



**HAL**  
open science

# Spoken Language Translation Graphs Re-decoding using Automatic Quality Assessment

Laurent Besacier, Benjamin Lecouteux, Ngoc Quang Luong, Ngoc Tien Le

## ► To cite this version:

Laurent Besacier, Benjamin Lecouteux, Ngoc Quang Luong, Ngoc Tien Le. Spoken Language Translation Graphs Re-decoding using Automatic Quality Assessment. ASRU, 2015, Scotsdale, United States. hal-02095256

**HAL Id: hal-02095256**

**<https://hal.science/hal-02095256v1>**

Submitted on 10 Apr 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Abstract

This paper investigates how automatic quality assessment of spoken language translation (SLT) can help re-decoding SLT output graphs and improving the overall speech translation performance.

Using robust word confidence measures (from both ASR and MT) to re-decode the SLT graph leads to a significant BLEU improvement (more than 2 points) compared to our SLT baseline (French-English task).

## Introduction

### Context

- Automatic quality assessment of spoken language translation (SLT)
- Also named confidence estimation (CE)
- Pointing out correct parts and errors in a speech translated output

### Useful for

- Interactive speech to speech translation
- Computer-assisted translation (from speech or text)

### Claim

- An accurate CE can also help to improve SLT itself through **search graph re-decoding**

## Formalisation

$x_f$  source signal,  $f = (f_1, f_2, \dots, f_M)$  transcription of  $x_f$ .

$\hat{e} = (e_1, e_2, \dots, e_N)$  translation of  $f$  and  $\hat{e} = \operatorname{argmax}_e \{p(e/x_f, f)\}$

Word Confidence Estimation (WCE) can be seen as finding sequence  $q$  where  $q = (q_1, q_2, \dots, q_N)$  and  $q_i \in \{good, bad\}$

$\hat{e} = \operatorname{argmax}_e \sum_q p(e, q/x_f, f)$   $\hat{e} = \operatorname{argmax}_e \sum_q p(q/x_f, f, e) * p(e/x_f, f)$

$\hat{e} \approx \operatorname{argmax}_e \{ \max_q \{ p(q/x_f, f, e) * p(e/x_f, f) \} \}$

$p(q/x_f, f, e)$  : WCE component

$p(e/x_f, f)$  : SLT component

## SLT Search Graph (SG) Re-decoding

For SLT N-best hypotheses  $e^N = \{e^1, e^2, \dots, e^N\}$ , each  $j$ -th word in the hypothesis  $e^i$ , denoted by  $e_{ij}$ , has a quality label,  $q_{ij}$ .

For all hypothesis  $H_k$  in the SG, the new transition cost is defined by:

$$\text{transition}'(H_k) = \text{transition}(H_k) + \begin{cases} \text{reward}(e_{ij}) & \text{if } q_{ij} = \text{good} \\ \text{penalty}(e_{ij}) & \text{otherwise} \end{cases} \quad (1)$$

with

$$\text{penalty}(e_{ij}) = -\text{reward}(e_{ij}) = \beta * \frac{\text{score}(H_k)}{\#\text{words}(H_k)} \quad (2)$$

## Analysis of SLT Hypotheses

example 1	
$f_{ref}$	une démobilisation des employés peut déboucher sur une démolisation <b>mortifère</b>
$f_{hyp}$	une démobilisation des employés peut déboucher sur une démolisation <b>mort y faire</b>
$e_{hyp}$ baseline	a <b>demobilisation employees</b> can lead to a <b>penalty demoralisation</b>
$e_{hyp}$ with re-decoding	a <b>demobilisation of employees</b> can lead to a <b>demoralization death</b>
$e_{ref}$	<b>demobilization of employees</b> can lead to a <b>deadly demoralization</b>
example 2	
$f_{ref}$	celui-ci a indiqué que l'intervention s'était parfaitement bien <b>déroulée</b> et que les examens post- <b>opérateurs</b> étaient normaux
$f_{hyp}$	celui-ci a indiqué que l'intervention c'était parfaitement bien <b>déroulés</b> , et que les examens post <b>opérateur</b> étaient normaux.
$e_{hyp}$ baseline	it has indicated that the speech <b>that was well</b> conducted, and that the tests were <b>normal post route</b>
$e_{hyp}$ with re-decoding	<b>he</b> indicated that the intervention is <b>very well done</b> , and that the tests <b>after operating were normal</b>
$e_{ref}$	<b>he</b> indicated that the operation went <b>perfectly well</b> and the <b>post-operative tests were normal</b>
example 3	
$f_{ref}$	general motors repousse jusqu'en janvier le plan pour <b>opel</b>
$f_{hyp}$	general motors repousse jusqu' en janvier le plan pour <b>open</b>
$e_{hyp}$ baseline	general motors postponed until january <b>the plan to open</b>
$e_{hyp}$ with re-decoding	general motors puts until january <b>terms to open</b>
$e_{ref}$	general motors postponed until january <b>the plan for opel</b>

Table : Exemples of French SLT output with and w/o re-decoding

## Building an Efficient WCE System

$$\hat{q} = \operatorname{argmax}_q \{p(q/x_f, f, e)\}$$

need training data with quadruplet  $(x_f, f, e, q)$  available,

so instead we compute:  $\hat{q} = \operatorname{argmax}_q \{p_{ASR}(q/x_f, f)^\alpha * p_{MT}(q/e, f)^{1-\alpha}\}$

### $p_{ASR}(q/x_f, f)$

- system described in [1]
- acoustic / graph / linguistic / lexical features
- boosting classifier

### $p_{MT}(q/e, f)$

- system described in [2]
- multiple features and CRF classifier
- using our open-source toolkit available online  
<http://github.com/besacier/WCE-LIG>

## Experimental Setting

### French ASR

- Kaldi toolkit [3]
- HMM/SGMM / 3-grams

### Test corpus

- 2643 French utterances
- 5h of speech
- Cross-validation for tuning / decoding
- Quality labels  $q_i \in \{good, bad\}$  obtained with TERp-A toolkit [4]

### French-English SMT

- Moses toolkit [5]
- 1.6M parallel sent.
- 48M monolingual sent.
- medium-size system
- WMT shared task

task	ASR (WER)	MT (BLEU)	% good	% bad
MT	0%	52.8%	82.5%	17.5%
SLT	26.6%	30.6%	65.5%	34.5%

Table : Baseline MT and SLT performance on 2643 utt.

## Re-decoding Results

### 2-pass (graph re-decoding)

- MT features only
- ASR features only
- joint MT+ASR features

system	baseline	redecoding ASR	redecoding MT	redecoding SLT
WCE feat.	none	$p(q/x_f, f)$	$p(q/f, e)$	$p(q/x_f, f, e)$
Perf.	30.60%	31.12%	31.89%	<b>32.82%</b>

Table : SLT perf. (BLEU) after 2d pass (2643 utt.)

## References

- Laurent Besacier, Benjamin Lecouteux, Ngoc Quang Luong, Kaing Hour, and Marwa Hadjsalah, "Word confidence estimation for speech translation," in *Proceedings of The International Workshop on Spoken Language Translation (IWSLT)*, Lake Tahoe, USA, December 2014.
- Ngoc Quang Luong, Laurent Besacier, and Benjamin Lecouteux, "An Efficient Two-Pass Decoder for SMT Using Word Confidence Estimation," in *European Association for Machine Translation (EAMT)*, Dubrovnik, Croatie, jun 2014.
- Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely, "The kaldi speech recognition toolkit," in *IEEE 2011 Workshop on Automatic Speech Recognition and Understanding*, Dec. 2011.
- Matthew Snover, Nitin Madnani, Bonnie Dorr, and Richard Schwartz, "Terp system description," in *MetricsMATR workshop at AMTA*, 2008.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst, "Moses: Open source toolkit for statistical machine translation," in *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, Prague, Czech Republic, June 2007, pp. 177–180.