# Bike: Bilingual Keyphrase Experiments

**David Nadeau, Caroline Barrière** and **George Foster**
Institute for Information Technology
National Research Council Canada
Gatineau, Quebec, Canada

National Research Council Canada  Conseil national de recherches Canada

Canada

## Introduction

- In this talk, I report <u>experiments</u> we conducted on the task of <u>translating lists of keyphrases</u>.

# Translating lists of keyphrases



# Translating lists of keyphrases

- Related to:
  - Sub-sentence-level translation (noun phrases, terminology, …)
  - Query translation (cross-lingual IR);
  - Cross-lingual summarization (Summary as a list of Keyphrases)

- BIKE (Bilingual Keyphrase Experiments):

  - Bidirectional French ←→ English keyphrase translation

  - Strategy pieces:
    - Statistical MT
    - Terminological resource
    - Inflectional morphology
    - WSD (list of keyphrases as a context)

  - Experiments in training and combining pieces.

- Collection of 3058 scientific papers from the
  Canada Institute for Scientific and Technical Information (**CISTI**)
  in ten domains (biochemistry, botany, chemistry, civil engineering,
  environment (x2), genomics, geotechnical, microbiology and pharmacology)

- Each document is a tuple $\{A_f, A_e, K_f, K_e, T\}$

# Task

$A_e$ (English abstract)

$K_e$ (English keyphrases)

$A_f$ (French abstract)

$K_f$ (French keyphrases)

T (Full text)

NRC·CNRC
Institute for Information Technology

---

# Task (performance metric)

NRC·CNRC
Institute for Information Technology

- English → French
- For each keyphrase, 1 candidate translation is produced

- Exact translation required (reproductibility, lowerbound exp.)

- **Accuracy**: ratio of correctly translated keyphrases to the total number of keyphrases

- Results are reported for all domains (30% held-out split)
- Significance testing done using 10-domain split

- Baseline (no translation)

Statistical MT:

- Phrase-based statistical MT [Koehn *et al*., 2003]

- We calculate the conditional Fr-given-En probabilities, and retain <u>only the most probable translation</u>

- Hansard Model: Canadian Hansard parallel corpora
- CISTI Model: 40% training split (abstracts and keyphrases)
  – CISTI « global »: all journals
  – CISTI « individual » one model per journal

- How to improve statistical MT translation?
- Problem:

| English | French |
|---|---|
| fiber | *de* fibre |
| population sizes | *des* dimensions des populations |

- Remove prefixes ("de la ", "le ", "la ", "les ", "l'", "du ", "de ", ",")
  and suffixes (" ," , " de", " du", " des") from French translation
  proposed by the MT system.

| | Accuracy (%) |
|---|---|
| baseline experiment | 20.21 |
| hansard | 26.11 |
| cisti individual | 36.50 |
| cisti global | 39.54 |
| cisti global + correction | 41.26 |

† statistically significant at the 95% level

**Institute for Information Technology** — NRC·CNRC

- Terminological resource:
  the *Grand Dictionnaire Terminologique* (GDT)

- Try « exact match » translation

- Use the first entry



---

**Institute for Information Technology** — NRC·CNRC

- **Important problem**: lot of author keyphrases are plural while the ressource contains only singular terms

- We handled limited inflectional morphology:
  - Detection of English plural (e.g.: word ending [^f]ves)
  - Singularization of English term (e.g.: ves → f)
  - Pluralization of French term (e.g.: al → aux)

  - Simple heuristics for multi-word expressions
    ex.:     pomme de terre → pomme**s** de terre
             pomme de terre frite → pomme**s** de terre frite**s**

- How to choose a candidate in a terminological resource?
- Problem: different *domain,* different *translation*

  e.g: marché = market (finance)
       marché = market place (commerce)
       marché = contract (law)

- We introduce the Minimal Domain Set (MDS) algorithm

- **Idea**: find the minimal set of *coherent* domains covering all keyphrases.

1. Calculate the frequency of each domain *F(D)*
2. Calculate the number of domain per keyphrase $|D_{Ki}|$
3. Sort Keyphrases in ascending order of $|D_{Ki}|$

   *For each keyphrase:*

4. **Likeliness**: From the list of domains $D_{Ki}$, build a reduced list containing only the domains with the highest frequency F(D)
5. **Coherence**: From this reduced list, select the domain which has the highest coherence with a member of MDS. Add this domain to MDS.

bar (metallurgy, textile, law, automotive, …30 other!)

iron (metallurgy, shoe)

cement (shoe, textile)

rubber (shoe, rubber, leasure, graphic)

lace (clothing, shoe, brewing)

| (Step 1) F(shoe) | = 4 | (Step 2) $|D_{iron}|$ | = 2 |
|---|---|---|---|
| F(metallurgy) | = 2 | $|D_{cement}|$ | = 2 |
| F(textile) | = 2 | $|D_{lace}|$ | = 3 |
| F(*other*) | = 1 | $|D_{rubber}|$ | = 4 |
| | | $|D_{bar}|$ | = 34 |

(Step 3) L := iron, cement, lace, rubber, bar.

(first loop)  **Iron (**metallurgy, shoe**)**

(Step 4)  Domain likeliness: F(shoe) = 4

MDS = {shoe}

…

(last loop)  **bar** (metallurgy, textile, law, automotive, …30 other!)

(Step 4)  Domain likeliness: F(metallurgy) = F(textile) = 2

(Step 5)  Domain coherence:
C(shoe, metallurgy) = 0.1, C(shoe, textile) = 0.3

MDS = {shoe, textile}

|  | Accuracy (%) |
|---|---|
| GDT | 35.66 |
| GDT + morphology | 38.30 |
| GDT + morph + MDS | 39.22 |

† statistically significant at the 99% level
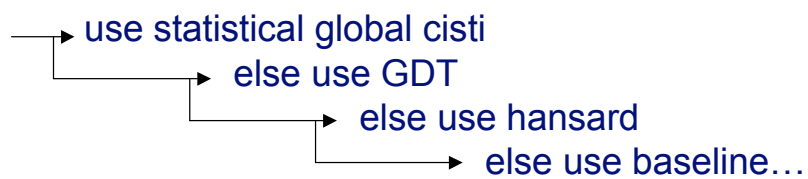
‡ statistically significant at the 95% level

**Experiment 4**
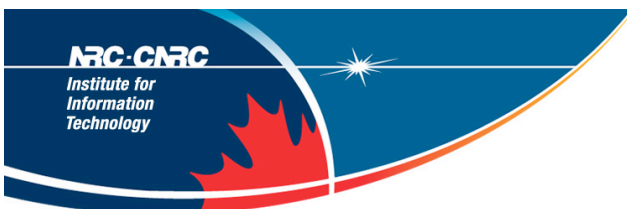
Combining all previous modules:

- Cascading modules

  – For instance:

use statistical global cisti
else use GDT
else use hansard
else use baseline…

- Combine modules (And optimize by genetic search)
- Use **all candidates** from phrase tables and GDT
- We used a weighting model involving six features. Features modify the original probabilities.

    – [1.071] candidate proposed by MT cisti model;
    – [0.227] candidate proposed by MT hansard model;
    – [0.477] candidate proposed by GDT;
    – [1.464] candidate proposed by more than one source;
    – [0.853] candidate using morphology;
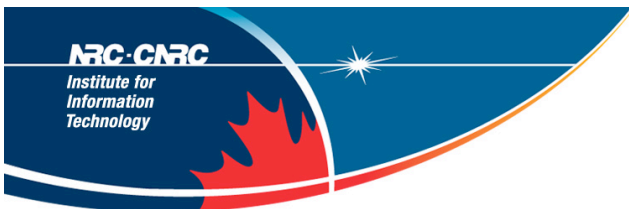    – [1.257] GDT candidate(s) in MDS;

**Experiment 4 results**

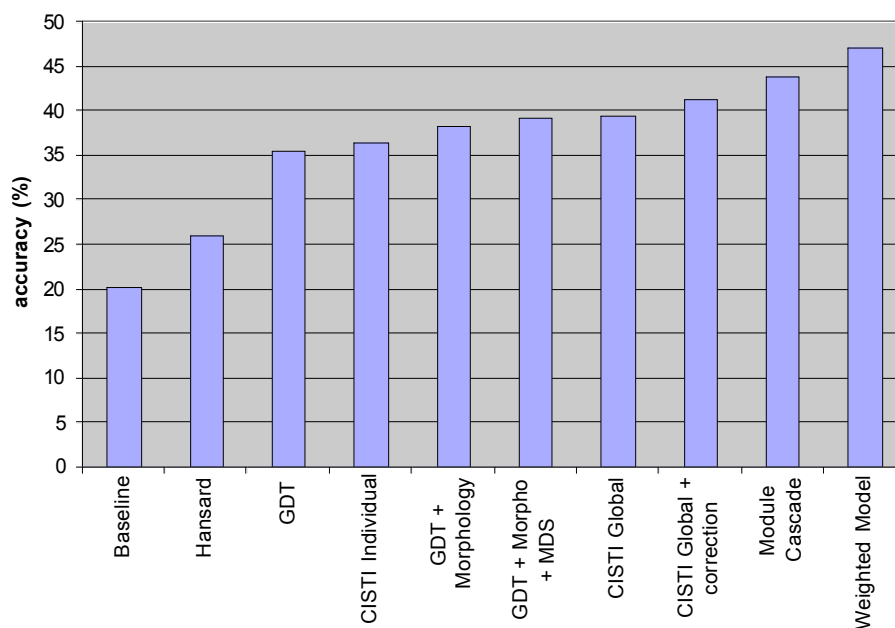| | Accuracy (%) |
|---|---|
| Module cascade | 43.83 |
| Best weighted model | 47.16 |

- Overall best solution found by genetic search and 6 features

- Statistical MT and terminological resource <u>are complementary</u>

- We show how to use the keyphrase list to « disambiguate the sense » of a keyphrase

**All experiments**

Two promising avenues:

- Include Termium, another Fr-En terminological ressource
- Follow (Jayaraman and Lavie, 2005) idea to combine module outputs and *create* candidate translations