

A Large-Scale Evaluation of Pre-editing Strategies for Improving User-Generated Content Translation

Violeta Seretan, Pierrette Bouillon, Johanna Gerlach, *Université de Genève*

The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007-2013) under grant agreement n° 288769

In the **ACCEPT** project...

Consortium	Project Description on Website	Facts
	<p>The ACCEPT project is part of the seventh European framework programme and works on enabling machine translation for the emerging community content paradigm, allowing citizens across the EU better access to communities in both commercial and non-profit environments. There are five countries participating in the project: Switzerland, United Kingdom, Germany, Ireland, and France.</p> <p>ACCEPT: Automated Community Content Editing PoTat Seventh Framework Programme: THEME ICT-2011.4.2(a) Language Technologies</p>	<p>Duration: 36 months Start date: 01/01/2012 Finish date: 31/12/2014 EU funding: €1,825,000 EU project officer: Pierre-Paul Sondag Website: www.accept-project.eu</p>

We want to bring machine translation closer to user communities.

Mission *Help communities share knowledge more effectively across the language barrier.*

Use Cases

<p>support forum community (commercial use case)</p>	<p>volunteer translator community (NGO use case)</p>
--	--

Aim: Help companies engage with their customers across the language barrier. Help NGOs deliver critical information in the right language at the right time.

R&D Axes

<p>Develop minimally-intrusive strategies for pre-editing for statistical machine translation (SMT).</p>	<p>Improve SMT for community content: <ul style="list-style-type: none"> - domain adaptation - linguistic back-off - text analytics. </p>	<p>Develop strategies for post-editing to learn corrections and automate them (in loops).</p>
--	--	---

We developed technology to pre-edit and post-edit community content:

1 Pre-editing plug-in

2 Post-editing plug-in

3 APIs

4 Documentation

5 ... and a lot more

In particular, rules for content checking based on shallow parsing.

Case	<i>norton/Norton</i>
Punctuation	<i>When .../When ...,</i>
Hyphenation	<i>system-based</i>
Spelling	<i>instalation</i>
Spaces	<i>environ.Hier</i>
Redundant words	<i>the the</i>
Homophones	<i>sur 'on' / sûr 'sure'</i>
Word choice	<i>participate to/in</i>
Word form	<i>aimera/aimerait bien</i>
Agreement	<i>problem occur</i>
Ungram. sequences	<i>data the to</i>
Style	<i>C'est nickel/parfait</i>
SMT-specific	<i>have to/must</i>

For details: [1]

English	~50 rules
French	~90 rules

automatic manual

Rule Example acrolinx

```
//example: a dogs
TRIGGER(80) == @det_sg^1 [{{@mod}@noun}}]*! @noun_p1^2
-> { $det_sg, $noun_p1 }
-> { mark : $det_sg, $noun_p1; }

//example: a dogs -> a dog
SUGGEST(10) == $det_sg []* $noun_p1
-> { suggest: $det_sg -> $det_sg, $noun_p1 -> $noun_p1; generateInflections({number="singular"}); }
```

http://atempto.fi.uzh.ch/site/cni2012/slides/lehmannetal_forumdata.pdf

We evaluated the rules:

We conducted a large-scale evaluation campaign with Master's students in Translation to assess the impact of pre-editing on translation quality.

1 Pre-editing

Re : Erreur 3048,3
Tu as le dernier correctif.
Attends 48h. Fais de temps en temps un LiveUpdate manuel et on verra à chaque redémarrage du PC. (...) N'aites crainte, tu es toujours protégé.
Merci de tenir au courant si le problème persiste.

Re : Erreur 3048,3
Vous avez le dernier correctif.
Attendez 48h. Faites de temps en temps un LiveUpdate manuel et on verra à chaque redémarrage du PC. (...) N'ayez crainte, vous êtes toujours protégé.
Veuillez me tenir au courant si le problème persiste.

2 Translation



For details: [2]

3 Evaluation



For details: [3]

Example 8 of 20 in set

Look at the original sentence and two possible translations (all differences are marked in blue). Choose which of the translations is better.

Original Re : Erreur 3048,3
Tu as le dernier correctif.
Attends 48h. Fais de temps en temps un LiveUpdate manuel et on verra à chaque redémarrage du PC.
Il n'est pas impossible que cela vienne des serveurs Symantec qui sont surchargés ou qui sont en maintenance pour LiveUpdate. Il y a des problèmes en ce moment pour dispatcher le dernier correctif.
N'aites crainte, tu es toujours protégé.
Merci de tenir au courant si le problème persiste.

First translation Re : Erreur 3048,3
You have the latest patch.
Please wait 48 hours. Made from time to time a manual LiveUpdate and we will see to each restart of the PC.
It is not impossible that this comes from the Symantec servers that are overloaded or which are in maintenance for LiveUpdate. There are problems in this moment to dispatcher the latest patch.
Have no fear, you are still protected.
Please keep me informed if the problem persists.

Second translation Re : Erreur 3048,3
You have the latest patch.
Look 48 hours. Am from time to time a manual LiveUpdate and we will see to each restart of the PC.
It is not impossible that this comes from the Symantec servers that are overloaded or which are in maintenance for LiveUpdate. There are problems in this moment to dispatcher the latest patch.
Did fear, you are always protected.
Thank you to take aware if the problem persists.

Categories for hypothesis testing: better, same, worse

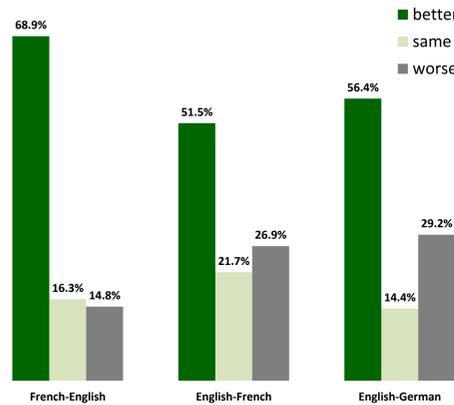
Fleiss k: French-English 0.43, English-French 0.20, English-German 0.38 (up to moderate)

1000 posts avg. length (words) English 93.7 French 78.6 time spent for judging/feedback avg. seconds/post - outliers removed French-English 43 19, English-French 30 20, English-German 52 31

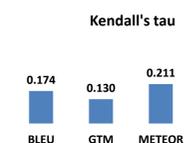
We built reference translations for a subset of the data and computed automatic metric scores.

50 posts French-English BLEU GTM METEOR TER

And this is what we found:



But the gain in quality is not reflected by automatic metric scores (non-significant difference).



There is weak/weak or no correlation with human judgments.

Pre-editing improves translation quality significantly (McNemar test, $p < 0.0001$).

Error analysis showed we need to improve named entity recognition and avoid auto spell check.

Pre-editing is worth the effort.

Pre-editing is relatively fast

English	238
French	261

Words per minute (wpm)

and has a positive impact on post-editing, according to a related study [4]:

- The post-editing time was roughly reduced by half thanks to pre-editing.
- The processing speed increased from 28 to 52 wpm (or to 37 wpm, if we consider the time spent pre-editing).

So, pre-edit your message to reach a wider audience!

The ACCEPT plug-in & middleware can be freely tested and downloaded: www.accept-portal.eu

