



UNIVERSITÉ DE BELGRADE
FACULTÉ DE PHILOLOGIE
Département d'études romanes

avec le concours de



Journées serbes des DICTIONNAIRES



Ministère de la recherche, du développement technologique
et des innovations de la République de Serbie
& Agence universitaire de la Francophonie (AUF)

COLLOQUE INTERNATIONAL FRANCOPHONE

(Méta)phraséographie et phraséomatique : exploration et innovation dans le projet DiCoP

CHEN Lian 陈恋

loselychen@gmail.com

LLL, Université d'Orléans



Membre associé aux laboratoires LT2D-Centre Jean Pruvost, Cergy Paris Université



Phraséologie et phraséoculturologie

Unité phraséologique (UP)

shúyǔ 熟语

Unité phraséologique : « [...] **séquence polylexicale** constituée de deux ou plusieurs mots graphiques catégoriellement liés, contigus ou non. Les UP se caractérisent linguistiquement par : (i) **un certain degré de fixité syntaxique** (blocage des propriétés transformationnelles et ordre des constituants inaltérable) ; et/ou (ii) un certain degré **de figement sémantique** (non-compositionnalité au moins partielle) ; et/ou (iii) un certain degré **de figement lexical** (restriction paradigmatic) ; et / ou (iv) une contrainte sur l'emploi en situation de communication. » (BOLLY, 2011 : 28)

Ex : collocations, proverbes,
expressions idiomatiques

Ex: 谚语 *yànyǔ*, 成语 *chéngyǔ*, 惯用语
guànyòng yǔ, 歇后语 *xiēhòuyǔ*, etc.



L'étude de ce phénomène culturel idiomatique spécifique en phraséologie s'appelle :
« **phraséoculturologie** » 熟语文化学 *shú yǔwén huàxué* (L. Chen 2022).



(Méta)phraséographie et phraséomatique

Métaphraséographie (L. Chen 2022)

une (sous-)discipline dont l'objectif est l'étude des types de dictionnaires UP et des méthodes par lesquelles ils sont construits, ainsi qu'un objet de réflexion et de recherche sur la conception phraséographique.

→ théorique

Phraséographie (L. Chen 2022, 2023)

une branche de la phraséologie appliquée dont l'objet d'étude est le développement de collections, de glossaires et de dictionnaires d'UP

→ appliqué

Phraséomatique (L. Chen 2023)

À l'ère numérique et de l'intelligence artificielle (IA) contemporaine, le domaine de la phraséologie informatisée se trouve confronté à des défis importants, mettant en lumière l'émergence d'une discipline particulière : **la phraséomatique, ou la phraséologie informatique** (Chen 2023)

Unités phraséologiques (UP)

VS

Shúyǔ 熟语

Défigements /Jeux de mots
tout feu tout **femme**, défigé de [être]
tout feu tout **flamme**

熟语活用 *Shúyǔ huóyòng*

默默无**蚊** *mòmòwúwén* défigé de 默默无**闻** *mòmòwúwén*
faire quelque chose sans être entendu/rester inconnu/tomber dans l'oubli.
Le caractère « 蚊 wén » (moustique) a remplacé le « 闻 wén » (entendre).

Collocations

gravement handicapé
 déficientaire
 brûlé
 préjudiciable
 intoxiqué
 blessé

惯用语 *Guànyòng yǔ* (expression usuelle)

戴高帽 *dàigāomào* (porter, haut, chapeau) : porter un chapeau haut/
flatter quelqu'un/lécher les bottes, encenser

歇后语 *Xiēhòuyǔ* (calembour à tiroir) :

和尚打伞 – 无发(法)无天 *héshàngdǎsǎn – wúfà(fǎ)wú tiān* : Un moine
tient un parapluie – il n'a ni cheveu (discipline) ni ciel/ni dieu/il est
sans respect ; sans foi ni loi

Proverbes

Le chat sauvage ne sait pas où il a pris
une poule.

谚语 *Yànyǔ* (proverbe)

虚心使人进步，骄傲使人落后 *xūxīn shǐ rén jìnbù, jiāo'ào shǐ rén luòhòu* : L'orgueil provoque l'échec l'humilité engendre le progrès.

Expressions idiomatiques

avoir le coeur sur la main

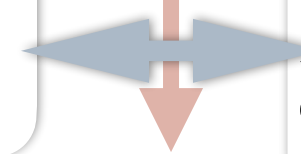
成语 *Chéngyǔ* (expression toute faite)

龙腾虎跃 *lóngténg-hǔyuè* (le dragon, sauter, le tigre, sauter): le dragon
et le tigre bondissent/être plein de dynamisme, travailler ferme et viser
haut.

D
e
g
r
é

d
e

f
i
g
e
m
e
n
t



2. Contexte du projet

Les ressources numériques propres aux UP restent majoritairement monolingues

在线成语词典

提供成语的汉语拼音、释义、出处、示例的在线查询

ABCDEFGHIJKLMNOPQRSTUVWXYZ [按拼音查]

输入查询内容: 词目 ▾

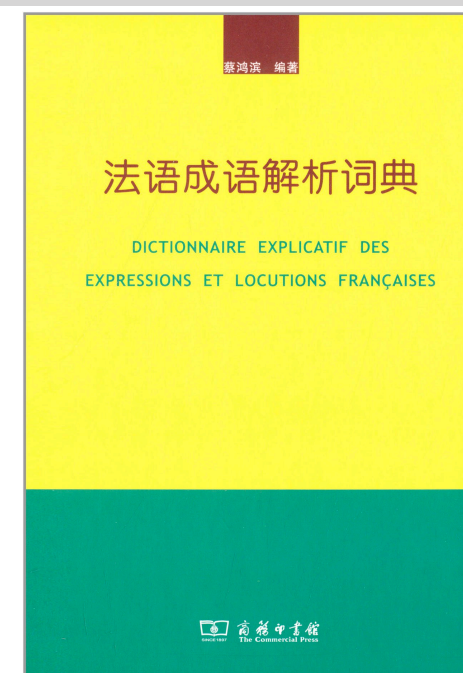
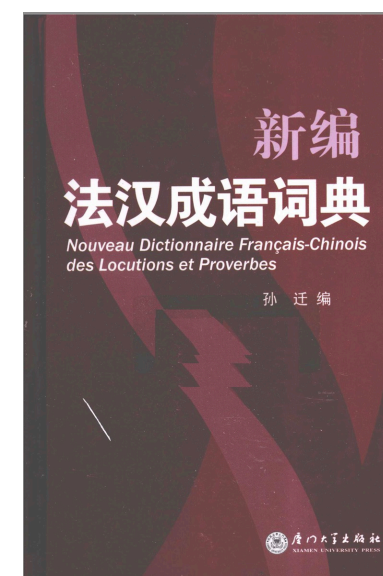
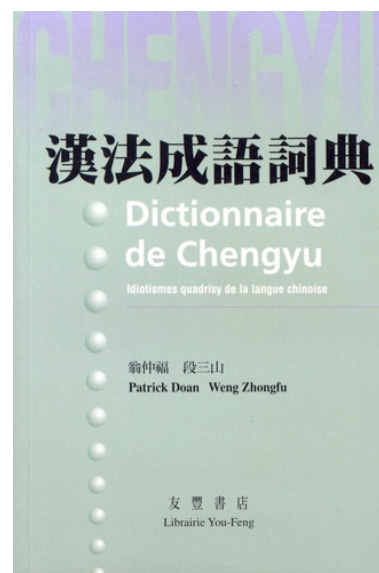


首页 成语词典 成语谜语 歇后语 成语故事 成语接龙 成语文章 看图猜成语

Expressio.fr

lintern@ute

Dictionnaires papiers bilingues chinois-français et français-chinois



2. Contexte du projet

Français- Portugais

CNRTL Centre National de Ressources Textuelles et Lexicales

Ortolang Outils et Ressources pour un Traitement Optimisé de la LANGue

■ Accueil ■ Portail lexical ■ Corpus ■ Lexiques ■ Dictionnaires ■ Métalexicographie ■ Outils ■ Contact



Dictionnaire d'expressions idiomatiques

Français - Portugais - Français



accueil

parcourir

rechercher

classement par concepts

télécharger

aide

■ Recherche d'une expression ou d'un terme

en français : un coup de main

rechercher

en portugais :

rechercher

Error connecting to mysql

https://www.cnrtl.fr/dictionnaires/expressions_idiomatiques/recherche.php

Français- Chinois ?
Chinois -Français ?

DiCoP

AccueilDiCoP

3. DiCoP

My DiCoPDéconnecterAide

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

www.phraseologia.com/DiCoP.php

AccueilDiCoP

My DiCoPDisconnectHelp

Phraséographie numérique : macrostructure et microstructure élaborées à l'aide d'outils numériques

auto-apprentissage des UP pour les débutants

base de données pour améliorer la traduction automatique et le Traitement Automatique de langue (TAL).

DiCoP

Dictionary and Corpus of Phraseology

phraseologia.com

Dictionary

DiCoP-Learning

DiCoP-Text

Fixed expressions

Collocations

Defrosting

Chinese metalexicography

Conferences

un coup de main

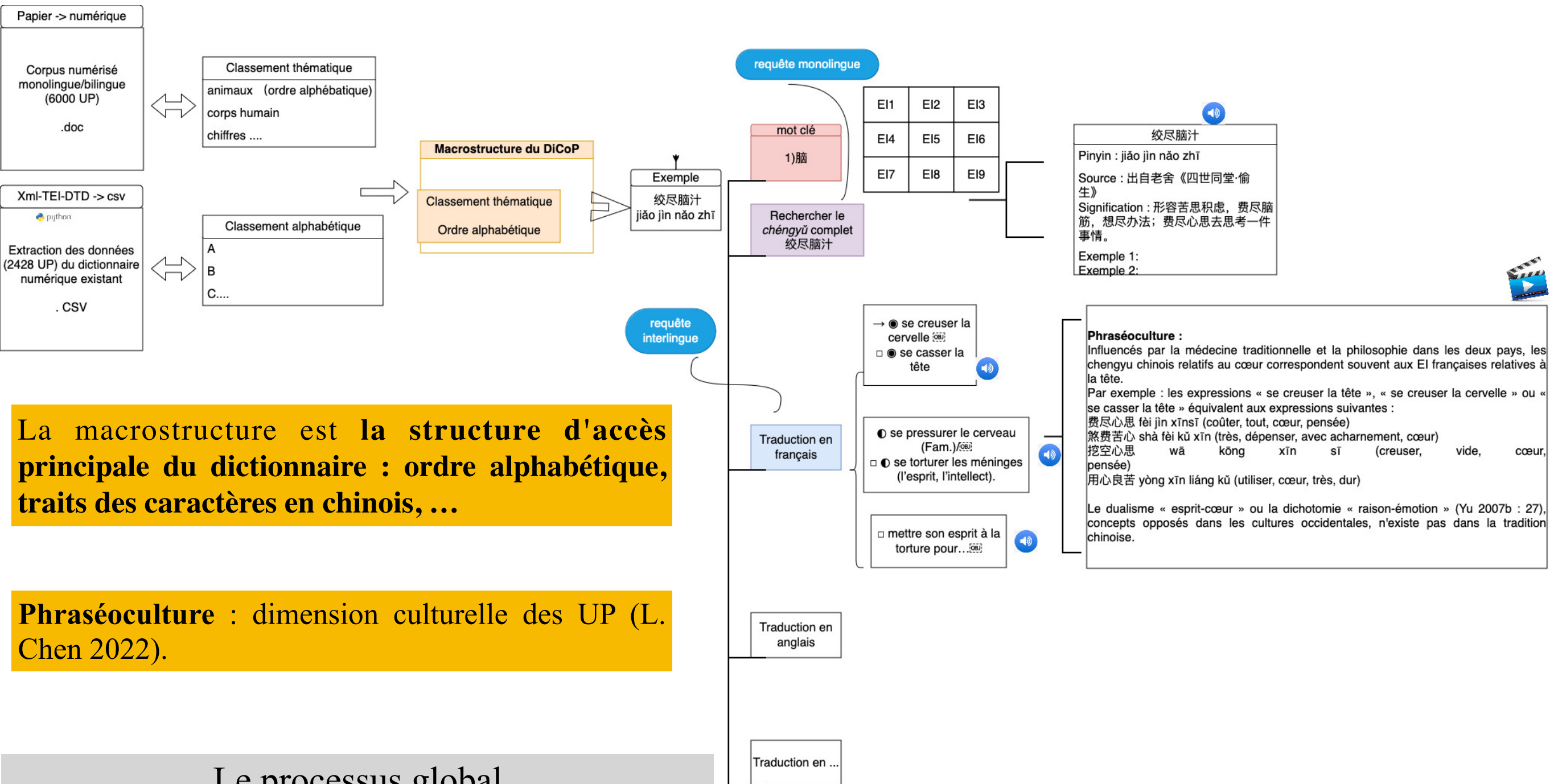
Search »

中文

Objectifs :

À l'ère numérique et de l'intelligence artificielle (IA) contemporaine, le projet DiCoP (**Dictionnaire et Corpus de la Phraséologie** (disponible sur phraseologia.com, en cours de développement) vise à développer **un dictionnaire numérique de phraséologie multilingue** (dans un premier temps français-chinois/chinois-français, et à terme multilingue), **basé sur un corpus d'unités phraséologiques et des bases de données associées pour déterminer leur fréquence d'utilisation** (dans les journaux, les œuvres littéraires, les manuels scolaires, etc.) **dans la pratique et donc leur vitalité**, pour améliorer leur traduction automatique, et dans le but de faciliter l'accès aux unités phraséologiques.

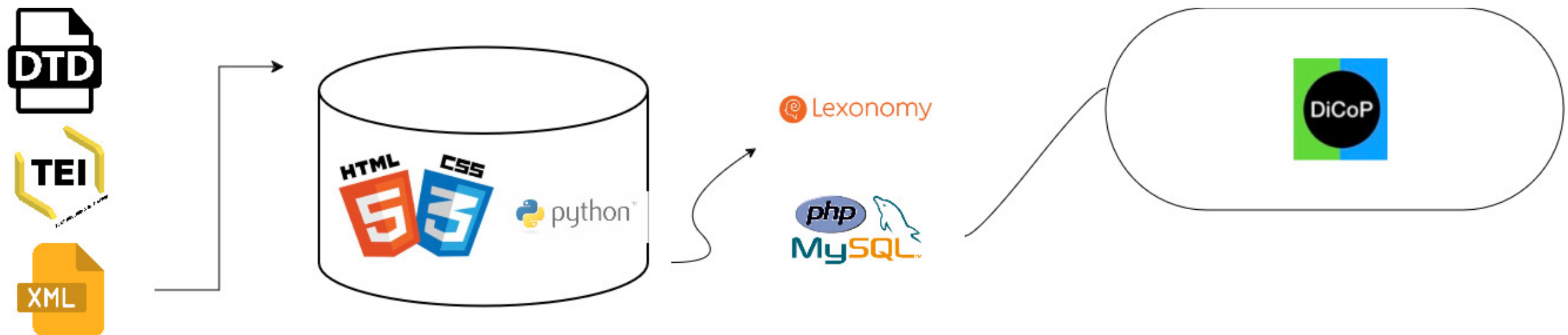
3.1 Macrostructure du DiCoP



La macrostructure est la structure d'accès principale du dictionnaire : ordre alphabétique, traits des caractères en chinois, ...

Phraséoculture : dimension culturelle des UP (L. Chen 2022).

Le processus global



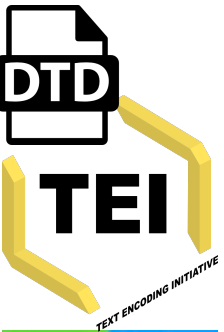
Microstructure : l'ensemble des informations ordonnées de chaque article, réalisant un programme d'informations constant pour tous les articles, et qui se lisent horizontalement à la suite de l'entrée

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE product SYSTEM "product.dtd">

<entry>
  <word>
    <simplified> <b> 回头是岸 </b></simplified><br style="display:block;"/>
    <traditional><b> 回頭是岸 </b></traditional><br style="display:block;"/>
    <translation>retourner, tête, être, rivage</translation><br/>
  </word><br/>
  <pinyin font="italic">huí tóu shì àn</pinyin> <audio src="audio.mp3" /><br/>
  <definition><br/>
    <literal> si on tourne la tête, on retrouve le rivage</literal><br/>
    <implicit>il suffit de revenir sur son erreur pour recouvrer le droit chemin/□ revenir dans
    le droit chemin/il n'est jamais trop tard pour se corriger (pour s'amender, pour bien
    faire).
  </implicit><br/>
  </definition><br/>
  <source> Premier acte de 《度翠柳》 (Du Cuiliu), dynastie des Yuan, anonyme </source><br/>
  <history>Le texte originel est «<0xa0>苦海无边, 回头是岸<0xa0>Kǔ hǎi wú biān, huí tóu shì àn » :
  La mer de souffrances est immense, mais il suffit de reculer pour trouver terre ferme. Cette
  expression vient de la conception bouddhiste, selon laquelle une personne coupable qui
  accepte de se repentir, va pouvoir arriver de « l'autre côté » et libérer son âme en expiant
  ses péchés. Cette croyance est assez semblable à celle de la religion chrétienne qui insiste
  sur la nécessité du repentir pour « sauver son âme ».</history><br/>
  <video src="video.mp4" /> <br/>
  <example> 奉劝那些赌徒们, 赌海无边, 回头是岸, 还是早些醒悟吧! Je voudrais dire à ces joueurs que la
  mer des jeux d'argent est sans limites, et que le retour en arrière est le rivage ; ils
  feraient donc mieux de se réveiller plus tôt !</example><br/>
  <example> 爱上一个这样的人多不幸, 回头是岸 。C'est malheureux de tomber amoureux d'une telle
  personne, mais il n'est jamais trop tard pour se corriger et il est facile de se retourner.</
  example><br/>
</entry>
```



L'Extensible Markup Language, généralement appelé XML, « langage de balisage extensible » en français, est un métalangage informatique de balisage générique qui est un sous-ensemble du Standard Generalized Markup Language.



3.2 Microstructure informatisée du DiCoP

La norme Extensible Markup Language (XML) basée sur la Document Type Definition - Text Encoding Initiative (DTD-TEI) en XML pour les dictionnaires a été retenue, notamment dans l'optique d'assurer la pérennité de la ressource.

[Accueil](#)[DiCoP](#)[My DiCoP](#)[Déconnecter](#)[Aide](#)[A](#) [B](#) [C](#) [D](#) [E](#) [F](#) [G](#) [H](#) [I](#) [J](#) [K](#) [L](#) [M](#) [N](#) [O](#) [P](#) [Q](#) [R](#) [S](#) [T](#) [U](#) [V](#) [W](#) [X](#) [Y](#) [Z](#)

```
1 <!ELEMENT woxinchangdan DEF_CONTENU>
2 <!ATTLIST woxinchangdan TYPE OBLIGATION VALEUR_DEFAULT>
3
4 <!ELEMENT entry (word, pinyin, audio, definition, source, history, video, examples)>
5 <!ELEMENT word (simplified, traditional, translation)>
6 <!ELEMENT definition (literal, implicit)>
7 <!ELEMENT examples (example+)>
8
9 <!ELEMENT simplified (#PCDATA)>
10 <!ELEMENT traditional (#PCDATA)>
11 <!ELEMENT translation (#PCDATA)>
12 <!ELEMENT pinyin (#PCDATA)>
13 <!ELEMENT audio EMPTY>
14 <!ELEMENT literal (#PCDATA)>
15 <!ELEMENT implicit (#PCDATA)>
16 <!ELEMENT source (#PCDATA)>
17 <!ELEMENT history (#PCDATA)>
18 <!ELEMENT video EMPTY>
19 <!ELEMENT example (#PCDATA)>
20
21 <!ATTLIST audio src CDATA #REQUIRED>
22 <!ATTLIST video src CDATA #REQUIRED>
```


3.2 Microstructure informatisée de DiCoP

```
pf.xml
1 <?xml version="1.0" encoding="UTF-8"?>
2 <!DOCTYPE TEI SYSTEM "dico_tei.dtd">
3 <TEI>
4   <teiHeader>
5     <fileDesc>
6       <titleStmt>
7         <title>Dictionnaire électronique d'expressions idiomatiques</title>
8         <author>Claudia Maria Xatara</author>
9         <respStmt>
10          <name>Etienne Petitjean</name>
11        </respStmt>
12      </titleStmt>
13      <publicationStmt>
14        <distributor>
15          <name>ATILF UMR7118 - CNRS/UN2/UHP</name>
16        </distributor>
17      </publicationStmt>
18      <sourceDesc>
19        <p/>
20      </sourceDesc>
21    </fileDesc>
22  </teiHeader>
```

```
dico_tei.dtd
1 <?xml version="1.0" encoding="UTF-8"?>
2
3
4 <!ENTITY % _type "(reg| hint | gram | syn | ant | freq | dom | style)">
5
6 <!ELEMENT TEI (teiHeader, text)>
7 <!ELEMENT teiHeader (fileDesc)>
8 <!ELEMENT fileDesc (titleStmt, publicationStmt, sourceDesc)>
9 <!ELEMENT titleStmt (title, author, respStmt) >
10 <!ELEMENT respStmt (name) >
11 <!ELEMENT title (#PCDATA) >
12 <!ELEMENT author (#PCDATA) >
13 <!ELEMENT publicationStmt (distributor) >
14 <!ELEMENT distributor (name) >
15 <!ELEMENT name (#PCDATA) >
16 <!ELEMENT p (#PCDATA) >
17 <!ELEMENT sourceDesc (p) >
18 <!ELEMENT text (body) >
19 <!ELEMENT body (entry+) >
20 <!ELEMENT entry (form, gramGrp, sense+, usg*)>
21 <!ELEMENT form (orth) >
22 <!ELEMENT orth (#PCDATA) >
23 <!ELEMENT gramGrp (pos, gen*, number*) >
24 <!ELEMENT pos (#PCDATA)>
25 <!ELEMENT gen (#PCDATA) >
26 <!ELEMENT number (#PCDATA) >
27 <!ELEMENT sense (def, usg*, cit, usg, xr*, cit+) >
28 <!ATTLIST sense n CDATA #IMPLIED >
29 <!ELEMENT def (#PCDATA) >
30 <!ELEMENT cit ((quote, bibl) | (tr, ref*, usg*)) >
31 <!ATTLIST cit type (example|translation) #REQUIRED >
32 <!ELEMENT quote (#PCDATA |oVar)* >
33 <!ELEMENT oVar (#PCDATA) >
34 <!ELEMENT bibl (ref, date) >
35 <!ELEMENT ref (#PCDATA) >
36 <!ATTLIST ref target CDATA #IMPLIED >
37 <!ELEMENT date (#PCDATA) >
38 <!ELEMENT usg (#PCDATA)>
39 <!ATTLIST usg type %_type; #REQUIRED>
40 <!ELEMENT xr (ref, usg*, ref?) >
41 <!ATTLIST xr type (ant | syn) #REQUIRED >
42 <!ELEMENT tr (#PCDATA) >
43
```

Enfin, nous utilisons le format TEI-DTD pour faciliter l'utilisation dans d'autres langages, car il s'agit d'une norme universelle qui s'applique à différents types de textes et à des langues différentes.



<https://lexonomy.elex.is/>

Un système de rédaction de dictionnaires en ligne pour rédiger et publier des dictionnaires.

un maximum de 500 entrées



PHP : Hypertext Preprocessor, plus connu sous son sigle PHP, est **un langage de programmation libre, principalement utilisé pour produire des pages Web dynamiques via un serveur web.**

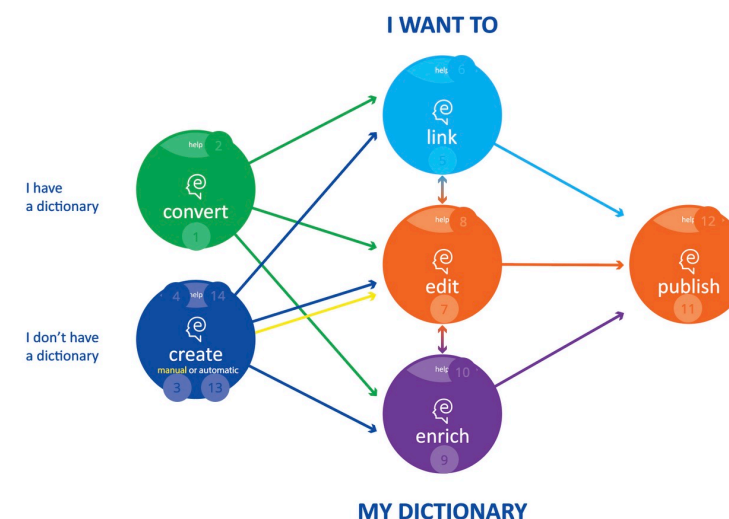
MySQL : un système de gestion de bases de données relationnelles.



langage de programmation



<https://elex.is/>



Dictionary and Corpus of Phraseology (DiCoP)

NEW

ID

13



EDIT

CLONE

DELETE



grue (n.f.) Faire le pied de grue « attendre debout » [exp.] v. 鹄立, 久等

- 1 该表达式来源于17世纪的固定表达faire la jambe de grue和 faire de la grue (和16世纪的faire la grue)。在 Faire de la grue中, grue这个名词有动词“等待”的含义 (如在诗人 Maurice Scève作品中), 但是 Bonaventure des Périers认为, 在16世纪的语境中faire de la grue中的grue一词也含有鸟类:鹤 的比喻义。
- 2 在 16 世纪, 鹤的比喻用途通常是贬义的(être grue 愚蠢, suivre la multitude comme les grues [Calvin], s'en aller comme des grues [ibid.] 像鹤一样跟随众人, 没有自己的想法。以及 Le Roux 在 1752 年给出的变体:être planté comme une grue, 与 être planté comme un sot同义:一动 不动地站着, 等待很久, 久等。最后, “妓女”的意思来自faire le pied de grue, 人行道上妓女等 待客人的形象比喻。

Et d'indiquer des yeux les dizaines de Noires et de mûlatresses qui font, tout près, le pied de grue en se pavanant lascivement. (Le Point, 18-03-1995)

并使眼色表明附近几十个黑人妇女和黑白混血女人正在作出淫荡姿态等客上门呢。

Actuellement, le site Web (PHP+MySQL) n'a pas encore été rendu public.



[Accueil](#)
[DiCoP](#)

[My DiCoP](#)
[Déconnecter](#)
[Aide](#)

A B C D E F G H I J K L M N O P Q R **S** T U V W X Y Z

localhost:8888/AFPC2/sangumaolu3.php



Photo

三顾茅庐 三顧茅廬
trois, se rendre, chaumière
sāngùmáolú
Aller visite qqn à trois reprises dans sa chaumière
rendre plusieurs visites à un homme de talent pour le prier de sortir d'une vie recluse ou de la retraite/inviter qqn plusieurs fois de suite à prendre une fonction importante.
「出师表 Chūshībiǎo」 (mémoires) écrits par Liang Zhuge (227-228)
LIU Bei (161 - 223), puissant seigneur chinois de la fin de la dynastie Han et du début de la période des Trois Royaumes, veut que le célèbre stratège ZHU Geliang devienne son conseiller militaire. Lors de ses deux premières visites, ZHU Geliang est intentionnellement absent. La troisième fois, LIU Bei réussit à le rencontrer et à le persuader d'accepter ce poste.
这一年冬天有人三顾茅庐，感恩知己，一夕倾谈遂相许以驱驰了。（唐弢《桥》）

UP + Phraséoculture



Vidéo

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

www.phraseologia.com/DiCoP.php

Accueil

DiCoP

My DiCoP

Disconnect

Help



Dictionary

DiCoP-Learning

DiCoP-Text

Fixed expressions

Collocations

Defrosting

Chinese metalexicography

Conferences

un coup de main

Search »

中文

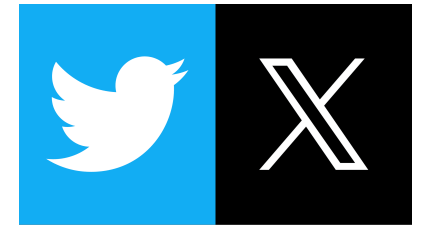
DiCoP-Text : La phraséologie pose particulièrement des difficultés en matière de traduction, car influencée tant par la linguistique que la culture ou la stylistique, et reflète fortement les choix et la technique de traduction du traducteur. **DiCoP-Text** sera **une base de données** (corpus monolingue, parallèle, contrastive, multilingue) **permettant de connaître la fréquence d'emploi des UP**, afin d'en vérifier la vitalité dans la pratique. Il doit permettre de scruter sans peine l'utilisation des UP dans les traductions, pour des études lexicométriques et contribuer aux recherche scientifiques en matière de la **phraséotraduction automatique** et du TAL (traitement automatique de la langue).

5. Corpus de phraséotraductologie : DiCoP-Text

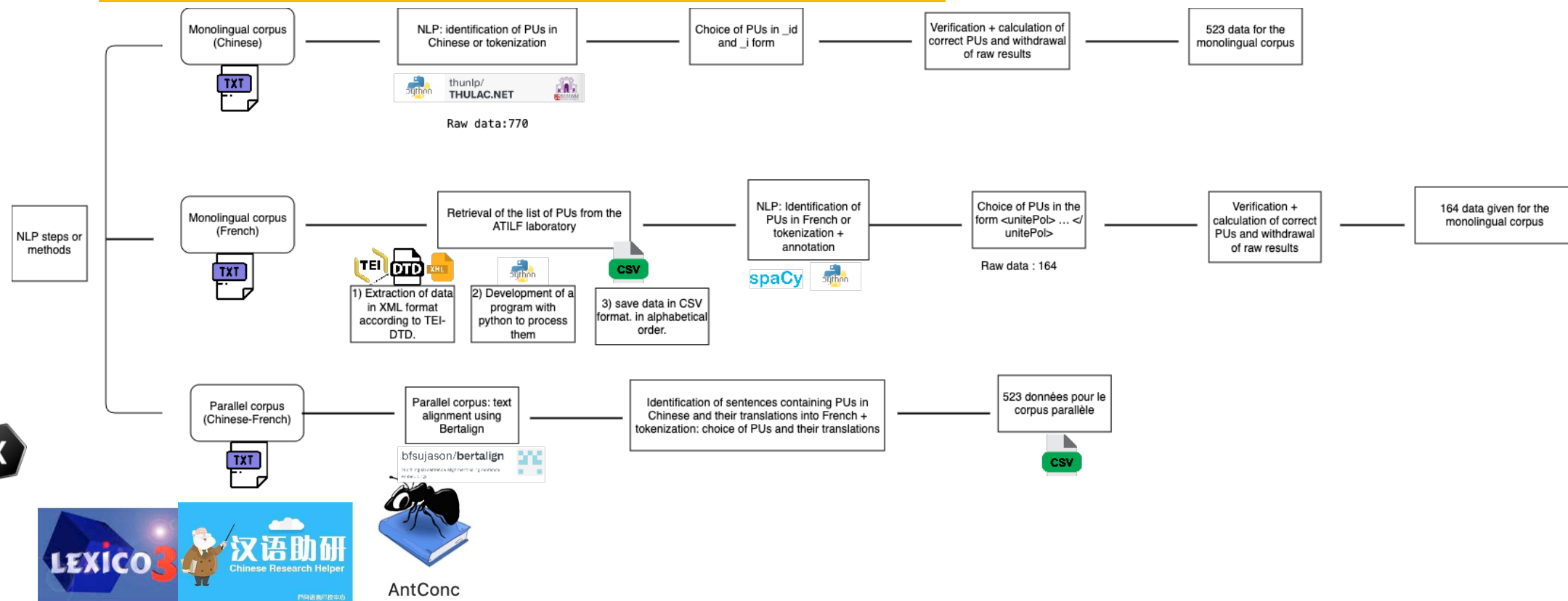
La collecte des textes et le traitement numérique sont actuellement en cours. Le choix du corpus doit refléter la diversité linguistique de la langue et donc être suffisamment large pour garantir une représentation adéquate des UP et améliorer la précision et la fiabilité du DiCoP.

- **représenter une variété de genres (littérature, poésie, discours, roman, science fiction, etc.)** pour trouver un équilibre entre le langage formel et informel (textes officiels, académiques, ainsi que dialogues, conversations quotidiennes, etc.).
- ressources modernes (20e siècle et suivant)
- références déjà accessibles sous forme numérique (ou déjà ocerisées).

Journaux : le « défigement » (活用 huóyòng en chinois), phénomène que l'on rencontre par exemple à travers les jeux de mots contenus **dans les titres, les slogans et les publicités**.



Corpus en cours de développement + TAL



Méthode de TAL suivie dans le projet DiCoP et la base de données DiCoP-Text



TAL : Identifier les UP en chinois ou token

```
1 ##终端命令: cd ~/Documents。 然后python3 import\ thulac.py
2
3
4
5 import thulac
6 thu1 = thulac.thulac()
7 #t = '''“红色联合”对“四·二八兵团”总部大楼的攻击已持续了两天，他们的旗帜在
8 #大楼顶上出现了一个娇小的身影，那个美丽的女孩子挥动着一面“四·二八”的大旗，
9 #“红色联合”的战士们欢呼起来，几个人冲到楼下，掀开四·二八的旗帜，抬起下面
10 #'''
11
12 f=open("texte_ch.txt","r")
13 f_out=open("text_ch_tok.txt","w")
14 for line in f:
15     text = thu1.cut(line, text=True)
16     f_out.write(text+"\n")
17     #print(text)
18
19 f.close()
20 f_out.close()
21
```



thunlp/
THULAC.NET



_f (_w 这个_r 绝密_a 信息_n 是_v 基地_n 许多_m 中层_f 干部_n 都_d 不_d 知道_v 的_u) _w , _w 她_r 把_p 与_p 外界_n 精神_n 上_f 的_u 联系_v 也_d 斩断_v 了_u , _w 只是_d 埋头_v 于_p 工作_v 。 _w 这_r 以后_f , _w 她_r 更_d 深_a 地_u 介入_v 到_v 红岸_n 系统_n 的_u 技术_n 接心_v , _w 开始_v 承担_v 比较_d 重要_a 的_u 研究_v 课题_n 。 _w 对于_p 杨卫宁_np 给_v 予_g 叶文洁_np 的_u 信任_v , _w 雷志成_np 一直_d 耿耿于怀_id , _w 但_c 他_r 还是_d 很_d 愿意_v 将_p 重要_a 课题_n 交_v 到_v 叶文洁_np 手上_s —_w 以_p 叶文洁_np 的_u 身份_n , _w 她_r 对_p 自己_r 的_u 研究_v 成果_n 没有_v 任何_r 权利_n ; _w 而_c 基地_n 中_f , _w 只有_c 雷志成_np 是_v 天体_n 物理_n 专业_n 出身_v 的_u , _w 是_v 当时_t 少见_a 的_u 知识分子_n 政委_n ; _w 这样_r , _w 叶文洁_np 的_u 成果_n 和_c 论文_n 最后_f 都_d 被_p 他_r 占_v 去_v , _w 使_v 他_r 成_v 了_u 部队_n 政工_j 干部_n 中_f 又红又专_id 的_u 典型_n 。 _w 调_v 叶文洁_np 进入_v 红岸_n 基地_n 的_u 最初_a 缘由_n , _w 是_v 她_r 读_v 研究生_n 时_g 发表_v 在_p 《_w 天文学报_n 》_w 上_f 的_u 那_r 篇_q 试图_v 建立_v 太阳_n 数学_n 模型_n 的_u 论文_n 。 _w 其实_d , _w 与_p 地球_n 相比_v , _w 太阳_n 是_v 一个_m 更_d 简单_a 的_u 物理_n 系统_n , _w 只是_d 由_p 氢_n 和_c 氦_n 这_r 两_m 种_q 很_d 简单_a 的_u 元素_n 构成_v 。 _w 它_r 的_u 物理_n 过程_n 虽然_c 剧烈_a , _w 但_c 十分_m 单纯_a 。 _w 只是_d 氢_v 至_v 氦_n 的_u 聚变_v , _w 所以_c , _w 有_v 可能_n 建立_v 一个_m 数学_n 模型_n 来_v 对_p 太阳_n 进行_v 较为_d 准确_a 的_u 描述_v 。 _w 那_r 论文_n 本_d 来_v 是_v 一_m 篇_q 很_d 基础_n 的_u 东西_n , _w 但_c 杨卫宁_np 和_c 雷志成_np 却_d 从中_d 看到_v 了_u 解决_v 红岸_n 监听_v 系统_n 一个_m 技术_n 难题_n 的_u 希望_n 。 _w 凌日_t 干扰_v 问题_n 一直_d 困扰_v 着_u 红岸_n 的_u 监听_v 操作_v 。 _w 这个_r 名词_n 是_v 从_p 刚_d 出现_v 的_u 通信卫星_n 技术_n 中_f 借_v 来_v 的_u , _w 当_p 地球_n 、_w 卫星_n 和_c 太阳_n 处于_v 同_v 一_m 条_q 直线_n 时_g , _w 地面_n 接收_v 天线_n 对准_v 的_u 卫星_n 是_v 以_p 太阳_n 为_v 背景_n 的_u 。 _w 太阳_n 是_v 一个_m 巨大_a 的_u 电磁_n 发射源_n , _w 这时_r 地面_n 接收_v 的_u 卫星_n 微波_n 就_d 会_v 受到_v 太阳_n 电磁_n 辐射_v 强烈_a 干扰_v 。 _w 这个_r 问题_n 后_f 来_v 直到_v 二十一_m 世纪_n 都_d 无法_v 解决_v 。 _w 红岸_n 所_u 受到_v 的_u 日凌_n 干扰_v 与_p 此_r 类似_v , _w 不_d 同_v 的_u 是_v 干扰_v 源_g (_w 太阳_n) _w 位于_v 发射源_n (_w 外_f 太空_s) _w 和_c 接收器_n 之间_f , _w 与_p 通信卫星_n 相比_v , _w 红岸_n 所_u 受_v 的_u 凌日_t 干扰_v 出现_v 的_u 时间_n 更_d 频繁_a , _w 也_d 更_d 严重_a 。 _w 实际_a 的_u 红岸_n 系统_n 又_d 比_p 原_a 设计_v 缩水_v 了_u 许多_m , _w 监听_v 和_c 发射_v 系统_n 共_d 用_v 一个_m 天线_n 。 _w 这_r 使得_v 监听_v 的_u 时间_n 较为_d 珍贵_a , _w 日凌_nz 干扰_v 也_d 就_d 成为_v 一个_m 严重_a 问题_n 了_u 。 _w 杨卫宁_np 和_c 雷志成_np 的_u 想法_n 很_d 简单_a : _w 搞_v 清_a 太阳_n 发射_v 的_u 电磁波_n 在_p 监测_v 波段_n 上_f 的_u 频谱_n 规律_n 和_c 特征_n , _w 用_p 数字_n 滤波_n 滤掉_v 它_r , _w 就_d 可_v 排除_v 干扰_v 。 _w 两_m 人_n 都_d 是_v 技术_n 专_np 家_n , _w 在_p 这_r 外行_n 领导_n 内行_n 的_u 年代_n , _w 这_r 是_v 难能可贵_id 的_u 。 _w 但_c 杨卫宁_np 不_d 是_v 天体_n 物理_n 专业_n 的_u , _w 雷志成_np 则_d 是_v 走_v 政工_j 道路_n 的_u 人_n , _w 在_p 专业_n 上_f 不_d 可能_v 知道_v 得_u 太_d 深_a 。 _w 其实_d 太阳_n 电磁_n 辐射_n 的_u 稳定_a 只_d 局限_v 于_p 包括_v 可见光_n 在内_u 的_u 从_p 近紫外_s 到_v 中红_nz 外_f 波段_n , _w 在_p 其他_r 的_u 波段_n 上_f , _w 它_r 的_u 辐射_n 是_v 动荡不定_i 的_u 。 _w 叶文洁_np 首先_d 明智_a 地_u 在_p 第一_m 份_q 研究_v 报告_n 中_f 明确_a 一点_m : _w 在_p 太阳_n 黑子_n 、_w 耀斑_n 、_w 日冕_n 物质_n 抛射_v 等_u 太阳_n 剧烈_a 爆发性_n 活动_v 期间_f , _w 日_g 凌_g 干扰_v 无法_v 排除_v 。 _w 于是_c , _w 研究_v 对象_n 只_d 局限_v 于_p 太阳_n 正常_a 活动_v 时_g 红岸_s 监测_v 波段_n 内_f 的_u 电磁_n 辐射_v 。 _w 基地_n 内_f 的_u 研究_v 条件_n 还是_d 不_d 错_v 的_u , _w 资料室_n 可以_v 按_p 课题_n 内容_n 调_v 来_v 较_d 全_a 的_u 外文_n 资料_n , _w 还有_v 很_d 及时_a 的_u 欧_j 美_j 学术_n 期刊_n , _w 在_p 那个_r 年代_n 这_r 是_v 件_q 很_d 不_d 容易_a 的_u 事_n 。 _w 叶文洁_np 还_d 可以_v 通过_p 军线_n , _w 与_p 中科院_j 两_m 家_q 研究_v 太阳_n 的_u 科研_n 单位_n 联系_v , _w 通过_p 传真_n 得到_v 他们_r 的_u 实时_a 观测_v 数据_n 。 _w 叶文洁_np 的_u 研究_v 持续_v 了_u 半_m 年_q , _w 丝毫_m 看_v 不_d 到_v 成功_a 的_u 希望_n 。 _w 她_r 很快_d 发现_v , _w 在_p 红岸_n 的_u 观测_v 频率_n 范围_n 内_f , _w 太阳_n 的_u 辐射_v 变幻莫测_id 。 _w 通过_p 对_p 大量_m 观测_v 数据_n 的_u 分析_v , _w 叶文洁_np 发现_v 了_u 令_v 她_r 迷惑_a 的_u 神秘_a 之_u 外_n : _w 有时_d , _w 上述_a 某_r 一_m 频段_n 辐射_v 发生_v 突变_v 时_q , _w 太阳_n

	UP chinoises identifiées par Thulac			
Total	771			
Taux de réussite	65.43% (504/771)			
Taux d'erreur	34.57% (267/771)			
	Sous forme <u>id</u> (567 UP)		Sous forme <u>i</u> (204 UP)	
	Correct	Erreur	Correct	Erreur
Nombre d'EI	321	246	183	21
Pourcentage	56.61 %	43.39%	89.71%	10.29%

TAL : Identifier les UP en français ou token

ce ne fut qu'après de longues minutes qu'elle relâcha ses bras encore tendus lentement vers la tribune, s'assit à côté du cadavre et serra la main déjà les yeux vides se perdaient dans le lointain.

Lorsque quelqu'un vint pour emporter le corps, Ye Wenjie sortit de la poche qu'elle glissa dans la main de son père : sa pipe.

Quand Wenjie quitta silencieusement les lieux, il n'y avait déjà plus persc terrain de sport.

Elle prit le chemin qui la ramenait chez elle.

Lorsqu'elle arriva devant le dortoir des enseignants, elle entendit au pren provenir de leur maison : ce rire émanait de la bouche de la femme qu'elle mère.

Tourner les talons

Wenjie <unitePol>tourna les talons</unitePol> sans bruit, laissant ses piec semblait.

Elle s'aperçut que ceux-ci l'avaient transportée jusque devant la porte de principal lorsqu'elle était en quatrième année de licence, et aussi son ami

À compter des deux années de maîtrise d'astrophysique, puis ensuite, lorsqu avec le lancement de la révolution et jusqu'à aujourd'hui, le Pr Ruan avait personne dont elle s'était sentie le plus proche.

Ruan Wen était jadis partie étudier à Cambridge.

```
import spacy
from spacy.matcher import PhraseMatcher

nlp = spacy.load("fr_core_news_sm")
matcher = PhraseMatcher(nlp.vocab, attr="LEMMA")
terms = ["écorcher le renard", "retourner à ses moutons", "tourner les truies au foin"]
# Only run nlp.make_doc to speed things up
patterns = [nlp(text) for text in terms]
matcher.add("TerminologyList", patterns)

f = open("rabelais.txt", "r")
for ligne in f:
    doc = nlp(ligne)
    matches = matcher(doc)
    for match_id, start, end in matches:
        span = doc[start:end]
        print(span.text)
f.close()
```



identifié 164 EI

par la balise «
<unitePol> </
unitePol> ».



[7] src = ""中国，1967年。

“红色联合”对“四·二八兵团”总部大楼的攻击已持续了两天，他们的旗帜在大楼周围躁动地飘扬着，仿佛渴望干柴的火种。“红色联合”的大楼顶上出现了一个娇小的身影，那个美丽的女孩子挥动着一面“四·二八”的大旗，她的出现立刻招来了一阵杂乱的枪声，射击的武器五花八门。“红色联合”的战士们欢呼起来，几个人冲到楼下，掀开四·二八的旗帜，抬起下面纤小的遗体，作为一个战利品炫耀地举了一段，然后将

tgt = ""

Chine, 1967.

L'assaut de l'Union rouge contre le quartier général de la brigade du 28 Avril durait depuis déjà de
Tout autour de l'édifice, les drapeaux de la brigade claquaient au vent, telles des torches attendan
L'officier en chef de l'Union rouge était terriblement inquiet, non qu'il ait peur des miliciens en
Ce qu'il craignait, c'étaient les quelques dizaines de fours de métal en fusion entreposés à l'intér
Ceux-ci, reliés par des détonateurs électriques, étaient pleins à ras bord d'explosifs.

Il ne pouvait les voir, mais il pouvait sentir leur présence magnétique.

Un doigt sur un bouton et tous seraient réduits en cendres.

Les jeunes gardes rouges de la brigade du 28 Avril étaient capables d'une telle folie.

Par comparaison avec les vétérans qui avaient mûri à mesure qu'ils affrontaient les tempêtes, cette
Plus fous que la folie elle-même.

Une silhouette gracile émergea au sommet de la tour, celle d'une belle jeune fille brandissant un gr
Son irruption fut accueillie par une volée de coups de feu.

C'était un cocktail hétérogène d'armes : des vieilles carabines à l'américaine, en passant par des Z
On trouvait même des armes blanches comme des coutelas et des lances, qui formaient avec les armes à
Les miliciens de la brigade du 28 Avril connaissaient la chanson par cœur.



aligner.print_sents()

terminées

中国，1967年。
Chine, 1967.

“红色联合”对“四·二八兵团”总部大楼的攻击已持续了两天，他们的旗帜在大楼周围躁动地飘扬着，仿佛渴望干柴的火种。

L'assaut de l'Union rouge contre le quartier général de la brigade du 28 Avril durait depuis déjà deux jours.

“红色联合”的指挥官心急如焚，他并不惧怕大楼的守卫者，那二百多名“四·二八”战士，与诞生于1966年初、经历过大规模检阅和大串联的“红色联合”。

L'officier en chef de l'Union rouge était terriblement inquiet, non qu'il ait peur des miliciens en faction.

他怕的是大楼中那十几个大铁炉子，里面塞满了烈性炸药，用电雷管串联起来，他看不到它们，但能感觉到它们磁石般的存在，开关一合，玉石俱焚。

Ce qu'il craignait, c'étaient les quelques dizaines de fours de métal en fusion entreposés à l'intérieur du bâtiment.

Les jeunes gardes rouges de la brigade du 28 Avril étaient capables d'une telle folie.

比起已经在风雨中成熟了许多的第一代红卫兵，新生的造反派们像火炭上的狼群，除了疯狂还是疯狂。

Par comparaison avec les vétérans qui avaient mûri à mesure qu'ils affrontaient les tempêtes, cette génération était encore plus sauvage.

大楼顶上出现了一个娇小的身影，那个美丽的女孩子挥动着一面“四·二八”的大旗，她的出现立刻招来了一阵杂乱的枪声，射击的武器五花八门，有轻武器，有重武器。

Une silhouette gracile émergea au sommet de la tour, celle d'une belle jeune fille brandissant un grand étendard.

——连同那些梭标和大刀等冷兵器，构成了一部浓缩的近现代史.....“四·二八”的人在前面多次玩过这个游戏，在楼顶上站出来的人，除了挥舞旗帜外，还会做一些其他的事情。

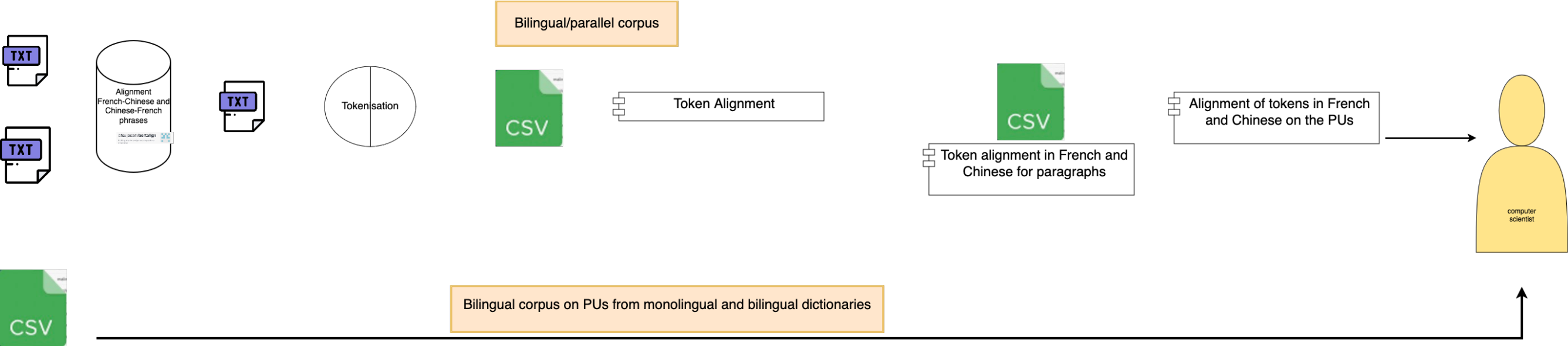
On trouvait même des armes blanches comme des coutelas et des lances, qui formaient avec les armes à feu une collection impressionnante.

La fille qui venait d'apparaître sur le toit croyait sans doute qu'elle connaîtrait aujourd'hui le même destin que ses camarades.

这次出来的女孩儿显然也相信自己还有那样的幸运她挥舞着战旗，挥动着自己燃烧的青春，敌人将在这火焰中化为灰烬，理想世界明天就会在她那面旗帜下诞生。

Tableau : Phrases comportant les 549 UP, et leur traduction

mot_chinois	sent_chinois	sent_fr	mot_fr
心急如焚	“红色联合”的指挥官心急如焚.	L’officier en chef de l’Union rouge était terriblement inquiet.	était terriblement inquiet
玉石俱焚	开关一合，玉石俱焚.	Un doigt sur un bouton et tous seraient réduits en cendres.	tous seraient réduits en cendres
五花八门	射击的武器五花八门.	C’était un cocktail hétérogène d’armes .	un cocktail hétérogène
愈演愈烈	射击的武器五花八门：有陈旧的美式卡宾枪、捷克式	C’était un cocktail hétérogène d’armes : des vieilles carabines à l’américaine, en passant par des	qui avait déclenché une explosion des conflits
全身而退	每次他们都能在弹雨中全身而退，为自己挣到了崇高	Chaque fois, ceux qui parvenaient à sortir indemnes de la pluie de projectiles étaient couverts de	ceux qui parvenaient à sortir indemnes de la pluie de projecti
化为灰烬	敌人将在这火焰中化为灰烬，理想世界明天就会在她	L’ennemi serait bientôt consumé par les flammes et demain, un monde idéal jaillirait de son sang	consumé par les flammes
错综复杂	任何一处都有 错综复杂 的对立派别在格斗。	En ces temps où coexistaient de multiples factions, des conflits complexes éclataient partout entre différentes obédiences.	complexes
牛鬼蛇神	与其他的 牛鬼蛇神 相比，反动学术权威有他们的特点	Les membres de l’Autorité académique réactionnaire possédaient des particularités qui les distinguaient des autres adversaires malfaisants de la révolution.	des autres adversaires malfaisants
肃然起敬	都自己结束了他们那曾经让人 肃然起敬 的生命。	Hai Mo et tant d’autres décidèrent de mettre eux-mêmes un terme à des vies autrefois très estimées .	très estimées
旷日持久	旷日持久 的批判将鲜明的政治图像如水银般注入了他们的意识	après une séance de critique ayant duré jusqu’au soir , la glorieuse idéologie politique de la révolution se déversait comme du mercure dans leur conscience, ébranlant les fondations d’un édifice intellectuel jadis façonné par le savoir et la raison.	e ayant duré jusqu’au soir
痛哭流涕	并为此 痛哭流涕 ，	et ils pleuraient de douleur.	et ils pleuraient de douleur
牛鬼蛇神	他们的忏悔往往比那此非知识分子的 牛鬼蛇神 要深刻得多	Les confessions des intellectuels étaient toujours beaucoup plus profondes et sincères que celles des autres ennemis de la révolution.	des autres ennemis
牛鬼蛇神	只有处于第一阶段的 牛鬼蛇神 才能对他们那早已过度兴奋的神经产生有效的刺激	Seuls les dangereux ennemis qui restaient bloqués à la première étape étaient capables de susciter une ferveur dans leurs esprits depuis longtemps fanatisés.	les dangereux ennemis
死路一条	顽固下去是 死路一条 ！	, vous devez bien vous rendre compte que la force unie contre laquelle vous résistez est si puissante que si vous continuez à vous obstiner, elle vous conduira à la mort !	elle vous conduira à la mort
有奶便是	，他有奶便是娘	Et, comme les mauvais esprits se rencontrent	Et, comme les mauvais esprits se rencontrent
有备而来	并且显然 有备而来	Elle avait manifestement préparé sa tirade	manifestement préparé sa tirade
彻头彻尾	更是 彻头彻尾 的反动唯心主义……”	Elle prétend que l'univers a des limites, ce qui n'est autre chose que de l'idéalisme réactionnaire...	autre chose que
滔滔不绝	听着妻子 滔滔不绝 的演讲	En entendant son épouse déverser son éloquent discours,	déverser son éloquent discours
两手空空	你 两手空空 地上来，	Tu es venue les mains vides	les mains vides
迫不及待	绍琳 迫不及待 地要继续下去了	Shao Lin n'attendait que cette occasion	n'attendait que
摇摇欲坠	以维持自己那 摇摇欲坠 的精神免于彻底垮掉	elle devait continuer à parler sans relâche pour préserver son esprit d'un effritement total.	d'un effritement total
不知所措	他这突然的反击令批判者们一时 不知所措 。	Cette soudaine contre-attaque plongea un moment ses accusateurs dans l'embarras.	dans l'embarras
实践出真知	这等于说正确的哲学是从天上掉下来的，这反对 实践出真知	C'est donc remettre en cause l'idée que les connaissances véritables proviennent de la pratique,	les connaissances véritables proviennent de la pratique,
无言以对	绍琳和两名大学红卫兵 无言以对	Ni Shao Lin ni ses deux étudiants gardes rouges n'étaient en mesure de justifier convenablement ce paradoxe.	n'étaient en mesure de justifier convenablement ce paradoxe
无坚不摧	但来自附中的四位小将自有她们“ 无坚不摧 ”的革命方式	Mais ce n'était pas le cas des quatre jeunes filles, persuadées que la révolution ne devait souffrir aucune attaque.	ne devait souffrir aucune attaque



Conclusion et perspectives