# A Graphical Interface for MT Evaluation and Error Analysis

Meritxell Gonzàlez, Jesús Giménez, Lluís Màrquez
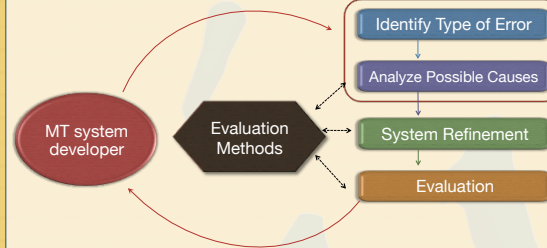TALP Research Center, Univesitat Politècnica de Catalunya

## Goals

This work presents an online graphical interface to access ASIYA, an open source toolkit to evaluate automatic translations using an heterogeneous set of metrics and meta-metrics.

1. To allow MT developers to evaluate their test beds using a large set of metric scores
2. To detect and analyze the errors of the MT systems using just their Internet browsers
3. To help developers to understand the strengths and weaknesses of the evaluation measures

## The ASIYA Toolkit

- MT system developer
- Evaluation Methods
  - Identify Type of Error
  - Analyze Possible Causes
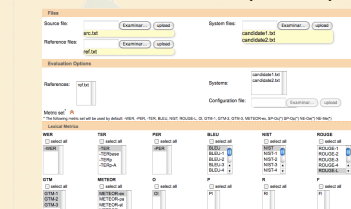  - System Refinement
  - Evaluation

- ⚙ More than 500 metric variants for MT evaluation
- ⚙ Various similarity principles: precision, recall and overlap
- ⚙ Different linguistic layers:
  - ⚙ **Lexical similarity:** based on n-gram similarity and edit distance based on word form,
    *e.g., PER, TER, WER, BLEU, NIST, GTM, METEOR*
  - ⚙ **Syntactic similarity:** based on part-of-speech tags, base phrase chunks, and dependency and constituency trees
    *e.g., SP-Overlap-POS, DP-HWCM, CP-STM*
  - ⚙ **Semantic similarity:** based on named entities, semantic roles and discourse representation
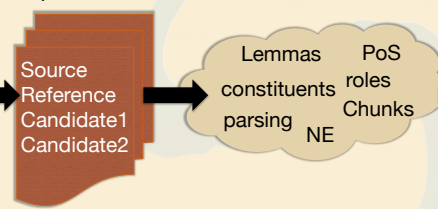    *e.g., NE-Overlap, SR-Overlap, DRS-Overlap*

## The Online Graphical Interface     (available at http://asiya.lsi.upc.edu/demo)

1. The online form allows to upload files and select the required options.

2. Asiya generates a number of metric-dependent information…

Source Reference Candidate1 Candidate2

Lemmas constituents parsing PoS roles Chunks NE

3. …that produce an interactive evaluation report.

| Systems | Segment | -WER | -TER | -PER | BLEU | NIST | ROUGE-L | GTM-1 | GTM-2 | OI |
|---|---|---|---|---|---|---|---|---|---|---|
| candidate1 | 1 | -0.3333 | -0.3333 | -0.3333 | 0.2104 | 3.0157 | | 0.6286 | 0.6667 | 0.3043 | 0.5 |
| candidate1 | 2 | | | | | | | | 1 | 1 | 1 |
| candidate1 | 3 | -0.642 | | | | | | 2357 | 0.4118 | | |
| candidate2 | 1 | -0.666 | | | | | | 354 | 0.5926 | | |
| candidate2 | 2 | -0.411 | | | | | | 444 | 0.6364 | | |
| candidate2 | 3 | -0.2143 | -0.1429 | -0.0714 | 0.6046 | 5.1749 | | 0.88 | 0.963 | 0.6242 | 0.9286 |

**candidate1:** The master ruled that the technician leaves the Barca bench by problems with some of his players. **src:** El capitán descarta que el técnico abandone el banquillo del Barça por problemas con algunos de sus jugadores.

### Metric values

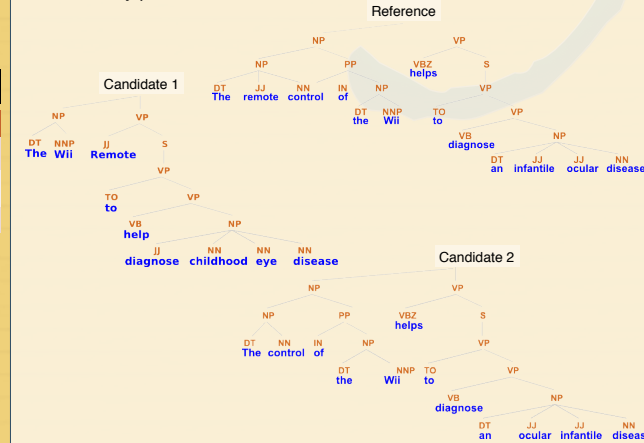apertium  babelfish  bing  candidate1  candidate2

4. The visualization of the linguistic information is useful for MT developers, such as interactive annotations and parse trees:

| Source | El | mando | de | la | Wii | ayuda | a | diagnosticar | una | enfermedad | ocular | infantil | . |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ref | The | remote | control | of | the | Wii | helps | to | diagnose | an | infantile | ocular | disease | . |
| | DT | JJ | NN | IN | DT | NNP | VBZ | TO | VB | DT | JJ | J | NN | . |
| | B-NP | I-NP | I-NP | B-PP | B-NP | I-NP | B-VP | I-VP | I-VP | B-NP | I-NP | I-NP | I-NP | O |
| | O | O | O | O | O | B-ORG | O | O | O | O | O | O | O | O |
| Cand. 1 | The | Wii | Remote | to | help | diagnose | childhood | eye | disease | . |
| | DT | NPP | NPP | TO | VB | VB | NN | NN | NN | . |
| | B-NP | I-NP | I-NP | B-VP | I-VP | I-VP | B-NP | I-NP | -NP | O |
| | O | B-ORG | I-ORG | O | O | O | O | O | O | O |
| Cand. 2 | The | control | of | the | Wii | helps | to | diagnose | an | ocular | infantile | disease | . |
| | DT | NN | IN | DT | NNP | VBZ | TO | VB | DT | JJ | JJ | NN | . |
| | B-NP | I-NP | B-PP | B-NP | I-NP | B-VP | I-VP | I-VP | B-NP | I-NP | I-NP | I-NP | O |
| | O | O | O | O | B-ORG | O | O | O | O | O | O | O | O |

### Constituency parse trees:

Reference

Candidate 1

Candidate 2

## Future Work

- ⚙ Improve usability of the interface, e.g., allow input texts, dote the parse trees with more interactions.

- ⚙ Create a database to save test sets and results.

- ⚙ Show word alignments and use them to calculate metrics under a new principle.

- ⚙ Detect and classify errors automatically.

- ⚙ Create a search engine to filter results and obtain specific good/bad examples from the test set.