Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

# Towards a learning approach
# for abbreviation detection and resolution

Klaar Vanopstal, Bart Desmet, Véronique Hoste

LT[3], Language and Translation Technology Team
University College Ghent
{klaar.vanopstal,bart.desmet,véronique.hoste}@hogent.be

Department of Applied Mathematics & Computer Science
Ghent University
Krijgslaan 281 (S9), 9000 Gent, Belgium

May 19, 2010

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

1. Background

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

1 Background

2 Annotation

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

1. Background

2. Annotation

3. Pattern-based approach

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

1. Background

2. Annotation

3. Pattern-based approach

4. Learning-based approach

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

Problem
Use

## Problem

Information explosion $\Rightarrow$ growing number of (bio)medical
abbreviations.
New abbreviations are created; not always known to the reader.
$\Rightarrow$ automatic detection and resolution

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

Problem
Use

# Use

- information retrieval
- information extraction
- NER
- anaphora resolution

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

## Corpus

- English
    - AbbRE: reliable standard but limited size
    - Medstract: publicly available and commonly used

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

## Corpus

- English
    - AbbRE: reliable standard but limited size
    - Medstract: publicly available and commonly used
- Dutch: no resources available

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

## Corpus

- English
  - AbbRE: reliable standard but limited size
  - Medstract: publicly available and commonly used
- Dutch: no resources available
- Abstracts from 2 medical journals:
  - *Nederlands Tijdschrift voor Geneeskunde* (NTvG); 29,978 words
  - *Belgisch Tijdschrift voor Geneeskunde* (TvG); 36,757 words

  $\Rightarrow$ total of 66,739 words

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

Different **types** of abbreviations included in annotations:

- **Truncation**

### Example

*adm* for *adm*inistration

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

Different **types** of abbreviations included in annotations:

- **Truncation**

Example

*adm* for *adm*inistration

- **First letter initialization**

Example

*AAA* for *a*bdominal *a*ortic *a*neurysm

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

Different **types** of abbreviations included in annotations:

- **Truncation**

### Example

*adm* for <u>adm</u>inistration

- **First letter initialization**

### Example

*AAA* for <u>a</u>bdominal <u>a</u>ortic <u>a</u>neurysm

- **Opening letter initialization**

### Example

*HeLa* for <u>He</u>nrietta <u>La</u>cks

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

- **Syllabic initialization**

### Example

*BZD* for *benzodiazepine*

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

- **Syllabic initialization**

### Example

*BZD* for *benzodiazepine*

- **Substitution initialization**

### Example

*Fe* for *iron*

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

**Corpus**
Labels

- **Syllabic initialization**

### Example

*BZD* for <u>ben</u><u>zo</u><u>di</u>azepine

- **Substitution initialization**

### Example

*Fe* for *iron*

- **Combination of letters and numbers**

### Example

*CXCR4* for *chemokine receptor fusin*

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

## Labels

1. **ABBR**: Dutch abbreviations which have a full form in their local context

### Example

Hoge-resolutie-computertomografie (**HRCT**)
<u>EN</u>: High resolution computed tomography (HRCT)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

## Labels

1. **ABBR**: Dutch abbreviations which have a full form in their local context

### Example

Hoge-resolutie-computertomografie (**HRCT**)
<u>EN</u>: High resolution computed tomography (HRCT)

2. **ABBR_DE**: Dutch abbreviations with full form in abstract (not in local context)

### Example

de pathofysiologie van het **CFS**
<u>EN</u>: the pathophysiology of CFS

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

3. **DEF**: Dutch full forms which define an abbreviation in their local context

### Example

**Hoge-resolutie-computertomografie** (HRCT)
<u>EN</u>: High resolution computed tomography (HRCT)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

3. **DEF**: Dutch full forms which define an abbreviation in their local context

### Example

**Hoge-resolutie-computertomografie** (HRCT)
<u>EN</u>: High resolution computed tomography (HRCT)

4. **ABBR_IN_COMP**: part of a compound word; no definition in the abstract

### Example

**HIV**-patiënten
(EN: HIV patients)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

5. **ABBR_IN_COMP_DE**: part of a compound word; full form in abstract

### Example

ernstige *reumatoïde artritis* (RA)-vasculitis. Bij de ziekte van Wegener en **RA**-vasculitis...
EN: severe rheumatoid arthritis (RA) vasculitis. Wegener's disease and RA vasculitis...)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

5. **ABBR_IN_COMP_DE**: part of a compound word; full form
   in abstract

### Example

ernstige *reumatoïde artritis* (RA)-vasculitis. Bij de ziekte van
Wegener en **RA**-vasculitis...
EN: severe rheumatoid arthritis (RA) vasculitis. Wegener's disease
and RA vasculitis...)

6. **ABBR_NO_DEF**: abbreviations without full form

### Example

AIDS, HIV

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

7. **ABBR_EN**: English abbreviation with Dutch/English definition in local context

### Example

endosonografie (**EUS**)
<u>EN</u>: endoscopic ultrasound (EUS)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

7. **ABBR_EN**: English abbreviation with Dutch/English definition in local context

### Example

endosonografie (**EUS**)
<u>EN</u>: endoscopic ultrasound (EUS)

8. **DEF_EN**: English full form which accompanies an English abbreviation

### Example

Mini Mental State Examination (**MMSE**)

$\Rightarrow$ Kappa score: 0.89

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

|                  | NTvG | TvG   |
|------------------|------|-------|
| ABBR             | 11.60| 14.25 |
| ABBR_DE          | 30.62| 22.55 |
| ABBR_IN_COMP     | 7.14 | 22.43 |
| ABBR_IN_COMP_DE  | 16.85| 4.96  |
| ABBR_NO_DEF      | 27.65| 29.12 |
| ABBR_EN          | 6.14 | 6.69  |
| TOTAL %          | 3.36 | 2.19  |

Table: Labels and their frequencies in the corpus (%)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

|                | NTvG    | TvG      |
|----------------|---------|----------|
| def: loc       | 17.74%  | 20.94 %  |
| def: broad     | 47.47%  | 27.50%   |
| def: loc/broad | 65.21%  | 48.45%   |

Table: Abbreviations and defined abbreviations in the corpus

$\Rightarrow$ Between 45% and 52% of the abbreviations are undefined

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

## Challenges

- English abbreviations with Dutch full form: no match

### Example

HAART = **k**rachtige **a**ntiretrovirale **t**herapie

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

## Challenges

- English abbreviations with Dutch full form: no match

### Example

HAART = **k**rachtige **a**ntiretrovirale **t**herapie

- Parenthetical patterns

### Example

gunstige uitkomst (**score 5**)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

## Challenges

- English abbreviations with Dutch full form: no match

### Example

HAART = **k**rachtige **a**ntiretrovirale **t**herapie

- Parenthetical patterns

### Example

gunstige uitkomst (**score 5**)

- Syllabic initialization

### Example

CVS = **c**hronische-**v**ermoeidheids**s**yndroom
<u>EN</u>: CFS = **c**hronic **f**atigue **s**yndrome)

Background
**Annotation**
Pattern-based approach
Learning-based approach
Conclusions and future work

Corpus
**Labels**

## Challenges

- English abbreviations with Dutch full form: no match

### Example

HAART = **k**rachtige **a**ntiretrovirale **t**herapie

- Parenthetical patterns

### Example

gunstige uitkomst (**score 5**)

- Syllabic initialization

### Example

CVS = **c**hronische-**v**ermoeidheids**s**yndroom
<u>EN</u>: CFS = **c**hronic **f**atigue **s**yndrome)

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
Results

## Pattern-based approach - Related research

$\Rightarrow$ Use of patterns to detect abbreviations:

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

**Related research**
Own approach
Results

## Pattern-based approach - Related research

$\Rightarrow$ Use of patterns to detect abbreviations:

- short uppercase words
- typical patterns: "long form (short form)" or "short form (long form)"

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

**Related research**
Own approach
Results

# Pattern-based approach - Related research

$\Rightarrow$ Use of patterns to detect abbreviations:

- short uppercase words

- typical patterns: "long form (short form)" or "short form (long form)"

- <u>identification of definitions:</u>

    - window of **2\*N** (Taghva & Gilbreth, 1999)
      or **3\*N** words (Stanford Medical Abbreviation Method (Chang & Schütze, 2006))
    - **text markers**: () " =
    - **linguistic cues**: "short", "or" (Park & Byrd, 2001)

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

**Related research**
Own approach
Results

- + use of NLP tools to refine the search space of the definitions (Pustojevski et al., 2001) and/or to tackle the problem of function word matching

### Example

**ADL** = **a**ctiviteiten <u>van het</u> **d**agelijkse **l**even
<u>EN</u>: daily life activities

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

<u>2 steps:</u>

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

**2 steps:**

- Abbreviation detection

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

**2 steps:**

- Abbreviation detection
- Definition matching

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

## Step 1: abbreviation detection:

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

### Step 1: abbreviation detection:

- capital letters / combinations of capital letters with 1-3 lowercased letters or numbers

#### Example

QSRL
pANCA
CDG1A

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

### Step 1: abbreviation detection:

- capital letters / combinations of capital letters with 1-3 lowercased letters or numbers

### Example

QSRL
pANCA
CDG1A

- window of 3*N words

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

### Step 1: abbreviation detection:

- capital letters / combinations of capital letters with 1-3 lowercased letters or numbers

#### Example

QSRL
pANCA
CDG1A

- window of 3*N words
- text markers () = " ' $\Rightarrow$ list of candidate definitions

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

## Step 2: definition matching:

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

### Step 2: definition matching:

- list of candidate definitions

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

### Step 2: definition matching:

- list of candidate definitions
- matching: first letter of abbreviation - words in candidate definition
  $\Rightarrow$ matching word + rest of the 3*N sequence = definition

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

| Abbreviations | | | |
|---|---|---|---|
| | precision | recall | FB1 |
| TvG | 83.89 | 78.64 | 81.18 |
| NTvG | 82.05 | 83.07 | 82.56 |
| Definitions | | | |
| | precision | recall | FB1 |
| TvG | 74.49 | 83.36 | 78.68 |
| NTvG | 68.03 | 85.50 | 75.77 |

Table: Results of the pattern-based approach

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

# Error Analysis

- Errors in abbreviation detection step

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

# Error Analysis

- Errors in abbreviation detection step
  - **Titles** printed in capital letters

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

# Error Analysis

- Errors in abbreviation detection step
    - **Titles** printed in capital letters
    - **Roman numerals** confused with capitalized i, v or x

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

## Error Analysis

- Errors in abbreviation detection step
    - **Titles** printed in capital letters
    - **Roman numerals** confused with capitalized i, v or x
    - **single letters** which are not abbreviations (e.g. hepatitis **A**)

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

## Error Analysis

- Errors in abbreviation detection step
    - **Titles** printed in capital letters
    - **Roman numerals** confused with capitalized i, v or x
    - **single letters** which are not abbreviations (e.g. hepatitis **A**)
    - abbreviations with **word-internal capital letters** (e.g. mmHg (EN: Torr))

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

## Error Analysis

- Errors in abbreviation detection step
    - **Titles** printed in capital letters
    - **Roman numerals** confused with capitalized i, v or x
    - **single letters** which are not abbreviations (e.g. hepatitis **A**)
    - abbreviations with **word-internal capital letters** (e.g. mmHg (EN: Torr))
    - abbreviations with **no typical orthographical characteristics** (e.g. min)

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

- Errors in definition matching step

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

- Errors in definition matching step
  - **error percolation**

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

- Errors in definition matching step
    - **error percolation**
    - **mislinked words** (e.g. **h**et **h**epatitis-**A**-**v**irus (HAV))

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

- Errors in definition matching step
    - **error percolation**
    - **mislinked words** (e.g. **h**et **h**epatitis-**A**-**v**irus (HAV))
    - **function words** (e.g. **op** evidentie gebaseerde zorg (EBZ)
      (EN: evidence-based medicine (EBM))

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

- Errors in definition matching step
    - **error percolation**
    - **mislinked words** (e.g. **h**et **h**epatitis-**A**-**v**irus (HAV))
    - **function words** (e.g. **op** evidentie gebaseerde zorg (EBZ)
      (EN: evidence-based medicine (EBM)))
    - **English** abbreviations with a **Dutch** definition

Background
Annotation
**Pattern-based approach**
Learning-based approach
Conclusions and future work

Related research
Own approach
**Results**

- Errors in definition matching step
  - **error percolation**
  - **mislinked words** (e.g. **h**et **h**epatitis-**A**-**v**irus (HAV))
  - **function words** (e.g. **op** evidentie gebaseerde zorg (EBZ)
    (EN: evidence-based medicine (EBM))
  - **English** abbreviations with a **Dutch** definition
  - **contractions** (e.g. **therapiegebonden** secundaire
    myelodysplasie (**t** - MDS) en acute leukemie (**t** - AL).
    (EN: the incidence of therapy-related secondary myelodysplasia
    (t-MDS) and acute leukemia (t-AL).)

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

**Related research**
Own approach
Results

## Learning-based approach - Related research

- Often in combination with pattern-based techniques, e.g.
  Stanford Medical Abbreviation Method (2006), Chang et al.
  (2002)

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

**Related research**
Own approach
Results

# Learning-based approach - Related research

- Often in combination with pattern-based techniques, e.g. Stanford Medical Abbreviation Method (2006), Chang et al. (2002)
- Pattern-based detection of abbreviations + learning-based matching with definitions

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

**Related research**
Own approach
Results

# Learning-based approach - Related research

- Often in combination with pattern-based techniques, e.g. Stanford Medical Abbreviation Method (2006), Chang et al. (2002)
- Pattern-based detection of abbreviations + learning-based matching with definitions
- examples of features:
  - % of characters aligned at beginning of word
  - % of characters aligned on syllable boundary
  - number of words that were skipped (negative weight)

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

**Related research**
Own approach
Results

# Learning-based approach - Related research

- Often in combination with pattern-based techniques, e.g. Stanford Medical Abbreviation Method (2006), Chang et al. (2002)
- Pattern-based detection of abbreviations + learning-based matching with definitions
- examples of features:
  - % of characters aligned at beginning of word
  - % of characters aligned on syllable boundary
  - number of words that were skipped (negative weight)

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

Related research
**Own approach**
Results

## Own approach

- Preprocessing steps:

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

Related research
**Own approach**
Results

## Own approach

- Preprocessing steps:
  - tokenization

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

## Own approach

- Preprocessing steps:
  - tokenization
  - POS tagging + NP chunking (Daelemans & van den Bosch, 2005)

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

Related research
**Own approach**
Results

- Learning experiments

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

Related research
**Own approach**
Results

- Learning experiments
  - YamCha (Kudo & Matsumoto, 2003): open source sequence tagger using SVM

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

Related research
**Own approach**
Results

- Learning experiments
  - YamCha (Kudo & Matsumoto, 2003): open source sequence tagger using SVM
  - 10-fold cross-validation

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

Related research
**Own approach**
Results

- Feature vector:

Background
Annotation
Pattern-based approach
Learning-based approach
Conclusions and future work

Related research
**Own approach**
Results

- Feature vector:
    - token
    - POS
    - name initials
    - sentence-initial position
    - morphological features (initial capital letter, completely capitalized, internal capital letters, lowercased, roman number, punctuation, hyphens, exclusively consonants)
    - prefix and suffix information
    - symbolic word shape feature: all morphological (binary) features
    - feature to match 1st letter of abbreviation against words in 3*N sequence

Background
Annotation
Pattern-based approach
**Learning-based approach**
Conclusions and future work

Related research
Own approach
**Results**

# Results

| Abbreviations | | | |
|---|---|---|---|
| | precision | recall | FB1 |
| TvG | 95.31 | 92.26 | 93.76 |
| NTvG | 96.76 | 90.97 | 93.78 |
| Definitions | | | |
| | precision | recall | FB1 |
| TvG | 86.92 | 78.18 | 82.32 |
| NTvG | 87.19 | 78.00 | 82.34 |

Table: Ten-fold cross-validation results of the learning experiments.

# Conclusions

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Conclusions

- **annotated dataset** of $+/-$ 67,000 words (Dutch, medical)

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Conclusions

- **annotated dataset** of +/- 67,000 words (Dutch, medical)
- **2 approaches**: **pattern-based** and **classification-based**

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Conclusions

- **annotated dataset** of $+/-$ 67,000 words (Dutch, medical)
- **2 approaches**: **pattern**-based and **classification**-based
- classification-based approach **outperforms** the pattern-based approach on both tasks:

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Conclusions

- **annotated dataset** of $+/-$ 67,000 words (Dutch, medical)
- **2 approaches**: **pattern-based** and **classification-based**
- classification-based approach **outperforms** the pattern-based approach on both tasks:
  - abbreviation detection: 93% F-score

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Conclusions

- **annotated dataset** of +/- 67,000 words (Dutch, medical)
- **2 approaches**: **pattern**-based and **classification**-based
- classification-based approach **outperforms** the pattern-based approach on both tasks:
    - abbreviation detection: 93% F-score
    - definition matching: 82% F-score

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

# Future work

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Future work

- incorporate information from **error analysis** into learning approach

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Future work

- incorporate information from **error analysis** into learning approach
- apply **decompounding** techniques (syllabic initializations)

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Future work

- incorporate information from **error analysis** into learning approach
- apply **decompounding** techniques (syllabic initializations)
- **cross-lingual matching**: external sources $+$ MT techniques

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Future work

- incorporate information from **error analysis** into learning approach
- apply **decompounding** techniques (syllabic initializations)
- **cross-lingual matching**: external sources $+$ MT techniques
- **undefined** abbreviations: external sources

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Future work

- incorporate information from **error analysis** into learning approach
- apply **decompounding** techniques (syllabic initializations)
- **cross-lingual matching**: external sources + MT techniques
- **undefined** abbreviations: external sources
- **F-scores** per label (now focus on abbreviations and definitions)

Background
Annotation
Pattern-based approach
Learning-based approach
**Conclusions and future work**

## Future work

- incorporate information from **error analysis** into learning approach
- apply **decompounding** techniques (syllabic initializations)
- **cross-lingual matching**: external sources + MT techniques
- **undefined** abbreviations: external sources
- **F-scores** per label (now focus on abbreviations and definitions)
- **English** corpus