# Automatic Recognition and Translation of Tunisian Dialect Named Entities into Modern Standard Arabic

Roua Torjmen and Kais Haddar

rouatorjmen@gmail.com          kais.haddar@yahoo.fr

Miracl Laboratory, University of Sfax, Tunisia

Outline

- Introduction
- Previous work
- Linguistic study
- Proposed method
- Experimentation and evaluation
- Conclusion and perspectives

# Introduction

Tunisian Dialect (TD) corpus contains enormously NEs

ANER helps in several fields:

Information retrieval

Information Extraction

Automatic indexing

Classification of documents

# Introduction

TD does not have a standard spelling

TD words not only of Arab origin but of different origins

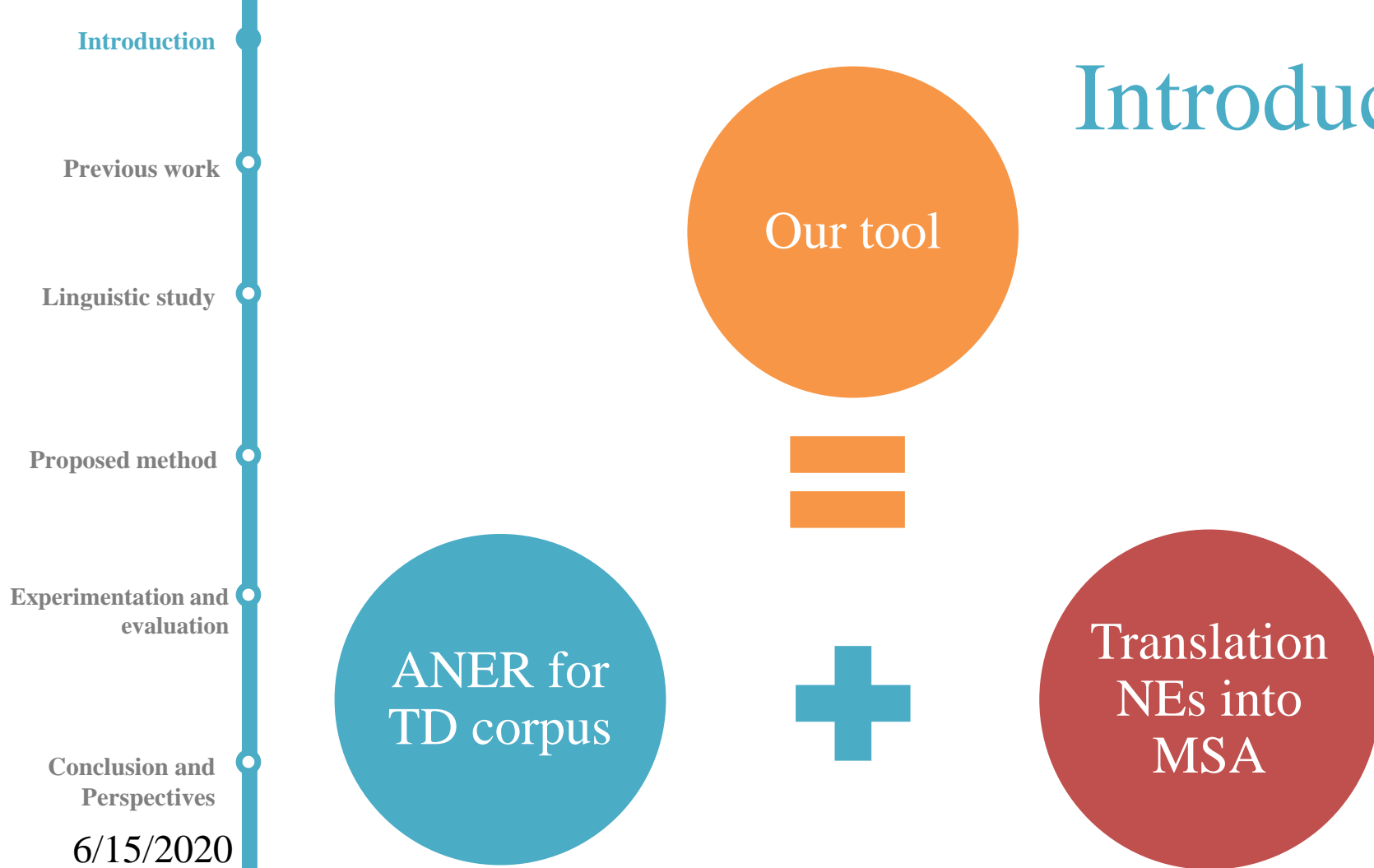No capital letters at the beginning of a word for TD proper noun

TD proper nouns can be written in different manners

# Introduction

Our tool

=

ANER for TD corpus

+

Translation NEs into MSA

# Introduction

Provide a linguistic study for TD

Reuse a bilingual dictionary TD-MSA

Establish a set of syntactic grammars that offers NEs recognition and their translation

Implement them in NooJ linguistic platform

# Previous work

|  | Approach | Tool |
|---|---|---|
| El Bazi et al., 2019 | Statistic | Bidirectional Long Short-Term Memory (LSTM) |
| Hkiri et al., 2017 | Hybrid | ANNIE/GATE and CRF |
| Zaghouani, 2012 | Linguistic | Language-independent rules |
| Fehri et al., 2011 Hasni et al., 2009 Khalfallah et al., 2009 Mesfar, 2007 | Linguistic | Finite state transducers |

6/15/2020

# Previous work

No work deals with ANER for TD corpus

Finite transducers have given encouraging results in:

- MSA ANER (Mesfar)

- MSA-French (Fehri)

# Linguistic study

ENAMEX: People,
organizations,
locations

TIMEX: Date,
time. . .

NUMEX: Money,
percentages. . .

# Linguistic study

**Percentage**

سبعين بالمية

*sab'In bilmiyah*
(seventy percent)

سبعين في المية

*sab'In filmiyah*
(seventy percent)

70 %

**Rule:**
**(number|digit⁺)(preposition definite_noun|%)**

**Rule:**
**(number|digit⁺) weight_unit**

**weight**

70 كيلو

*70kIlO*
( 70 kilo)

سبعين كيلو

*sab'iin kIlO*
(seventy kilo)

6/15/2020

11

# Linguistic study

## Measure

سبعين كتاب

*sab'iin ktAb*
(seventy books)

6 كتب

*6 ktob*
(6 books)

**Rule:**

**(number|digit⁺)(devise|symbol)**

## Money

70 دولار

*70 dOlAr*
(70 dollars)

$ 70

سبعين دولار

*sab'iin dOlAr*
(seventy dollars)

**Rule:**

**(number|digit⁺) indefinite_noun**

Translation keeps the same syntax structure for NUMEX NEs

➔ Word for word translation

6/15/2020

12

# Linguistic study

**TIMEX NE**

# Linguistic study

**Period**

أربعة سوايع

*'arb'ah swAyi'*
(four hours)

مدة أربعة سوايع

*moddet 'arb'ah swAyi'*
(four hours duration)

عندي أربعة سنين

*'andI 'arb'ah snin*
(four years duration)

عندي نهارين

*'andI nhArIn*
(two days duration)

**Rule:**
**(prep suffix|indefinite_noun)?**
**(number|digit⁺) * time_unit**

# Linguistic study

Hour

نهار التلاث مع الثمنية

*nhar iltlAth m'A ilthmanya*
(Tuesday at 8 am)

غدوة الصباح/غدا صباحا

*ghodwah ilsbAh*
(tomorrow morning)

التسعة متاع الصباح/
تاسعة صباحا

*iltes'ah mtA' ilsbAh*
(at 8 am)

**Rule:
((indefinite_noun day_name)? (prep)? (number|digit⁺) (prep part_day)?)|(noun part_day)**

# Linguistic study

Age

عمري 5 سنين

*'omrI 5 snIn* (my age is 5 years)

قفلت ستة سنين

*qfalt sittah snIn* (my age is 6 years)

**Rule:**
(**قفل_conjugated|suffix عمر**)
(**number|digit$^+$**) **time_unit**

Translation keeps the same syntax structure for TIMEX NEs except for hours

➔ Word for word translation & adjustment translation

# Linguistic study

**ENAMEX NE**

# Linguistic study

# Linguistic study

**Name**

| عصام الشوالي |
| --- |
| Issam Alchawali |

عصام

Issam

أنا عصام

*ana issAm*
(i am Issam)

اسمي عصام

*ismI issAm*
(my name is Issam)

سي عصام الشوالي

*sI issAm alchawAlI*
(Mr Issam Alchawali)

**Rule:**
**(Pronoun| اسم**
**suffix|civility_noun)?**
**first_name last_name?**

6/15/2020

19

Introduction

Previous work

Linguistic study

Proposed method

Experimentation and evaluation

Conclusion and Perspectives

# Linguistic study

**Profes-sion**

الفنان عاصي الحلاني

*ilfannAn 'AssI alhallAnI*
(the singer Assi AlHallani)

وزير الثقافة محمد النزالي

*wazIr ilthakAfah muhammad ilnzAli*
(Culture  Minister Mohamed AlNzali)

نخدم رئيس تحرير

*Nikhdim ra'Is tahrIr*
(I work Chief Editor)

**Rule: (Profession_noun definite_noun? first_name last_name?) | (خدم_conjugacated Profession_noun definite_noun?)**

Introduction

Previous work

Linguistic study

Proposed method

Experimentation and
evaluation

Conclusion and
Perspectives

# Linguistic study

Nationality

أحنا التوانسة

Ahna twAnsah
(we Tunisians)

هند صبري  التونسية

hind sabrI iltUnsiyah
(Tunisian Hind Sabri)

الفرنساوي جاك

IfransAwI jAk
(French Jack)

**Rule:**

**((first_name last_name?| pronoun) nationality_noun) | (nationality_noun first_name last_name?)**

Translation keeps the same syntax structure for Person ENAMEX NEs

➔ Word for word translation

6/15/2020

21

# Linguistic study

**Localisation ENAMEX NE**

# Linguistic study

**Country**

تونس

*tUnis*
(Tunisia)

بلادي تونس

*blAdI tUnis*
(My country Tunisia)

**Rule:**
(suffix)? Country_name (بلاد)
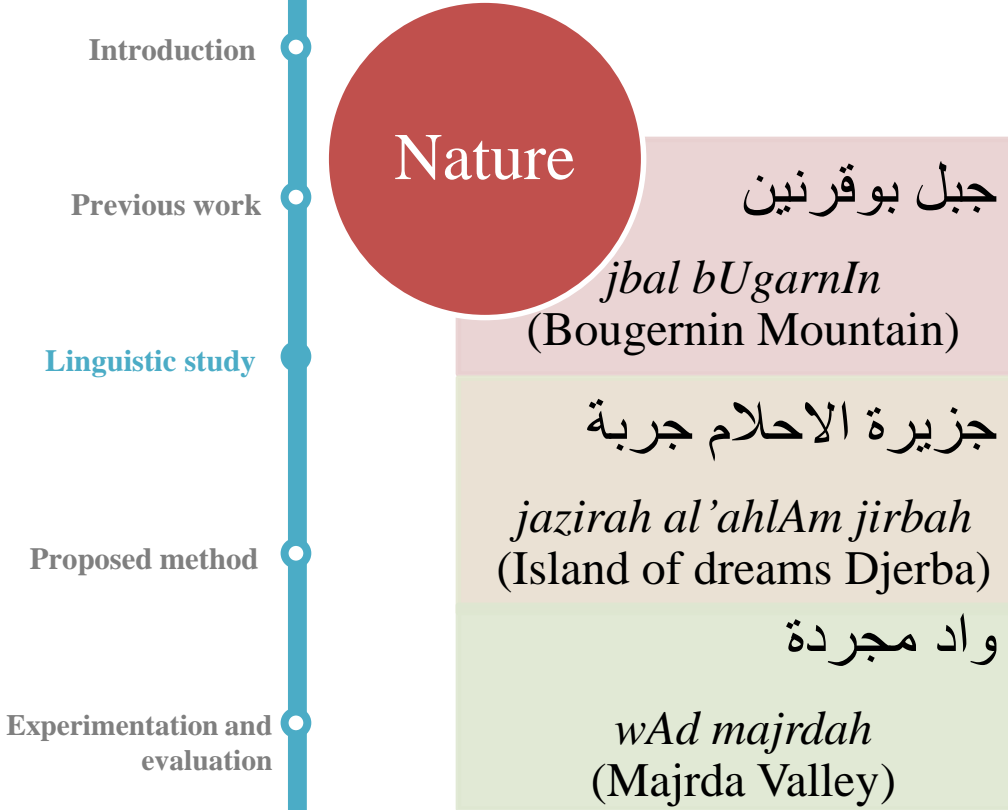
**City**

صفاقس

*sfAqis*
(Sfax)

ريف صفاقس

*rIf sfAqis*
(Sfax countryside)

تونس العاصمة

tUnis il'Asmah
(Tunis)

**Rule:**
(العاصمة)? City_name (ريف)?

# Linguistic study

Nature

| |
|---|
| جبل بوقرنين |
| *jbal bUgarnIn* (Bougernin Mountain) |
| جزيرة الاحلام جربة |
| *jazirah al'ahlAm jirbah* (Island of dreams Djerba) |
| واد مجردة |
| *wAd majrdah* (Majrda Valley) |

**Rule:**
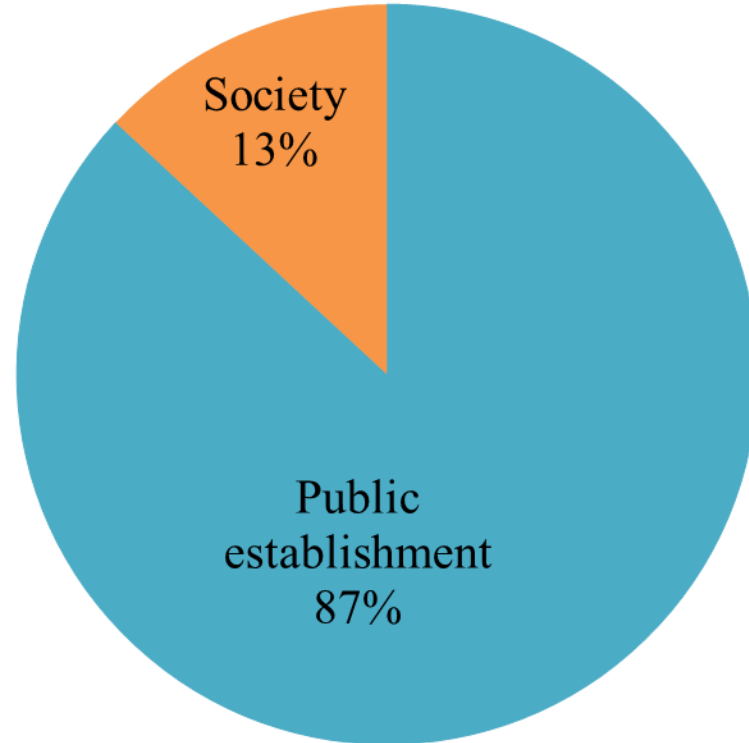**City_name (جبل|جزيرة|واد)**

Translation keeps the same syntax structure for location ENAMEX NEs

➔ Word for word translation

# Linguistic study

**Organisation ENAMEX NE**

Introduction

Previous work

Linguistic study

Proposed method

Experimentation and evaluation

Conclusion and Perspectives

# Linguistic study

**Compa-ny**

بنك الامان

*Bank il'amAn*
(Amen bank)

مقاولات شعبان

*muqAWlAt cha'bAn*
( Chaaben Contracting)

مغازة كارفور

*maghAzah kArfUr*
(carrefour market)

**Rule:
company_noun (first_name?
Last_name) definite noun))**

Translation keeps the same syntax structure for organization ENAMEX NEs
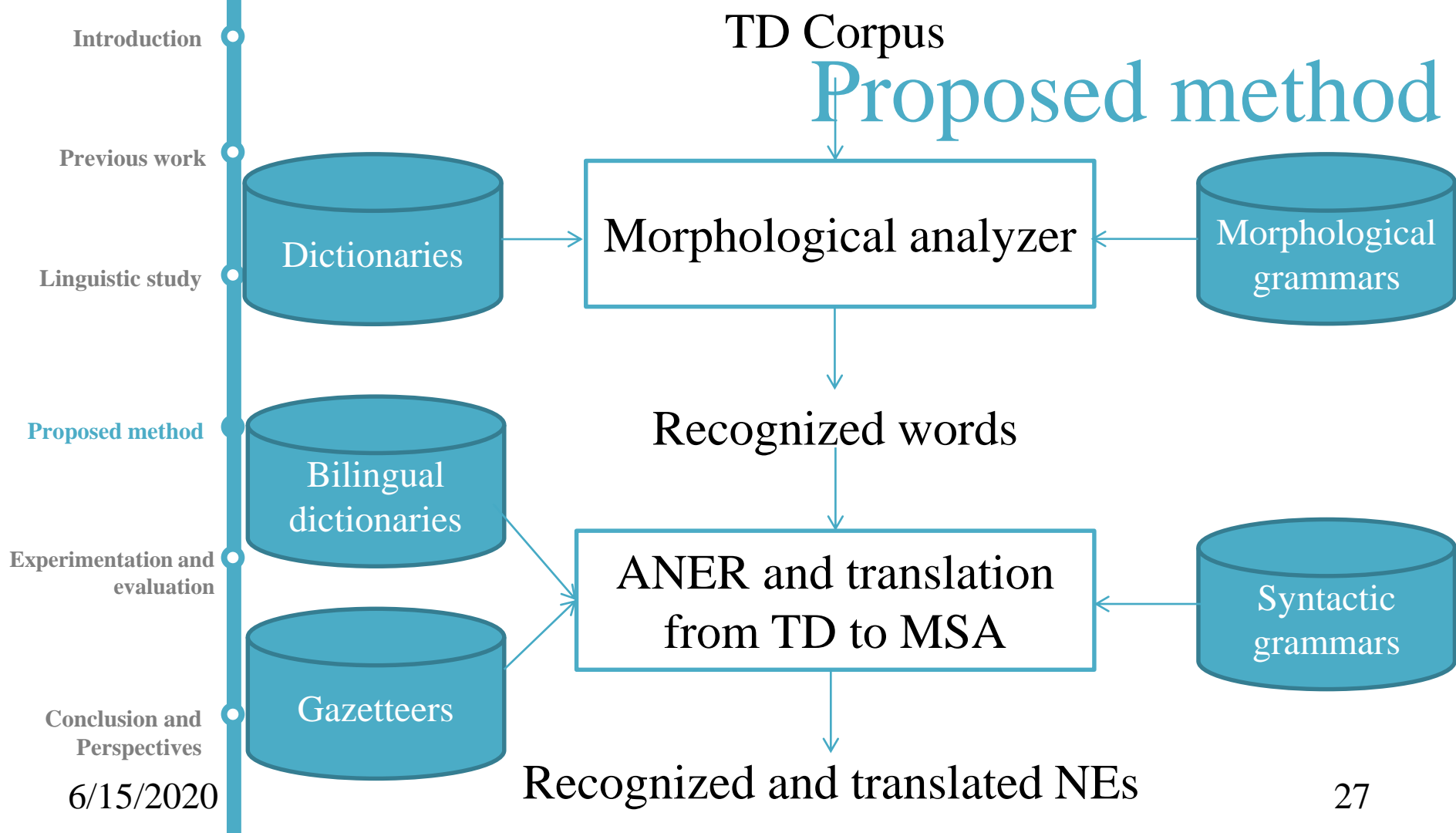
➔ Word for word translation

6/15/2020

26

## Proposed method

TD Corpus

Dictionaries → Morphological analyzer ← Morphological grammars

Recognized words

Bilingual dictionaries

Gazetteers → ANER and translation from TD to MSA ← Syntactic grammars

Recognized and translated NEs

# Proposed method

Bilingual dictionary thanks to + AR

"يَزِيد"=N+Prenom+m+AR,يَزِيدْ
"يَسري"=N+Prenom+m+AR,يُسْرِي
"يعقوب"=N+Prenom+m+AR,يَعْقُوبْ
"يُوسف"=N+Prenom+m+AR,يُوسِفْ
"يونس"=N+Prenom+m+AR,يُونِسْ

"سعيد"=N+Nomfamille+s+AR,سْعِيدْ
"بورقيبة"=N+Nomfamille+s+AR,بُورْقِيبَةْ
"قروي"=N+Nomfamille+s+AR,قَرْويِ

6/15/2020

28

Introduction

Previous work

Linguistic study

Proposed method

Experimentation and
evaluation

Conclusion and
Perspectives

# Proposed method

## Grammar of ANER and translation from TD to MSA

نبيل القروي رئيس حزب قلب تونس

# Proposed method

نبيل القروي رئيس حزب قلب تونس

# Experimentation and evaluation

Introduction

Previous work

Linguistic study

Proposed method

Experimentation and
evaluation

Conclusion and
Perspectives

# Experimentation and evaluation
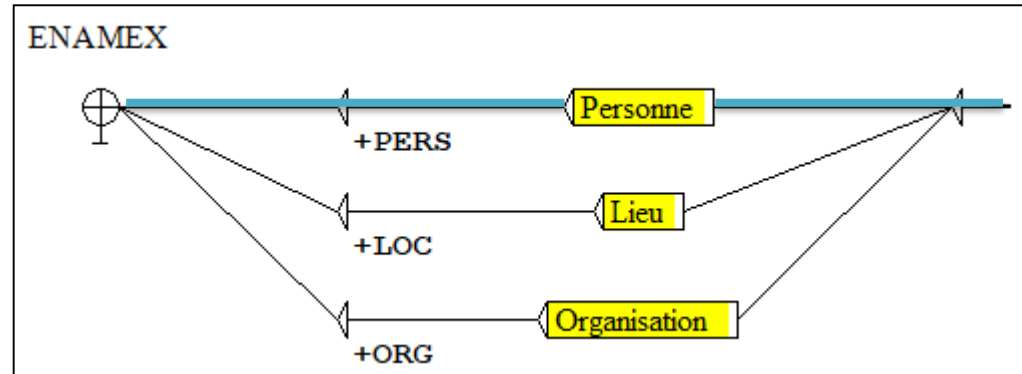
Pattern is:
- ○ a string of characters:
- ○ a PERL regular expression:

الفطايري علي يضحك

| 14 | | 9 | | 0,01 |

e=s | علي,N+Prenom+Genre=m+AR="علي" | فطائري,N+Nombre=s+Genre=m+Métier+AR="فطائري" | ال,PREF

في كرهبتو و جايك انا HEURE+> صباح الغد<TIMEX/غدوة الصباح
عيشة بلاهي كان عندك نومرو NOM+PERS+> يوسف<ENAMEX/يوسف

Index
- ○ Shortest matches
- ● Longest matches
- ○ All matches

Limitation
- ● All occurrences
- ○ Only: [100] occ.
- ☐ 1 occ. per match

☑ Reset Concordance    N o o J

# Experimentation and evaluation

Test corpus contains around 20000 words

|  | **Precision** | **Recall** |
|---|---|---|
| TIMEX | 96% | 94% |
| NUMEX | 98% | 91% |
| ENAMEX | 91% | 83% |

Our dictionary treat words of different origin

Our syntactic grammars use adjustment translation

Noise is present due to lack of disambiguation rules

# Conclusion and Perspectives

NE Linguistic study

ANER and Translation grammars based on transducers and implemented in NooJ

Word for word and adjustment translation

# Conclusion and Perspectives

Increased dictionary coverage

Building translation grammars for other NEs

Disambiguation of words

Thank you for your attention