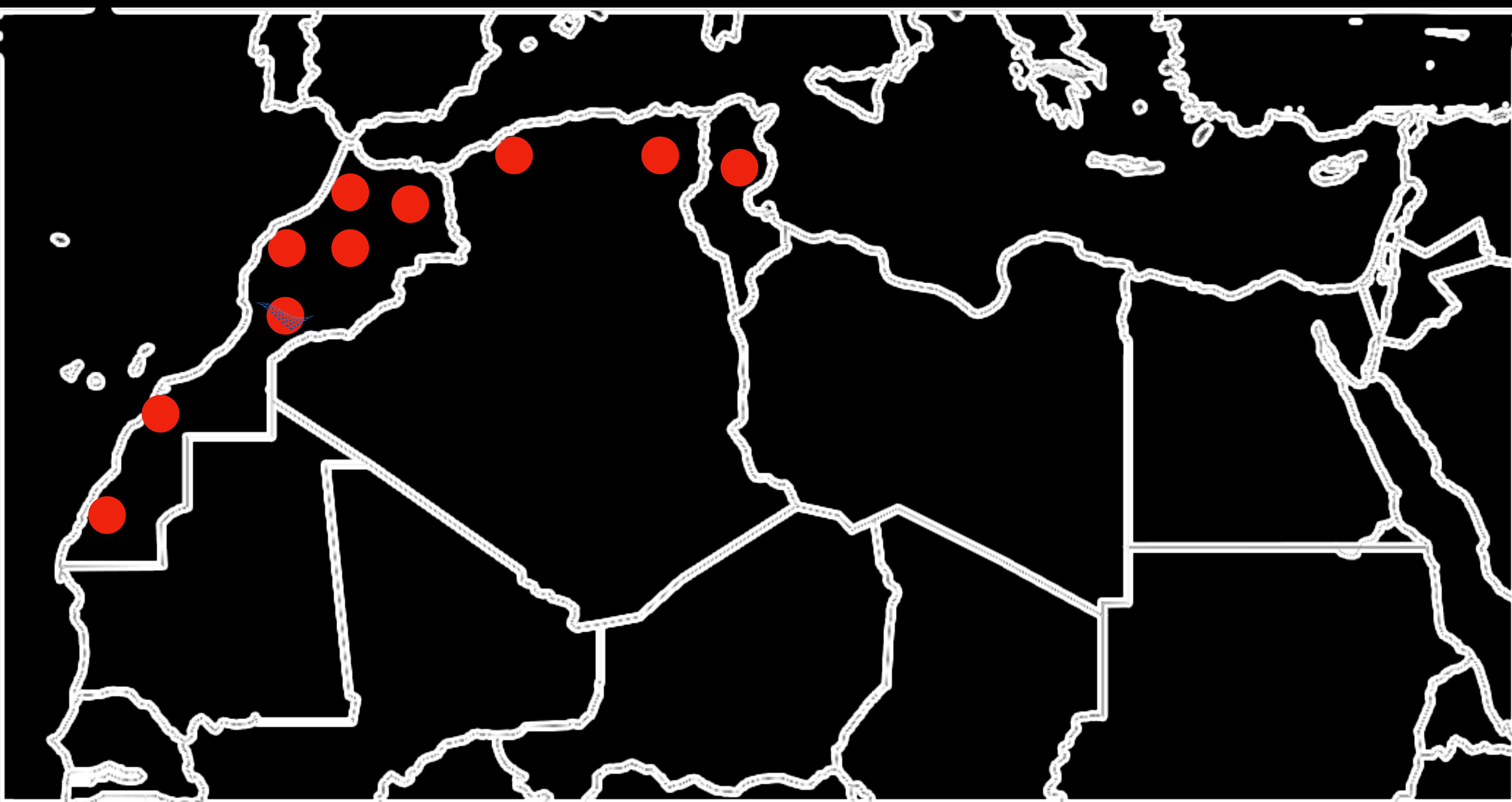


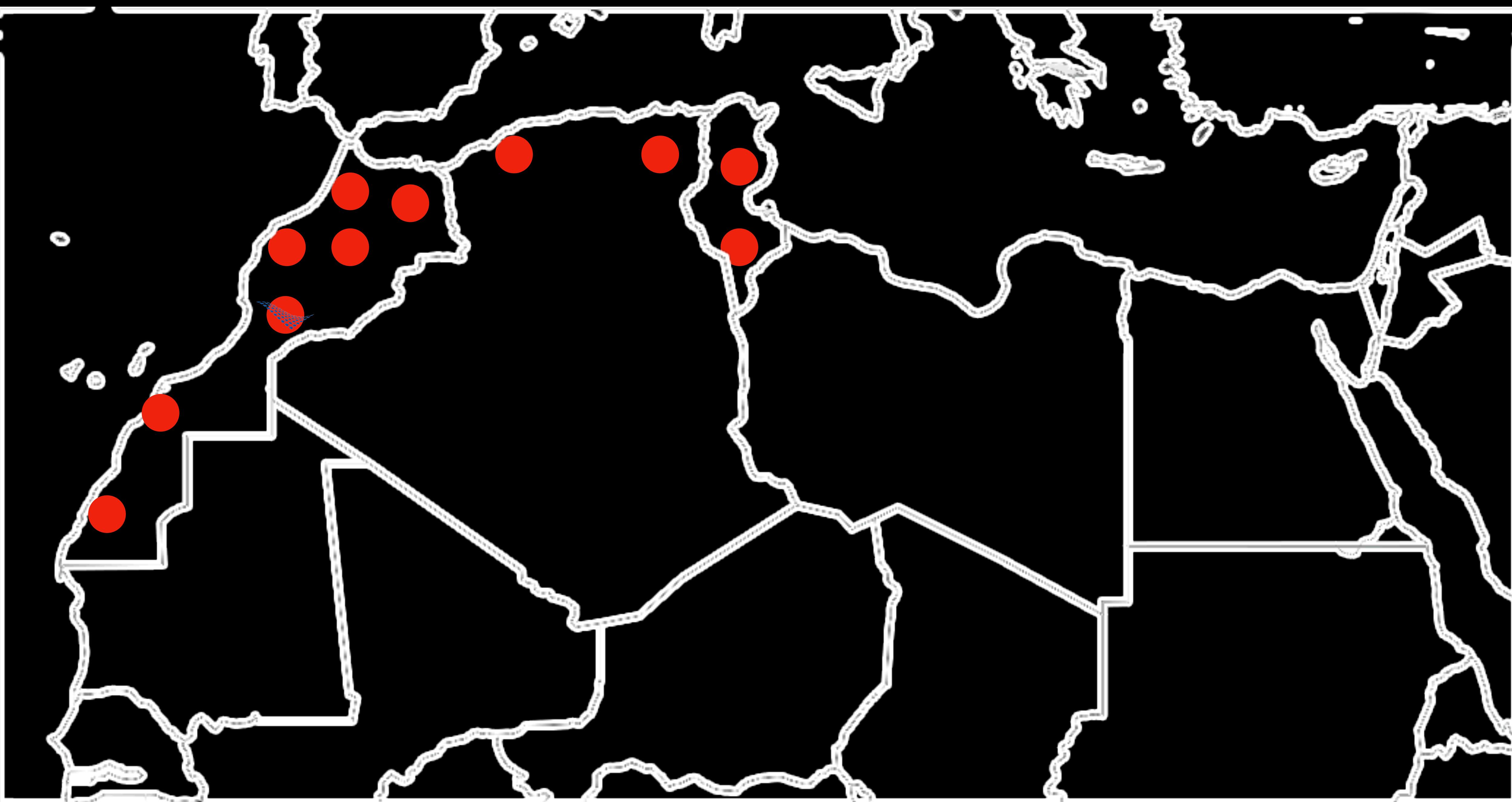


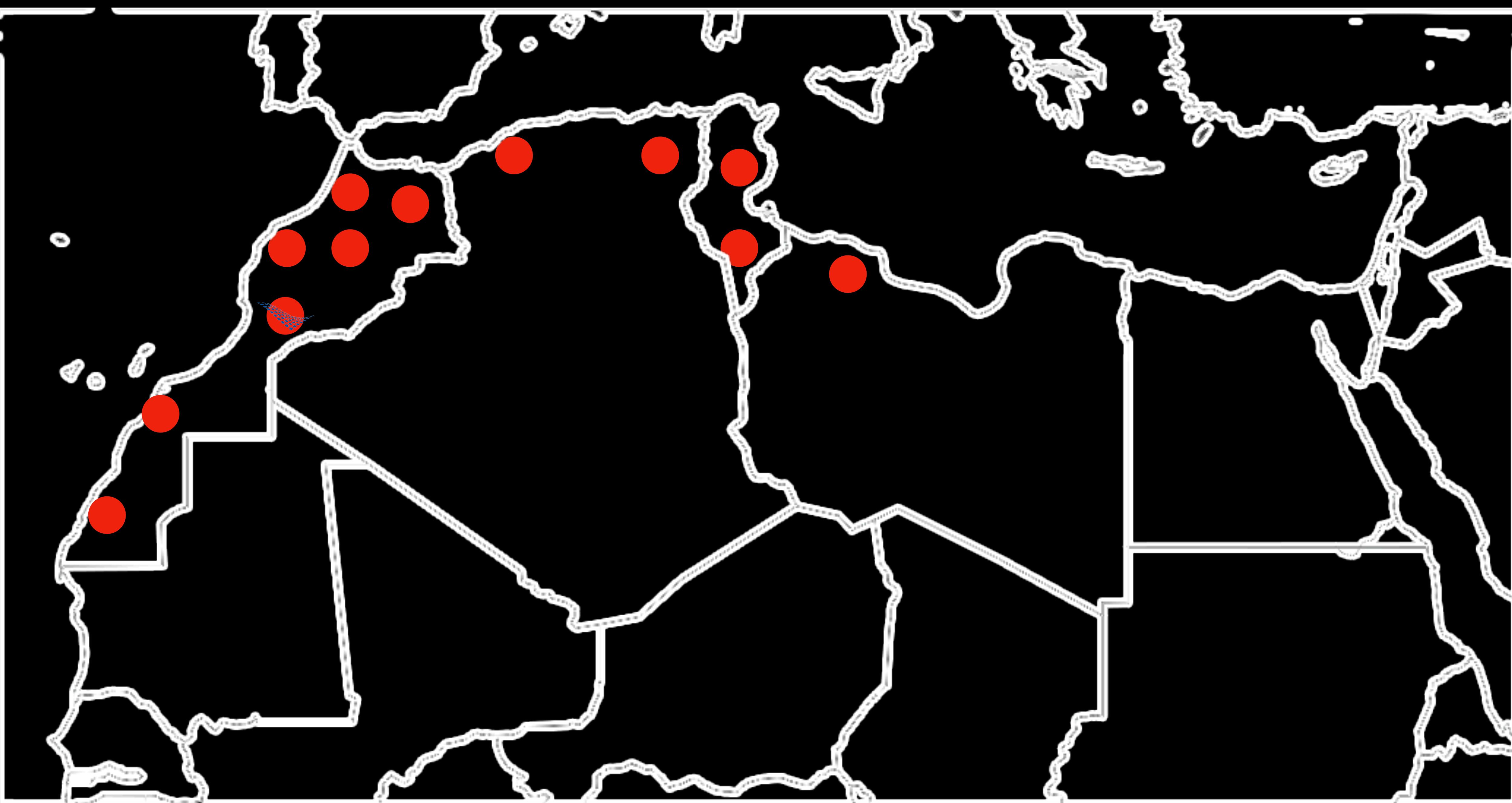


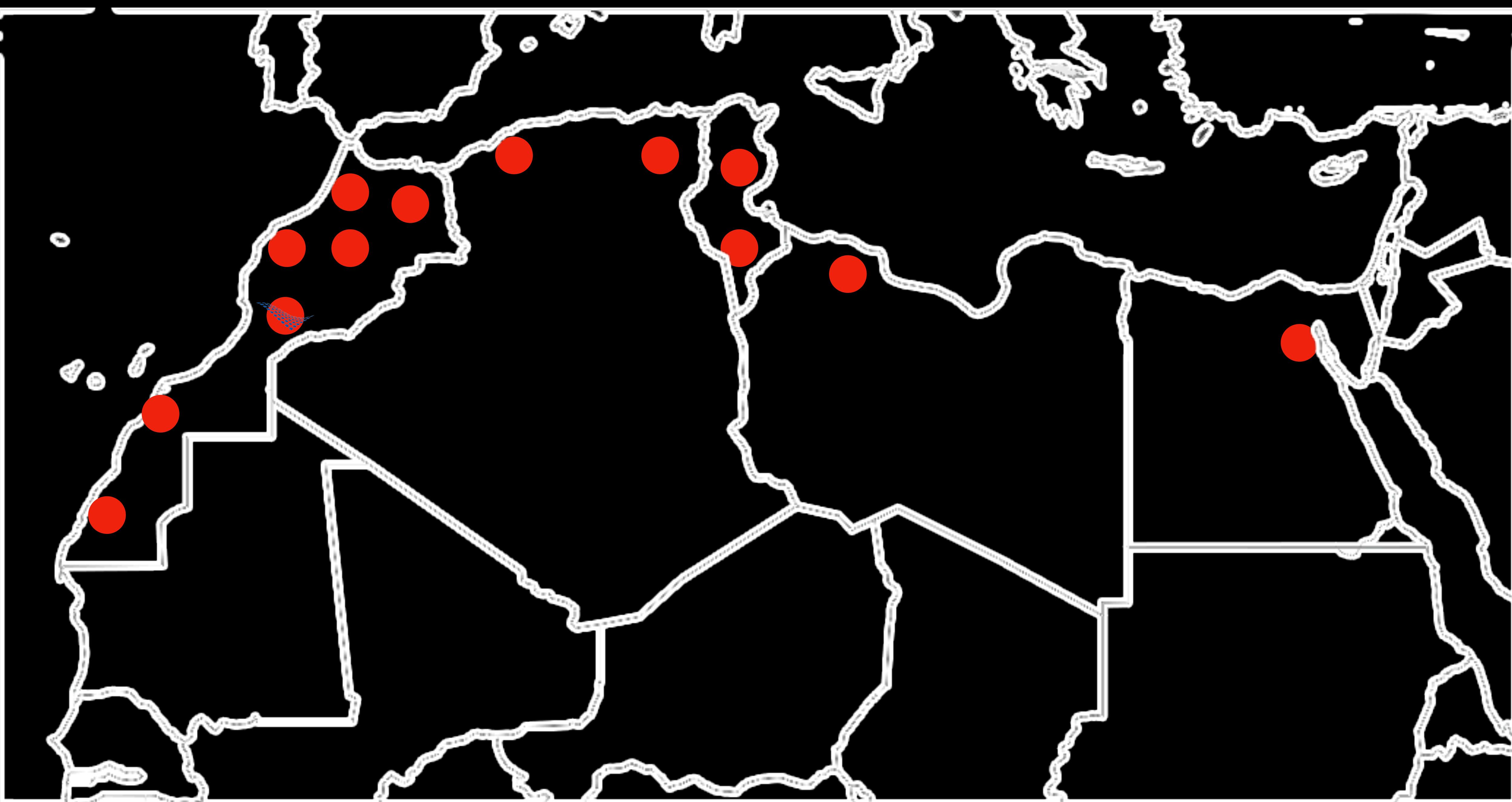


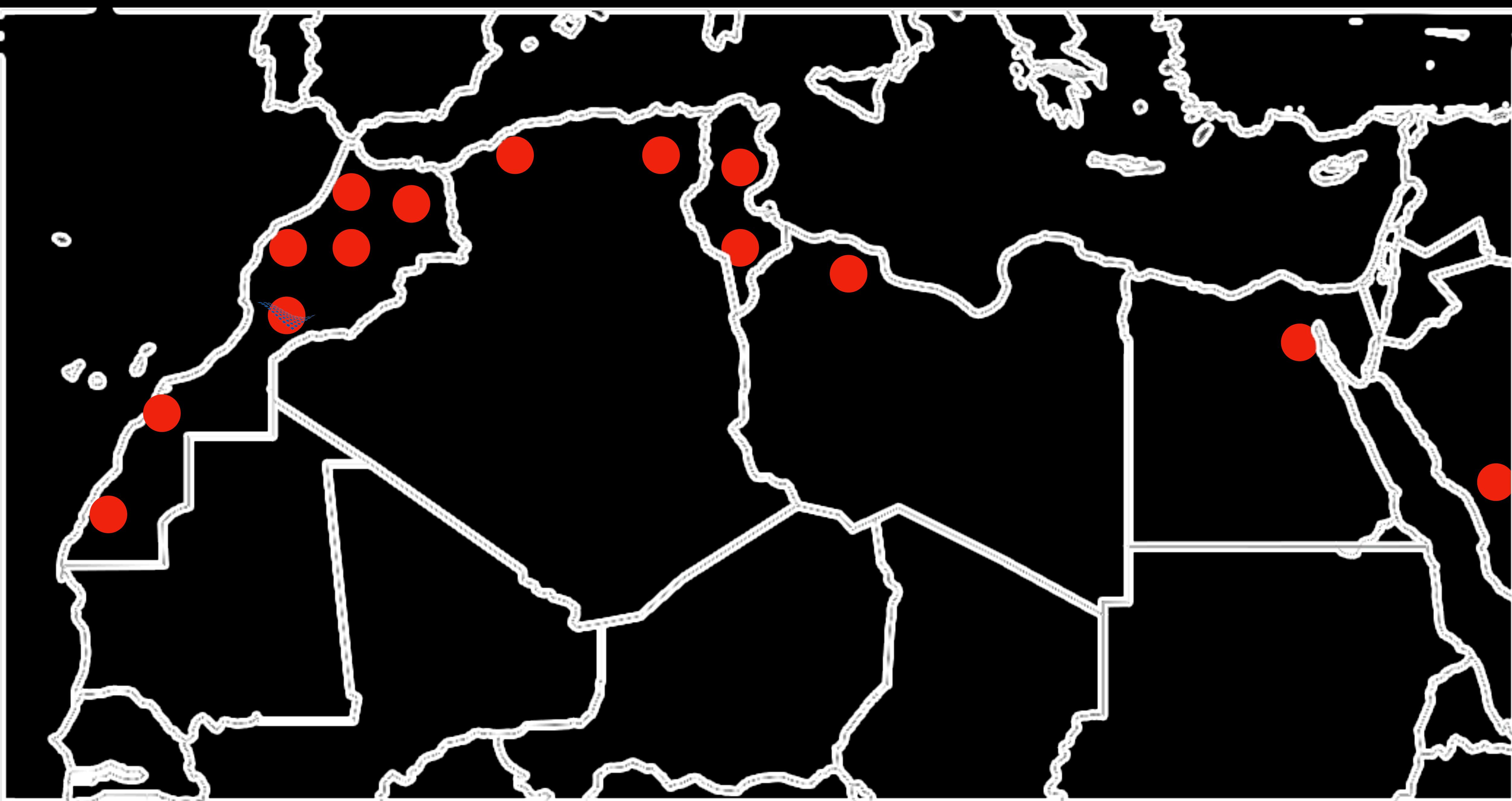


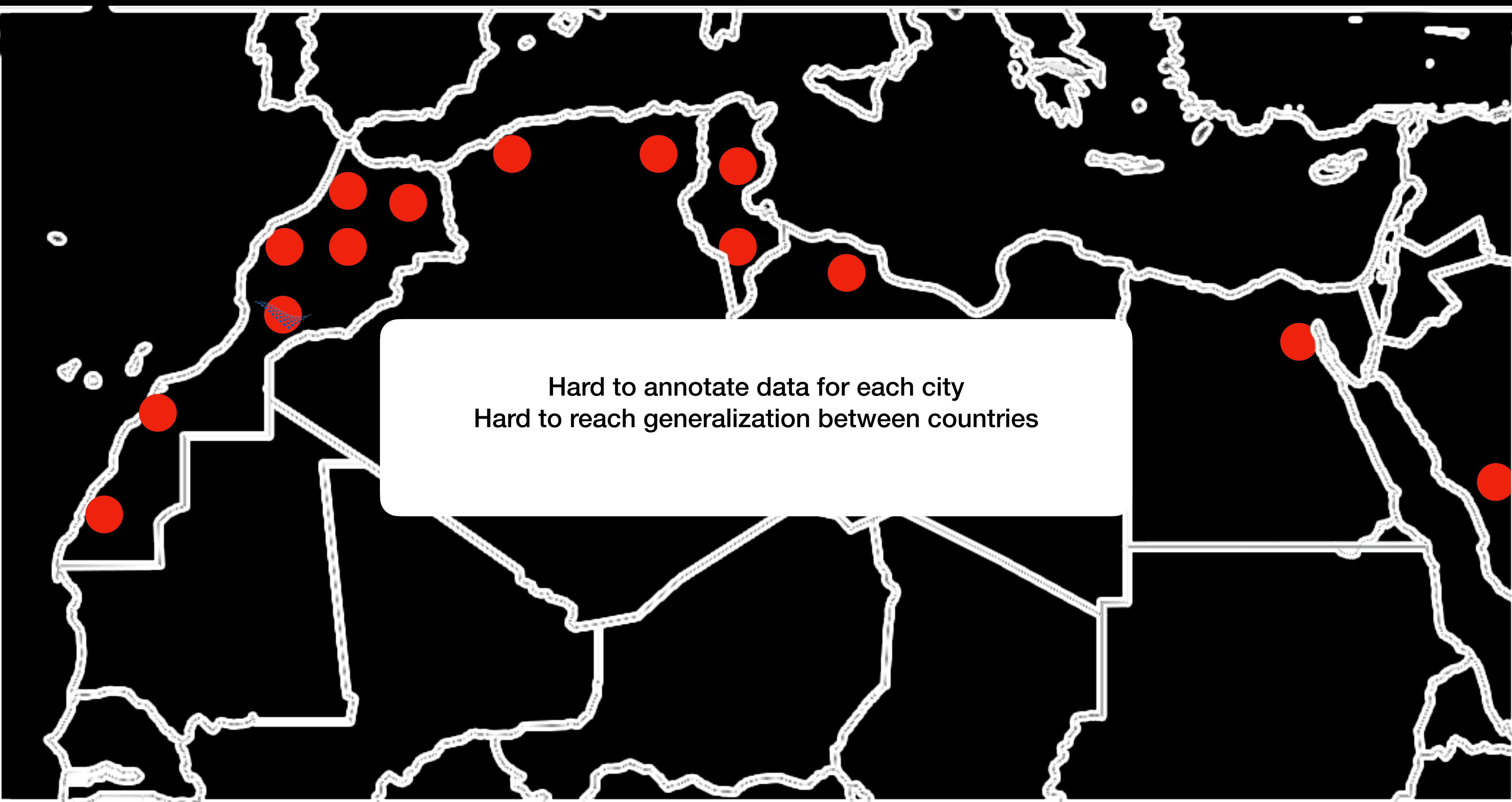














Everybody wants
Dialectal Arabic NLP
systems.

Hard to annotate data for each city
Hard to reach generalization between countries



Everybody wants
Dialectal Arabic NLP
systems.

Hard to annotate data for each city
Hard to reach generalization between countries



Nobody can put
the work in
annotating data.



Everybody wants
Dialectal Arabic NLP
systems.

Hard to annotate data for each city

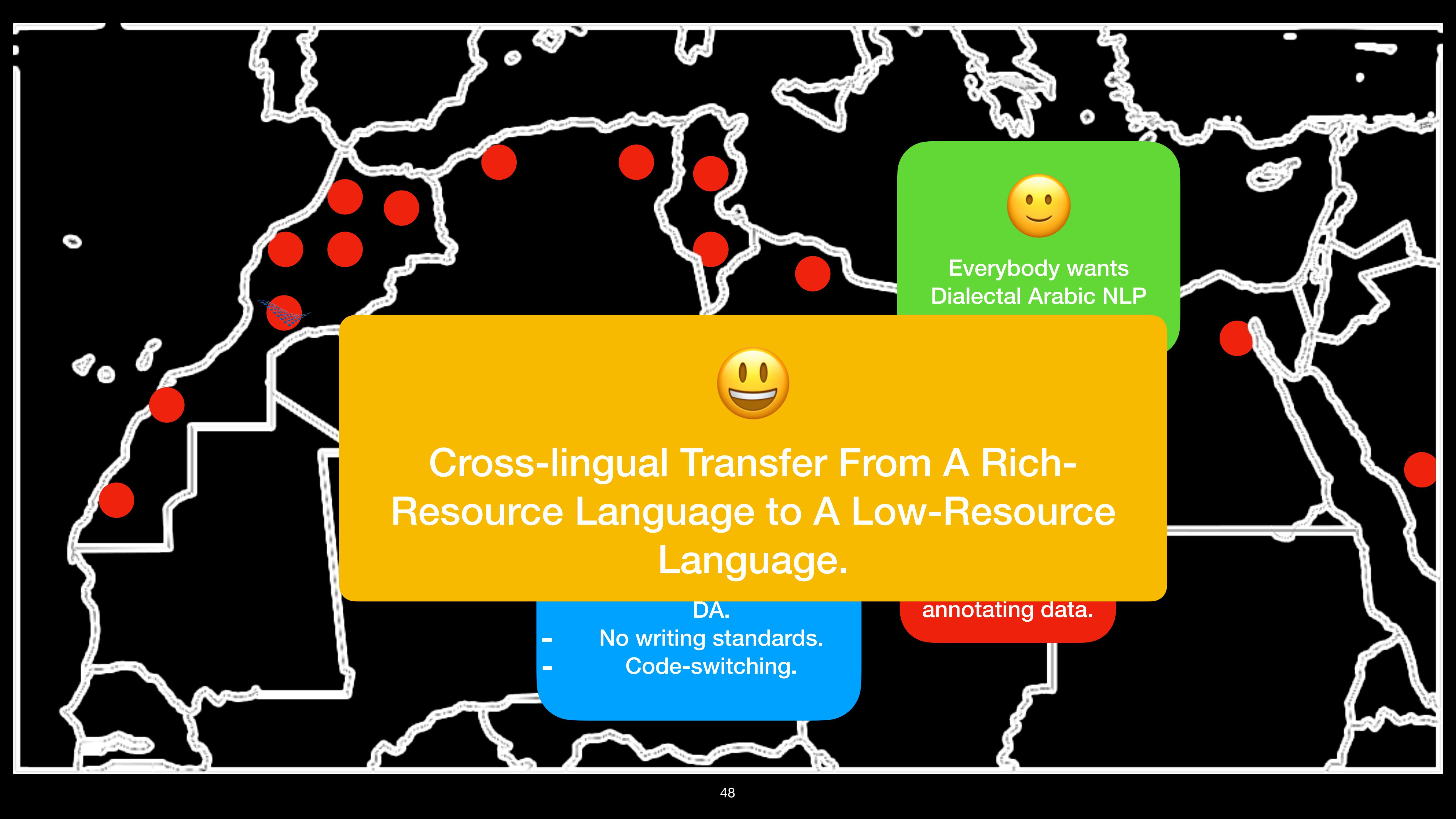
Hard to reach generalization between countries



- More than 100 variants of DA.
- No writing standards.
- Code-switching.



Nobody can put
the work in
annotating data.



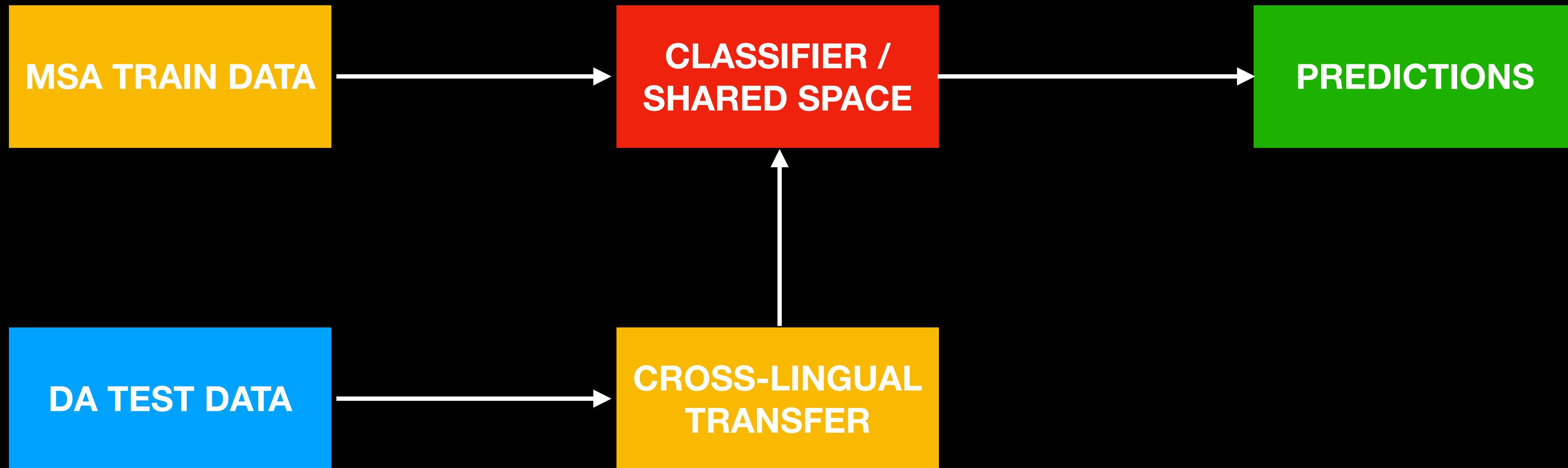
Cross-lingual Transfer From A Rich-Resource Language to A Low-Resource Language.

DA.

- No writing standards.
- Code-switching.

annotating data.

Cross-lingual from MSA to DA



Motivations

- Modern Standard Arabic and Dialectal Arabic share similarities on multiple linguistic levels (orthography, morphology, phonology, etc).
- Availability of large Modern Standard Arabic pre-trained language models and annotated datasets.

Our research

- Word-to-word translation between Arabic dialects
- Unsupervised Sentence Classification for cross-domain & cross-dialects Arabic
- Unsupervised Sequence Labeling framework for Arabic dialects

Word-to-word translation between Arabic dialects



Word Embeddings

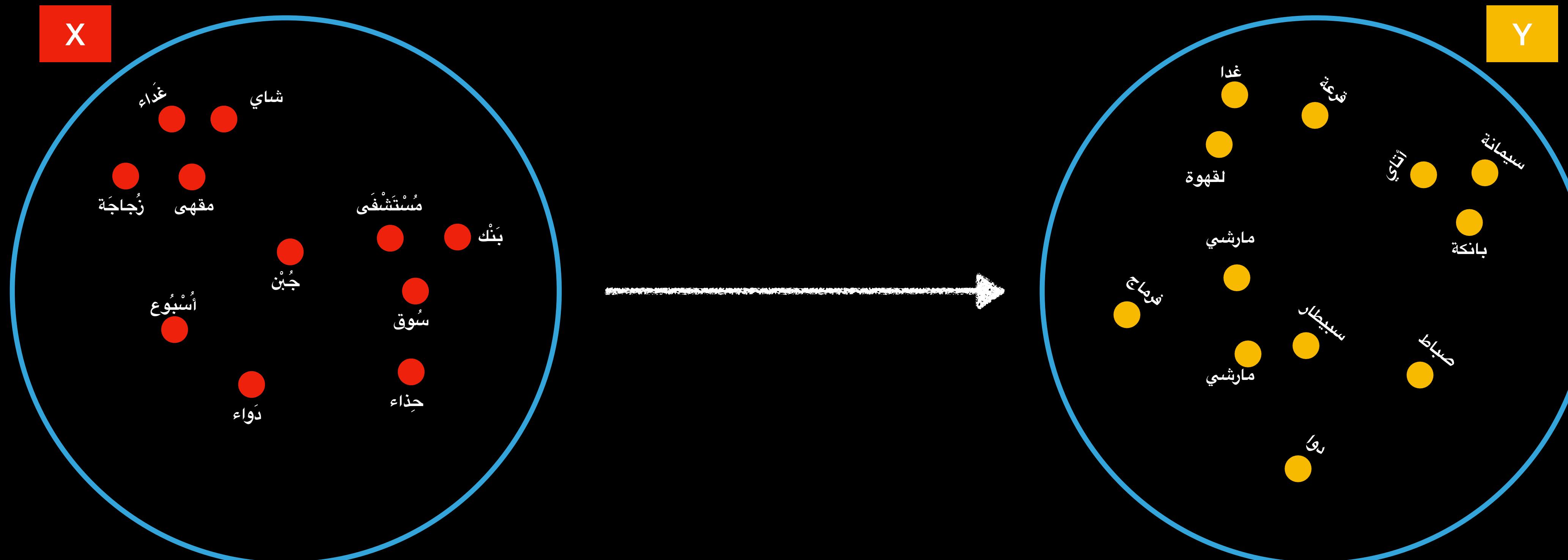
Word2vec (Mikolov, 2013)
FastText (Bojanowski, 2016)
BERT (Devlin, 2018)
XLM-RoBERTa (Conneau, 2019)

...

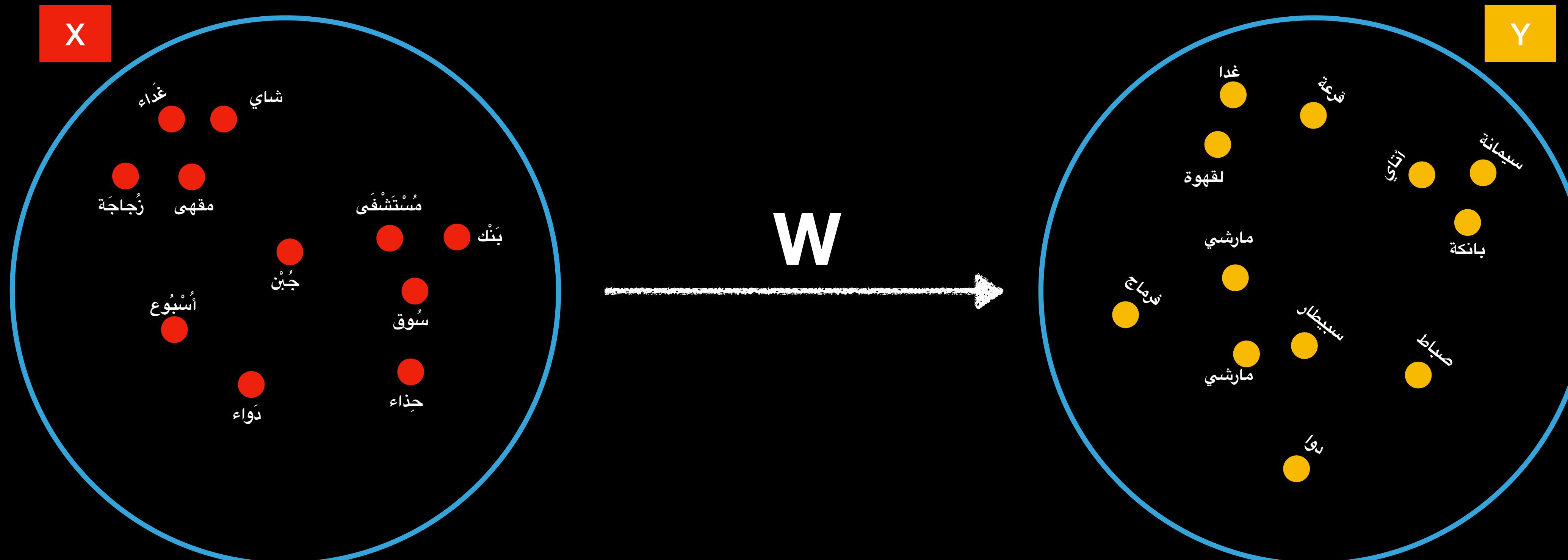
Multi-dialectal Arabic Word embedding Approach



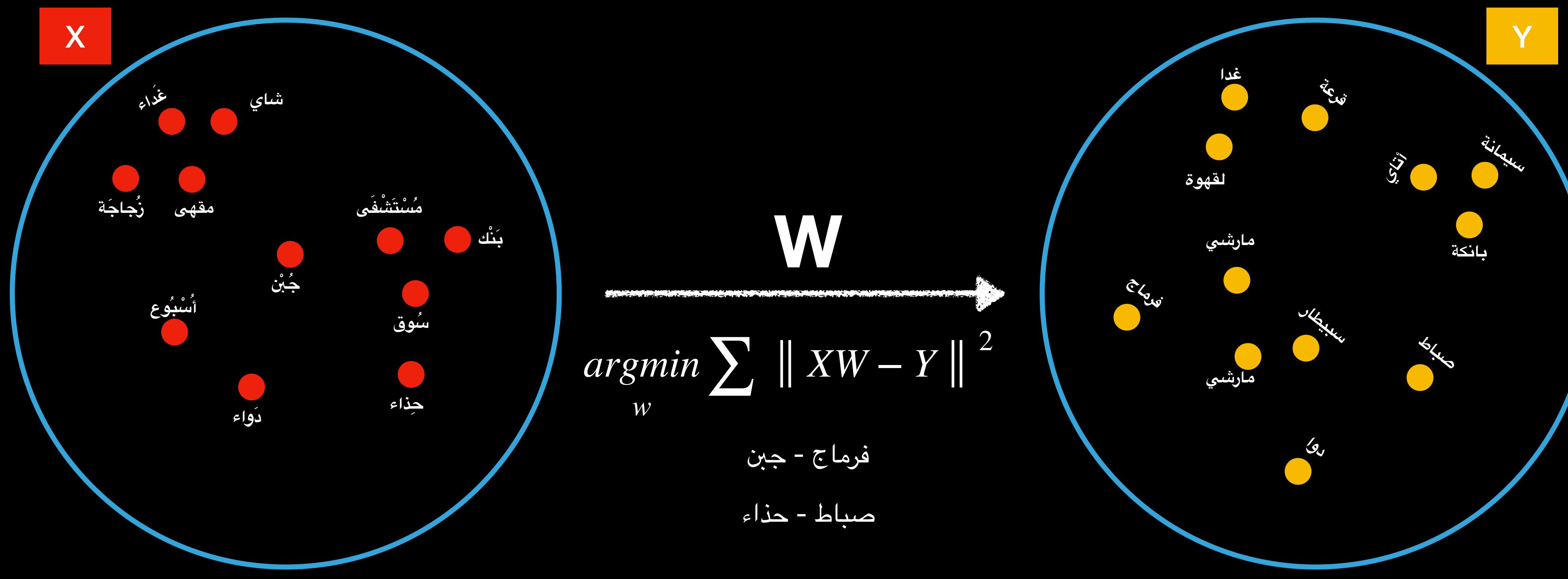
Multi-dialectal Arabic Word embedding Approach



Multi-dialectal Arabic Word embedding Approach



Multi-dialectal Arabic Word embedding Approach



Multi-dialectal Arabic Word embedding

Approach

- Orthographic extension of the word embeddings.
- The encoding of the orthographic variations between Arabic dialects.

Extension of embeddings using orthographic features

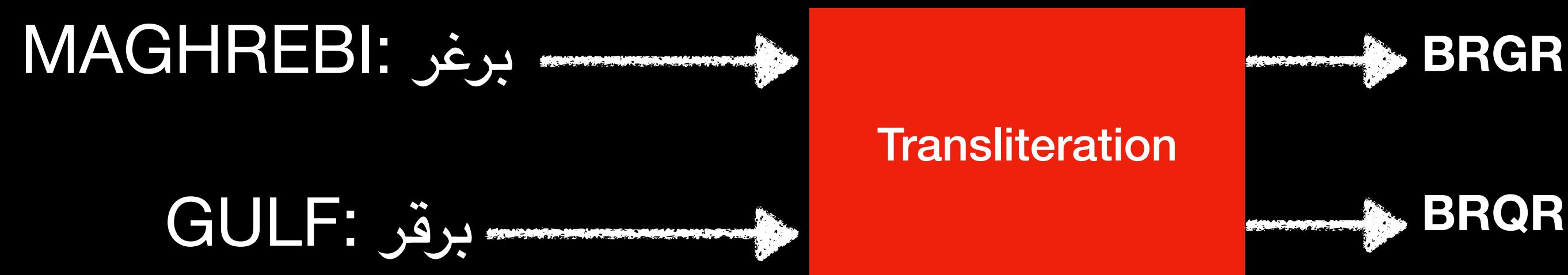
Approach

MAGHREBI: بِرْغَرْ

GULF: بِرْقَرْ

Extension of embeddings using orthographic features

Approach



Vocabulary = {BRGR, BRQR}
Alphabet = {B, R, G, Q}
Embeddings dimensions: 100

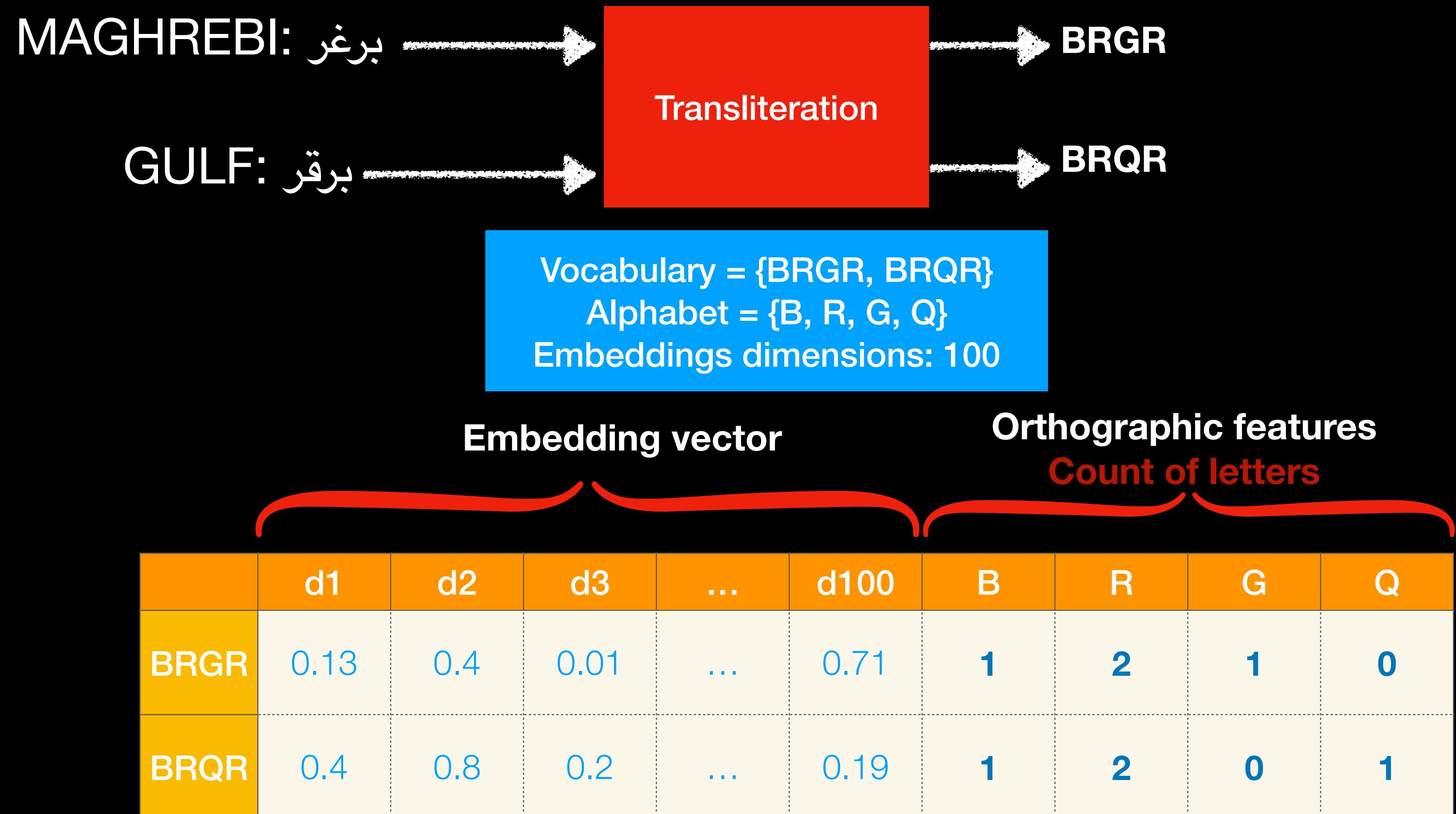
Embedding vector

A red bracket is drawn above the first 100 dimensions of the embedding vector, spanning from d1 to d100.

	d1	d2	d3	...	d100
BRGR	0.13	0.4	0.01	...	0.71
BRQR	0.4	0.8	0.2	...	0.19

Extension of embeddings using orthographic features

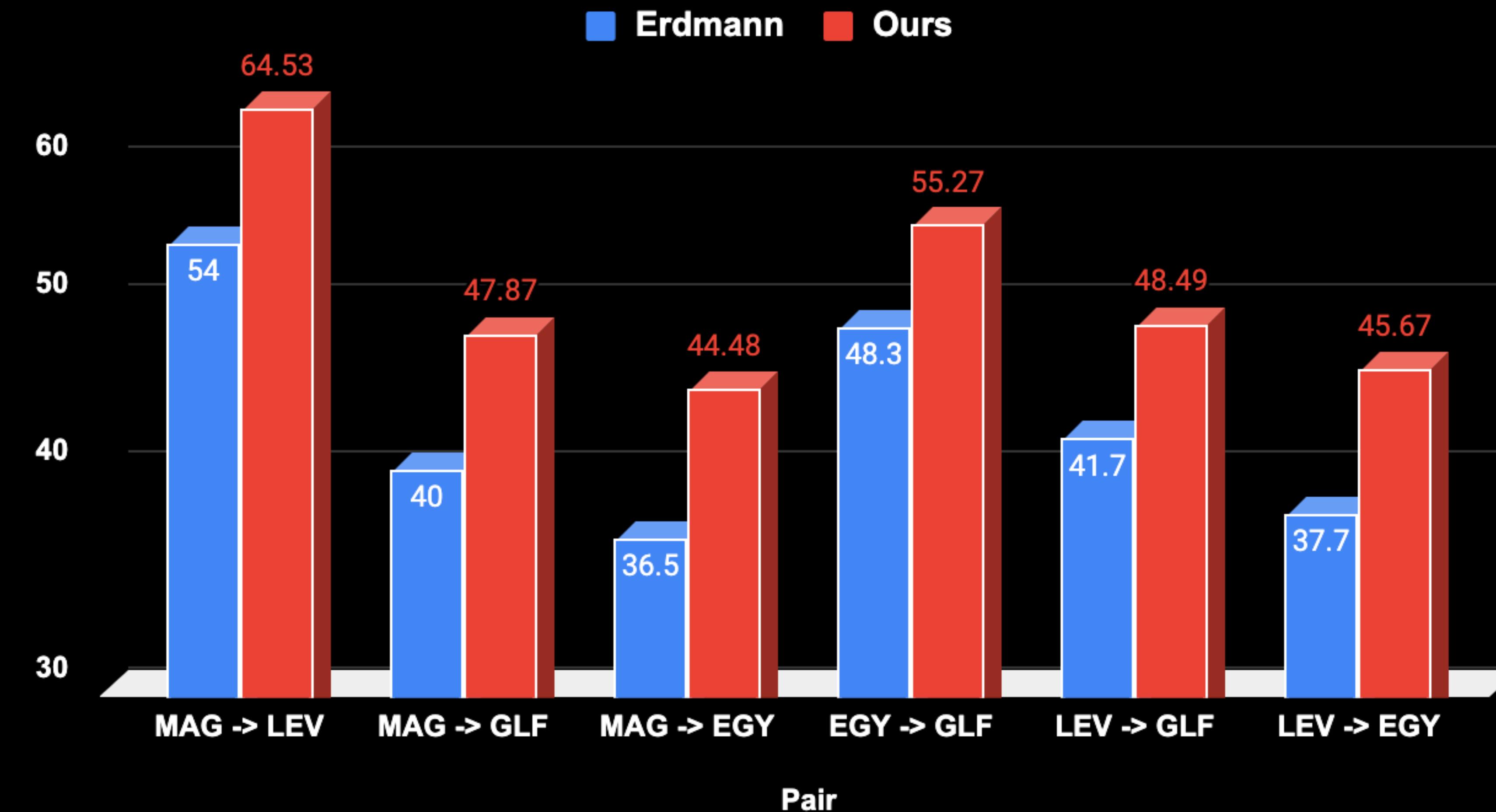
Approach



Multi-dialectal Arabic Word embedding

Results

Accuracy of bilingual translations



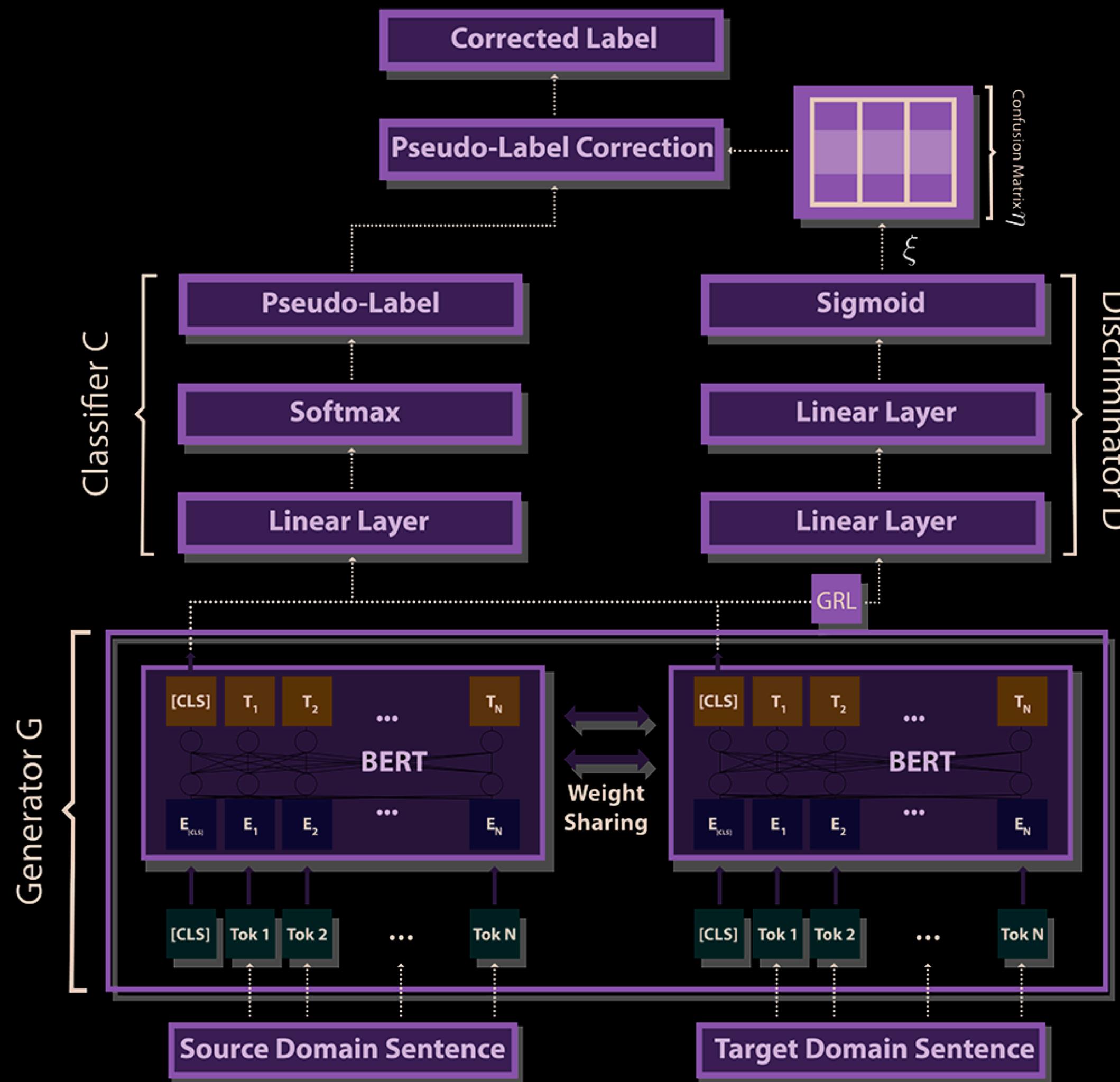
Multi-dialectal Arabic Word embedding

Limitations

- Performance is good in the word-to-word translation task. 
- Performance is low on downstream tasks such as sentence classification. 
- Hard and complex to extend to the sequence labeling tasks. 

Unsupervised Sentence classification for cross-domain cross-dialect Arabic dialects

Unsupervised sentence classification for DA Approach

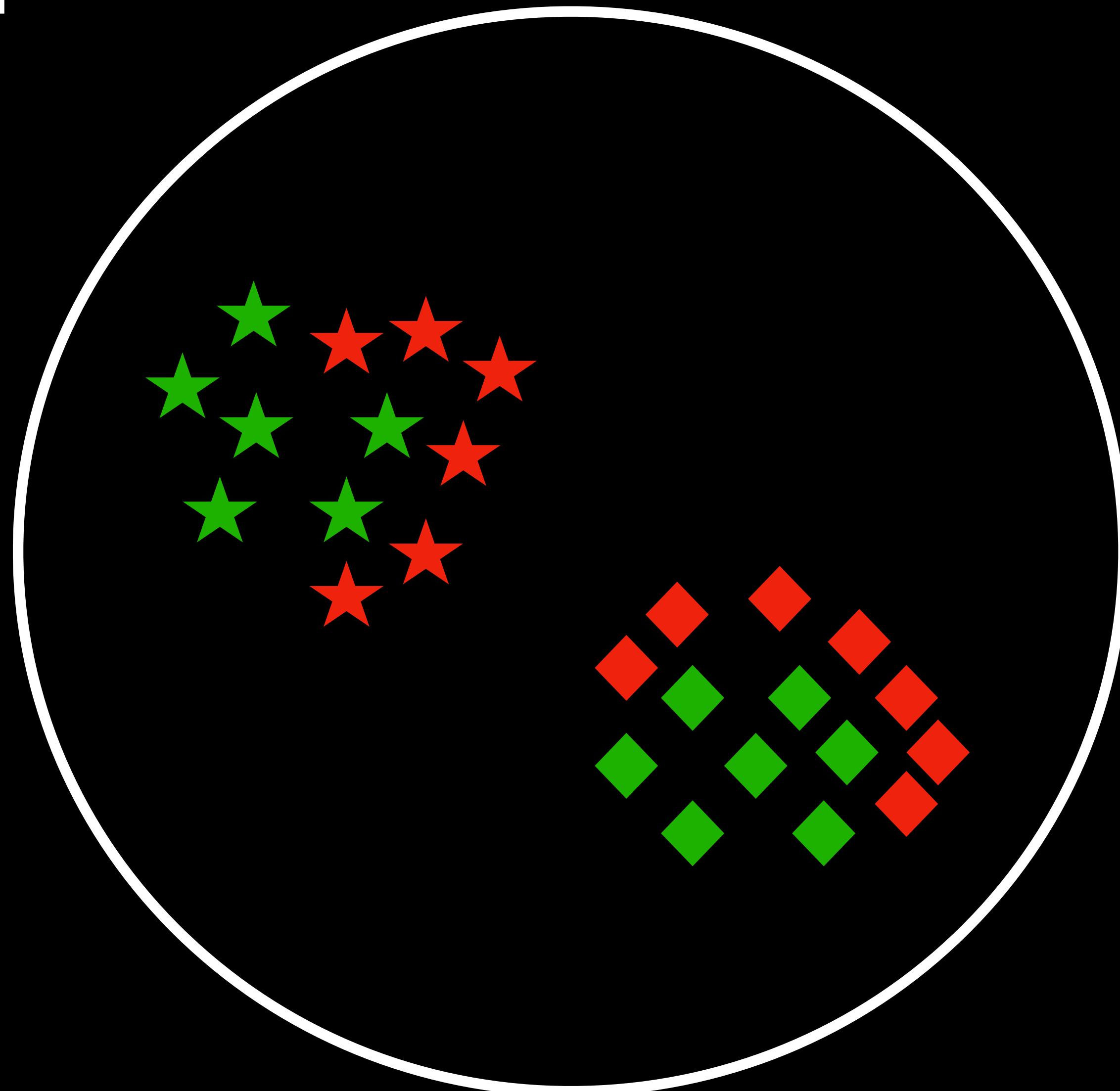


1. Self-training
2. Domain adaptation
 1. Alignment between the the sentences' embeddings of the source and target languages.
 2. Class-wise alignment.

El Mekki, A., El Mahdaouy, A., Berrada, I., and Khoumsi, A. 2021. Domain Adaptation for Arabic Cross-Domain and Cross-Dialect Sentiment Analysis from Contextualized Word Embedding. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (pp. 2824–2837).

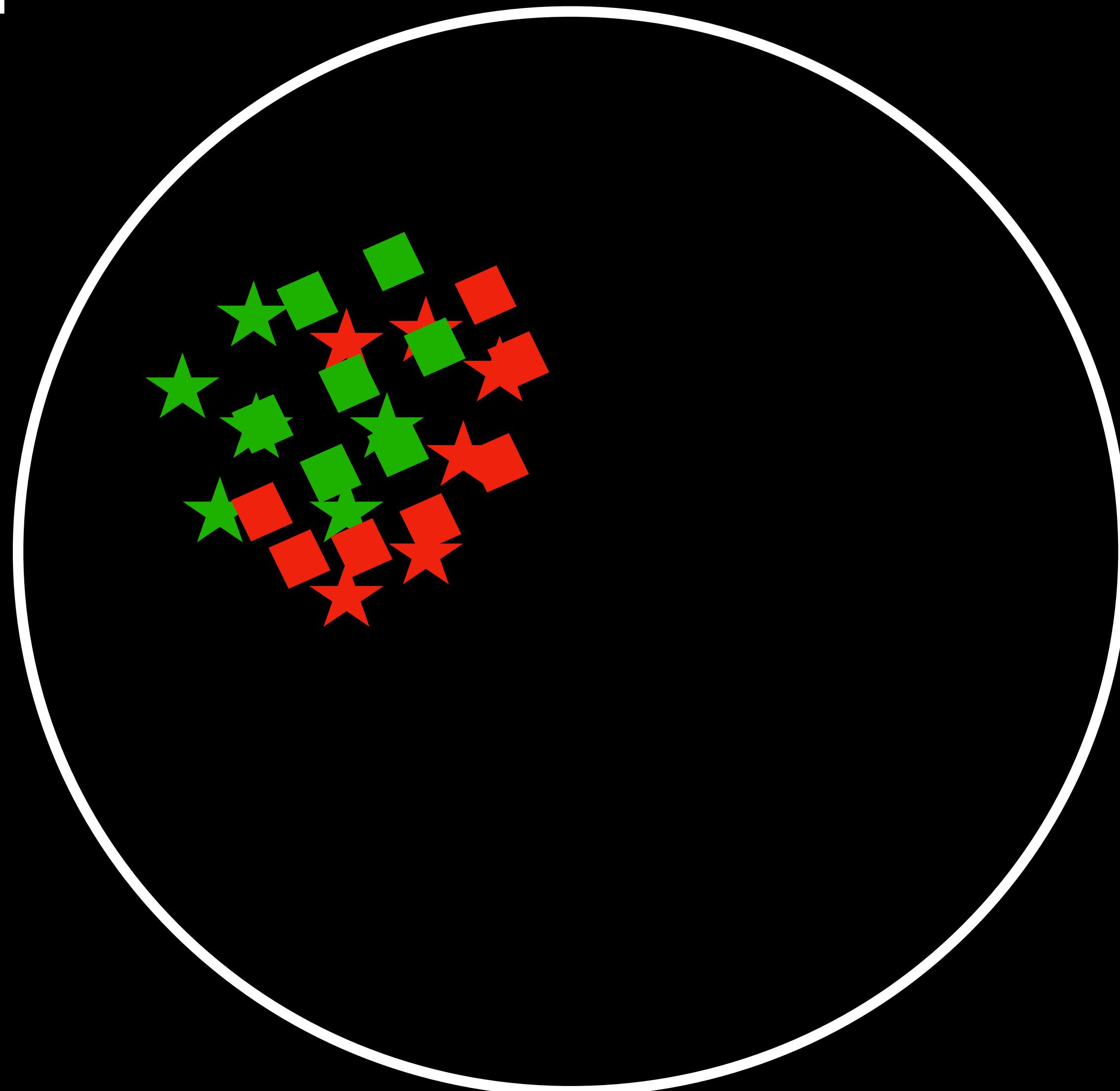
Domain Adaptation

- Positive** (Green)
- Negative** (Red)
- ★ Standard Arabic
- ◆ Dialectal Arabic

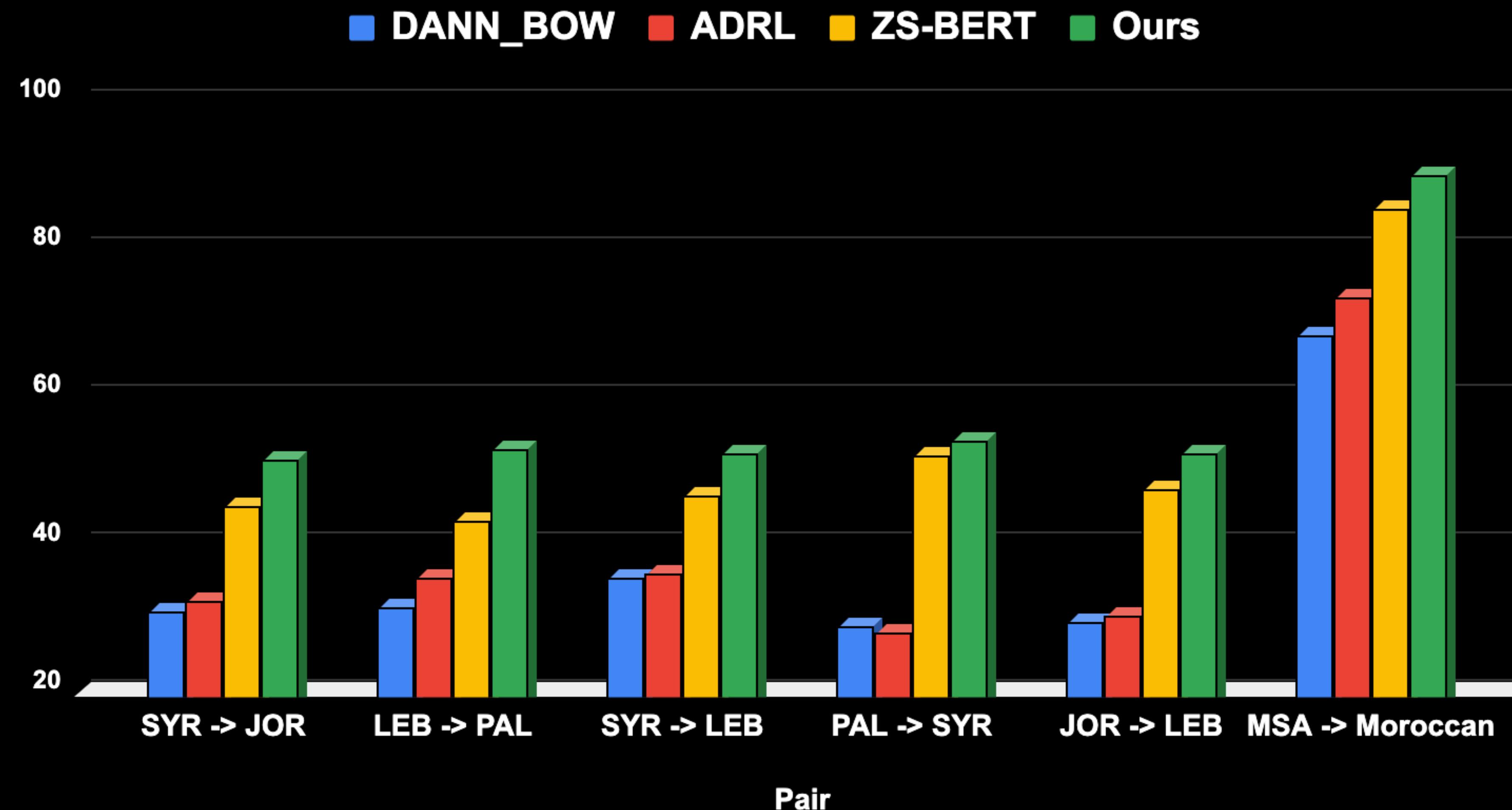


Domain Adaptation

- Positive**
- Negative**
- ★ Standard Arabic
- ◆ Dialectal Arabic



Unsupervised sentence classification for Dialectal Arabic Results



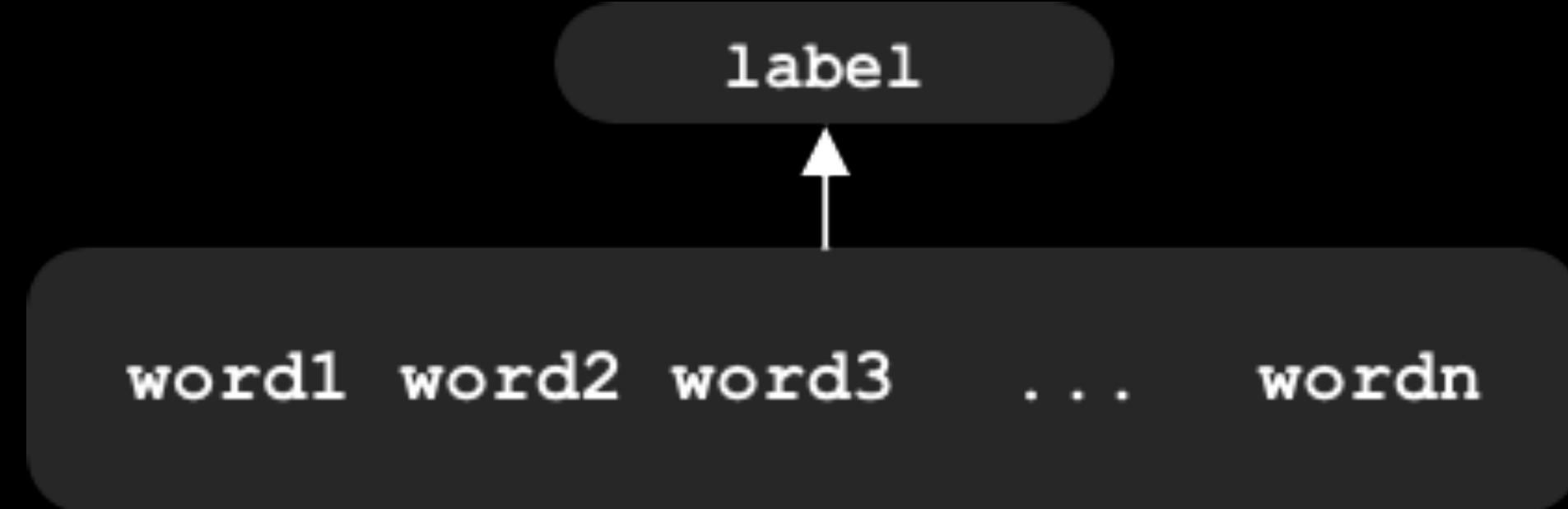
Unsupervised sentence classification for DA

Limitations

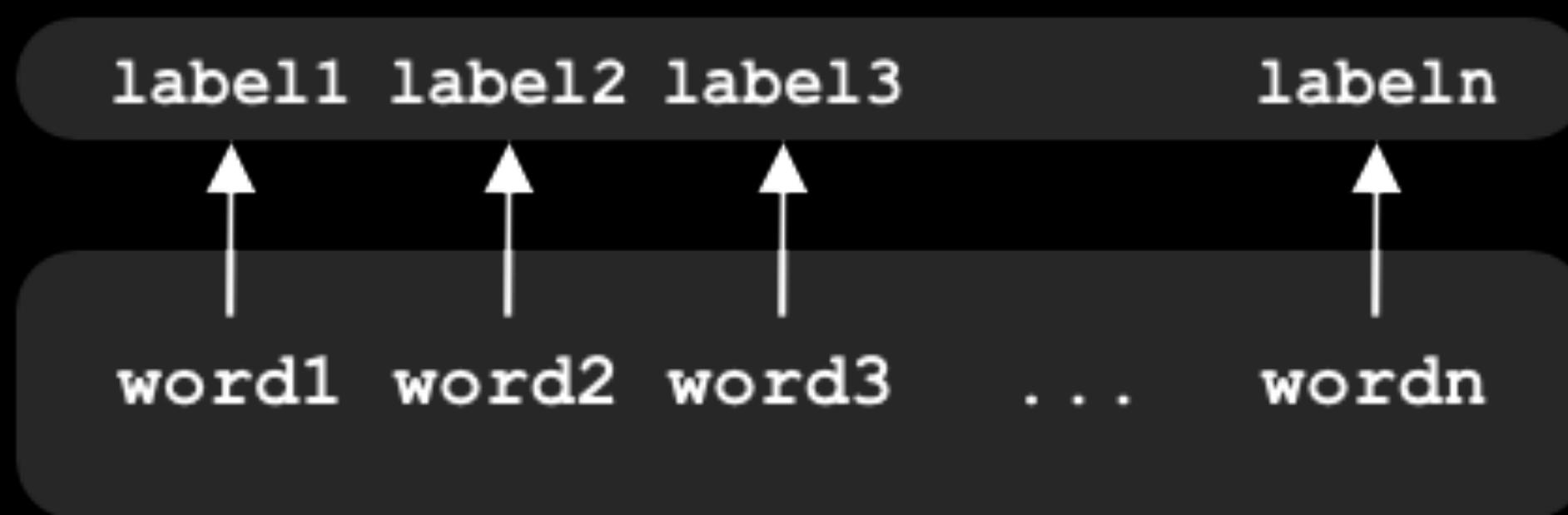
- Good performance on sentence classification tasks. 
- Hard and complex to extend to the sequence labeling tasks. 

AdaSL: An Unsupervised Domain Adaptation Framework For Arabic Multi- dialectal Sequence Labeling

Sentence classification

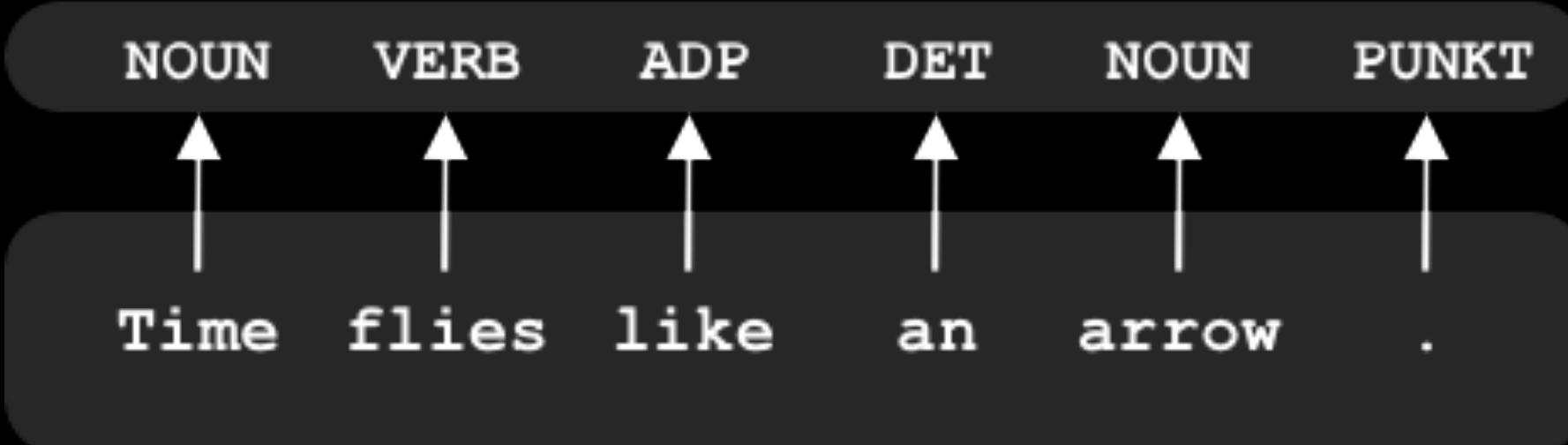


Sequence labeling

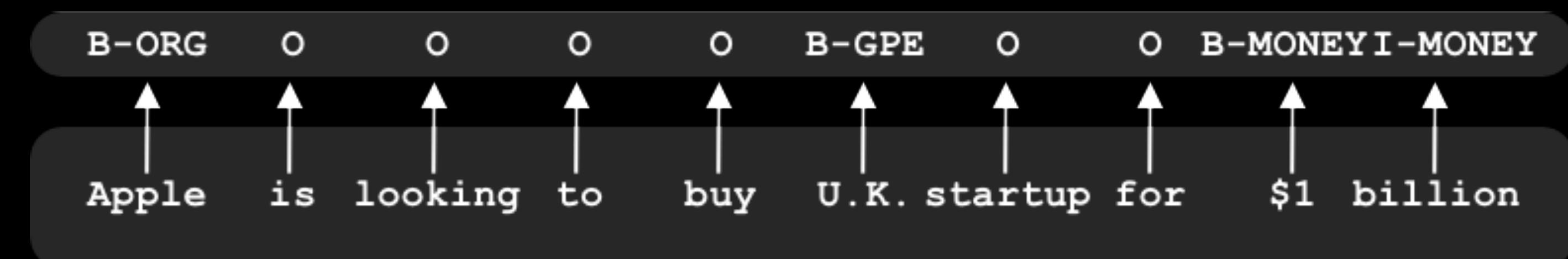


Examples of sequence labeling tasks:

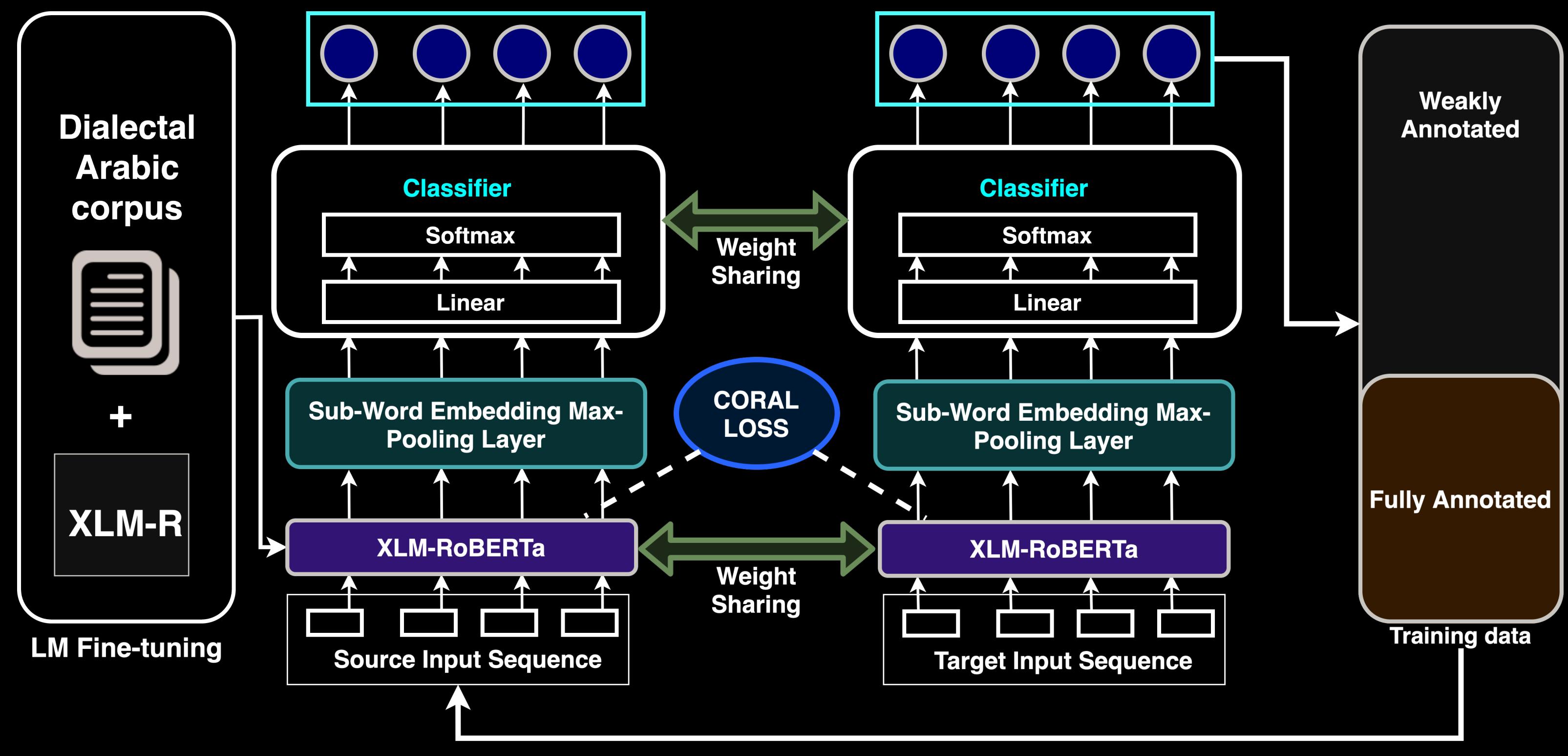
Part-of-Speech (POS) tagging



Named-entity recognition (NER)



Unsupervised Sequence Labeling for Dialectal Arabic Framework

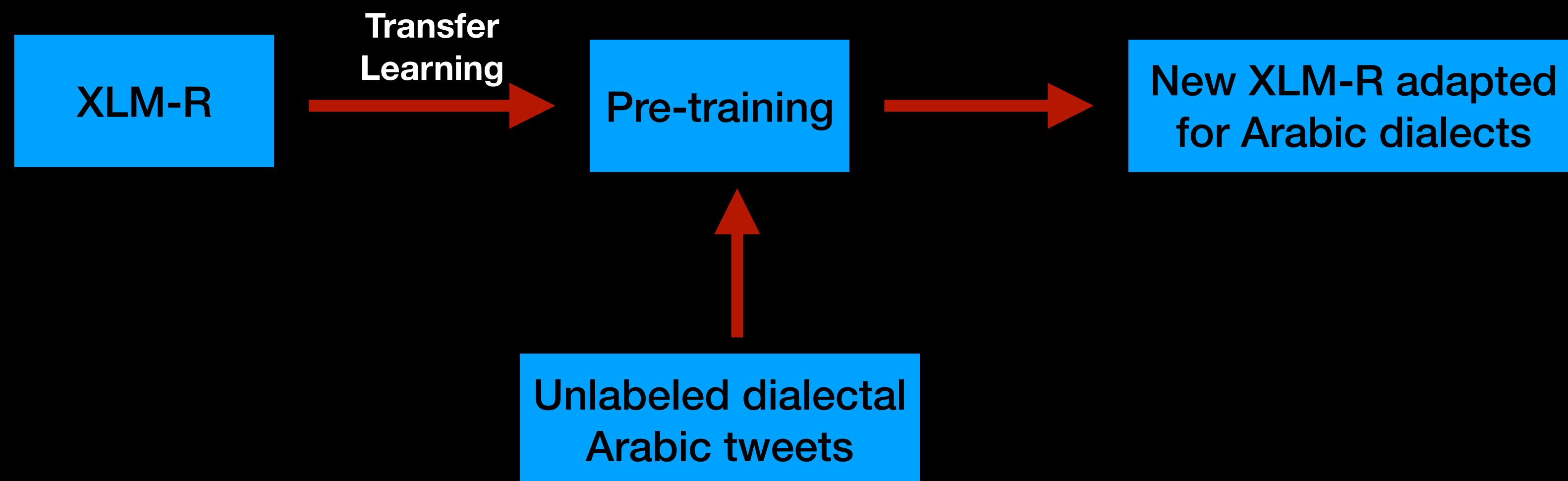


- 1) Language model domain adaptive fine-tuning
- 2) Sub-word embedding pooling
- 3) Unsupervised MSA-DA distributions alignment
- 4) Iterative self-training

Unsupervised Sequence Labeling for Dialectal Arabic

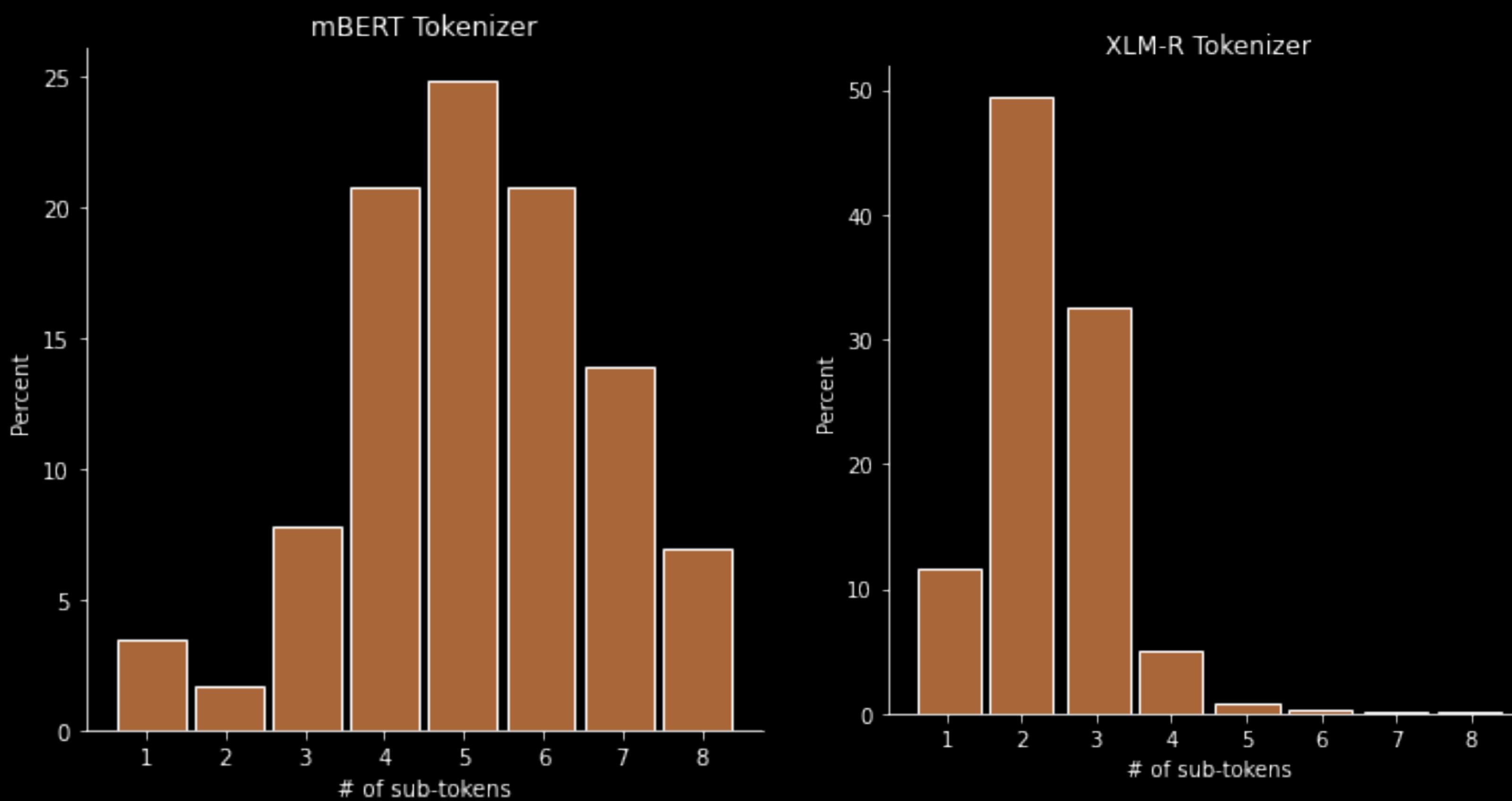
1) Domain adaptive fine-tuning (Masked-language modeling on the target dialect)

- 1) Collected 5 Million Dialectal Arabic Tweets.
- 2) Fine-tuning the multilingual language model on the collected tweets using masked-language modeling objective.



Unsupervised Sequence Labeling for Dialectal Arabic

2) Sub-word embedding pooling



- 1) Multilingual Language models tokenize Arabic words into multiple subwords.
- 2) Arabic is morphologically complex.
- 3) State-of-the-art sequence labeling methods use the 1st sub-token to label the word.

Tokenizers

- MSA tokenizer
- Dialectal Arabic tokenizer

```
tokenizer.tokenize("الجو جميل اليوم")  
['الجو', 'جميل', 'اليوم']
```

```
tokenizer.tokenize("الجو زوين ليوما")  
['الجو', 'زوين', 'ليوما']
```

- GPT-3 Tokenizer on English
- GPT-3 Tokenizer on Arabic

GPT-3 Codex

```
Mohammed VI University
```

Clear Show example

Tokens	Characters
4	22

Mohammed VI University

GPT-3 Codex

```
جامعة محمد السادس
```

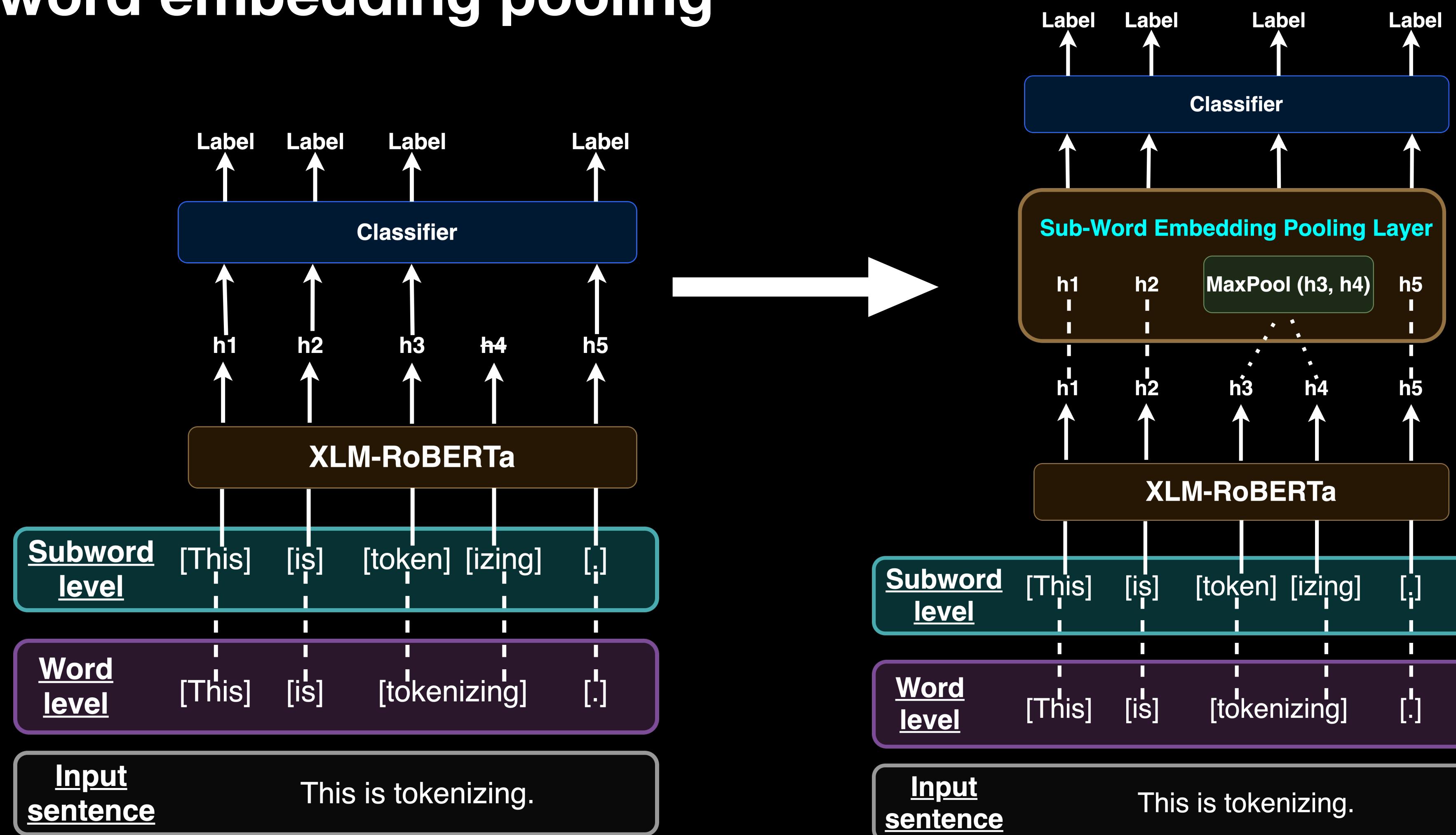
Clear Show example

Tokens	Characters
16	17

سُلْطَانِيَّةُ مَرْكَازِيَّةُ دُمَّالِيَّةُ اَنْسَالِيَّةُ

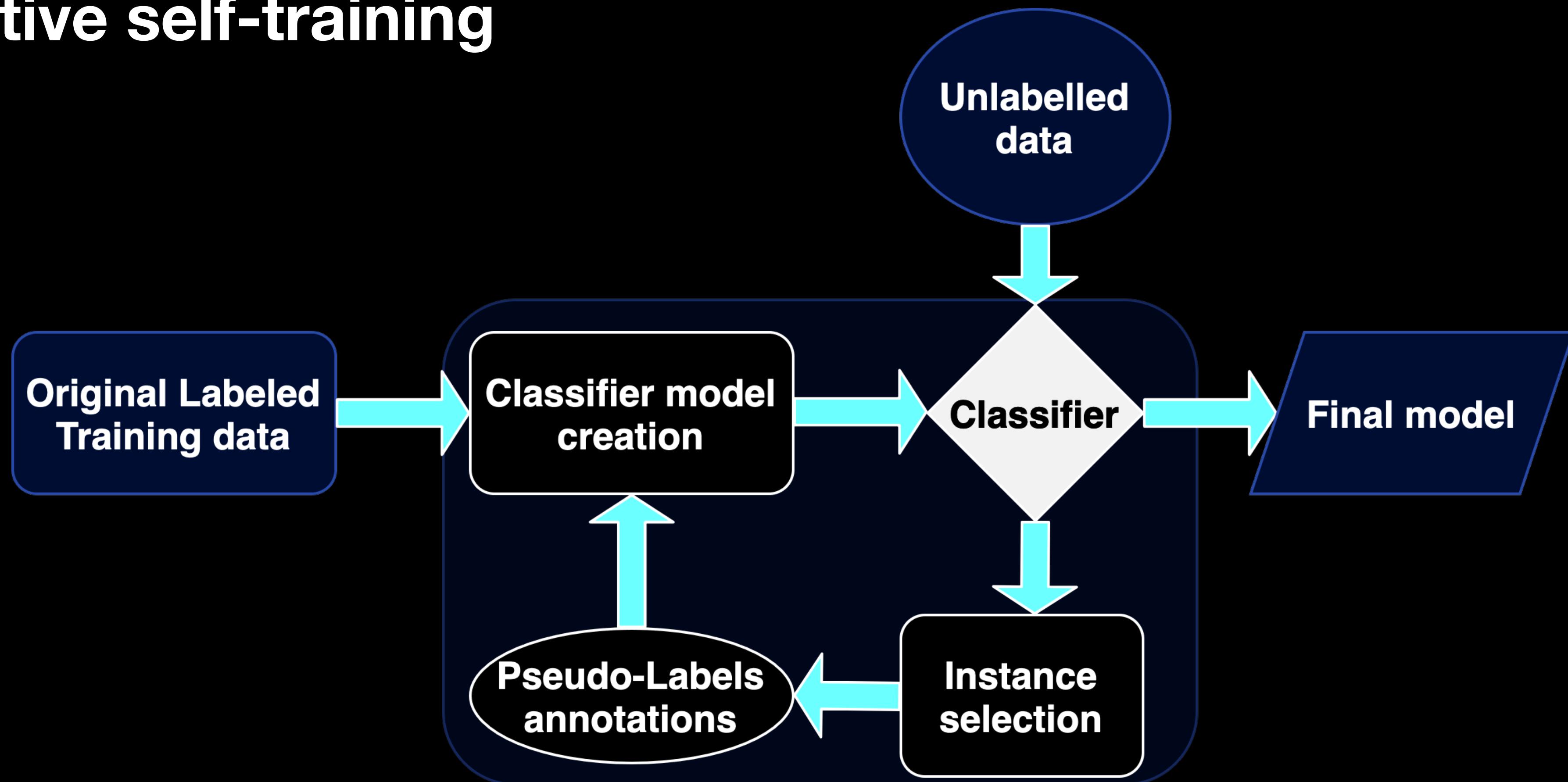
Unsupervised Sequence Labeling for Dialectal Arabic

2) Sub-word embedding pooling



Unsupervised Sequence Labeling for Dialectal Arabic

4) Iterative self-training



Unsupervised Sequence Labeling for Dialectal Arabic

Evaluation scenarios

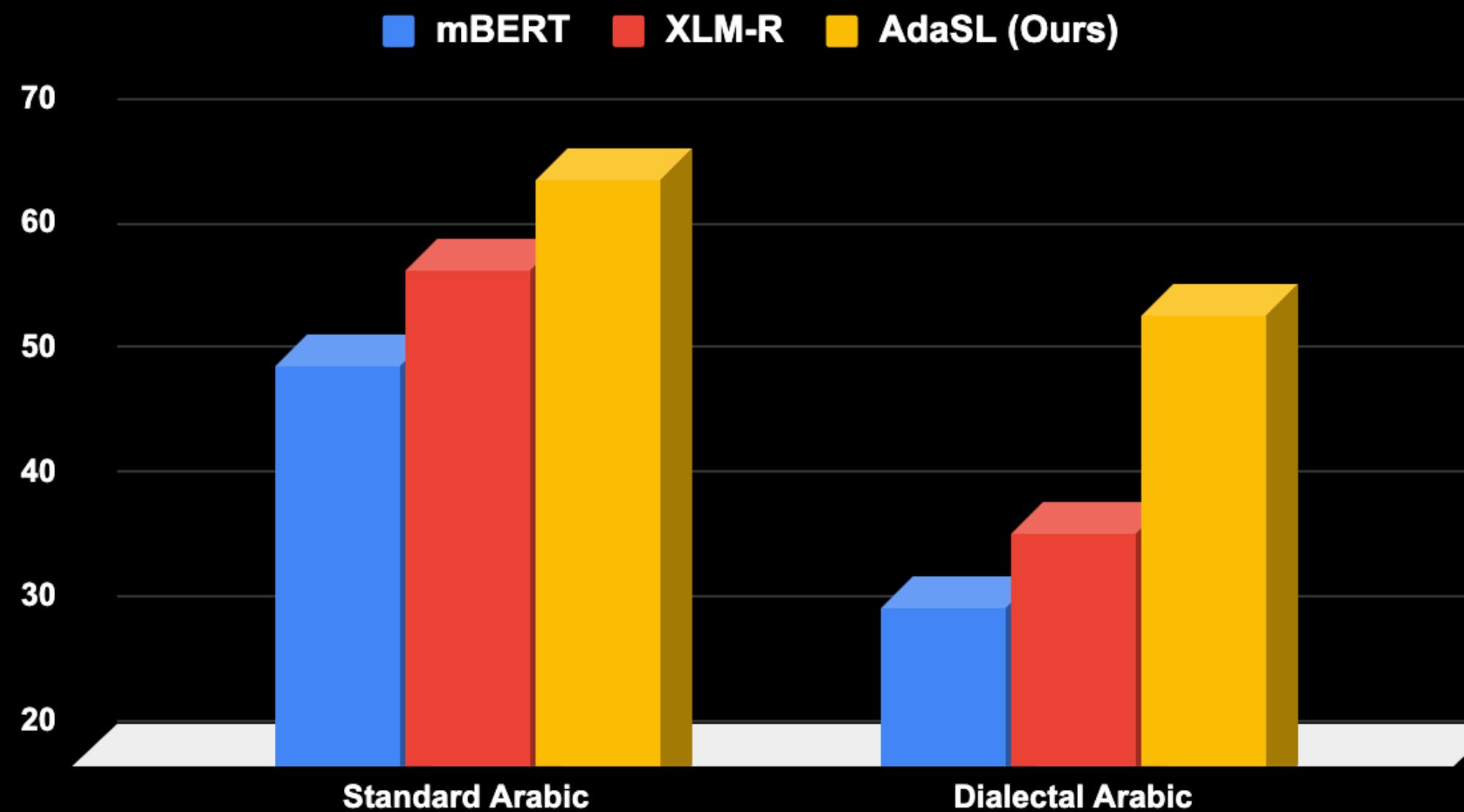


- **Named Entity recognition task:**
 - Training dataset: MSA
 - Test dataset: MSA (Twitter), DA (Twitter)
 - Evaluation metric: Macro-F1 score
- **Part-of-speech tagging task:**
 - Training dataset: MSA
 - Test dataset: Maghreb (MAG), Egyptian (EGY), Leventine (LEV) and Gulf (GLF)
 - Evaluation metric: Accuracy

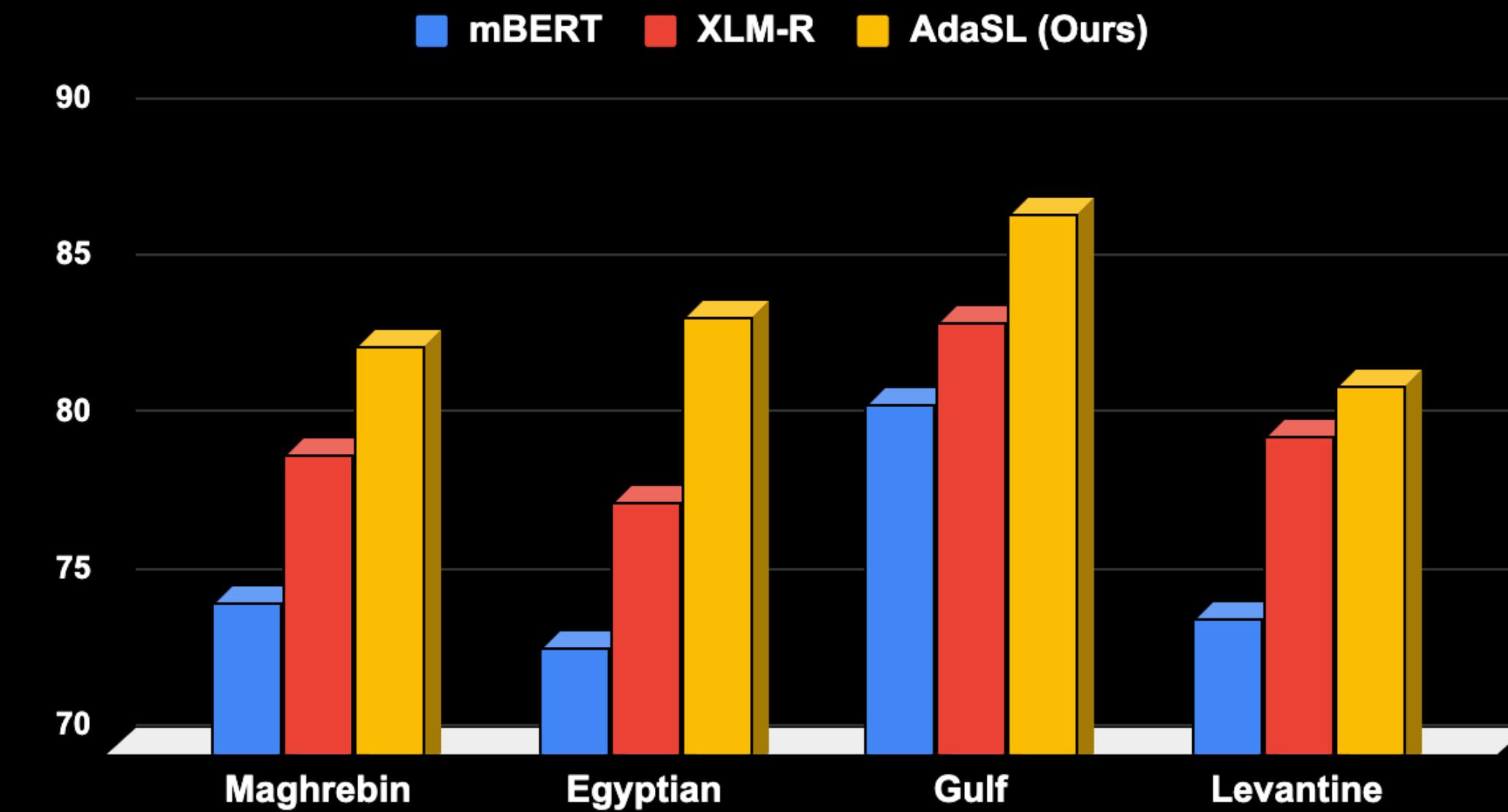
جامعة محمد السادس متعددة التخصصات التقنية مؤسسة دولية للتعليم العالي متخصصة في العلوم، والتكنولوجيا، والعلوم الاجتماعية، والأعمال. المقر الرئيسي ديالها كайн فمدينة ابن جرير لي جات حدا مدينة
برهان الدين والملك محمد السادس ودشنها الملك محمد السادس نهار 12 يناير 2017. مؤسسة المكتب الشريف للفوسفاط OCP، ORG وفالعيون. صاوباتها عندها فاراباط GPE مراكش GPE بالمغرب،

Unsupervised Sequence Labeling for Dialectal Arabic Results

Named Entity Recognition Results



Part-Of-Speech Tagging Results



جامعة محمد السادس التقنية مؤسسة دولية للتعليم العالي متخصصة في العلوم، والتكنولوجيا، والعلوم الاجتماعية، والأعمال. المقر الرئيسي ديالها كاين فمدينة ابن جرير GPE لي جات حدا مدينة

مراكش GPE بالمغرب، فالرباط GPE وفالعيون. صاوباتها مؤسسة المكتب الشريف للفوسفاط ORG نهار 12 يناير 2017. DATE OCP، ORG ودشنها الملك محمد السادس PERS

Conclusions

- A good performance on Arabic dialects NLP systems can be achieved with zero labeled data.
- Semi-supervised approaches can lead to better results (few-shot learning).
- More effort should be performed to make labeled data for underrepresented languages available (crowdsourcing, more use of local languages on the web rather than using 2nd languages).



Resources

- Stanford NLP course:
 - [https://www.youtube.com/playlist?
list=PLoROMvodv4rOSH4v6133s9LFPRHjEmbJ](https://www.youtube.com/playlist?list=PLoROMvodv4rOSH4v6133s9LFPRHjEmbJ)
- Jay Alammar Blog:
 - <https://jalammar.github.io/>
- Coding:
 - Python
 - Pytorch
 - Huggingface

THANK YOU FOR YOUR ATTENTION!

Any questions?

abdelah.elmekki@um6p.ma