



# DiCoP project and DiCoP-Text corpus for the enrichment of Language Models and Automatic Translation

Lian Chen, Wenjun Sun, Flora Badin

## ► To cite this version:

Lian Chen, Wenjun Sun, Flora Badin. DiCoP project and DiCoP-Text corpus for the enrichment of Language Models and Automatic Translation. XXI EURALEX International Congress, Oct 2024, Cavtat, Croatia. . hal-04725787

HAL Id: hal-04725787

<https://hal.science/hal-04725787v1>

Submitted on 8 Oct 2024

**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.

# DiCoP project and DiCoP-Text corpus for the enrichment of Language Models and Automatic Translation

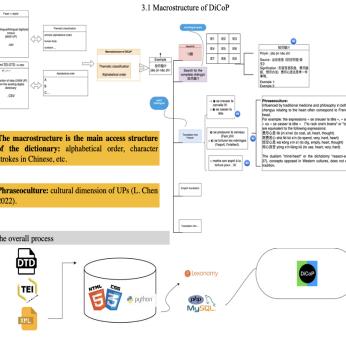
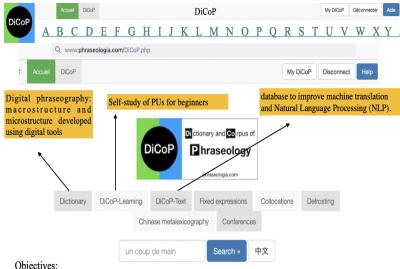


CHEN Lian 陈恋  
loselychen@gmail.com

SUN Wenjun 孙文君  
wenjun.sun@univ-lr.fr

BANDIN Flora  
flora.badin@univ-orleans.fr

La Rochelle  
Université

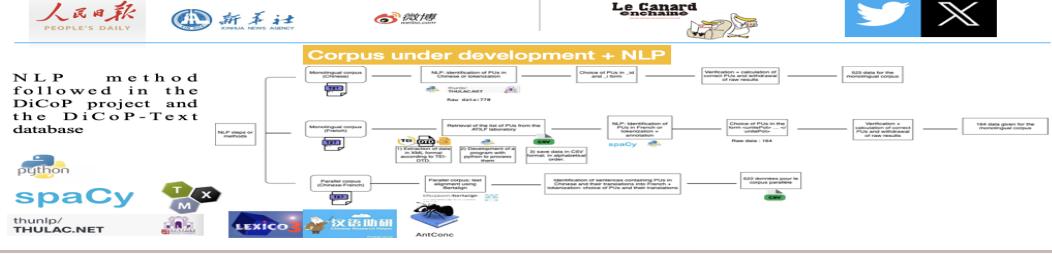


**Objectives:**  
In the digital age and contemporary artificial intelligence (AI), the DiCoP project (Dictionary and Corpus of Phraseology available on phraseologia.com, under development) aims to develop a digital dictionary of multilingual phrases, including French-Chinese-Chinese-French, and eventually multilingual, based on a corpus of phraseological units and associated databases to determine their frequency of use (in newspapers, literary works, school textbooks, etc.) in practice and therefore their vitality, to improve their automatic translation, and with the aim of facilitating access to phraseological units.

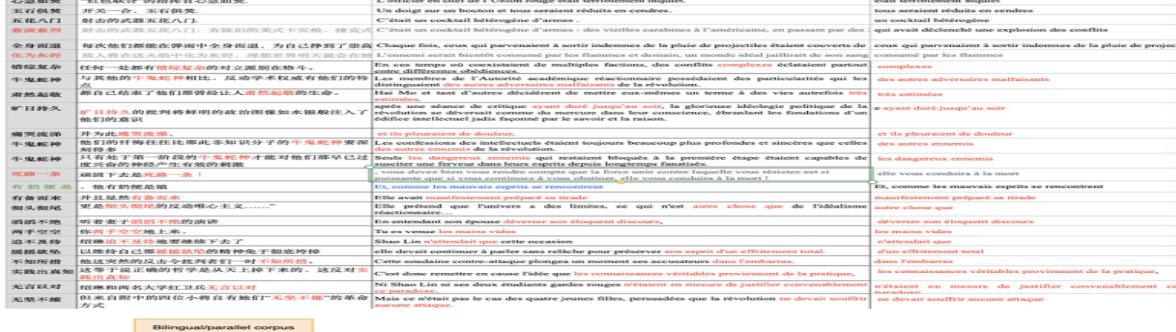
**The collection of texts and digital processing are currently underway. The choice of the corpus should reflect the linguistic diversity of the language and therefore be broad enough to ensure adequate representation of the PUs and improve the accuracy and reliability of the DiCoP.**

- represent a variety of genres (literature, poetry, speeches, novels, science fiction, etc.) to find a balance between formal and informal language (official academic texts, as well as dialogues, everyday conversations, etc.)
- modern resources (20th century and later)
- references already accessible in digital form (or already OCRed) : HTML, XML, TEI, etc.

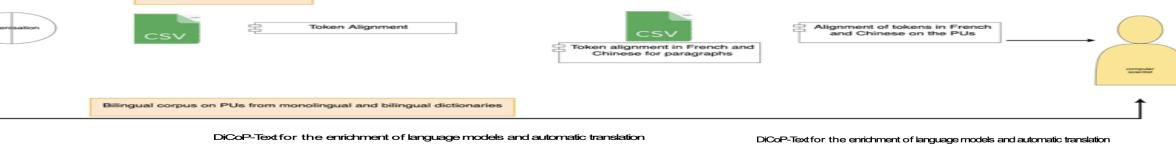
**Newspapers: “défigement” (活用 huóyòng in Chinese), a phenomenon that can be seen, for example, in wordplay in headlines, slogans and advertisements.**



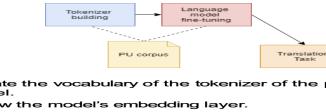
**DiCoP-Text:** Phraseology is particularly difficult in translation, as it is influenced by linguistics, culture and stylistics, and strongly reflects the translator's choices and translation technique. DiCoP-Text will be a database (monolingual, parallel, contrastive, multilingual corpus) to determine the frequency of use of PUs, in order to verify their vitality in practice. It should make it possible to easily scrutinize the use of PUs in translations, for lexicometric studies and contribute to scientific research in the field of automatic phraseological translation and NLP (natural language processing).



**Table: Sentences containing the 549 PUs, and their translation**



DiCoP-Text for the enrichment of language models and automatic translation



- ▶ Update the vocabulary of the tokenizer of the pre-trained model.
- ▶ Renew the model's embedding layer.
- ▶ Fine-tune the model with all PUs and sentence training set.
- ▶ Validate the model with sentence test set.

DiCoP-Text for the enrichment of language models and automatic translation

- ▶ **Basic language models:** Mbart (Tang et al., 2020), M2m100 (Fan et al., 2021), Nllb (Costa-Jussà, Marta R. et al., 2022), and Mrebel (Cabot et al., 2023).
- ▶ **Metrics:** SentiBLE (Postolache, 2018).
- ▶ **Enriched corpora:** 50k data items, divided the corpus into training, validation, and test sets at a ratio of 6:2:2.
- ▶ **Tokenizers:** 409 new tokens were added for each to tokenizer.

<sup>1</sup>Model weights are from: <https://huggingface.co/Babelscape/mrebel-large>, <https://huggingface.co/facebook/mbart-large-50-many-to-many-mmmt>, <https://huggingface.co/facebook/nllb-200-distilled-600M>, and <https://huggingface.co/facebook/m2m100-418Bm>.

DiCoP-Text for the enrichment of language models and automatic translation

**Table:** The experimental results.

Model	Original Model	Fine-tuned Model	Improved Value
Mbart	2,0686	16,7191	14,6505
M2m100	16,7198	16,2120	11,4122
Nllb	4,6345	17,3355	12,6040
Mrebel	0,0131	7,1798	7,1667

- ▶ Introduction of PUs can help translation models better translate PUs.
- ▶ Although each fine-tuned model outperformed its original model, the speed was lower.
- ▶ Volume of the training data.
- ▶ Improving the model.
- ▶ Expanding corpus size.
- ▶ Improving the Chinese-French translation ability of models with the Chinese-French corpus of ordinary texts.

## Conclusion

The first evaluation of the NLP tools in the DiCoP-Text project provided a detailed overview of the effectiveness of the DiCoP-Text corpus and the improved LMs. Our proposal aimed to improve LMs by integrating more fixed expressions and refining linguistic models for more accurate identification and translation of PUs. However, room for improvement exists:

1) In the future, we envision broader applicability of our DiCoP project. Indeed, expansion to other languages would strengthen its relevance and applicability across the IT language community. We also aim to provide more information, including details concerning user interfaces, accessibility, and integrating user feedback into ongoing development.

2) We studied the effect of fine-tuning the translation model using a PU corpus. This approach involved updating the tokenizer, training the model to integrate these PUs, refining the model from sentences containing them, and testing based on sentence-level data. The results indicated that fine-tuning the model with PUs could improve its translation capacity. However, given the limitations of the model and the corpus volume, additional efforts are necessary to refine its translation capacity. Thus, our future research will focus on expanding the corpus and improving the model's Chinese-French translation capability.