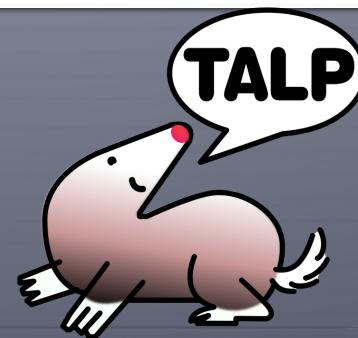




Meritxell Gonzàlez

TALP Research Center – Universitat Politècnica de Catalunya  
LREC 2014 – May 31st, Reykjavik, Iceland

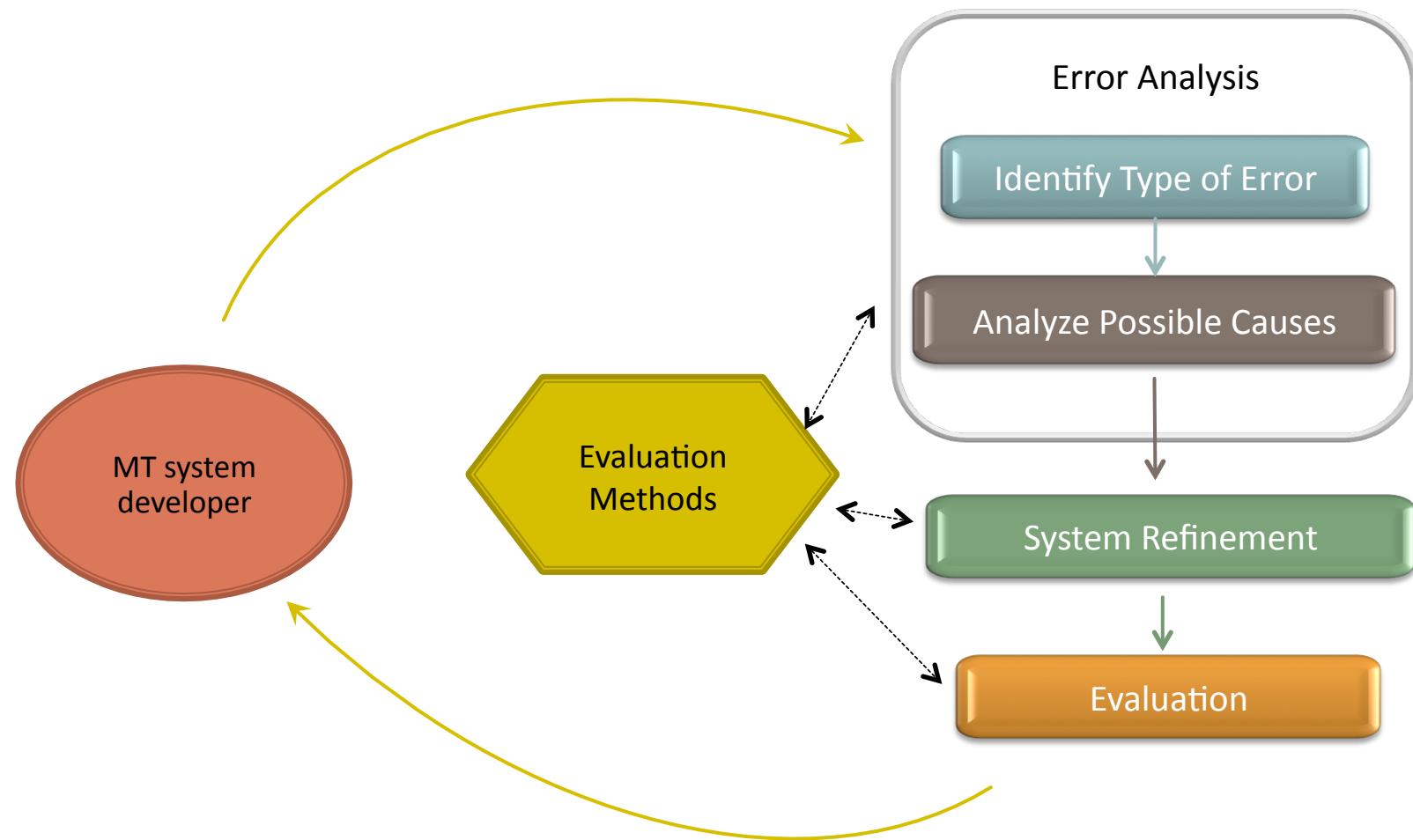
## Automatic MT Evaluation



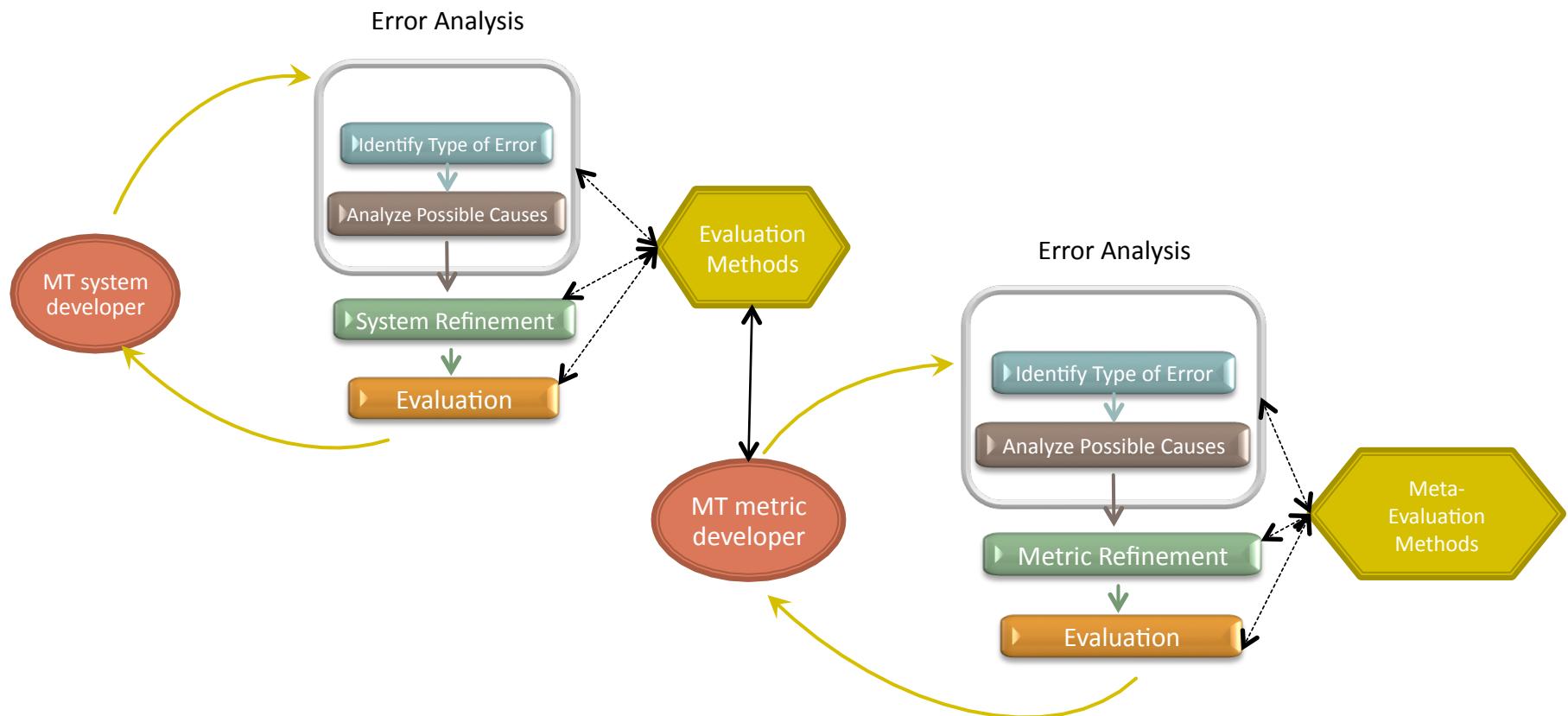
# Overview

- Introduction
- Automatic MT evaluation
- Linguistically motivated evaluation measures
- Quality estimation
- The Asiya toolkit

# MT Development cycle (1)



# MT Development cycle (2)



# Difficulties of the MT evaluation (1)

- Machine Translation is an open NLP problem
  - The correct translation is not unique.
  - The set of valid translations is not small.
  - The quality of a translation is a fuzzy concept.

## Difficulties of the MT evaluation (2)

- Quality aspects are heterogeneous:
  - **Adequacy (or Fidelity)**: Does the output convey the same meaning as the input sentence? Is part of the message lost, added, or distorted?
  - **Fluency (or Intelligibility)**: Is the output fluent? This involves both grammatical correctness and idiomatic word choices.
  - **Post-edition** effort: time required to repair the translation, number of key strokes, etc.

# Example

- El mando de la Wii ayuda a diagnosticar una enfermedad ocular infantil.
- The remote control of the Wii helps to diagnose an infantile ocular disease.
- The control of the Wii help to diagnose an ocular illness childish.
- The control of the Wii helps to diagnose an infantile ocular disease.
- The Wii remote helps diagnose a childhood eye disease.
- The Wii Remote to help diagnose childhood eye disease.
- The control of the Wii helps to diagnose an ocular infantile disease.
- The mando of the Wii helps to diagnose an infantile ocular disease.

# Manual vs. automatic evaluation

- Categorisation problem for human annotations.
  - 5-point likert scale [LDC05]
  - 4-point likert scale [TAUS13]

Adequacy		Fluency	
4	All	4	Flawless
3	Most	3	Good
2	Little	2	Disfluent
1	None	1	Incomprehensible

- Ranking problem for human annotations [Cal12]
- Regression problem for automatic metrics.

# Meta-Evaluation

- Correlation with human assessments
  - Pearson (system level)
  - Spearman
  - Kendall's tau (segment level)
- Consistency (ranking)
- AvgDelta [Cal12]

# Human Annotation Tools

- BLAST [Sty11] - annotate
- Appraise [Fed12] - rank
- DQF [Tau12] – best practices
- Costa MT Evaluation Tool [Chat13] – error classification

# Appraise [Fed12]

1/3

Sentence #402

Russian → English

Глава МИД Турции Муаммер Гюлер сначала опроверг факт строительства, затем заявил, что стена строится для обеспечения безопасности границ. "Мэр Нусайбина Айше Гекхан, член ПМД, объявила голодовку, таким образом превратив свой протест против строительства стены в смертельную борьбу", - заявили в пресс-службе партии.

— Source

Best ← Rank 1 ⚡ Rank 2 ⚡ Rank 3 ⚡ Rank 4 ⚡ Rank 5 ⚡ → Worst

The Foreign Minister of Turkey Muammer Gyuler at first disproved the construction fact, then declared that the wall is under construction for safety of borders.

— Translation 1

Best ← Rank 1 ⚡ Rank 2 ⚡ Rank 3 ⚡ Rank 4 ⚡ Rank 5 ⚡ → Worst

Turkey " s Foreign Minister Guler initially denied the fact of construction, then said that the wall is being built to ensure border security.

— Translation 2

Best ← Rank 1 ⚡ Rank 2 ⚡ Rank 3 ⚡ Rank 4 ⚡ Rank 5 ⚡ → Worst

Turkey's foreign Minister Muammer Guler first denied the fact of construction, then said that the wall is being built for the border security.

— Translation 3

Best ← Rank 1 ⚡ Rank 2 ⚡ Rank 3 ⚡ Rank 4 ⚡ Rank 5 ⚡ → Worst

The foreign minister of Turkey Muammer Guler initially denied that construction was taking place, then stated that the wall is being built to ensure border security.

— Translation 4

Best ← Rank 1 ⚡ Rank 2 ⚡ Rank 3 ⚡ Rank 4 ⚡ Rank 5 ⚡ → Worst

The head of the Ministry of Foreign Affairs of Turkey Muammer Gyuler first refuted the fact of building, then he stated that the wall is built for providing of a security of frontiers.

— Translation 5

Submit

Reset

Skip Item

# Interannotator Agreement

- Cohen's kappa coefficient [Coh60]

- WMT13 [Boj13]

- Kappa interpretation [Lan77]

- 0.0–0.2 slight
- 0.2–0.4 fair
- 0.4–0.6 moderate
- 0.6–0.8 substantial
- 0.8–1.0 almost perfect

Pair	Inter- $\kappa$	Intra- $\kappa$
CZ-EN	0.244	0.479
EN-CZ	0.168	0.290
DE-EN	0.299	0.535
EN-DE	0.267	0.498
ES-EN	0.277	0.575
EN-ES	0.206	0.492
FR-EN	0.275	0.578
EN-FR	0.231	0.495
RU-EN	0.278	0.450
EN-RU	0.243	0.513

# Benefits of Automatic Evaluation (1)

- Compared to manual evaluation, automatic measures are:
  - Cheap (vs. costly)
  - Objective (vs. subjective)
  - Reusable (vs. not-reusable)

## Benefits of Automatic Evaluation (2)

- Automatic evaluation metrics have notably accelerated the development cycle of MT systems
  - Error analysis
    - Identify and analyze weak points
  - System optimization
    - Ranking of N-best list and parameter estimation
  - System comparison
    - Phrase- or system-based combination

# Active Topic of Research

- Annual metrics competition organized by the WMT workshop series and supported by the EC
  - <http://www.statmt.org/wmt14/>
  - Both Evaluation Measure and Confidence Estimation
- Biannual OpenMT metric competition organized by NIST and supported by DARPA
  - <http://www.nist.gov/itl/iad/mig/openmt.cfm>
  - Evaluation Measures for informal data genres and speech translations
- 1st Workshop on Asian Translation, Tokyo, October 2014
  - <http://orchid.kuee.kyoto-u.ac.jp/WAT/>
  - Japanese-Chinese, test data is prepared using paragraph as a unit

# Overview

- Introduction
- **Automatic MT evaluation**
- Linguistically motivated evaluation measures
- Quality estimation
- The Asiya toolkit

# MT Automatic Evaluation (1)

- Setting:
  - Compute the similarity between a **system's output** and one or several **reference translations**.
- Challenge:
  - The similarity measure should be able to discriminate whether the two sentences convey the same meaning (**semantic equivalence**).

# MT Automatic Evaluation (2)

- Goals:
  - **Low cost**
  - **Tunable**
  - **Meaningful**
  - **Coherent**
  - **Consistent**

# First Approaches

- Lexical similarity as a measure of quality
  - Edit Distance: WER [Nieoo], PER [Til97], TER [Snoo6]
  - Precision: BLEU [Papo1], NIST [Dodo2]
  - Recall: ROUGE [Lino4a]
  - Precision/Recall: GTM [Melo3], METEOR [Bano5,Den10]

# Precision and Recall of Words (1)

- The remote control of the Wii helps to diagnose an infantile ocular disease.
- The Wii remote helps diagnose a childhood eye disease.

# Precision and Recall of Words (2)

- The remote control of the Wii helps to diagnose an infantile ocular disease .
- The Wii remote helps diagnose a childhood eye disease .

▪ **Precision:** 
$$\frac{correct}{output\_length} = \frac{7}{10} = 0.7$$

▪ **Recall:** 
$$\frac{correct}{reference\_length} = \frac{7}{14} = 0.5$$

▪ **F-measure:** 
$$\frac{precision * recall}{(precision + recall)/2} = \frac{0.35}{0.6} = 0.583$$

# Precision and Recall of Words (3)

- The remote control of the Wii helps to diagnose an infantile ocular disease .
- Wii the control of the remote to diagnose disease helps an ocular infantile .

- **Precision:** 
$$\frac{correct}{output\_length} = \frac{14}{14} = 1.00$$

No Penalty for  
reordering!

- **Recall:** 
$$\frac{correct}{reference\_length} = \frac{14}{14} = 1.00$$

- **F-measure:** 
$$\frac{precision * recall}{(precision + recall)/2} = \frac{1.00}{1.00} = 1.00$$

# IBM BLEU (1)

- “The main idea is to use a **weighted average of variable length phrase matches** against the reference translations. This view gives rise to a family of metrics using various weighting schemes. We have selected a promising baseline metric from this family.” [Papoi]

# IBM BLEU (2)

- **Modified N-gram** precision between machine translation output and reference translation.
  - Usually with  $n$ -grams of size 1 to 4
- Modified  $n$ -gram precision on the entire corpus

$$P_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n-gram \in C} count_{clip}(ngram)}{\sum_{C' \in \{Candidates\}} \sum_{n-gram' \in C'} count_{clip}(ngram')}$$

- Brevity penalty for too short translations.

$$BP = \begin{cases} 1 & \text{if } c > r \\ e^{1-\frac{r}{c}} & \text{otherwise} \end{cases}$$

- Typically computed over the entire corpus, not single sentences.

# IBM BLEU (3)

- The remote control of the Wii helps to diagnose an infantile ocular disease .
- The control of the Wii helps to diagnose an ocular infantile disease .

■  $w_n = 1/4$

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^N w_n \log P_n\right)$$

	$p_n$	$BP * p_n$	$w_n \log P_n$
1-gram precision	$13/13 = 1.0$	0.926	0
2-gram precision	$8/12 = 0.667$	0.617	-0.405
3-gram precision	$6/11 = 0.545$	0.505	-0.606
4-gram precision	$5/10 = 0.5$	0.463	-0.693
Brevity penalty	0.926		
$P_n$	Same, only one sentence		
<b>BLEU score</b>	<b>0.6046</b>		

# Problems of lexical similarities (1)

- The **reliability** of lexical metrics depends strongly on **the heterogeneity/representativity** of reference translations.
- Actually, human translations tend to score low on BLEU.
- Underlying Cause:
  - Lexical similarity is neither a **sufficient** nor a **necessary** condition so that two sentences convey the same meaning.

## Problems of lexical similarities (2)

- Statistical MT systems heavily rely on the training data.
- Testsets tend to be similar (domain, register, sublanguage) to training materials.
- N-gram based metrics favour MT systems which closely replicate the lexical realization of the references.
- Statistical MT systems tend to share the reference sublanguage and be favoured by  $n$ -gram-based measures.

# Overview

- Introduction
- Automatic MT evaluation
- **Linguistically motivated evaluation measures**
- Quality estimation
- The Asiya toolkit

# Linguistically motivated measures (1)

- Extending Lexical Similarity Measures to increase robustness [Gim09]
  - Lexical variants:
    - Morphological information (i.e., stemming )  
**ROUGE** and **METEOR**
    - Synonymy lookup : **METEOR** (based on WordNet)
  - Paraphrasing support:
    - Extended versions of **METEOR**, **TER**
  - Equivalent reference translation graph:
    - **HyTER** [Dre12]

# METEOR (1) [Bano5]

- Parameterized harmonic mean of word P and R

$$F_{mean} = \frac{P \cdot R}{\alpha \cdot P + (1 - \alpha) \cdot R}$$

- Matching algorithm

- Exact matching
  - Partial credit for matching stems
  - Partial credit for matching synonyms

$$Pen = \gamma \cdot \left(\frac{ch}{m}\right)^\beta$$

- N-gram **penalty** based on the number of **chunks** with longer length of adjacent words **matched** in both strings
- Final score:

$$METEOR = (1 - Pen) \cdot F_{mean}$$

# METEOR (2) [Den10]

- Extensions
  - METEOR-NEXT
    - Weighted matches depending on the type
    - Phrase-level matches
    - new matching algorithm accounting for start-positions distance
  - Paraphrasing
    - Paraphrase tables from parallel corpora
    - Used by the paraphrase matcher
  - $\delta$  parameter: content vs. function words discrimination

# More linguistically-motivated measures

- Features capturing **syntactic** and **semantic** information
  - Shallow parsing, constituency and dependency parsing, named entities, semantic roles, textual entailment, discourse representation, error categories, ...
- Some linguistically-motivated measures:
  - IQmt [Gim09] – syntactic and semantics
  - MaxSim [Chao8] - syntactic
  - RTE [Pado9] – textual entailment
  - VERTa [Com14] – syntactic and semantics

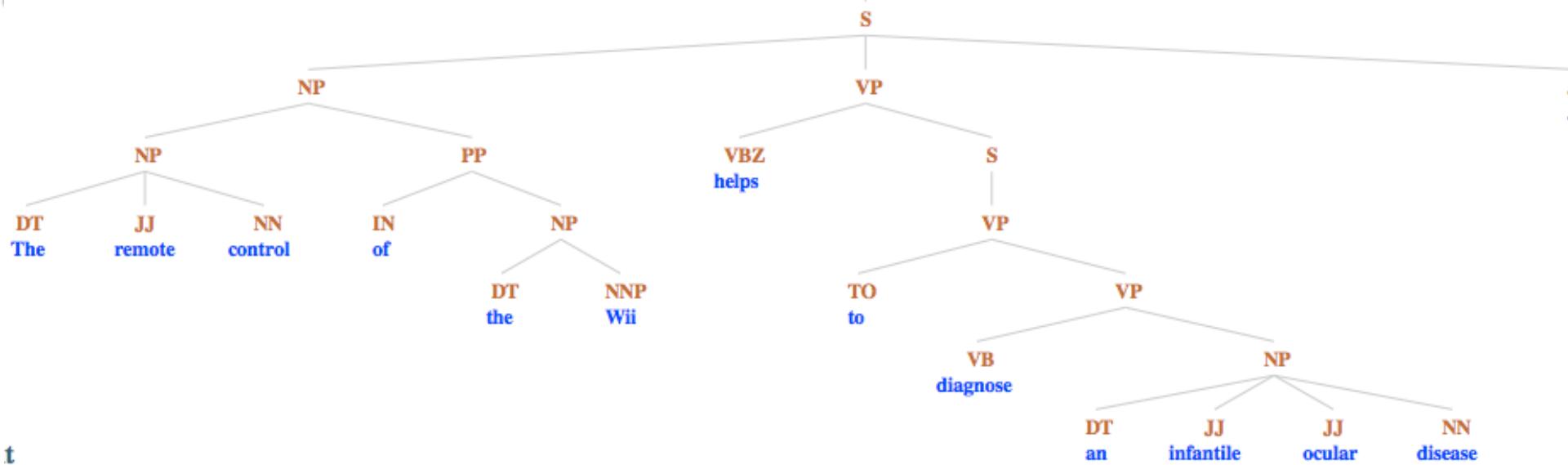
# Example 1: Structural Similarity (1)

- Rather than comparing sentences at lexical level:  
**Compare the linguistic structures** and the words within them [Gim10]
- Compare different linguistic-level elements
  - Words, lemmas, POS, Chunks
  - Parsing Trees
  - Named entities and semantic roles
  - Discourse representation (logical forms)

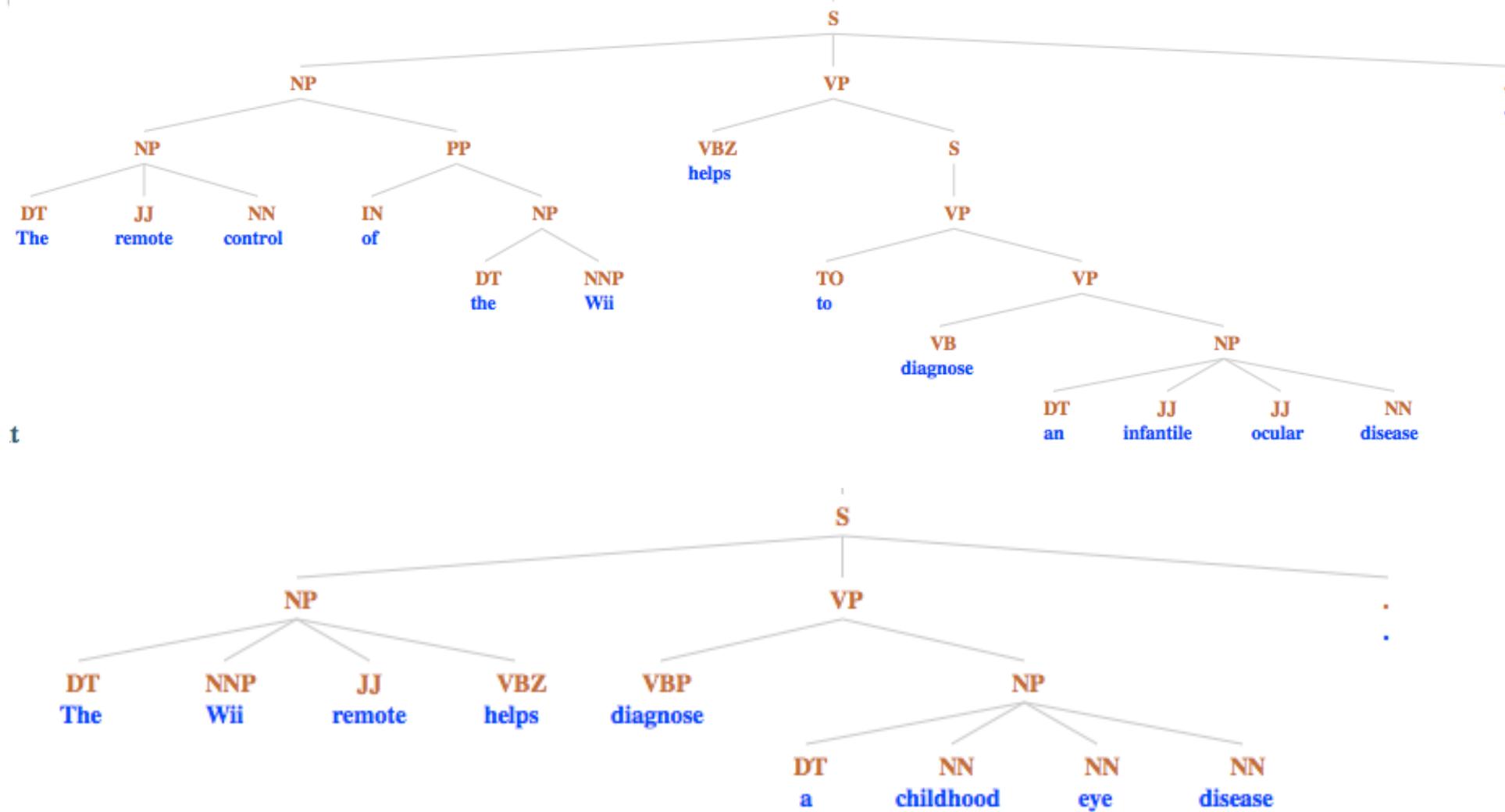
# Example 1: Structural Similarity (2)

- The remote control of the Wii helps to diagnose an infantile ocular disease.
- The Wii remote **helps diagnose** a childhood eye disease.

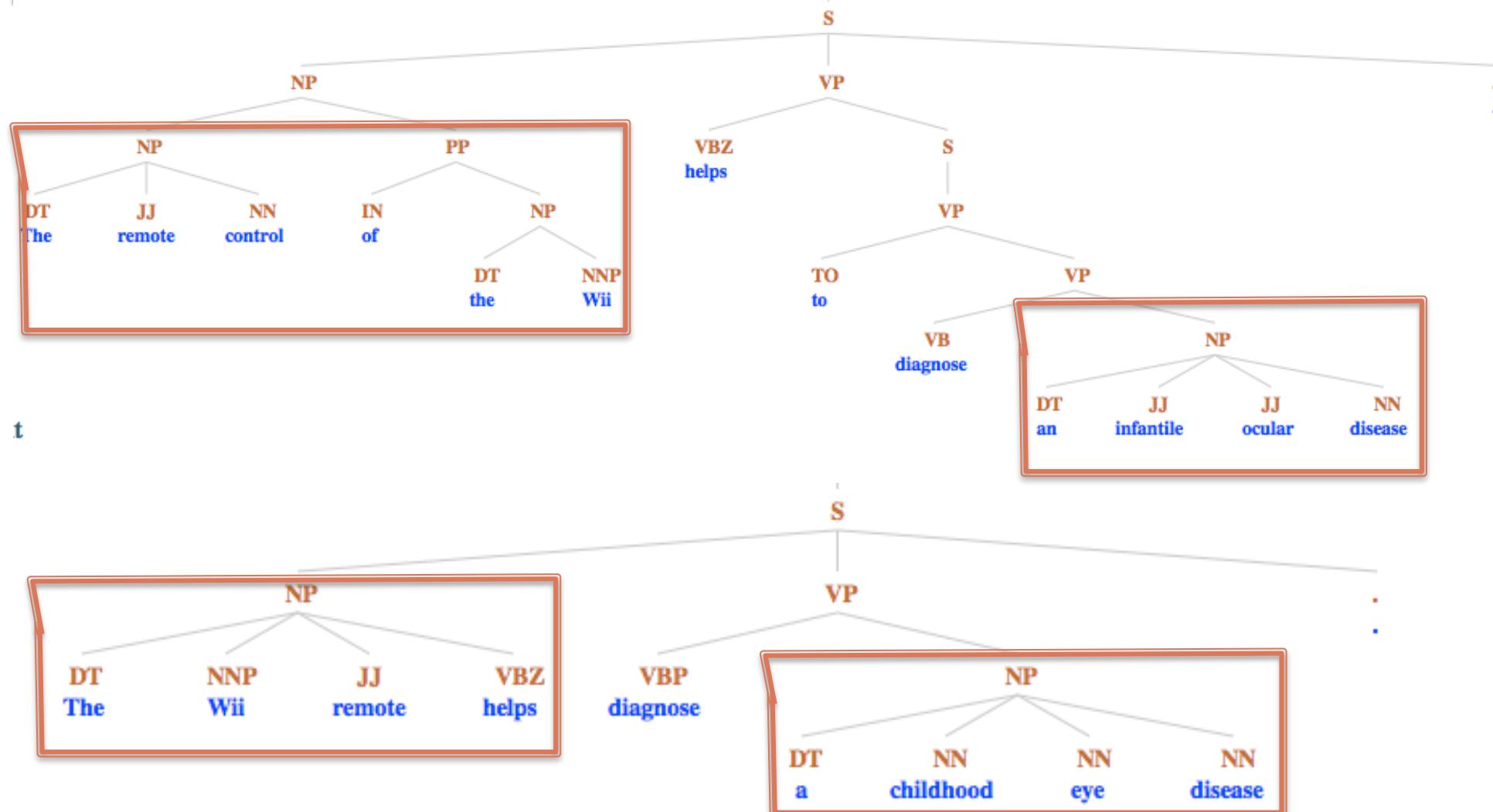
# Example 1: Structural Similarity (3)



## Example 1: Structural Similarity (4)



# Example 1: Structural Similarity (4)



# Measuring structural similarity (1)

- Linguistic Element (LE): abstract reference to any possible type of linguistic unit, structure, or relationship among them.
  - For instance: POS tags, word lemmas, NPs, semantic roles, dependency relations, etc.
- A sentence can be seen as a bag (or a sequence) of LEs of a certain type

# Measuring structural similarity (2)

- **OVERLAP** [Gim07]: generic similarity measure among linguistic elements inspired by the Jaccard coefficient [Jac1901]
- SEMPOS [Maco8] is a MT evaluation measure that considers several **overlapping variations**
- **MATCHING** is a more strict variant [Gim10]
  - All items inside an element are considered the same unit.
  - Computes the proportion of fully translated LEs according to their types.

# Overlap (1)

$$O(t) = \frac{\sum_{i \in (items_t(cand) \cap items_t(ref))} count_{cand}(i,t)}{\sum_{i \in (items_t(cand) \cup items_t(ref))} \max(count_{cand}(i,t), count_{ref}(i,t))}$$

$$O(*) = \frac{\sum_{t \in T} \sum_{i \in (items_t(cand) \cap items_t(ref))} count_{cand}(i,t)}{\sum_{t \in T} \sum_{i \in (items_t(cand) \cup items_t(ref))} \max(count_{cand}(i,t), count_{ref}(i,t))}$$

# Overlap (2)

- The remote control of the Wii helps to diagnose an infantile ocular disease.
- The Wii remote helps diagnose a childhood eye disease.
- Overlap:
  - Intersection: 13
  - Union: 25
  - $O_l = 13/25 = 0.52$

Words	Reference	Candidate
the	2	1
remote	1	1
control	1	
of	1	
wii	1	1
helps	1	1
to	1	
diagnose	1	1
an	1	
a		1
infantile	1	
childhood		1
ocular	1	
eye		1
disease	1	1
.	1	1

# Overlap (3)

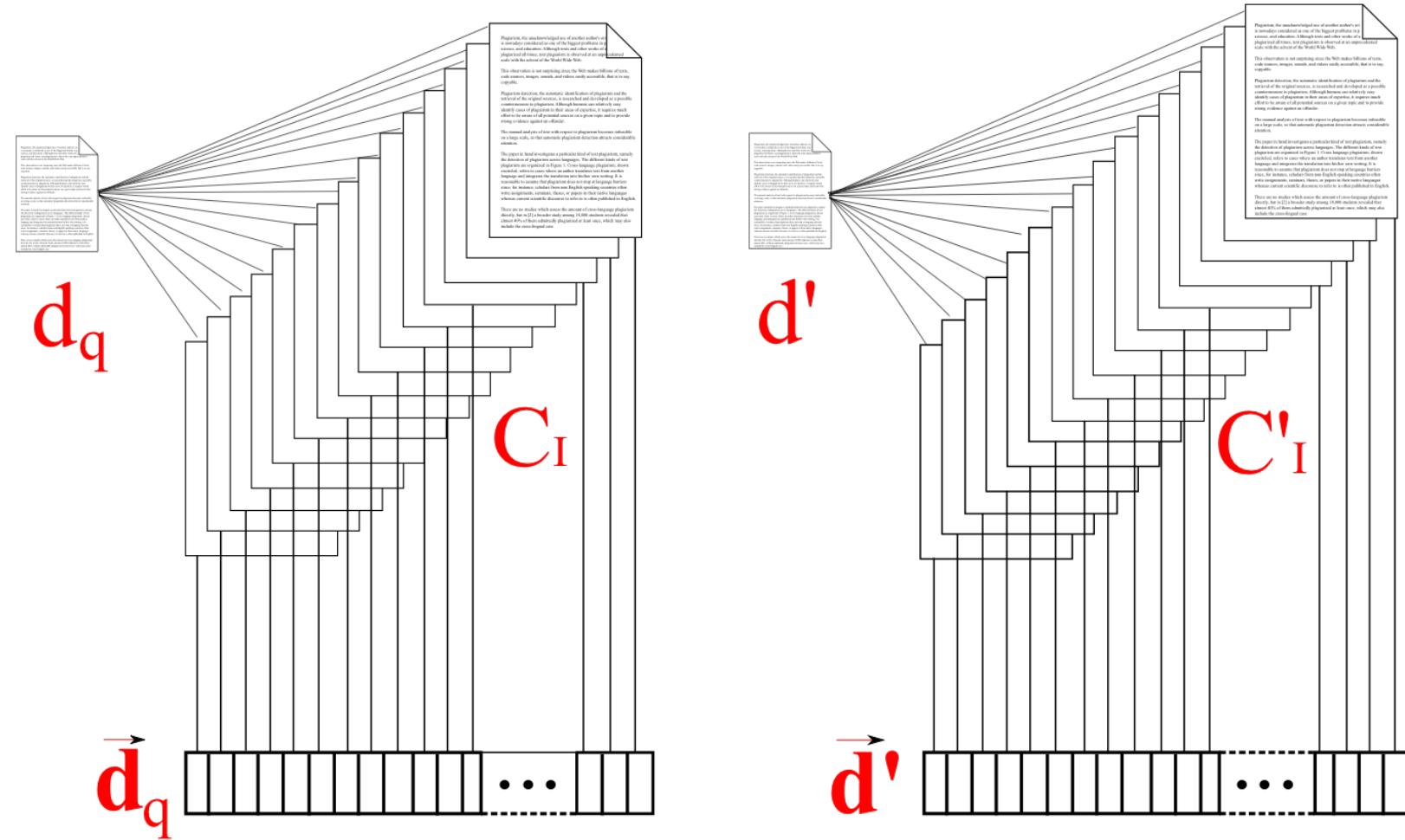
- The remote control of the Wii helps to diagnose an infantile ocular disease.
- DT JJ NN IN DT NNP VBZ TO VB DT JJ JJ NN .
- The Wii remote helps diagnose a childhood eye disease.
- DT NNP JJ VBZ VB DT NN NN NN .
- Overlap:
  - Intersection: 9
  - Union: 15
  - $O_l = 9/15 = 0.6$

Words	Reference	Candidate
DT	3	2
JJ	3	1
NN	2	3
IN	1	
NNP	1	1
VBZ	1	1
TO	1	
VB	1	1
.	1	1

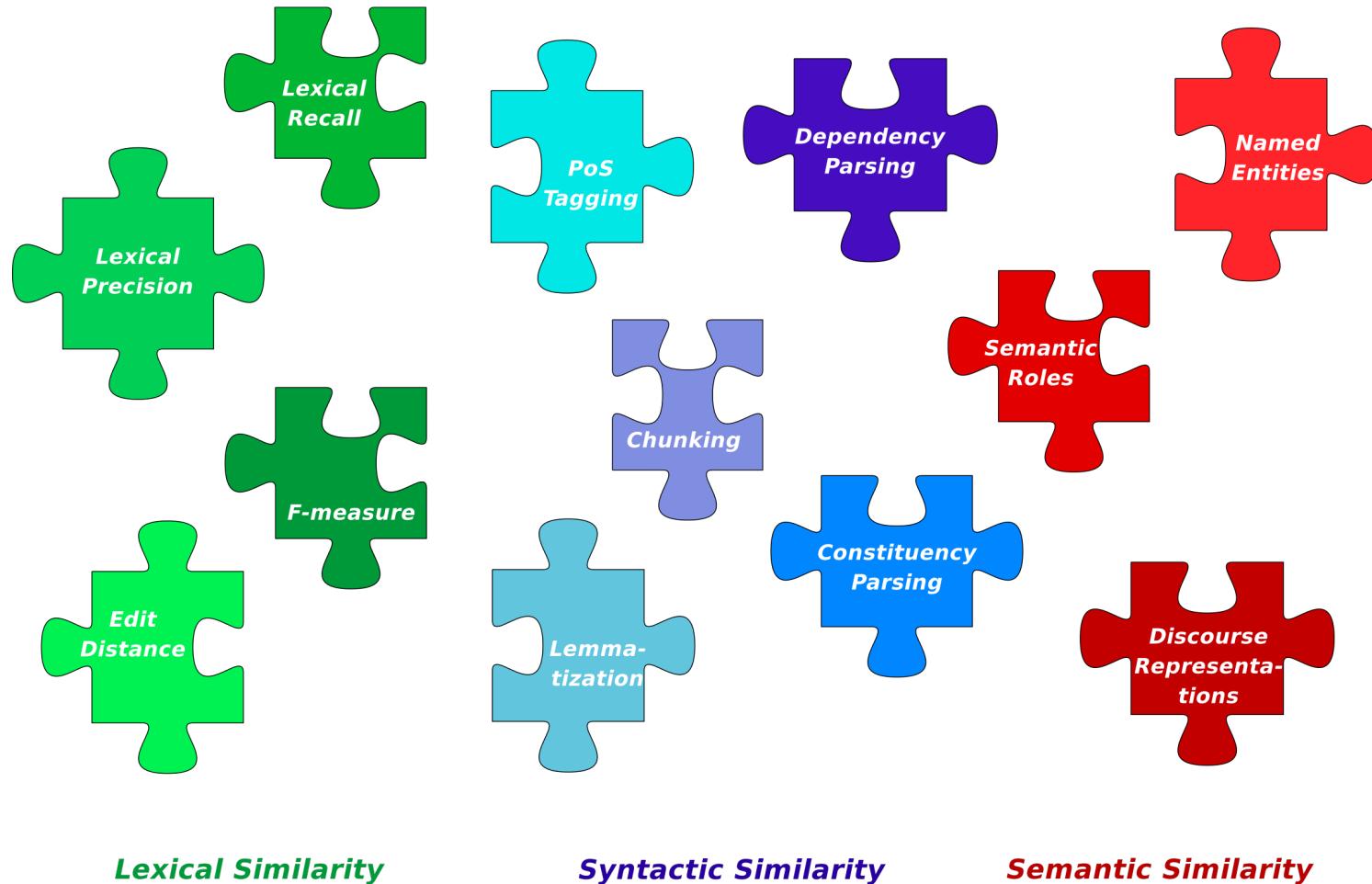
# More linguistically-motivated measure (2)

- CL-Explicit Semantic Analysis
  - CL-ESA requires a significant comparable corpus  $C_I$
  - $d_q \in L$  ( $d' \in L'$ ) is represented as a vector of relations to the index collection  $C_I$  ( $C_I'$ )
  - Monolingual similarities are computed over the VSM (e.g., the cosine of the vocabulary) [Poto8]

# Example 2: Semantic Analysis (2)



# Towards Heterogeneous MT Evaluation

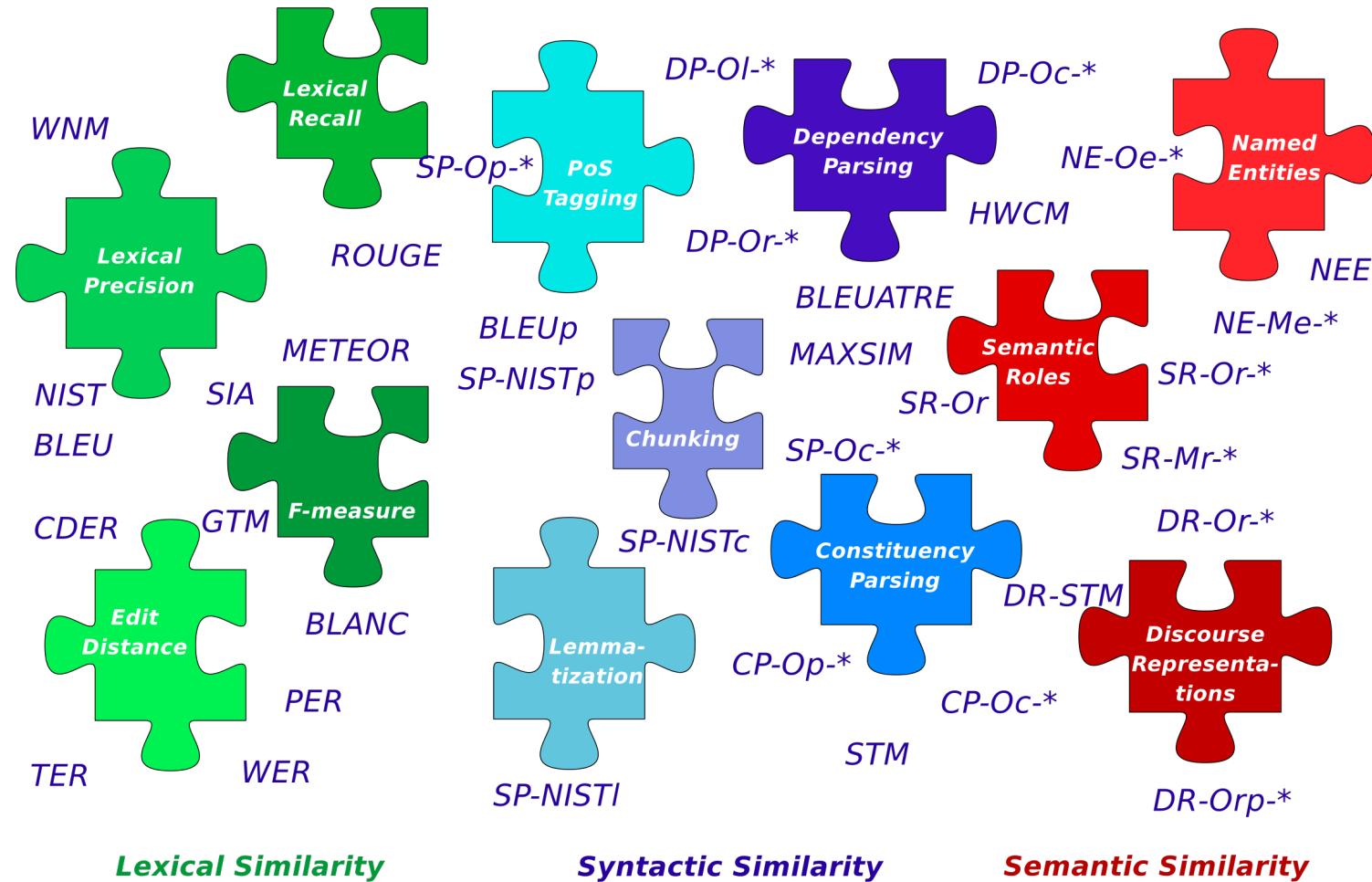


**Lexical Similarity**

**Syntactic Similarity**

**Semantic Similarity**

# Towards Heterogeneous MT Evaluation



# Metric Combination

- Different measures capture different aspects of similarity
- Simple Approach
  - ULC: Uniformly-averaged linear combination of measures
- But, which ones?
  - Simple hill climbing approach to find the best subset of measures M on a development corpus
    - $M = \{ROUGE_w, METEOR, DP-HWC_r, DP-O_c(*), DP-O_l(*), DP-O_r(*), CP-STM_4, SR-O_r(*), SR-O_{rv}, DR-O_{rp}(*)\}$

# Estimate Models

- The goal is to **combine the scores** conferred by different evaluation measures **into a single measure of quality** such that their relative contribution is adjusted on the basis of human feedback (i.e. from human assessments).
- Examples:
  - AMBER [Che12] – downhill simplex
  - SIMBLEU (ROSE) [Son11] – SVM
  - SPEDE [Wan12] – pFSM for regression
  - TERRORCAT [Fis12] - SVM on error categories

# Overview

- Introduction
- Automatic MT evaluation
- Linguistically motivated evaluation measures
- **Quality estimation**
- The Asiya toolkit

# Quality Estimation (1)

- Setting:
  - Quality assessment without reference translations
- Information available:
  - Source sentence, candidate translation(s) and, possibly, MT system information
- Motivation:
  - System ranking (system selection)
  - Hypotheses re-ranking (parameter optimization)
  - Feedback filtering (especially end-users)
  - Post-edition effort (industry pricing)

# Quality Estimation (2)

- Relevant Work:
  - Johns Hopkins University Summer Workshop, 2003.  
“Confidence Estimation for Machine Translation”. [Bla04]
- Recent work:
  - (Specia et al., 2009;2010), (Soricut and Echihabi, 2010),  
(Giménez and Specia 2010), (Pighin et al., 2011), (Avramidis, 2012)
- WMT shared task on Quality Estimation
  - [Cal12] WMT12 – 11 participants
  - [Boj13] WMT13 – 14 participants
  - (3d edition at WMT 2014)

# Quality Estimation Features (1)

- System-dependent
  - internal system probabilities/scores (automatic score)
  - features over **n-best translation hypotheses**
    - language modelling
    - candidates rank
    - score ratio
    - average candidates length
    - length ratio
    - ...

# Quality Estimation Features (2)

- System-independent
  - source (translation difficulty)
    - Source sentence length
    - Ambiguity → dictionary/alignment/WordNet-based
      - e.g, number of candidate translations per word or phrase
  - target (fluency)
    - OOV
    - Language models: perplexity, log probability

# Quality Estimation Features (3)

- System-independent
  - source-target (adequacy)
    - length factor
    - punctuation and symbols concurrency
    - candidate matching → dictionary-/alignment-based
    - character  $n$ -grams [McNo4]
    - pseudo-cognates [Sim92]
    - word alignments [Gon14]

# QE Challenges

- QE is a difficult task
  - Few corpus available
  - Too domain-oriented

DE-EN, task 1.2, QE2013	Kendall's $\tau$ ties ignored
DFKI-logregFss33	0.31
DFKI-logregFss24	0.28
UPC-1	0.27
UPC-2	0.24
DCU-CCG	0.18
CNGL-SVRPLSF1	0.17
CNGL-SVRF1	0.17
Baseline	0.08
Oracle BLEU	0.22
Oracula METEOR-ex	0.20

# Overview

- Introduction
- Automatic MT evaluation
- Linguistically motivated evaluation measures
- Quality estimation
- **The Asiya toolkit**

# Asiya

- **Asiya** is an Open Toolkit for Automatic Machine Translation and (Meta-)Evaluation  
<http://asiya.lsi.upc.edu>
- Asiya provides:
  - Automatic evaluation measures using several linguistic layers for a variety of languages
  - Quality Estimation measures
  - Meta-evaluation metrics
  - Learning schemes

# Asiya

- Languages:
  - English, Spanish, Catalan
  - Czech, French, German and Russian with limited resources
- Similarity principles
  - Precision, Recall, Overlap, Matching, ...
- Linguistic layers:
  - Lexical, Syntactic, Semantic, Discourse

# Metrics and Meta-metrics

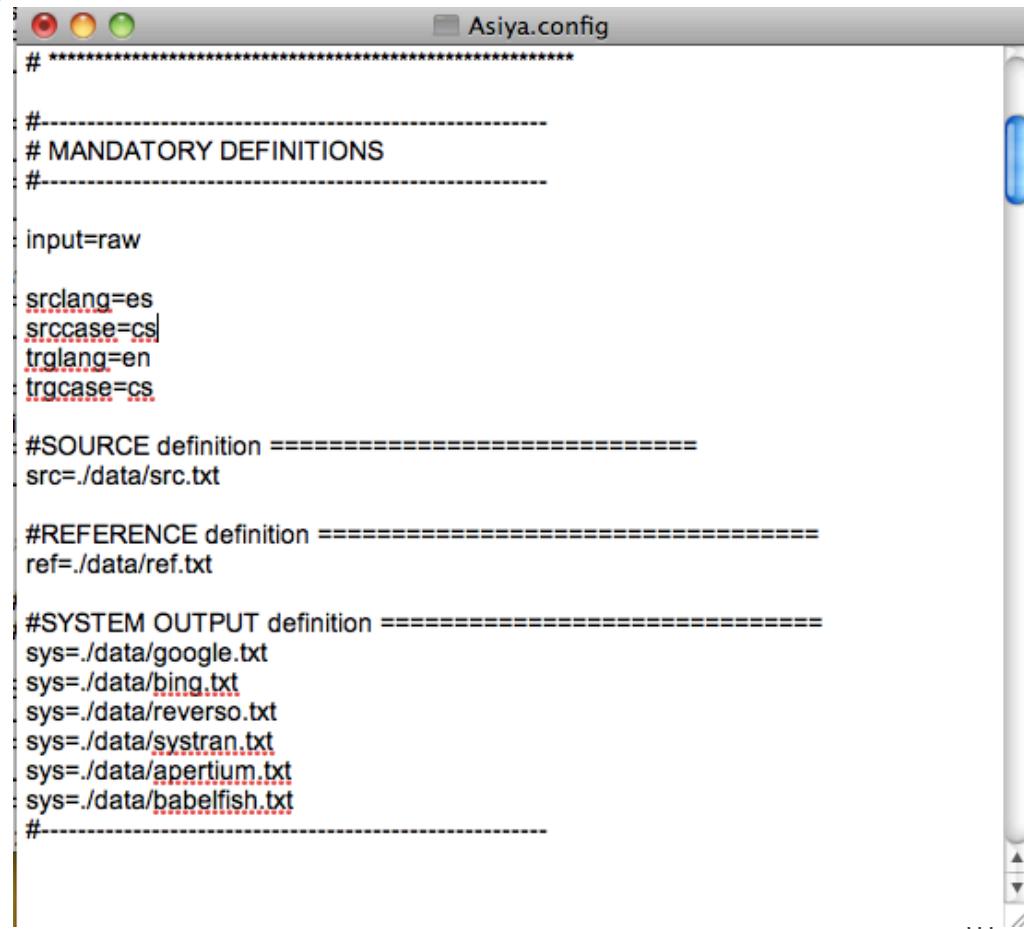
813 metrics are available for language 'es' -> 'en'

# Asiya how to (1)

- Asiya operates over testbeds (or test suites).
  - a testbed is a collection of test cases:
    - Source segment
    - Candidate translation(s)
    - Reference translation(s)

# Asiya how to (2)

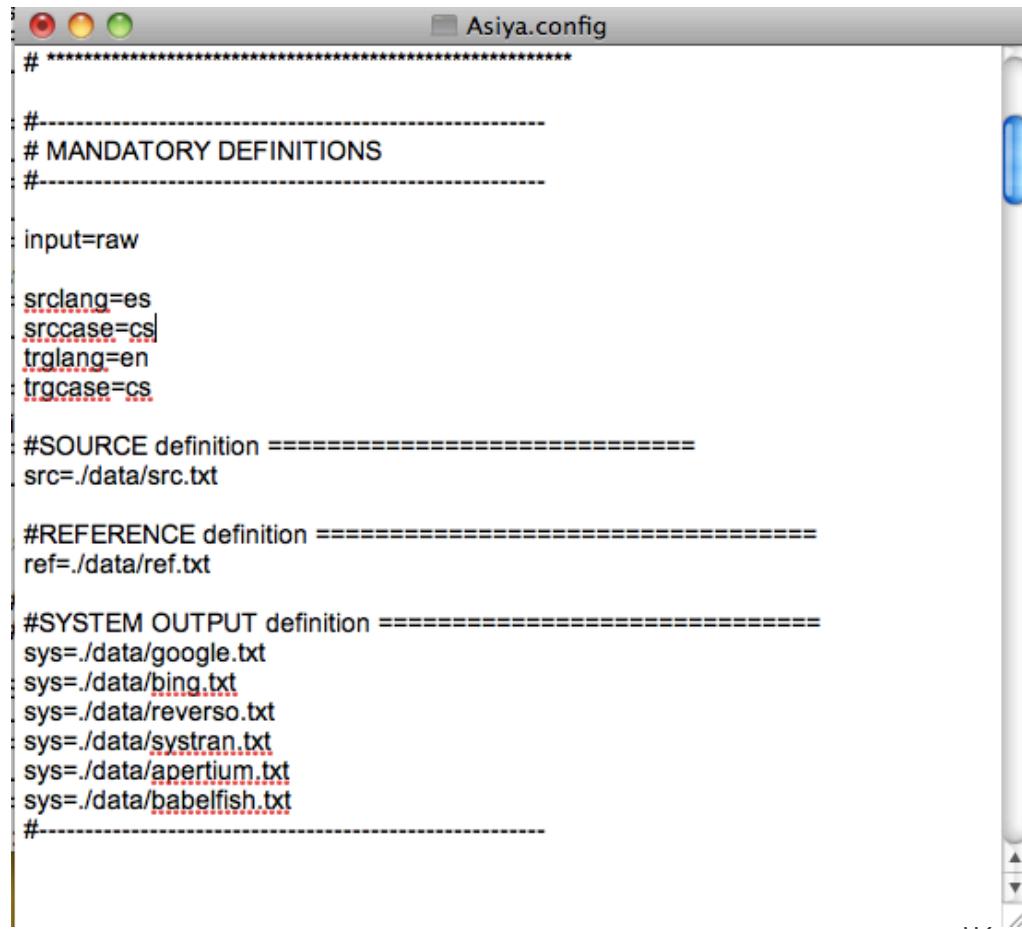
- Asiya.pl Asiya.config
- Asiya.config:



```
# *****
#
# MANDATORY DEFINITIONS
#
input=raw
srclang=es
srcase=cs
trglang=en
trgcase=cs
#
#SOURCE definition =====
src=./data/src.txt
#
#REFERENCE definition =====
ref=./data/ref.txt
#
#SYSTEM OUTPUT definition =====
sys=./data/google.txt
sys=./data/bing.txt
sys=./data/reverso.txt
sys=./data/systran.txt
sys=./data/apertium.txt
sys=./data/babelfish.txt
#
```

# Asiya how to (3)

- General Options
  - Input format
    - Raw
    - Nist
  - Language pair
    - Srclang
    - Trglang
  - Predefined sets of metrics, systems and references



```
# *****
#
# MANDATORY DEFINITIONS
#
input=raw

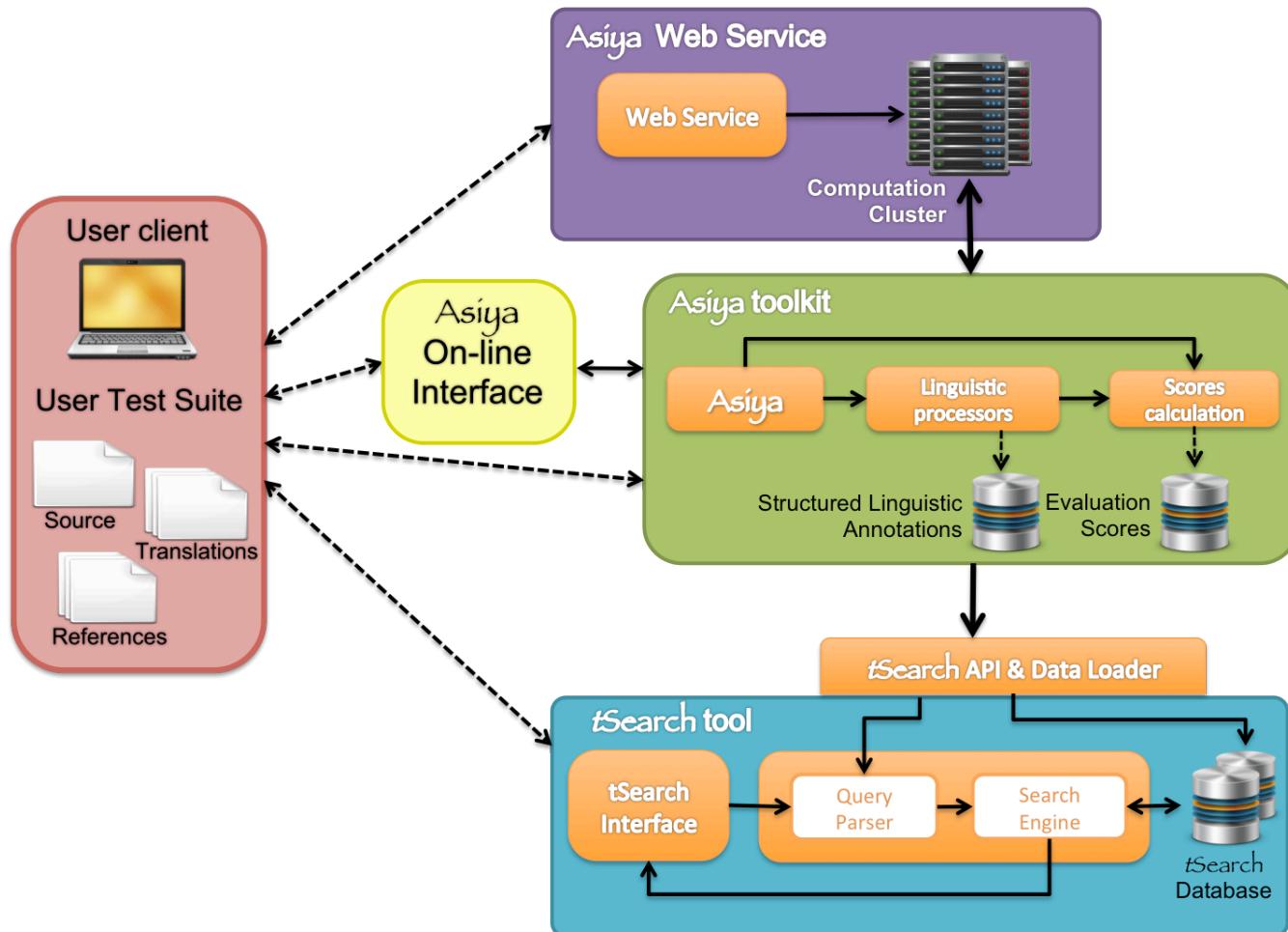
srclang=es
srcase=cs
trglang=en
trgcase=cs

#SOURCE definition =====
src=./data/src.txt

#REFERENCE definition =====
ref=./data/ref.txt

#SYSTEM OUTPUT definition =====
sys=./data/google.txt
sys=./data/bing.txt
sys=./data/reverso.txt
sys=./data/systran.txt
sys=./data/apertium.txt
sys=./data/babelfish.txt
# -----
```

# Asiya Interfaces



# Hands-on

<http://asiya.lsi.upc.edu>

- Choose the languages
- Write some sentences or upload a SMALL file. Try to introduce several errors:
  - lexical disagreement, missing prepositions,
- Use some linguistic measures in addition to the lexical ones:
  - *BLEU*, *NIST*, *ROUGE<sub>w</sub>*, *METEOR-pa*
  - *SP-Op*(\*), *DP-HWC<sub>r</sub>*, *DP-O<sub>r</sub>*(\*), *CP-STM<sub>4</sub>*
- Run it and look how the segment level scores identify the errors in each sentence
- Look at the parse trees
- Use the tSearch interface to find interesting sentences according to the scores and the parse trees

# References

# References

- [LDC05] NIST Multimodal Information Group. NIST 2005 Open Machine Translation (OpenMT)
- [TAUS13] TAUS. Quality Evaluation using Adequacy and/or Fluency Approaches.  
<https://evaluation.taus.net/resources/adequacy-fluency-guidelines>
- [Tau12] *Nora Aranberri and Rahzeb Choudhury. Advancing Best Practices in Machine Translation Quality Evaluation.* TAUS 2012.
- [Sty11] Sara Stymne. BLAST: A Tool for Error Analysis of Machine Translation Output. 2011.
- [Fed12] Christian Federmann. Appraise: An Open-Source Toolkit for Manual Phrase-Based Evaluation of Translations. LREC 2012.
- [Cha13] Konstantinos Chatzitheodorou, Stamatis Chatzistamatis. COSTA MT Evaluation Tool: An Open Toolkit for Human Machine Translation Evaluation. The Prague Bulletin of Mathematical Linguistics No. 100, 2013, pp. 83–89.

# References

- [Coh60] Cohen, Jacob (1960). "A coefficient of agreement for nominal scales". *Educational and Psychological Measurement* 20 (1): 37–46.
- [Boj13] Ondřej Bojar, Christian Buck, Chris Callison-Burch, Christian Federmann, Barry Haddow, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut and Lucia Specia. *Findings of the 2013 Workshop on Statistical Machine Translation*. WMT13, 2013.
- [Lan77] Landis, J.R.; Koch, G.G. (1977). The measurement of observer agreement for categorical data. *Biometrics* 33 (1): 159–174.
- [Cal12] Chris Callison-Burch, Philipp Koehn, Christof Monz, Matt Post, Radu Soricut and Lucia Specia. *Findings of the 2012 Workshop on Statistical Machine Translation*. WMT13

# References

- [Nie00] Nie en, S., Och, F. J., Leusch, G., & Ney, H. An Evaluation Tool for Machine Translation: Fast Evaluation for MT Research. Proceedings of the 2nd International Conference on Language Resources and Evaluation. LREC, 2000.
- [Til97] Tillmann, C., Vogel, S., Ney, H., Zubiaga, A., & Sawaf, H. Accelerated DP based Search for Statistical Translation. Proceedings of European Conference on Speech Communication and Technology. 1997.
- [Snoo06] Snover, M., Dorr, B., Schwartz, R., Micciulla, L., & Makhoul, J. A Study of Translation Edit Rate with Targeted Human Annotation. Proceedings of the 7th Conference of the Association for Machine Translation in the Americas (AMTA) pp. 223-231, 2006.
- [Papo01] Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation, RC22176 (Technical Report). IBM T.J. Watson Research Center. 2001.

# References

- [Dodo2] Doddington, G. (2002). Automatic Evaluation of Machine Translation Quality Using N-gram Co-Occurrence Statistics. Proceedings of the 2nd International Conference on Human Language Technology(pp. 138-145). 2002.
- [Lino4a] Lin, C.-Y., & Och, F. J. Automatic Evaluation of Machine Translation Quality Using Longest Common Subsequence and Skip-Bigram Statics. Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL). 2004.
- [Melo3] Melamed, I. D., Green, R., & Turian, J. P. Precision and Recall of Machine Translation. Proceedings of the Joint Conference on Human Language Technology and the North American Chapter of the Association for Computational Linguistics (HLT-NAACL). 2003.
- [Bano5] Satanjeev Banerjee and Alon Lavie, "METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments", *Proceedings of the ACL 2005 Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization*, 2005.

# References

- [Den10] Michael Denkowski and Alon Lavie, "METEOR-NEXT and the METEOR Paraphrase Tables: Improved Evaluation Support For Five Target Languages", *Proceedings of the ACL 2010 Joint Workshop on Statistical Machine Translation and Metrics MATR*, 2010.
- [Dre12] Dreyer, Markus and Marcu, Daniel. HyTER: Meaning-equivalent Semantics for Translation Evaluation. In Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. NAACL HLT, 2012
- [Gim09] Jesús Giménez and Lluís Màrquez. On the Robustness of Syntactic and Semantic Features for Automatic MT Evaluation. In *Proceedings of the Fourth Workshop on Statistical Machine Translation, WMT 2009*, Athens, Greece, 2009.
- [Chao08] Chan YS, Ng HT. MAXSIM: a maximum similarity metric for machine translation evaluation. In Proceedings of ACL-08/HLT, pp 55-62. 2008.

# References

- [Pad09] S. Pado and M. Galley and D. Jurafsky and C. Manning. Robust Machine Translation Evaluation with Entailment Features. Proceedings of ACL, 2009.
- [Com14] Elisabet Comelles, Jordi Atserias, Victoria Arranz, Irene Castellon and Jordi Sesé. VERTa: Facing a Multilingual Experience of a Linguistic MT Evaluation. LREC, 2014.
- [Gim10] Jesús Giménez, Lluís Màrquez. Linguistic Measures for Automatic Machine Translation Evaluation. *Machine Translation*, Springer Netherlands, 2010.
- [Gim07] Jesús Giménez and Lluís Màrquez. Linguistic Features for Automatic Evaluation of Heterogeneous MT Systems. In *Proceedings of WMT 2007 (ACL'07)*, June 2007.
- [Jac01] Jaccard, Paul. Étude comparative de la distribution florale dans une portion des Alpes et des Jura", *Bulletin de la Société Vaudoise des Sciences Naturelles* 37: 547–579. 1901.

# References

- [Maco8] Machácek, Matous and Bojar, Ondrej. Approximating a Deep-syntactic Metric for MT Evaluation and Tuning. Proceedings of the Sixth Workshop on Statistical Machine Translation. WMT '11, 2011.
- [Poto8] Martin Potthast, Benno Stein, Maik Anderka. A Wikipedia-Based Multilingual Retrieval Model. In Advances in Information Retrieval, Vol. 4956, pp. 522-530. 2008.
- [Che12] Boxing Chen, Roland Kuhn, and George Foster. Improving amber, an MT evaluation metric. In Proceedings of the Seventh Workshop on Statistical Machine Translation. ACL 2012.
- [Son11] Xingyi Song and Trevor Cohn. Regression and ranking based optimisation for sentence level MT evaluation. In Proceedings of the Sixth Workshop on Statistical Machine Translation . 2011.
- [Wan12] Mengqiu Wang and Christopher Manning. SPEDE: Probabilistic edit distance metrics for MT evaluation. In Proceedings of the Seventh Workshop on Statistical Machine Translation. ACL 2012.

# References

- [Fis12] Mark Fishel, Rico Sennrich, Maja Popovic, and Ondrej Bojar. 2012. TerrorCat: a translation error categorization-based MT quality metric. In Proceedings of the Seventh Workshop on Statistical Machine Translation. ACL 2012.
- [Bla04] Blatz, John and Fitzgerald, Erin and Foster, George and Gandrabur, Simona and Goutte, Cyril and Kulesza, Alex and Sanchis, Alberto and Ueffing, Nicola. Confidence Estimation for Machine Translation. Proceedings of the 20th International Conference on Computational Linguistics, COLING 2004.
- Giménez and Specia, 2010. Lucia Specia and Jesús Giménez. Combining Confidence Estimation and Reference-based Metrics for Segment-level MT Evaluation. In *Ninth Conference of the Association for Machine Translation in the Americas*, AMTA 2010.
- [Specia et al., 2010] Lucia Specia, Dhwaj Raj, and Marco Turchi. Machine translation evaluation versus quality estimation. *Machine Translation*, 24(1): 39–50, Springer Netherlands, 2010.

# References

- Special et al., 2009. Lucia Specia, Marco Turchi, Zhuoran Wang, John Shawe-Taylor, and Craig Saunders. Improving the Confidence of Machine Translation Quality Estimates. In *Machine Translation Summit XII*, 2009.
- Soricut and Echihabi, 2010. Radu Soricut and Abdessamad Echihabi. TrustRank: Inducing Trust in Automatic Translations via Ranking. *Proceedings of the Association for Computational Linguistics Conference (ACL-2010)*. 2010.
- Pighin et al., 2012. Daniele Pighin and Meritxell González and Lluís Màrquez. The UPC Submission to the WMT 2012 Shared Task on Quality Estimation *Proceedings of the 7th Workshop on Statistical Machine Translation* pg. 127--132. ACL 2012.
- Avramidis 2012. Quality Estimation for Machine Translation output using linguistic analysis and decoding features  
E Avramidis. Seventh Workshop on Statistical Machine Translation, 2012.
- [McNo4] Paul McNamee and James Mayfield. Character N-Gram Tokenization for European Language Text Retrieval. *Information Retrieval* , 7(1-2):73–97. 2004.
- [Sim92] Michel Simard, George F. Foster, and Pierre Isabelle. Using Cognates to Align Sentences in Bilingual Corpora. In *Proceedings of the Fourth International Conference on Theoretical and Methodological Issues in Machine Translation*. 2012.
- [Gon14] Meritxell Gonzàlez and Alberto Barrón-Cedeño and Lluís Màrquez. IPA and STOUT: Leveraging Linguistic and Source-based Features for Machine Translation Evaluation. Ninth Workshop on Statistical Machine Translation (WMT2014).

# Additional Slides

# Evaluation of syntactic measures

- NIST 2005 Arabic-to-English Exercise

Level	Metric	pall	ρSMT
Lexical	BLEU	0.06	0.83
	METEOR	0.05	0.90
Syntactic	POS	0.42	0.89
	DP	0.88	0.86
Semantic	CP	0.74	0.95
	SR	0.72	0.96
	DR	0.92	0.92
	DR-POS	0.97	0.90