

## HIGHLIGHTS



- Subtask of **sentence translation from summaries**, English → French

In what circumstances do granulomatous and eosinophilic gastritis occur?  
What are the etiologies of dysphagia in gastroesophageal reflux disease?

- Successful approach that makes use of **two flexible translation systems**

## DATA SOURCES

| Corpus          | Tokens (en) | weight |
|-----------------|-------------|--------|
| COPPA           | 10M         | -3     |
| EMEA            | 6M          | 26     |
| PATTR-ABSTRACTS | 20M         | 22     |
| PATTR-CLAIMS    | 32M         | 6      |
| PATTR-TITLES    | 3M          | 4      |
| UMLS            | 8M          | -7     |
| WIKIPEDIA       | 17k         | -5     |
| NEWSCOMMENTARY  | 4M          | 6      |
| EUROPARL        | 54M         | -7     |
| GIGA            | 260M        | 27     |
| all             | 397M        | 33     |

- Combining both data sources** drastically boosts performance

|         | DEVEL     | TEST      |
|---------|-----------|-----------|
| medical | 42.2± 0.1 | 39.6± 0.1 |
| WMT'13  | 43.0± 0.1 | 41.0± 0.0 |
| both    | 48.3± 0.1 | 45.4± 0.0 |

BLEU scores obtained by NCODE

## PART-OF-SPEECH TAGGING

- Medical data exhibit **different syntactic constructions and a specific vocabulary**
- We use a specific model trained on medical data

| PoS tagging | DEVEL     | TEST      |
|-------------|-----------|-----------|
| Standard    | 47.9± 0.0 | 44.8± 0.1 |
| Specialized | 48.3± 0.1 | 45.4± 0.0 |

## PROXY TEST SET

- Only a small development set is available (500 sentences)
- This makes both **system design and tuning challenging**
- We created an internal dev/test set (LMTEST) by extracting sentences from PATTR-ABSTRACTS

| DEVEL     | LMTEST    | NEWTST12  | TEST      |
|-----------|-----------|-----------|-----------|
| 48.3± 0.1 | 46.8± 0.1 | 26.2± 0.1 | 45.4± 0.0 |
| 41.8± 0.2 | 48.9± 0.1 | 18.5± 0.1 | 40.1± 0.1 |
| 39.8± 0.1 | 37.4± 0.2 | 29.0± 0.1 | 39.0± 0.3 |

## ERROR ANALYSIS

|              | <i>extra</i> |         | <i>missing</i> |         | <i>incorrect</i> |       |      |       | <i>unknown</i> |      | <b>all</b> |
|--------------|--------------|---------|----------------|---------|------------------|-------|------|-------|----------------|------|------------|
|              | word         | content | filler         | disamb. | form             | style | term | order | word           | term |            |
| SysComb      | 4            | 13      | 20             | 47      | 62               | 8     | 18   | 21    | 1              | 11   | <b>205</b> |
| OTF+VSM+Soul | 4            | 4       | 31             | 44      | 82               | 6     | 20   | 42    | 3              | 12   | <b>248</b> |

Manual error analysis following (Vilar et al., 2006) for the first 100 test sentences.

## SYSTEMS

MIRA

**NCODE** — bilingual *n*-gram approach to SMT

**OTF** — on-the-fly estimation of the parameters of a standart phrase-based model

**VSM** — Vector space model to perform domain adaptation

**Soul** — Continuous space models working on top of conventional language models (reranking); adapted language model (LM\*)

**SysComb** — Combination of both systems (reranking)

|                | DEVEL | TEST |
|----------------|-------|------|
| NCODE          | 48.5  | 45.2 |
| + Soul LM*     | 49.8  | 45.9 |
| + Soul LM*+ TM | 50.1  | 47.0 |
| OTF            | 46.6  | 42.5 |
| + VSM          | 46.9  | 42.8 |
| + Soul LM*     | 48.4  | 44.2 |
| + Soul LM*+ TM | 49.7  | 44.9 |
| SysComb        | 50.7  | 46.5 |

- NCODE outperforms OTF by 2.8 BLEU points
- Vector space model does not yield here any improvement
- Continuous space language models yield gains** of up to 2 BLEU points
- System combination **gain does not transfer to the test set**

## CONCLUSIONS

- Moderate to high-quality** translations
- Lack of an internal test** challenging
- More careful integration of medical terminology** proved necessary

