

# ComPair

Compare documents in any format, in any language

Lingua et Machina

Corporate Multilingual Content Management

September 2018

# Presentation

- ComPair is a universal document comparer, to be installed on your computer:
  - It compares two versions of a document,
  - If needed in two different file formats,
  - In the same language or in two different languages.
- All common office file formats are accepted, including scanned documents.
- All languages are processed.

# Pricing

- The first month of use is free for any language(s), please feel free to test.
- Thereafter, the cost is 5 € (excl. tax) per language and per month.
- The user can select any number of languages, starting with one.
- Selection of language(s) is done during the license request procedure.
- New language(s) can be added at any time.

# Installation

- Download the software here:

<http://lingua-et-machina.eu/ComPair-setup.exe>

- Install it on your computer,
- Fill the free license request, send it, receive the license, install it.
- You are ready to start, the following slides show the results for:
  - A comparison of two versions of a document,
  - A comparison of a document and its translation.

# Comparing versions

Lingua et Machina ComPair 2.0.2



The translation community has seen a major change over the last five years:  
machine translation has become good enough so that it has become advantageous for translators to post-edit it rather than translate from scratch.  
This is due to recent progress in statistical machine translation, that is, the training of a translation engine with a corpus of existing translations.  
Current enhancement of machine translation (MT) systems from human post-edition (PE) of raw outputs are somewhat efficient yet rather basic:  
the post-edited output is added to the training corpus and the translation model and language model are re-trained, with no clear view of how much has been improved and how much is left to be improved.  
In this approach, only the final PE result is used, no other user feedback on the raw MT quality is provided, such as the cognitive processes of the post-editor or the logging of the post-edition actions he has performed.  
The KEHATH project intends to address these issues in two ways:  
Firstly, leverage advanced machine learning (ML) techniques in the MT+PE loop.  
Our goal is to boost the impact of PE, that is, reach the same performance with less PE or better performance with the same amount of PE.  
In other words, we want to improve machine translation learning curves.  
For this purpose, active learning and reinforcement learning techniques will be proposed and evaluated.  
In the industrial context of KEHATH, we will have to face challenges such as MT systems heterogeneity (statistical and/or rule-based), and ML algorithms scalability to improve a domain-specific MT.  
Secondly, quality prediction (QP) on MT outputs is crucial for translation project managers.  
We have developed over the years a number of confidence estimation and error detection techniques in the laboratory and we will implement and evaluate them in real-world conditions.  
A shared concern will be to work on continuous domain-specific data flows to improve both MT and the performance of indicators for quality prediction.

0 The translation community has seen a major change over the last five years.  
1 This is explained by recent developments in statistical machine translation, that is to say the training of a translation engine from a body of existing translated texts.  
2 Thanks to progress in the training of statistical machine translation engines on corpora of existing translations, machine translation has become good enough so that it has become advantageous for translators to post-edit machine outputs rather than translate from scratch.  
3 However, current enhancement of machine translation (MT) systems from human post-edition (PE) are rather basic:  
4 the post-edited output is added to the training corpus and the translation model and language model are re-trained, with no clear view of how much has been improved and how much is left to be improved.  
5 Moreover, the final PE result is the only feedback used:  
6 available technologies do not take advantages of logged sequences of post-edition actions, which inform on the cognitive processes of the post-editor.  
7 The KEHATH project intends to address these issues in two ways.  
8 Firstly, we will optimise advanced machine learning techniques in the MT+PE loop.  
9 Our goal is to boost the impact of PE, that is, reach the same performance with less PE or better performance with the same amount of PE.  
10 In other words, we want to improve machine translation learning curves.  
11 For this purpose, active learning and reinforcement learning techniques will be proposed and evaluated.  
12 Along with this, we will have to face challenges such as MT systems heterogeneity (statistical and/or rule-based), and ML scalability so as to improve domain-specific MT.  
13  
14 Secondly, since quality prediction (QP) on MT outputs is crucial for translation project managers, we will implement and evaluate in real-world conditions several confidence estimation and error detection techniques previously developed at a laboratory scale.  
15 A shared concern will be to work on continuous domain-specific data flows to improve both MT and the performance of indicators for quality prediction.

# Comparing translations

Lingua et Machina ComPair 2.0.2



The translation community has seen a major change over the last five years.	0	La communauté de la traduction a vu un changement majeur au cours des cinq dernières années :
This is explained by recent developments in statistical machine translation, that is to say the training of a translation engine from a body of existing translated texts.	1	la traduction automatique est devenue suffisamment bonne pour qu'il devienne plus avantageux pour des traducteurs de post-éditer une traduction machine plutôt que de traduire directement.
Thanks to progress in the training of statistical machine translation engines on corpora of existing translations, machine translation has become good enough so that it has become advantageous for translators to post-edit machine outputs rather than translate from scratch.	2	Ceci s'explique par les évolutions récentes en traduction automatique statistique, c'est-à-dire l'entraînement d'un moteur de traduction à partir d'un corpus de textes traduits existants.
However, current enhancement of machine translation (MT) systems from human post-edition (PE) are rather basic:	3	Les systèmes actuels d'amélioration de traduction automatique à partir de la post-édition de sorties brutes sont relativement efficaces mais assez frustrés :
the post-edited output is added to the training corpus and the translation model and language model are re-trained, with no clear view of how much has been improved and how much is left to be improved.	4	le texte post-édité est ajouté au corpus d'entraînement et le modèle de traduction et le modèle linguistique sont entraînés à nouveau, sans vision précise de ce qui a été amélioré ni de ce qui reste à améliorer.
Moreover, the final PE result is the only feedback used, available technologies do not take advantages of logged sequences of post-edition actions, which inform on the cognitive processes of the post-editor.	5	Dans cette démarche, seul le résultat brut de la post-édition est utilisé, aucune autre information n'est mise à profit, comme par exemple les processus cognitifs du post-éditeur ou l'enregistrement des actions qu'il a effectuées.
The KEHATH project intends to address these issues in two ways:	6	Le projet KEHATH propose de revoir la boucle traduction automatique / post-édition de deux façons :
Firstly, we will optimise advanced machine learning techniques in the MT+PE loop.	7	D'une part, mettre en œuvre des techniques avancées en apprentissage automatique.
Our goal is to boost the impact of PE, that is, reach the same performance with less PE or better performance with the same amount of PE.	8	Notre objectif est de renforcer l'impact de la post-édition, c'est-à-dire atteindre les mêmes performances avec moins de post-édition ou alors atteindre de meilleures performances avec la même quantité de post-édition.
In other words, we want to improve machine translation learning curves.	9	En d'autres termes, nous souhaitons améliorer la courbe d'apprentissage des systèmes de traduction automatique spécialisés par domaine.
For this purpose, active learning and reinforcement learning techniques will be proposed and evaluated.	10	Pour cela, des techniques d'apprentissage actif ou par renforcement seront proposées et évaluées.
Along with this, we will have to face challenges such as MT systems heterogeneity (statistical and/or rule-based), and ML scalability so as to improve domain-specific MT.	11	Le contexte industriel de KEHATH nous confrontera par ailleurs aux défis de l'hétérogénéité des systèmes (statistiques ou par règles) et du passage à l'échelle des algorithmes d'apprentissage automatique.
Secondly, since quality prediction (QP) on MT outputs is crucial for translation project managers, we will implement and evaluate in real-world conditions several confidence estimation and error detection techniques previously developed at a laboratory scale.	12	D'autre part, la prédiction de qualité de traduction automatique est d'une utilité cruciale pour les chefs de projets de traduction, nous avons développé au fil du temps plusieurs techniques d'estimation de confiance et de détection d'erreurs en laboratoire, nous comptons les mettre en œuvre et les évaluer en conditions réelles.
A shared concern will be to work on continuous domain-specific data flows to improve both MT and the performance of indicators for quality prediction.	13	Nous partageons la conviction que ce travail doit s'appliquer à un flot continu de textes spécialisés par domaine, de façon à démontrer clairement l'amélioration de la traduction automatique et la performance des indicateurs de prédiction de qualité.
The overall goal of the KEHATH project is straightforward:	14	Le but du projet KEHATH est simple :
gain additional machine translation performance as fast as possible in each and every new industrial translation project, so that post-	15	gagner de la qualité de traduction automatique le plus vite possible pour chaque nouveau projet industriel de traduction, de façon à ce

# FAQ

- I received printed documents, how do I compare them?
  - Download and install the free software tesseract-ocr:
  - <https://digi.bib.uni-mannheim.de/tesseract/tesseract-ocr-setup-3.05.02-20180621.exe>
  - Scan the document(s) and save as PDF file(s),
  - Ask ComPair to process the scan(s), it calls tesseract-ocr automatically.
- Can I export the table in Word or Excel?
  - No, but there is an icon for copying the table to your clipboard,
  - You can then paste it in a Word or Excel file.
- Are there other versions of ComPair?
  - Yes, it is available online in SaaS mode (Software as a Service),
  - Or we can install it as a dedicated server for a company.

# FAQ: Translator's corner

- I'm a translator, how do I align translated documents?
  - Just select the languages and launch ComPair with your documents,
  - Then click on the "TMX" icon to export it as a standard memory file.
- Can I edit the alignments?
  - Yes, use is the command icon left of the segment number shown in slide 6.
- Can I merge several alignments in one memory?
  - Yes, please call us.



# More info?

- Contact : François Brown de Colstoun
  - +33 6 80 95 94 39
  - Skype : f\_brown\_de\_colstoun
  - [fbc@lingua-et-machina.com](mailto:fbc@lingua-et-machina.com)
- <http://www.lingua-et-machina.com/>