

# Towards Heterogeneous Automatic MT Error Analysis

*(6th LREC)*

Jesús Giménez and Lluís Màrquez

—

TALP Research Center  
Technical University of Catalonia

May 29, 2008



- 1 Introduction
- 2 Our Proposal
- 3 Applicability
- 4 Discussion

# Outline

## 1 Introduction

- The Role of Evaluation Methods
- Recent Advances in Automatic MT Evaluation

## 2 Our Proposal

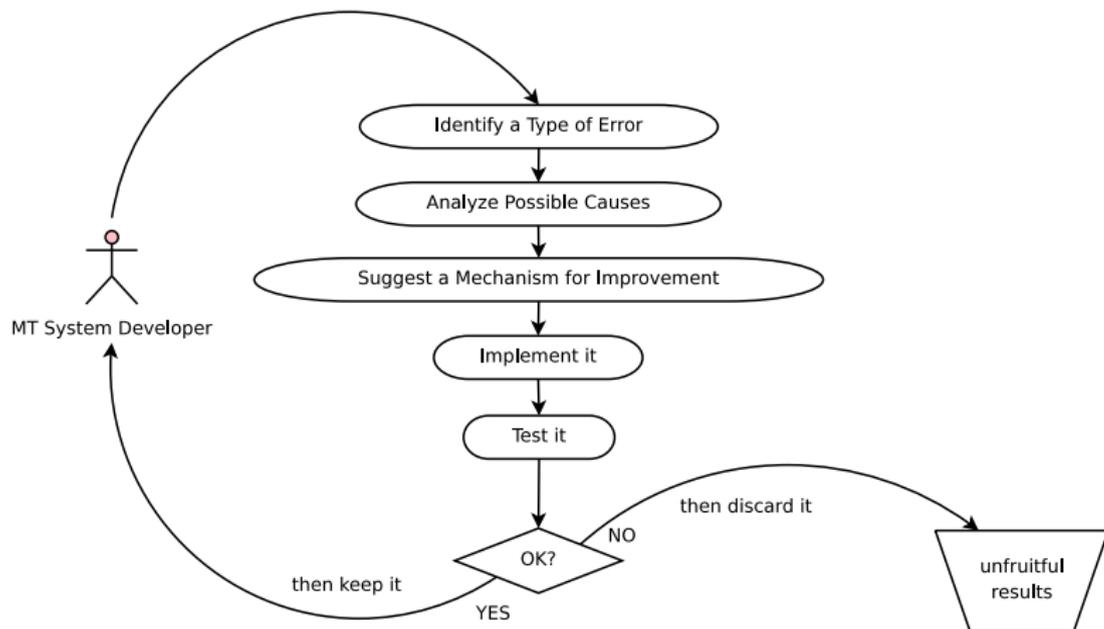
## 3 Applicability

## 4 Discussion

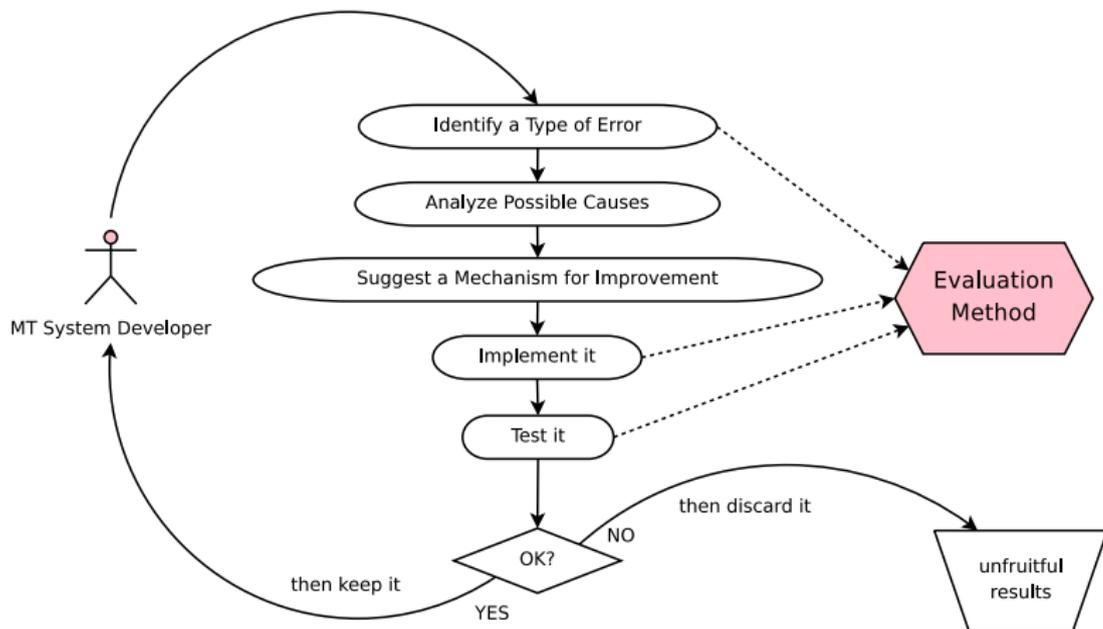
# Outline

- 1 Introduction
  - The Role of Evaluation Methods
  - Recent Advances in Automatic MT Evaluation
- 2 Our Proposal
- 3 Applicability
- 4 Discussion

# Development Cycle of MT systems



# Development Cycle of MT systems



# Error Analysis Today

- Error analyses are conducted manually
  - **low-level** analysis related to the linguistic analysis of translation quality (i.e., what?)
  - **high-level** analysis involving knowledge about the system architecture (i.e., why?)
- Error analyses require intensive human labor
- Automatic metrics are used only as quantitative evaluation measures
  - to identify high/low quality translations



# Error Analysis Today

- Error analyses are conducted manually
  - **low-level** analysis related to the linguistic analysis of translation quality (i.e., what?)
  - **high-level** analysis involving knowledge about the system architecture (i.e., why?)
- Error analyses require intensive human labor
- Automatic metrics are used only as quantitative evaluation measures
  - to identify high/low quality translations

# Error Analysis Today

- Error analyses are conducted manually
  - **low-level** analysis related to the linguistic analysis of translation quality (i.e., what?)
  - **high-level** analysis involving knowledge about the system architecture (i.e., why?)
- Error analyses require intensive human labor
- Automatic metrics are used only as quantitative evaluation measures
  - to identify high/low quality translations

# Error Analysis Today

- Error analyses are conducted manually
  - **low-level** analysis related to the linguistic analysis of translation quality (i.e., what?)
  - **high-level** analysis involving knowledge about the system architecture (i.e., why?)
- Error analyses require intensive human labor
- Automatic metrics are used only as quantitative evaluation measures
  - to identify high/low quality translations

# Error Analysis Today

- Error analyses are conducted manually
  - **low-level** analysis related to the linguistic analysis of translation quality (i.e., what?)
  - **high-level** analysis involving knowledge about the system architecture (i.e., why?)
- Error analyses require intensive human labor
- Automatic metrics are used only as quantitative evaluation measures
  - to identify high/low quality translations

# Metrics Based on Lexical Similarity

- **Edit Distance**  
WER, PER, TER
- **Precision**  
BLEU, NIST, WNM
- **Recall**  
ROUGE, CDER
- **Precision/Recall**  
GTM, METEOR, BLANC, SIA

# Outline

- 1** Introduction
  - The Role of Evaluation Methods
  - Recent Advances in Automatic MT Evaluation
- 2 Our Proposal
- 3 Applicability
- 4 Discussion



# Extending the Reference Lexicon

- Lexical variants
  - Morphological variations (i.e., stemming)  
→ ROUGE and METEOR
  - Synonymy lookup → METEOR (based on WordNet)
- Paraphrasing support
  - Zhou et al. [ZLH06]
  - Kauchak and Barzilay [KB06]
  - Owczarzak et al. [OGGW06]

# Beyond the Lexical Level

## Syntactic Similarity

- Shallow Parsing
  - Popovic and Ney [PN07]
  - Giménez and Màrquez [GM07]
- Constituency Parsing
  - Liu and Gildea [LG05]
- Dependency Parsing
  - Liu and Gildea [LG05]
  - Amigó et al. [AGGM06]
  - Mehay and Brew [MB07]
  - Owczarzak et al. [OvGW07]

# Beyond the Lexical Level

## Semantic Similarity

- Semantic Roles
  - Giménez and Màrquez [GM07]
- Named Entities
  - Reeder et al. [RMDW01]
  - Giménez and Màrquez [GM07]
- Discourse Representations
  - Giménez and Màrquez [GM08b]

# Outline

- 1 Introduction
- 2 Our Proposal**
  - A Smorgasbord of Features
- 3 Applicability
- 4 Discussion

# Rely on Automatic Metrics

**Idea:** Let automatic metrics do most of the *low-level* analysis, so system developers may concentrate on *high-level* analysis.



# Heterogeneous Error Analysis

- as automatic as possible
- as heterogeneous as possible
  - **Quality Aspects:** lexical, syntactic, semantic, etc.
  - **Granularity**
    - fine aspects → transfer of specific linguistic elements (e.g., what proportion of singular nouns are correctly translated?)
    - coarse aspects → overall linguistic structure (e.g., what proportion of the semantic role structure is correctly translated?)

# Heterogeneous Error Analysis

- as automatic as possible
- as heterogeneous as possible
  - **Quality Aspects:** lexical, syntactic, semantic, etc.
  - **Granularity**
    - fine aspects → transfer of specific linguistic elements (e.g., what proportion of singular nouns are correctly translated?)
    - coarse aspects → overall linguistic structure (e.g., what proportion of the semantic role structure is correctly translated?)

# Heterogeneous Error Analysis

- as automatic as possible
- as heterogeneous as possible
  - **Quality Aspects:** lexical, syntactic, semantic, etc.
  - **Granularity**
    - fine aspects → transfer of specific linguistic elements (e.g., what proportion of singular nouns are correctly translated?)
    - coarse aspects → overall linguistic structure (e.g., what proportion of the semantic role structure is correctly translated?)

# Heterogeneous Error Analysis

- as automatic as possible
- as heterogeneous as possible
  - **Quality Aspects:** lexical, syntactic, semantic, etc.
  - **Granularity**
    - fine aspects → transfer of specific linguistic elements (e.g., what proportion of singular nouns are correctly translated?)
    - coarse aspects → overall linguistic structure (e.g., what proportion of the semantic role structure is correctly translated?)

# Outline

- 1 Introduction
- 2 Our Proposal**
  - A Smorgasbord of Features
- 3 Applicability
- 4 Discussion

# Linguistic Similarities

- More than *500* metric variants operating at different *linguistic* levels:
  - Lexical
  - Shallow Syntactic (Lemmatization, PoS Tagging, and Base Phrase Chunking)
  - Syntactic (Constituency and Dependency Parsing)
  - Shallow Semantic (Semantic Roles and Named Entities)
  - Semantic (Discourse Representations)

# Shallow Syntactic Level

**SP-O<sub>p</sub>-★** Average overlapping between words belonging to the same PoS.

**SP-O<sub>c</sub>-★** Average overlapping between words belonging to the same phrase chunk type.

**SP-NIST<sub>l</sub>** NIST score over sequences of lemmas.

**SP-NIST<sub>p</sub>** NIST score over PoS sequences.

**SP-NIST<sub>ioB</sub>** NIST score over chunk IOB sequences.

**SP-NIST<sub>c</sub>** NIST score over sequences of chunks.

# Syntactic Level (i)

## ■ Dependency Overlapping

**DP-O<sub>l</sub>-\*** Average overlapping between words hanging at the same level.

**DP-O<sub>c</sub>-\*** Average overlapping between words hanging from terminal nodes (i.e., grammatical categories).

**DP-O<sub>r</sub>-\*** Average overlapping between words ruled by non-terminal nodes (i.e., grammatical relations).

# Syntactic Level (ii)

- **Head-word Chain Matching** (Liu and Gildea [LG05])
  - DP-HWC<sub>w</sub>** Average head-word chain matching up to length-4 word chains.
  - DP-HWC<sub>c</sub>** Average head-word chain matching up to length-4 category chains.
  - DP-HWC<sub>r</sub>** Average head-word chain matching up to length-4 relation chains.

# Syntactic Level (iii)

## ■ Syntactic Overlapping

**CP- $O_p$ -\*** Average overlapping between words belonging to the same PoS (similar to ' $SP-O_p$ -\*').

**CP- $O_c$ -\*** Average overlapping between words belonging to the same phrase type (similar to ' $SP-O_c$ -\*').

## ■ Syntactic Tree Matching (Liu and Gildea [LG05])

**CP-STM** Constituent tree matching averaged up to length-9 syntactic subpaths.

# Shallow Semantic Level (i)

## ■ Named Entity Overlapping/Matching

**NE-O<sub>e</sub>-★** Average lexical overlapping between named entities of the same type (excluding type 'O', i.e., Not-a-NE).

**NE-O<sub>e</sub>-★★** Average lexical overlapping between named entities of the same type (including 'O').

**NE-M<sub>e</sub>-★** Average lexical matching between named entities of the same type.

# Shallow Semantic Level (ii)

## ■ Semantic Role Overlapping/Matching

- SR-O<sub>r</sub>-\*** Average lexical overlapping between semantic roles (arguments and adjuncts) of the same type.
- SR-M<sub>r</sub>-\*** Average lexical matching between semantic roles of the same type.
- SR-O<sub>r</sub>** Role overlapping independently from the lexical realization.

# Semantic Level

## ■ Discourse Overlapping

**DR-O<sub>r</sub>-★** Average lexical overlapping between DR structures of the same type.

**DR-O<sub>rp</sub>-★** Average morphosyntactic overlapping between DR structures of the same type.

## ■ Semantic Tree Matching

**DR-STM** Matching between discourse representations averaged up to length-9 semantic subpaths.

# Linguistic Features at Work

ACL'07 MT Workshop (French/German/Spanish/Czech-to-English)

Metric	Adeq.	Fluen.	Rank	Const.	all
<b>SR-O<sub>r</sub>-★</b>	<b>.774</b>	<b>.839</b>	<b>.803</b>	.741	<b>.789</b>
ParaEval-Recall	.712	.742	.768	<b>.798</b>	.755
METEOR	.701	.719	.745	.669	.709
BLEU	.690	.722	.672	.602	.671
1-TER	.607	.538	.520	.514	.644
Max Adeq. Corr.	.651	.657	.659	.534	.626
Max Fluen. Corr.	.644	.653	.656	.512	.616
GTM	.655	.674	.616	.495	.610

# Outline

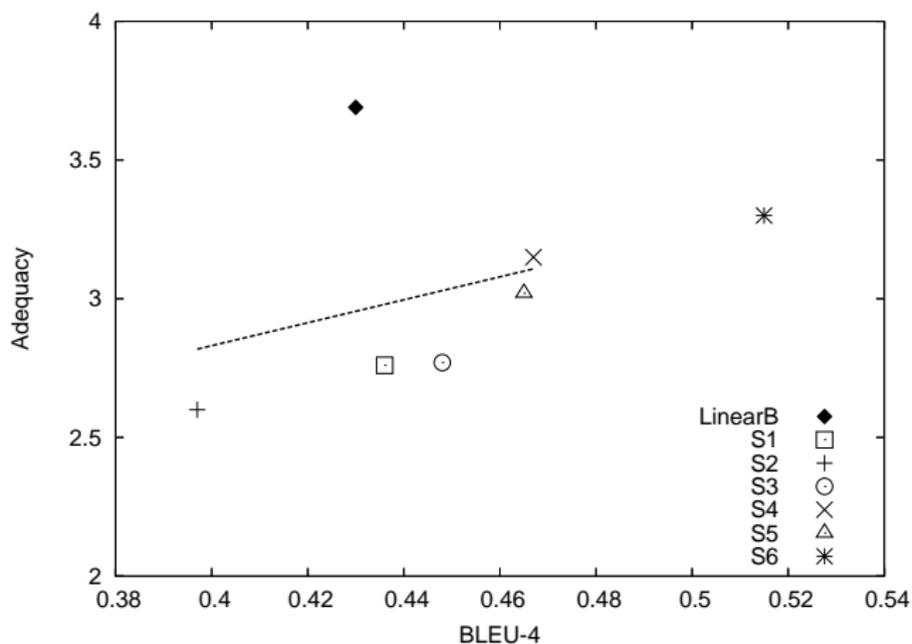
- 1 Introduction
- 2 Our Proposal
- 3 Applicability**
  - Settings
  - Document Level Error Analysis
  - Sentence Level Error Analysis
- 4 Discussion

# Outline

- 1 Introduction
- 2 Our Proposal
- 3 Applicability**
  - Settings
    - Document Level Error Analysis
    - Sentence Level Error Analysis
- 4 Discussion

# NIST 2005 MT Evaluation Puzzle

## ■ Arabic-to-English Translation Exercise [LP05]



# Linguistic Features Solved the Puzzle

- Giménez and Màrquez [GM07]

Feature	Metric	$R_{\text{sys}}$
<b>Lexical</b>	BLEU	0.06
	GTM	0.03
Syntactic	SP-NIST <sub>p</sub>	0.42
	DP-HWC <sub>r</sub>	0.88
	CP-STM	0.74
Semantic	SR- $O_r$ -*	0.61
	SR- $M_r$ -*	0.72
	DR- $O_r$ -*	0.92
	DR- $O_{rp}$ -*	0.97

# Linguistic Features Solved the Puzzle

- Giménez and Màrquez [GM07]

Feature	Metric	$R_{\text{sys}}$
<b>Lexical</b>	BLEU	0.06
	GTM	0.03
<b>Syntactic</b>	SP-NIST <sub>p</sub>	0.42
	DP-HWC <sub>r</sub>	<b>0.88</b>
	CP-STM	0.74
<b>Semantic</b>	SR- $O_r$ -*	0.61
	SR- $M_r$ -*	0.72
	DR- $O_r$ -*	0.92
	DR- $O_{rp}$ -*	<b>0.97</b>

# Outline

- 1 Introduction
- 2 Our Proposal
- 3 Applicability**
  - Settings
  - Document Level Error Analysis**
  - Sentence Level Error Analysis
- 4 Discussion

# A Note on Meta-Evaluation

- Metrics are automatically evaluated on the basis of **human likeness**, i.e., in terms of their ability to distinguish manual from automatic translations.
  - ORANGE, Lin and Och [LO04]
  - KING, Amigó et al. [AGPV05]
- We use the KING measure
  - *“A metric should never rank any reference translation lower in quality than any automatic translation.”*
- $\text{KING}(x)$  serves as an estimate of the impact on system performance of the quality aspects captured by metric  $x$

# Lexical Features

Feature	Metric	KING	LinearB	Best SMT
<b>Edit Distance</b>	1-PER	0.63	0.65	<b>0.70</b>
	1-TER	<b>0.70</b>	0.53	<b>0.58</b>
	1-WER	0.67	0.49	<b>0.54</b>
<b>Precision</b>	BLEU	0.65	0.47	<b>0.51</b>
	NIST	0.69	10.63	<b>11.27</b>
<b>Recall</b>	ROUGE <sub>W</sub>	0.68	0.31	<b>0.33</b>
<b>F-measure</b>	GTM ( $e = 1$ )	0.64	0.80	<b>0.85</b>
	GTM ( $e = 2$ )	0.66	0.31	<b>0.32</b>
	METEOR <sub>exact</sub>	0.68	0.60	<b>0.64</b>
	METEOR <sub>wnsyn</sub>	0.68	0.64	<b>0.68</b>

# Shallow Syntactic Features

Feature	Metric	KING	LinearB	Best SMT
<b>PoS Overlapping</b>	SP-O <sub>p</sub> -*	0.64	0.52	<b>0.55</b>
	SP-O <sub>p</sub> -J	0.26	0.53	<b>0.59</b>
	SP-O <sub>p</sub> -N	0.53	0.57	<b>0.63</b>
	SP-O <sub>p</sub> -V	0.43	0.39	<b>0.41</b>
<b>Chunk Overlapping</b>	SP-O <sub>c</sub> -*	0.63	0.54	<b>0.57</b>
	SP-O <sub>c</sub> -NP	0.60	0.59	<b>0.63</b>
	SP-O <sub>c</sub> -PP	0.38	0.63	<b>0.66</b>
	SP-O <sub>c</sub> -VP	0.41	0.49	<b>0.51</b>
<b>NIST<sub>x</sub></b>	SP-NIST <sub>I</sub> -5	0.69	10.78	<b>11.44</b>
	SP-NIST <sub>p</sub> -5	<b>0.71</b>	8.74	<b>9.04</b>
	SP-NIST <sub>job</sub> -5	0.65	6.81	<b>6.91</b>
	SP-NIST <sub>c</sub> -5	0.57	6.13	<b>6.18</b>

# Syntactic Features (i)

Feature	Metric	KING	LinearB	Best SMT
Dependency Parsing	DP-HWC <sub>w</sub> -4	0.59	0.14	0.14
	DP-HWC <sub>c</sub> -4	0.48	<b>0.42</b>	0.41
	DP-HWC <sub>r</sub> -4	0.52	<b>0.33</b>	0.31
	DP-O <sub>I</sub> -*	0.58	0.41	<b>0.43</b>
	DP-O <sub>c</sub> -*	0.60	0.50	<b>0.51</b>
	DP-O <sub>c</sub> -a	0.30	0.51	<b>0.57</b>
	DP-O <sub>c</sub> -aux	0.14	<b>0.56</b>	0.54
	DP-O <sub>c</sub> -det	0.35	<b>0.75</b>	0.73
	DP-O <sub>c</sub> -n	0.57	0.57	<b>0.59</b>
	DP-O <sub>c</sub> -v	0.37	0.43	<b>0.45</b>

# Syntactic Features (ii)

Feature	Metric	KING	LinearB	Best SMT
Dependency Parsing	DP-O <sub>r</sub> -*	<b>0.66</b>	0.36	0.36
	DP-O <sub>r</sub> -aux	0.14	<b>0.56</b>	0.54
	DP-O <sub>r</sub> -det	0.35	<b>0.75</b>	0.73
	DP-O <sub>r</sub> -fc	0.21	<b>0.26</b>	0.24
	DP-O <sub>r</sub> -i	<b>0.60</b>	<b>0.44</b>	0.43
	DP-O <sub>r</sub> -obj	0.43	<b>0.36</b>	0.35
	DP-O <sub>r</sub> -s	0.47	<b>0.52</b>	0,45
Constituency Parsing	CP-O <sub>p</sub> -*	<b>0.64</b>	0.52	<b>0.55</b>
	CP-O <sub>c</sub> -*	0.63	0.50	<b>0.53</b>
	CP-O <sub>c</sub> -NP	0.61	0.55	<b>0.58</b>
	CP-O <sub>c</sub> -PP	0.51	0.50	<b>0.53</b>
	CP-O <sub>c</sub> -SBAR	0.36	0.36	<b>0.38</b>
	CP-STM-9	0.58	0.35	0.35

# Shallow Semantic Features

Feature	Metric	KING	LinearB	Best SMT
<b>Named Entities</b>	NE-M <sub>e</sub> -*	0.32	0.53	<b>0.56</b>
	NE-M <sub>e</sub> -ORG	0.11	0.27	<b>0.29</b>
	NE-M <sub>e</sub> -PER	0.13	0.34	0.34
<b>Semantic Roles</b>	SR-M <sub>r</sub> -*	0.50	<b>0.19</b>	0.18
	SR-M <sub>r</sub> -A0	0.33	<b>0.31</b>	0.30
	SR-M <sub>r</sub> -A1	0.28	0.14	0.14
	SR-O <sub>r</sub>	0.41	<b>0.64</b>	0.63
	SR-O <sub>r</sub> -*	<b>0.53</b>	0.36	<b>0.37</b>
	SR-O <sub>r</sub> -AM-TMP	0.13	<b>0.39</b>	0.38

# Semantic Features

Feature	Metric	KING	LinearB	Best SMT
<b>Discourse Representations</b>	DR-O <sub>r</sub> -★	<b>0.59</b>	<b>0.36</b>	0.34
	DR-O <sub>r</sub> -card	0.12	<b>0.49</b>	0.45
	DR-O <sub>r</sub> -dr	0.56	<b>0.43</b>	0.40
	DR-O <sub>r</sub> -eq	0.12	<b>0.17</b>	0.16
	DR-O <sub>r</sub> -named	0.38	<b>0.48</b>	0.45
	DR-O <sub>r</sub> -pred	0.55	<b>0.38</b>	0.36
	DR-O <sub>r</sub> -prop	0.39	<b>0.27</b>	0.24
	DR-O <sub>r</sub> -rel	0.56	<b>0.38</b>	0.34
	DR-STM-9	0.40	0.26	0.26

# Outline

- 1 Introduction
- 2 Our Proposal
- 3 Applicability**
  - Settings
  - Document Level Error Analysis
  - Sentence Level Error Analysis**
- 4 Discussion

# Ex: Thousand Monks

<b>Ref 1:</b>	<p>Over 1000 monks and nuns , observers and scientists from over 30 countries and the host country attended the religious summit held for the first time in Myanmar which started today , Thursday .</p> <p>2: More than 1000 monks , nuns , observers and scholars from more than 30 countries , including the host country , participated in the religious summit which Myanmar hosted for the first time and which began on Thursday .</p> <p>3: The religious summit , staged by Myanmar for the first time and began on Thursday , was attended by over 1,000 monks an nuns , observers and scholars from more than 30 countries and host Myanmar .</p> <p>4: More than 1,000 monks , nuns , observers and scholars from more than 30 countries and the host country Myanmar participated in the religious summit , which is hosted by Myanmar for the first time and which began on Thursday .</p> <p>5: The religious summit , which started on Thursday and was hosted for the first time by Myanmar , was attended by over 1,000 monks and nuns , observers and scholars from more than 30 countries and the host country Myanmar .</p>
---------------	--

# Ex: Thousand Monks

<b>Info:</b>	<p>(1) → <b>subject:</b> over/more_than 1,000 monks and nuns, observers and scientists/scholars from over/more_than 30 countries , and/including the host country <b>action:</b> attended/participated_in</p> <p>(2) → <b>subject:</b> the religious summit <b>action:</b> began/started <b>temporal:</b> on Thursday</p> <p>(3) → <b>object:</b> the religious summit <b>action:</b> hosted <b>subject:</b> by Myanmar <b>mode:</b> for the first time</p>
<b>LinearB:</b>	1000 monks from more than 30 States and the host State Myanmar attended the Summit , which began on Thursday , hosted by Myanmar for the first time .
<b>Best SMT:</b>	Religious participated in the summit , hosted by Myanmar for the first time began on Thursday , as an observer and the world of the 1000 monk nun from more than 30 countries and the host state Myanmar .

## Ex: Thousand Monks - Lexical Features

Feature	Metric	LinearB	Best SMT
<b>Human</b>	Adequacy	<b>3</b>	2
	Fluency	<b>3.5</b>	2

<b>Edit Distance</b>	1-PER	0.64	<b>0.69</b>
	1-TER	<b>0.53</b>	0.51
	1-WER	0.40	<b>0.48</b>
<b>Precision</b>	BLEU	0.44	<b>0.45</b>
	NIST	9.04	<b>9.96</b>
<b>Recall</b>	ROUGE <sub>w</sub>	0.22	<b>0.23</b>
<b>F-measure</b>	GTM ( $e = 2$ )	0.30	<b>0.32</b>
	METEOR <sub>w<sub>nsyn</sub></sub>	0.59	<b>0.64</b>

## Ex: Thousand Monks - Shallow Syntactic Features

Feature	Metric	LinearB	Best SMT
<b>PoS Overlapping</b>	SP- $O_p$ -*	<b>0.52</b>	0.51
	SP- $O_p$ -IN	<b>0.71</b>	0.67
	SP- $O_p$ -NN	<b>0.67</b>	0.38
	SP- $O_p$ -NNP	0.60	<b>0.75</b>
	SP- $O_p$ -V	0.40	<b>0.75</b>
<b>Chunk Overlapping</b>	SP- $O_c$ -*	0.56	<b>0.60</b>
	SP- $O_c$ -NP	0.56	<b>0.60</b>
	SP- $O_c$ -PP	<b>0.80</b>	0.71
<b>NIST<sub>x</sub></b>	SP-NIST <sub>p</sub>	6.21	<b>8.36</b>
	SP-NIST <sub>c</sub>	<b>6.43</b>	6.25
	SP-NIST <sub>job</sub>	5.78	<b>6.41</b>

## Ex: Thousand Monks - Syntactic Features

Feature	Metric	LinearB	Best SMT
<b>Dependency Parsing</b>	DP-HWC <sub>w</sub> -4	<b>0.17</b>	0.16
	DP-O <sub>r</sub> -*	<b>0.46</b>	0.44
	DP-O <sub>r</sub> -mod	<b>0.62</b>	0.41
	DP-O <sub>r</sub> -obj	<b>0.29</b>	0.00
	DP-O <sub>r</sub> -pcomp-n	<b>0.71</b>	0.39
	DP-O <sub>r</sub> -rel	<b>0.33</b>	0.00
<b>Constituency Parsing</b>	CP-O <sub>c</sub> -*	<b>0.59</b>	0.48
	CP-O <sub>c</sub> -NP	<b>0.59</b>	0.55
	CP-O <sub>c</sub> -PP	<b>0.57</b>	0.54
	CP-O <sub>c</sub> -SB	<b>0.73</b>	0.00
	CP-O <sub>c</sub> -VP	<b>0.64</b>	0.42
	CP-STM-9	<b>0.34</b>	0.23

## Ex: Thousand Monks - Semantic Features

Feature	Metric	LinearB	Best SMT
<b>Semantic Roles</b>	SR- $O_r$	<b>0.84</b>	0.25
	SR- $O_r$ -*	<b>0.56</b>	0.18
	SR- $O_r$ -A0	<b>0.44</b>	0.10
	SR- $O_r$ -A1	<b>0.57</b>	0.28
<b>Discourse Representations</b>	DR- $O_r$ -*	<b>0.45</b>	0.34
	DR- $O_r$ -dr	<b>0.57</b>	0.40
	DR- $O_r$ -nam	<b>0.75</b>	0.24
	DR- $O_r$ -pred	0.44	<b>0.45</b>
	DR- $O_r$ -rel	<b>0.51</b>	0.32
	DR-STM-9	<b>0.32</b>	0.29

# Outline

- 1 Introduction
- 2 Our Proposal
- 3 Applicability
- 4 Discussion**
  - Conclusions
  - Future Work

# Outline

- 1 Introduction
- 2 Our Proposal
- 3 Applicability
- 4 Discussion**
  - Conclusions
  - Future Work

# Heterogeneous Automatic MT Error Analysis

- We have presented a valid path towards *heterogeneous automatic MT error analysis*:
  - Our approach allows developers to rapidly obtain detailed qualitative linguistic reports on their system's capabilities.
  - Human efforts may concentrate on high-level analysis.

# Hey! Linguistic Metrics are Not the Panacea<sup>1</sup>

- Linguistic metrics rely on:
  - 1 the representativity of the set of human references
    - lexicon
    - grammar
    - style...
  - 2 automatic linguistic processors are
    - domain-dependent
    - language-dependent
    - prone to error
    - slow

Sentence level performance must be improved!

---

<sup>1</sup>Panacea: a remedy for all ills or difficulties (see cure-all) 

# Hey! Linguistic Metrics are Not the Panacea<sup>1</sup>

- Linguistic metrics rely on:
  - 1 the representativity of the set of human references
    - lexicon
    - grammar
    - style...
  - 2 automatic linguistic processors are
    - domain-dependent
    - language-dependent
    - prone to error
    - slow

Sentence level performance must be improved!

---

<sup>1</sup>Panacea: a remedy for all ills or difficulties (see cure-all)

# Outline

- 1 Introduction
- 2 Our Proposal
- 3 Applicability
- 4 Discussion**
  - Conclusions
  - Future Work**

# Ongoing Steps...

- 1 Improving sentence level behavior:
  - Backing off to lexical similarity [GM08b]
  - Working on metric combinations [GM08a]
- 2 Porting metrics to languages other than English (e.g., Castilian Spanish, Catalan...)

# A New Interface

IQ<sub>MT</sub>
Heterogeneous MT Error Analysis

File Edit View Tools Help

Labels

NP

VP

PP

DT

NN

VBZ

Lexical

Syntactic

Semantic

Automatic Translation(s)

Human Reference(s)

# Thanks for your Attention

IQ<sub>MT</sub> v2.0 is publicly available at:

<http://www.lsi.upc.edu/~nlp/IQMT>



# Towards Heterogeneous Automatic MT Error Analysis

(6th LREC)

Jesús Giménez and Lluís Màrquez

—

TALP Research Center

Technical University of Catalonia

May 29, 2008



-  Enrique Amigó, Jesús Giménez, Julio Gonzalo, and Lluís Màrquez.  
MT Evaluation: Human-Like vs. Human Acceptable.  
*In Proceedings of COLING-ACL06, 2006.*
-  Enrique Amigó, Julio Gonzalo, Anselmo Peñas, and Felisa Verdejo.  
QARLA: a Framework for the Evaluation of Automatic Sumarization.  
*In Proceedings of the 43th Annual Meeting of the Association for Computational Linguistics, 2005.*
-  Jesús Giménez and Lluís Màrquez.  
Linguistic Features for Automatic Evaluation of Heterogeneous MT Systems.  
*In Proceedings of the ACL Workshop on Statistical Machine Translation, 2007.*

-  **Jesús Giménez and Lluís Màrquez.**  
Heterogeneous Automatic MT Evaluation Through  
Non-Parametric Metric Combinations.  
*In Proceedings of IJCNLP, 2008.*
-  **Jesús Giménez and Lluís Màrquez.**  
On the Robustness of Linguistic Features for Automatic  
MT Evaluation.  
*In Proceedings of the ELRA Workshop on Evaluation.  
Looking into the Future of Evaluation: when automatic  
metrics meet task-based and performance-based  
approaches, 2008.*
-  **David Kauchak and Regina Barzilay.**  
Paraphrasing for Automatic Evaluation.  
*In Proceedings of NLH-NAACL, 2006.*
-  **Ding Liu and Daniel Gildea.**

Syntactic Features for Evaluation of Machine Translation.  
*In Proceedings of ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, 2005.*



Chin-Yew Lin and Franz Josef Och.

ORANGE: a Method for Evaluating Automatic Evaluation Metrics for Machine Translation.  
*In Proceedings of COLING, 2004.*



Audrey Le and Mark Przybocki.

NIST 2005 machine translation evaluation official results.  
Technical report, NIST, August 2005.



Dennis Mehay and Chris Brew.

BLEUATRE: Flattening Syntactic Dependencies for MT Evaluation.

In *Proceedings of the 11th Conference on Theoretical and Methodological Issues in Machine Translation (TMI)*, 2007.



Karolina Owczarzak, Declan Groves, Josef Van Genabith, and Andy Way.

Contextual Bitext-Derived Paraphrases in Automatic MT Evaluation.

In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas (AMTA)*, 2006.



Karolina Owczarzak, Josef van Genabith, and Andy Way.  
Dependency-Based Automatic Evaluation for Machine Translation.

In *Proceedings of SSST, NAACL-HLT/AMTA Workshop on Syntax and Structure in Statistical Translation*, 2007.



Maja Popovic and Hermann Ney.

Word Error Rates: Decomposition over POS classes and Applications for Error Analysis.

*In Proceedings of the Second Workshop on Statistical Machine Translation*, pages 48–55, Prague, Czech Republic, June 2007. Association for Computational Linguistics.



Florence Reeder, Keith Miller, Jennifer Doyon, and John White.

The Naming of Things and the Confusion of Tongues: an MT Metric.

*In Proceedings of the Workshop on MT Evaluation "Who did what to whom?" at MT Summit VIII*, 2001.



Liang Zhou, Chin-Yew Lin, and Eduard Hovy.  
Re-evaluating Machine Translation Results with Paraphrase Support.

In *Proceedings of EMNLP*, 2006.

