# An Italian to Catalan RBMT system reusing data from existing language pairs

**Antonio Toral**, Mireia Ginestí-Rosell, Francis Tyers

2$^{nd}$ International Workshop on Free/Open-Source Rule-Based Machine Translation

2011/01/21

# Contents

- Two main approaches in Machine Translation: Rule-Based and Statistical

- Two main approaches in Machine Translation: Rule-Based and Statistical
- What is needed to build a system for a new language pair?

- Two main approaches in Machine Translation: Rule-Based and Statistical
- What is needed to build a system for a new language pair?
  - RBMT: dictionaries and rules

- Two main approaches in Machine Translation: Rule-Based and Statistical
- What is needed to build a system for a new language pair?
  - RBMT: dictionaries and rules
  - SMT: parallel corpus

- Two main approaches in Machine Translation: Rule-Based and Statistical
- What is needed to build a system for a new language pair?
    - RBMT: dictionaries and rules
    - SMT: parallel corpus
- Drawbacks:

- Two main approaches in Machine Translation: Rule-Based and Statistical
- What is needed to build a system for a new language pair?
    - RBMT: dictionaries and rules
    - SMT: parallel corpus
- Drawbacks:
    - RBMT: linguistic expertise on both languages, manual construction

- Two main approaches in Machine Translation: Rule-Based and Statistical
- What is needed to build a system for a new language pair?
  - RBMT: dictionaries and rules
  - SMT: parallel corpus
- Drawbacks:
  - RBMT: linguistic expertise on both languages, manual construction
  - SMT: only applicable to language pairs with big amounts of parallel data

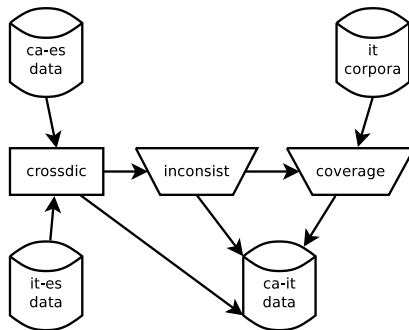- This paper: build RBMT system by exploiting data from existing pairs

- This paper: build RBMT system by exploiting data from existing pairs
- We build an MT system for pair $a$–$b$ given existing systems for pairs $a$–$c$ and $b$–$c$

- This paper: build RBMT system by exploiting data from existing pairs
- We build an MT system for pair $a$–$b$ given existing systems for pairs $a$–$c$ and $b$–$c$
- Italian→Catalan from Apertium's Italian–Spanish and Catalan–Spanish

- This paper: build RBMT system by exploiting data from existing pairs
- We build an MT system for pair $a$–$b$ given existing systems for pairs $a$–$c$ and $b$–$c$
- Italian→Catalan from Apertium's Italian–Spanish and Catalan–Spanish
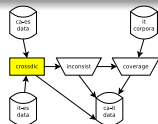- Motivation:

- This paper: build RBMT system by exploiting data from existing pairs
- We build an MT system for pair *a–b* given existing systems for pairs *a–c* and *b–c*
- Italian→Catalan from Apertium's Italian–Spanish and Catalan–Spanish
- Motivation:
    - RBMT competitive and useful for languages without parallel corpora

- This paper: build RBMT system by exploiting data from existing pairs
- We build an MT system for pair $a$–$b$ given existing systems for pairs $a$–$c$ and $b$–$c$
- Italian→Catalan from Apertium's Italian–Spanish and Catalan–Spanish
- Motivation:
    - RBMT competitive and useful for languages without parallel corpora
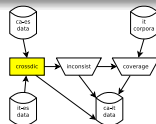    - Reusing data from similar pairs significantly reduces the amount of work

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
Inconsistencies
Coverage

Introduction
**Methodology**
Evaluation
Conclusions

**Crossdics**
Inconsistencies
Coverage

- Input dictionaries:

Introduction
**Methodology**
Evaluation
Conclusions

**Crossdics**
Inconsistencies
Coverage

- Input dictionaries:
  - es–it: mono es 11k, mono it 10k, bi 12k

Introduction
**Methodology**
Evaluation
Conclusions

**Crossdics**
Inconsistencies
Coverage

- Input dictionaries:
  - es–it: mono es 11k, mono it 10k, bi 12k
  - es–ca: mono es 44k, mono ca 40k, bi 51k

Introduction
**Methodology**
Evaluation
Conclusions

**Crossdics**
Inconsistencies
Coverage

- Input dictionaries:
  - es–it: mono es 11k, mono it 10k, bi 12k
  - es–ca: mono es 44k, mono ca 40k, bi 51k
- Output dictionaries:

Introduction
Methodology
Evaluation
Conclusions

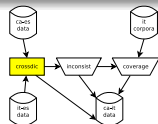Crossdics
Inconsistencies
Coverage

- Input dictionaries:
    - es–it: mono es 11k, mono it 10k, bi 12k
    - es–ca: mono es 44k, mono ca 40k, bi 51k
- Output dictionaries:
    - it–ca: mono it 7k, mono ca 8k, bi 9k

Introduction
Methodology
Evaluation
Conclusions

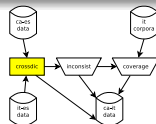Crossdics
Inconsistencies
Coverage

- Input dictionaries:
  - es–it: mono es 11k, mono it 10k, bi 12k
  - es–ca: mono es 44k, mono ca 40k, bi 51k
- Output dictionaries:
  - it–ca: mono it 7k, mono ca 8k, bi 9k
- Other linguistic data:

Introduction
**Methodology**
Evaluation
Conclusions

**Crossdics**
Inconsistencies
Coverage
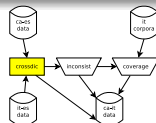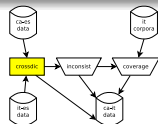
- Input dictionaries:
    - es–it: mono es 11k, mono it 10k, bi 12k
    - es–ca: mono es 44k, mono ca 40k, bi 51k
- Output dictionaries:
    - it–ca: mono it 7k, mono ca 8k, bi 9k
- Other linguistic data:
    - it tagger and disambiguation probabilities taken from it–es

Introduction
**Methodology**
Evaluation
Conclusions

**Crossdics**
Inconsistencies
Coverage

- Input dictionaries:
  - es–it: mono es 11k, mono it 10k, bi 12k
  - es–ca: mono es 44k, mono ca 40k, bi 51k
- Output dictionaries:
  - it–ca: mono it 7k, mono ca 8k, bi 9k
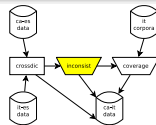- Other linguistic data:
  - it tagger and disambiguation probabilities taken from it–es
  - transfer rules: 35 taken from oc–ca (mainly noun phrases) + 9 manually created (verbs and clitic pronouns)

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
**Inconsistencies**
Coverage

- Reasons

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
**Inconsistencies**
Coverage

- Reasons
  - Differences of gender and number (it–ca)

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
**Inconsistencies**
Coverage

- Reasons
  - Differences of gender and number (it–ca)
  - Different ways of categorising lemmas and morphological features (es–it and es–ca)

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
**Inconsistencies**
Coverage

- Reasons
  - Differences of gender and number (it–ca)
  - Different ways of categorising lemmas and morphological features (es–it and es–ca)
- Solutions

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
**Inconsistencies**
Coverage

- Reasons
    - Differences of gender and number (it–ca)
    - Different ways of categorising lemmas and morphological features (es–it and es–ca)
- Solutions
    - Manually solve inconsistencies (identified automatically), 0.5 Person Months

Introduction
**Methodology**
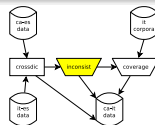Evaluation
Conclusions
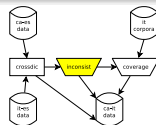
Crossdics
**Inconsistencies**
Coverage

- Reasons
    - Differences of gender and number (it–ca)
    - Different ways of categorising lemmas and morphological features (es–it and es–ca)
- Solutions
    - Manually solve inconsistencies (identified automatically), 0.5 Person Months
    - Substitute derived ca mono dictionary (8k) for that in es–ca (40k)

Introduction
Methodology
Evaluation
Conclusions

Crossdics
Inconsistencies
Coverage

- Coverage calculated on two Italian corpora: Europarl and Wikipedia

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
Inconsistencies
**Coverage**

- Coverage calculated on two Italian corpora: Europarl and Wikipedia
- 155 most frequent unknown words added to the system

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
Inconsistencies
**Coverage**

- Coverage calculated on two Italian corpora: Europarl and Wikipedia
- 155 most frequent unknown words added to the system

Introduction
**Methodology**
Evaluation
Conclusions

Crossdics
Inconsistencies
**Coverage**

- Coverage calculated on two Italian corpora: Europarl and Wikipedia
- 155 most frequent unknown words added to the system

|                  | Europarl   | Wikipedida  |
|------------------|------------|-------------|
| Number of words  | 46,569,602 | 241,563,615 |
| Initial coverage | 86.4%      | 75.5%       |
| Final coverage   | 88.9%      | 79.4%       |

Introduction
Methodology
Evaluation
Conclusions

Setting
Results

- Systems

Introduction
Methodology
Evaluation
Conclusions

Setting
Results

- Systems
  - Apertium (the it→ca system)

Introduction
Methodology
**Evaluation**
Conclusions

Setting
Results

- Systems
    - Apertium (the it→ca system)
    - Apertium-i (indirect translation using it→es and es→ca)

Introduction
Methodology
**Evaluation**
Conclusions

**Setting**
Results

- Systems
    - Apertium (the it→ca system)
    - Apertium-i (indirect translation using it→es and es→ca)
    - Google Translate

Introduction
Methodology
**Evaluation**
Conclusions

**Setting**
Results

- Systems
    - Apertium (the it→ca system)
    - Apertium-i (indirect translation using it→es and es→ca)
    - Google Translate
- Test set: 1k sentences from KDE4 (OPUS project)

Introduction
Methodology
**Evaluation**
Conclusions

**Setting**
Results

- Systems
    - Apertium (the it→ca system)
    - Apertium-i (indirect translation using it→es and es→ca)
    - Google Translate
- Test set: 1k sentences from KDE4 (OPUS project)
- Metrics: TER, GTM, BLEU, NIST

Introduction
Methodology
**Evaluation**
Conclusions

Setting
**Results**

| Metric | Apertium | Apertium-i | Google |
|--------|----------|------------|--------|
| TER    | **0.5703** | 0.6118   | 0.6785 |
| GTM    | **0.5162** | 0.4712   | 0.41637 |
| BLEU   | 0.2290   | 0.1492     | **0.2459** |
| NIST   | 5.6567   | 4.4753     | **6.1071** |

- GTM and TER: Apertium > Apertium-i > Google

Introduction
Methodology
**Evaluation**
Conclusions

Setting
**Results**

| Metric | Apertium | Apertium-i | Google |
|--------|----------|------------|--------|
| TER | **0.5703** | 0.6118 | 0.6785 |
| GTM | **0.5162** | 0.4712 | 0.41637 |
| BLEU | 0.2290 | 0.1492 | **0.2459** |
| NIST | 5.6567 | 4.4753 | **6.1071** |

- GTM and TER: Apertium > Apertium-i > Google
- BLEU*: Google $\simeq$ Apertium > Apertium-i

Introduction
Methodology
**Evaluation**
Conclusions

Setting
**Results**

| Metric | Apertium | Apertium-i | Google |
|--------|----------|------------|--------|
| TER | **0.5703** | 0.6118 | 0.6785 |
| GTM | **0.5162** | 0.4712 | 0.41637 |
| BLEU | 0.2290 | 0.1492 | **0.2459** |
| NIST | 5.6567 | 4.4753 | **6.1071** |

- GTM and TER: Apertium > Apertium-i > Google
- BLEU*: Google $\simeq$ Apertium > Apertium-i
- NIST*: Google > Apertium > Apertium-i

- it→ca RBMT system derived from es–it and es–ca systems

- it→ca RBMT system derived from es–it and es–ca systems
- Limited amount of manual work: correct inconsistencies, augment coverage and add some transfer rules

- it→ca RBMT system derived from es–it and es–ca systems
- Limited amount of manual work: correct inconsistencies, augment coverage and add some transfer rules
- Evaluated against RBMT indirect system and SMT system, yielding significant improvements over both

- it→ca RBMT system derived from es–it and es–ca systems
- Limited amount of manual work: correct inconsistencies, augment coverage and add some transfer rules
- Evaluated against RBMT indirect system and SMT system, yielding significant improvements over both
- System released as apertium-ca-it-0.1.0

Thanks! Questions?