



HAL
open science

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 SUN Wenjun 孙文君 BANDIN Flora
 loselychen@gmail.com wenjun.sun@univ-lr.fr flora.badin@univ-orleans.fr

The screenshot shows the DiCoP website with a search bar and various menu options like 'Dictionary', 'DiCoP-Learning', 'DiCoP-Text', 'Fixed expressions', 'Collocations', and 'Defining'. It also features a 'Photo' section and a 'Video' section.

5. Corpus of phraselogology: DiCoP-Text

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.

The diagram illustrates the 'Corpus under development + NLP' process. It shows a flow from source materials (newspapers, etc.) through NLP tools like spaCy, Thuluc.Net, and Lexico to a final corpus. It also mentions 'Le Canard Enchaîné' as a source.

DiCoP-Text: Phrasology 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 phrasological translation and NLP (natural language processing).

NLP: Identify PUs in Chinese or token

```

sans nom> EI-ATILF.py import thuluc.py
1 #终端命令: cd ~/Documents, 然后python3 import\ thuluc.py
2
3 import thuluc
4 # 1. 识别短语: 对“二八集团”总部大楼的攻击已持续了两天。他们的旗帜在
5 # 大楼顶上显示了一个新的身影。旗杆下“二八”的大旗
6 # 和“红色联合”的战士们欢呼起来。几个人冲到了楼下。楼下的
7 # 战士们
8
9 f=open("texte_ch.txt","r")
10 f_out=open("text_ch_tok.txt","w")
11 for line in f:
12     f_text = thuluc.cut(line, text=True)
13     f_out.write(f_text+"\n")
14
15 f.close()
16 f_out.close()
    
```

Liu Cixin's *Three Body Problem*. That's 186,079 words.

Total		Chinese PUs identified by Thuluc	
Success rates	771	65.43% (504/771)	
Error rate		54.57% (267/771)	
		Form_1 (204 PUs)	Correct Error
Number of hits	321	246	163 21
Percentage	56.61%	49.59%	69.71% 10.29%

The screenshot shows Python code for tokenization and a visualization of the corpus data, including a bar chart and a word cloud.

Table: Sentences containing the 549 PUs, and their translation

mot chinois	mot français
心无旁骛	心无旁骛
专心致志	专心致志
聚精会神	聚精会神
全神贯注	全神贯注
目不转睛	目不转睛
屏气凝神	屏气凝神
专心致志	专心致志
聚精会神	聚精会神
全神贯注	全神贯注
目不转睛	目不转睛
屏气凝神	屏气凝神
专心致志	专心致志
聚精会神	聚精会神
全神贯注	全神贯注
目不转睛	目不转睛
屏气凝神	屏气凝神

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

The diagram shows the pipeline for model fine-tuning: Bilingual/parallel corpus (CSV) -> Token Alignment -> Alignment of tokens in French and Chinese for paragraphs -> Translation Task. Below it is a table of experimental results.

Model	Original Model	Fine-tuned Model	Improved Value
Mbart	2,0685	16,7181	14,5505
M2m100	4,9078	16,2120	11,3042
Nllb	4,8345	17,1385	12,5040
Mnrel	0,0131	7,1798	7,1667

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.