

# Confidence measures for machine translation

Sylvain Raybaud, Caroline Lavecchia, David Langlois, Kamel  
Smaïli

LORIA (Nancy, France)/PAROLE group

ICAART 09  
21st January 2009

- 1 Overview of this work
  - Statistical Machine Translation
  - Overview of confidence measures
  - Contribution
- 2 Proposed confidence measures
  - Confidence measures based on mutual information
    - Intra-language MI
    - Inter-language MI
  - Language Model Based Confidence Measures
    - N-gram LM
    - Linguistic Features
  - Confidence Measures Combination
- 3 Conclusion and future work

- 1 Overview of this work
  - Statistical Machine Translation
    - Overview of confidence measures
    - Contribution
- 2 Proposed confidence measures
  - Confidence measures based on mutual information
    - Intra-language MI
    - Inter-language MI
  - Language Model Based Confidence Measures
    - N-gram LM
    - Linguistic Features
  - Confidence Measures Combination
- 3 Conclusion and future work

## A quick overview of SMT

source sentence:	$\mathbf{s} = s_1, \dots, s_l$
candidate translation:	$\mathbf{t} = t_1, \dots, t_j$

$$\hat{\mathbf{t}} = \arg \max_{\mathbf{t}} P(\mathbf{t}|\mathbf{s}) = \arg \max_{\mathbf{t}} P(\mathbf{s}|\mathbf{t}) \times P(\mathbf{t})$$

$P(\mathbf{s}|\mathbf{t})$ : translation model:

→ is the idea of the source sentence kept?

$P(\mathbf{t})$ : language model:

→ is the generated sentence correct?

- 1 Overview of this work
  - Statistical Machine Translation
  - Overview of confidence measures
  - Contribution
- 2 Proposed confidence measures
  - Confidence measures based on mutual information
    - Intra-language MI
    - Inter-language MI
  - Language Model Based Confidence Measures
    - N-gram LM
    - Linguistic Features
  - Confidence Measures Combination
- 3 Conclusion and future work

# Motivation

- Estimate **translation quality** without human references:
  - *Costly.*
  - *Generally not available.*
- Estimate **reliability** of each item in the translation (rather than the overall quality).
  - *Propose a better hypothesis.*
  - *Ask for user's verification.*

## Main approaches

- *A posteriori* relevance of a word:
  - estimate **word's probability of correctness** using the word-lattice or a translation table.
  - compute a **score** (predictive parameter), for example likelihood ratios.
  - **combine different measures** (perceptron, multi-layers perceptron, ...).

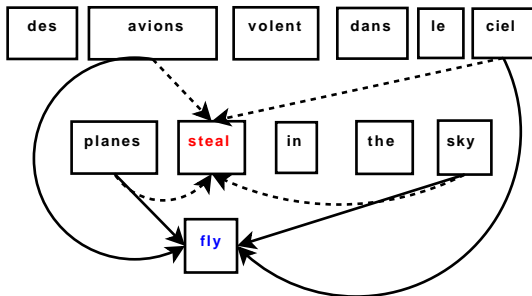
- 1 Overview of this work
  - Statistical Machine Translation
  - Overview of confidence measures
  - Contribution
- 2 Proposed confidence measures
  - Confidence measures based on mutual information
    - Intra-language MI
    - Inter-language MI
  - Language Model Based Confidence Measures
    - N-gram LM
    - Linguistic Features
  - Confidence Measures Combination
- 3 Conclusion and future work



## Our approach

Detect wrong or misplaced words by the mean of predictive parameters using the context.

*Example:*



## Proposed measures

- We compute different measures relying on:
  - **Mutual information** between words.
  - **N-gram** language model.
  - Words' **linguistic features**.
- Linear combination of the above measures.

## Performance estimation

- For every target word  $t$ :

$$class(t) = \begin{cases} \text{correct} & \text{if } C(t) > \delta \\ \text{incorrect} & \text{if } C(t) < \delta \end{cases}$$

- Machine classification compared to man-made one.
- Performance = reliability of the classification.

*Example:*

<i>source</i>	des	avions	volent	dans	le	ciel
<i>translation</i>	a	planes	steal	in	the	sky
<i>human</i>	0	1	0	1	1	1
<i>machine</i>	1	0	0	1	1	1

correct acceptances: 3

correct rejections: 1

# Evaluation metrics

- ROC curve, F-measure and Error Rate:

		AUTOMATIC	
		correct	incorrect
REFERENCE	correct	A	B
	incorrect	C	D

**correct acceptances:** A

**correct rejections:** D

**incorrect acceptances:** C

**incorrect rejections:** B

**Error Rate (ER):**  $\frac{B+C}{A+B+C+D}$

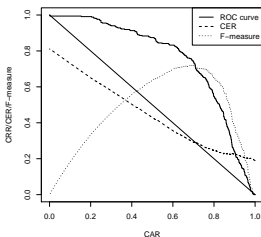
**Correct Acceptance Rate (CAR):**  $\frac{A}{A+B}$   $\left\{ \begin{array}{l} \text{proportion of correct words which} \\ \text{are actually accepted.} \end{array} \right.$

**Correct Rejection Rate (CRR):**  $\frac{D}{C+D}$   $\left\{ \begin{array}{l} \text{proportion of incorrect words which} \\ \text{are actually rejected.} \end{array} \right.$

**F-measure:**  $\frac{2 \times CAR \times CRR}{CAR + CRR}$

## Receiver Operating Characteristics curve

- Given threshold  $\rightarrow$  determines  $(CAR, CRR) \rightarrow$  a point.
- Testing different thresholds  $\rightarrow$  **ROC curve**.
- As the threshold increases, CRR goes from 0 to 1 and CAR from 1 to 0.



- Very easy interpretation: **the closer the curve is to (1, 1), the better.**

## Experimental settings

English to French phrase based translation	
<i>Language model</i>	SRILM toolkit, 1 to 5-grams models
<i>Translation model</i>	GIZA++, IBM-5 model
<i>Decoder</i>	Pharaoh
<i>Corpora</i>	EUROPARL 2005 (500k English-French sentences pairs)

- Confidence measures evaluated on the test corpus (around 5,000 words).

- 1 Overview of this work
  - Statistical Machine Translation
  - Overview of confidence measures
  - Contribution
- 2 Proposed confidence measures
  - Confidence measures based on mutual information
    - Intra-language MI
    - Inter-language MI
  - Language Model Based Confidence Measures
    - N-gram LM
    - Linguistic Features
  - Confidence Measures Combination
- 3 Conclusion and future work

## Mutual information

- Measures the *dependency* between two random variables.

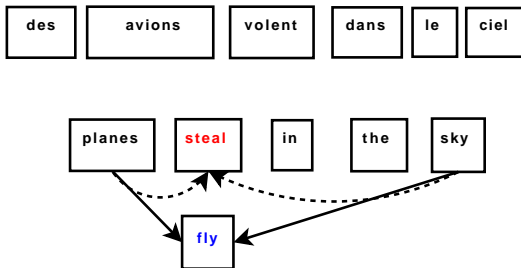
$$\mathcal{I}(X, Y) = \sum_x \sum_y P(X = x, Y = y) \log_2 \left( \frac{P(X = x, Y = y)}{P(X = x)P(Y = y)} \right)$$

- In statistical linguistic: detects words that co-occur.



# Intra-language MI

- Evaluates the relevance of a word given its context.



*How do we do that?*

# Training

$p(x)$  = probability that word  $x$  occurs in a sentence.

$p(x, y)$  = probability that words  $x$  and  $y$   
appear in the same sentence.

We consider the contribution of each word pair separately:

$$\mathcal{I}(x, y) = p(x, y) \log \left( \frac{p(x, y)}{p(x)p(y)} \right)$$

## Training (continued): Intra-language triggers

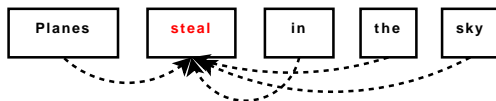
We get a so called **triggers list**:

word X triggers word Y	relationship is quantified by <i>mutual information</i>
------------------------	--

*Example:*

<i>security</i>	↔	<i>foreign</i>	$6.58 \cdot 10^{-3}$
<i>security</i>	↔	<i>policy</i>	$6.17 \cdot 10^{-3}$
<i>security</i>	↔	<i>social</i>	$4.33 \cdot 10^{-3}$
	...		
<i>policy</i>	↔	<i>common</i>	$1.12 \cdot 10^{-2}$
<i>policy</i>	↔	<i>foreign</i>	$1.05 \cdot 10^{-2}$
<i>policy</i>	↔	<i>agricultural</i>	$6.90 \cdot 10^{-3}$

## Word's confidence estimation



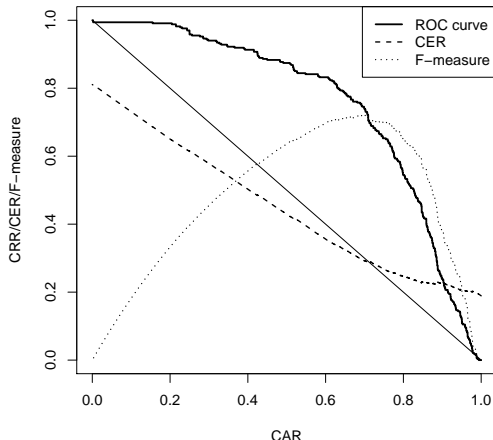
- Word's score: average mutual information between this word and the others in the sentence:

$$\mathbf{t} = t_1, \dots, t_J$$
$$C(t_j) = \frac{1}{J-1} \sum_{i=1, \dots, J, i \neq j} \mathcal{I}(t_i, t_j)$$

Possible tweaks:

- Take words' positions into account.
- Ignore function words (*the, of, ...*).

# Efficiency



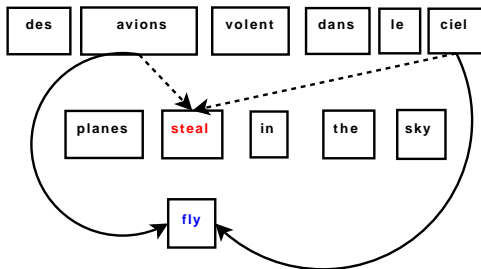
- Words' positions not taken into account.
- function words ignored.

Best classifier wrt.  
F-measure:

- F-measure: 0.72
- Error Rate: 0.30
- CAR: 0.69
- CRR: 0.75

## Inter-language MI

- Detects inconsistent translations.



*How do we do that?*

## Training

We are looking for couples  $(x, y)$  such that  $x \in \text{source}$  and  $y \in \text{translation}$

$p_S(x)$  = probability that  $\mathbf{x}$  occurs in a **source** sentence

$p_T(y)$  = probability that  $\mathbf{y}$  occurs in a **translated** sentence

$p(x, y)$  = probability that  $\mathbf{x}$  occurs in a **source**  
and  $\mathbf{y}$  in its **translation**

$$I(x, y) = p(x, y) \log \left( \frac{p(x, y)}{p_S(x)p_T(y)} \right)$$

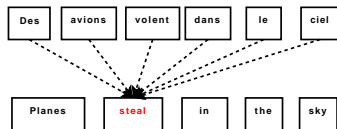
Remark: Can be used to build translation tables.

## Training (continued): Inter-languages triggers

<i>sécurité</i>	→	<i>security</i>	$4.39 \cdot 10^{-2}$
<i>sécurité</i>	→	<i>safety</i>	$3.76 \cdot 10^{-2}$
<i>sécurité</i>	→	<i>foreign</i>	$5.65 \cdot 10^{-3}$
	...		
<i>politique</i>	→	<i>policy</i>	$1.39 \cdot 10^{-1}$
<i>politique</i>	→	<i>political</i>	$5.75 \cdot 10^{-2}$
<i>politique</i>	→	<i>common</i>	$1.23 \cdot 10^{-2}$



## Word's confidence estimation



- A word's score is the average MI between this word and the words in the source sentence:

$$\mathbf{s} = s_1, \dots, s_l$$

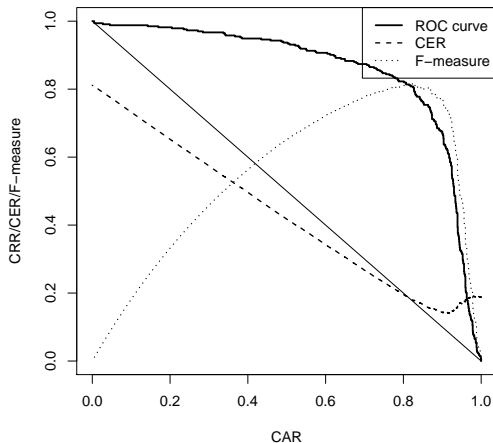
$$\mathbf{t} = t_1, \dots, t_j$$

$$C(t_j) = \frac{1}{l} \sum_{i=1, \dots, l} \mathcal{I}(s_i, t_j)$$

Same possible tweaks:

- Consider words positions.
- Ignore function words.

# Efficiency



- Limited distortion.
- function words ignored.

Best classifier wrt.  
F-measure:

- F-measure: 0.82
- Error Rate: 0.18
- CAR: 0.83
- CRR: 0.81

- 1 Overview of this work
  - Statistical Machine Translation
  - Overview of confidence measures
  - Contribution
- 2 Proposed confidence measures
  - Confidence measures based on mutual information
    - Intra-language MI
    - Inter-language MI
  - Language Model Based Confidence Measures
    - N-gram LM
    - Linguistic Features
  - Confidence Measures Combination
- 3 Conclusion and future work

## N-gram based measure

The decoder makes a tradeoff between  $P(\mathbf{t})$  and  $P(\mathbf{s}|\mathbf{t})$ ; unbalanced choices penalised:

$P(\mathbf{t})$  too low  $\rightarrow$   $\mathbf{t}$  often incorrect

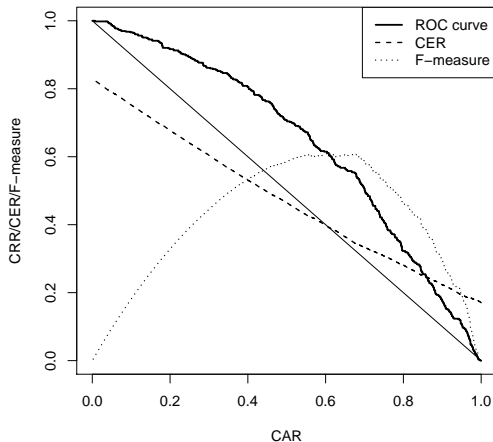
$$\mathbf{t} = t_1, \dots, t_j$$

$$C(t_j) = P(t_j | t_{j-n+1}, \dots, t_{j-1})$$

*Example:*

$$\begin{cases} P(\text{steal} | \text{the planes}) = 0.001 \\ P(\text{fly} | \text{the planes}) = 0.4 \end{cases}$$

# Efficiency



- Interesting performances but not as good as MI's.

Best classifier wrt.  
F-measure:

- F-measure: 0.61
- Error Rate: 0.39
- CAR: 0.61
- CRR: 0.61

## Linguistic Features

- Directly detect **grammatical errors**.
- Replace each word by a set of features: grammatical category, tense, person, gender, number.
- Features given by BDLEX.

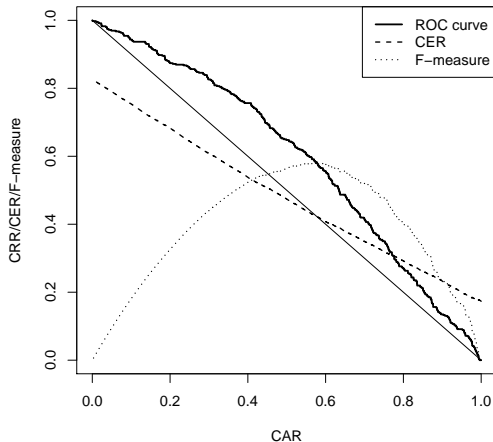
*Examples:*

$t = \text{was}$	$\rightarrow$	$\tilde{t} = \text{verb, indicative imperfect, 3rd person}$
$t = \text{word}$	$\rightarrow$	$\tilde{t} = \text{noun, masculine, singular}$

N-gram features model:

$$C(t_j) = P(\tilde{t}_j | \tilde{t}_{n-j+1}, \dots, \tilde{t}_{j-1})$$

# Efficiency



- Inefficient when used alone.
- Disambiguation is difficult.

Best classifier wrt.  
F-measure:

- F-measure: 0.58
- Error Rate: 0.43
- CAR: 0.57
- CRR: 0.59

- 1 Overview of this work
  - Statistical Machine Translation
  - Overview of confidence measures
  - Contribution
- 2 Proposed confidence measures
  - Confidence measures based on mutual information
    - Intra-language MI
    - Inter-language MI
  - Language Model Based Confidence Measures
    - N-gram LM
    - Linguistic Features
  - Confidence Measures Combination
- 3 Conclusion and future work



## Motivation

- Each measure detects a specific kind of error.
- Combining them yields more powerful classifiers.

# Perceptron

- Scores are **linearly combined**:

$$C(t) = \sum_{i=1}^k w_i C_i(t)$$

- **Weights are optimised** by comparing expected class  $c_{reference}$  to produced one  $c$ :

$$\forall t \forall i w_i \leftarrow w_i + \alpha (c_{reference}(t) - c(t)) C_i(t)$$

# Efficiency

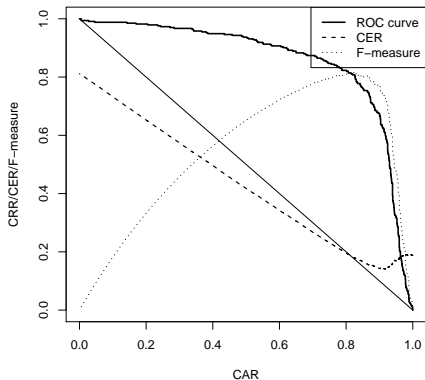


Figure: Inter-MI alone

F-measure: 0.82 - CAR: 0.83 - CRR: 0.81 - ER: 0.18

# Performances

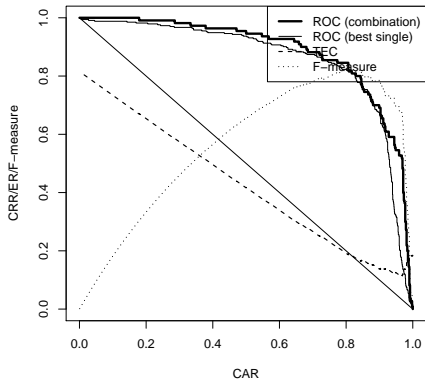


Figure: Intra- and inter-MI combined

F-measure: 0.83 - CAR: 0.81 - CRR: 0.85 - ER: 0.19

## Conclusion

- We propose four original confidence measures (plus combination).
- Some achieve results comparable to state-of-the-art.

## Future work

- Adapt more measures from ASR.
- Propose new original measures for MT.
- Come up with a smarter combination method.
- More thorough performance evaluation.

- Questions!