

# IQ<sub>MT</sub>: A Framework for Automatic MT Evaluation based on 'Human Likeness'

*A Tutorial*

Jesús Giménez (TALP-UPC)

October, 2006

- 1 Introduction
  - Current Metrics
  - What is next?
  
- 2 Approach
  - The QARLA Framework
  - QARLA for MT evaluation
  
- 3 A case of study: Europarl
  - Evaluation Outside QARLA
  - Evaluation Inside QARLA
  
- 4 Conclusions and Further Work

- 1 Introduction
  - Current Metrics
  - What is next?
- 2 Approach
- 3 A case of study: Europarl
- 4 Conclusions and Further Work

# Lexical Level

- Most metrics are developed on the basis of 'Human Acceptability'
- Most metrics consider only **lexical** similarities:
  - WER (Nießen et al., 2000)
  - BLEU (Papinenit et al., 2001)
  - NIST (Doddington, 2002)
  - PER (Leusch et al., 2003)
  - GTM (Melamed et al., 2003)
- Little effort has been devoted to:
  - introducing additional **linguistic information**
  - **combining** different metrics

# Remarkable Efforts

- Use of additional knowledge at the lexical level:
  - ROUGE (stemming)
  - METEOR (stemming + WordNet lookup)
  - WNM (word frequency weighting)
- Use of syntactic knowledge:
  - HPCM, STM, DSTM (Liu and Gildea, 2005)
- Combining different metrics:
  - Kulesza and Shieber (2004)

- 1 Introduction
  - Current Metrics
  - What is next?
- 2 Approach
- 3 A case of study: Europarl
- 4 Conclusions and Further Work

# (1) A new SuperMetric? or... (2) divide and conquer?

- ① Design a new metric XXX which considers information at different linguistic levels (lexical, syntactic, semantic)
- ② Design a set of *specialized metrics* which work at different levels and *combine* their scores into *a single measure of MT quality*?

## 1 Introduction

## 2 Approach

- The QARLA Framework
- QARLA for MT evaluation

## 3 A case of study: Europarl

## 4 Conclusions and Further Work



- 1 Introduction
- 2 Approach
  - The QARLA Framework
  - QARLA for MT evaluation
- 3 A case of study: Europarl
- 4 Conclusions and Further Work

# The QARLA assumption

“All manual references are equally optimal and the best similarity metric is the one that identifies and uses features that are common to all references, grouping them and separating them from the automatic candidates.”

- QARLA provides 3 measures:

**QUEEN** determines the quality of a set of systems

**KING** determines the quality of a set of metrics

**JACK** determines the quality of a test set

(Amigo et al., ACL 2005)

# The QUEEN assumption

“A good candidate should be similar to all models according to all metrics.”

Given a set of metrics  $X$ , a set of references  $R$ , and a candidate  $a$ :

$$\text{QUEEN}_{X,R}(a) = \text{Prob}(\forall x \in X : x(a, r) \geq x(r', r''))$$

# The QUEEN properties

- (i) it is able to combine different similarity metrics into a single evaluation measure.
- (ii) it is not affected by the scale properties of individual metrics, i.e. it does not require metric normalisation and it is not affected by metric weighting.
- (iii) Candidates which are very far from the set of models all receive QUEEN=0.
- (iv) The value of QUEEN is maximised for candidates that “merge” with the references according to all metrics in  $X$ .
- (v) The universal quantifier on the metric parameter  $x$  implies that adding redundant metrics does not bias the result of QUEEN.

# The KING measure

“A good metric should score human references higher than MT systems.”

Given sets of metrics  $X$ , references  $R$ , and systems  $A$ :

$$\text{KING}_{A,R}(X) =$$

$$\text{Prob}(\forall a \in A : \text{QUEEN}_{X,R-\{r\}}(r) \geq \text{QUEEN}_{X,R-\{r\}}(a))$$

KING measures the ability of a set of metrics to discern between automatic and human translations.

# The JACK measure

“References and systems should not be biased.”

Given sets of metrics  $X$ , references  $R$ , and systems  $A$ :

$$\text{JACK}(A, R, X) =$$

$$\text{Prob}(\exists a, a' \in A : \text{QUEEN}_{X,R}(a) > 0 \\ \wedge \text{QUEEN}_{X,R}(a') > 0 \wedge \forall x \in X : x(a, a') \leq x(a, r))$$

A high JACK value means that most references are closely and heterogeneously surrounded by automatic translations.

- 1 Introduction
- 2 Approach
  - The QARLA Framework
  - QARLA for MT evaluation
- 3 A case of study: Europarl
- 4 Conclusions and Further Work





# lQsetup <config-file>

'a priori' computing of similarities for given sets of metrics  $X$ , references  $R$ , and systems  $A$ :

```
#SRC =====  
source=./data/dev.src.spa.iso  
#REF =====  
ref=./data/dev.src.eng.iso  
ref=./data/dev.ref1.eng.iso  
ref=./data/dev.ref2.eng.iso  
#OUT =====  
sys=./data/sys_1.out  
...  
sys=./data/sys_n.out  
#-----
```

# IQsetup <config-file>

```
NAME=NIST_mt05_AE          /* testbed name (id) */
IQMT=/usr/local/IQMT-2.0   /* IQMT path */
LANG=ENG                   /* target language */
CASE=all                   /* case [lc|all]

doBLEU                     /* metric set selection*/
doNIST
doGTM
doMETEOR
doWER
doPER
doROUGE
```

# 'IQXML' similarity files

```
<IQ metric="BLEU-4" ref="R0"  
      score="0.3945" target="S0">  
  <S n="1">0.3033</S>  
  <S n="2">0.5833</S>  
  ...  
  <S n="1007">0.6852</S>  
  <S n="1008">0.8333</S>  
</IQ>
```

# IQeval

IQ methodology in 4 steps:

- 1 IQsetup
- 2 IQeval [-doKING | -optimizeKING]
- 3 IQeval -doQUEEN -M optimal\_set\_of\_metrics
- 4 IQeval -doJACK

- 1 Introduction
- 2 Approach
- 3 A case of study: Europarl**
  - Evaluation Outside QARLA
  - Evaluation Inside QARLA
- 4 Conclusions and Further Work

# Experimental Setting

- Openlab 2006 data (TC-STAR Consortium)
- Spanish-to-English
- 1,281,427 sentences for training
- 1,008 sentences for test
  - 3 human references
  - Outputs by 4 MT systems (WB, SYSTRAN, PB, PB++)
- 26 metrics from 7 different families  
(BLEU, NIST, GTM, WER, PER, ROUGE, METEOR)

- 1 Introduction
- 2 Approach
- 3 A case of study: Europarl
  - Evaluation Outside QARLA
  - Evaluation Inside QARLA
- 4 Conclusions and Further Work

# Standard Metrics Outside QARLA

METRIC	WB	SYSTRAN	PB	PB++
<b>1-PER</b>	0.66	0.70	<b>0.74</b>	<b>0.74</b>
<b>1-WER</b>	0.58	0.60	<b>0.64</b>	0.63
<b>BLEU-3</b>	0.50	0.56	<b>0.66</b>	<b>0.66</b>
<b>GTM-2</b>	0.33	0.36	<b>0.41</b>	<b>0.41</b>
<b>MTR-exact</b>	0.57	0.65	0.69	<b>0.70</b>
<b>NIST-3</b>	8.79	9.59	10.66	<b>10.72</b>
<b>RG-L</b>	0.56	0.63	0.66	<b>0.67</b>



- 1 Introduction
- 2 Approach
- 3 A case of study: Europarl
  - Evaluation Outside QARLA
  - Evaluation Inside QARLA
- 4 Conclusions and Further Work

# Evaluating with $IQ_{MT}$ . KING

Evaluation metric	KING
1-PER	0.30
1-WER	0.34
BLEU-3	0.32
GTM-2	0.32
MTR-exact	0.29
NIST-3	<b>0.37</b>
RG-L	0.33

$\mathbf{X} = \{\text{NIST-2, NIST-3, NIST-4, and 1-WER}\}$

$\text{KING}_{A,R}(\mathbf{X}) = \mathbf{0.38}$  (descriptive power)

# Evaluating with $IQ_{MT}$ . QUEEN (+JACK)

MT System	QUEEN <sub>X,R</sub>
WB	0.31
SYSTRAN	0.39
PB	0.45
PB++	<b>0.46</b>

JACK(A, R, X) = **0.77** (test set reliability)

# Error Analysis based on QUEEN

<b>source</b>	los ciudadanos esperan de nosotros algo más que la simple <b>gestión de las crisis ; esperan señales</b> y una política sostenible en estos ámbitos .
---------------	---

## systems

<b>WB</b>	the citizens expect of us something more than the simple <b>management of the crisis</b> and a sustainable policy in these areas . <b>expectantly signals</b>
<b>SYS-TRAN</b>	the citizens wait for of us something more than the simple <b>management of the crises; they wait for signals</b> and a sustainable policy in these scopes.
<b>PB</b>	the citizens expect us any more than simply <b>managing crises ; they hope signals</b> and a sustainable policy in these areas .
<b>PB++</b>	the citizens expect us something more than simply <b>crisis management ; they expect signs</b> and a sustainable policy in these areas .

## references

<b>R0</b>	the public expect more than just <b>crisis management ; they expect signs</b> , and a sustainable policy in these fields .
<b>R1</b>	citizens expect something more of us than just simple <b>crisis management ; they expect signs</b> and sustainable policies in these areas .
<b>R2</b>	the citizens expect from us something more than a simple <b>crisis management ; they expect signs</b> and a sustainable policy in these matters .

- 1 Introduction
- 2 Approach
- 3 A case of study: Europarl
- 4 Conclusions and Further Work**

# Conclusions and Further Work

- MT Evaluation based on 'Human Likeness'
- MT Evaluation over Metric Combinations
- Applications
  - MT evaluation of heterogenous MT systems (SMT, rule-based)
  - MT Development (tuning of parameters)
  - MT Error Analysis
- New metrics
  - Dependency Parsing
  - Shallow Parsing (PoS-tagging, lemma, chunking)
  - Full Parsing
  - Shallow Semantics (NERC, SRL)

# Thanks

IQ<sub>MT</sub> v1.3 is freely available at:

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