

Rapport du projet

Translittération et traitement de corpus multilingue : anglais, chinois, hébreu, birman, grec

Sous la direction de Madame Ilaine WANG

CHEN Lian 陈恋

30 décembre 2024

Sommaire

Introduction	2
I. Constitution du corpus	2
II. Analyse des caractères	3
1. Segmentation et analyse statistique	3
2) Relation graphème/phonème	9
III. Translittération vers l'alphabet latin	12
Conclusion	14

Translittération et traitement de corpus multilingue : anglais, chinois, hébreu, birman, grec

CHEN Lian 陈恋

30 Décembre 2024

Table des matières

1	Introduction	3
2	Constitution du corpus	3
3	Analyse des caractères	5
3.1	Segmentation et analyse statistique	5
3.2	Relation graphème/phonème	12
4	Translittération vers l' alphabet latin	17
5	Conclusion	19

1 Introduction

Ce projet se concentre sur le traitement d'un corpus multilingue axé sur le changement climatique. L'objectif principal est d'analyser les caractères de différents systèmes d'écriture en utilisant des articles provenant de Wikipédia, rédigés en cinq langues distinctes : birman, hébreu, chinois, anglais et grec. Pour établir ce corpus, nous avons appliqué des critères spécifiques, notamment l'inclusion d'au moins deux langues inconnues et la nécessité de fichiers contenant au minimum 2000 caractères.

Le projet se découpe en plusieurs étapes. Tout d'abord, nous avons constitué le corpus en récupérant le contenu des articles choisis, puis nous avons mis en place une analyse détaillée des caractères, en nous concentrant sur la segmentation, l'analyse statistique et la relation entre graphèmes et phonèmes. Enfin, la translittération des systèmes d'écriture vers l'alphabet latin a été réalisée pour faciliter l'interprétation des résultats. Ces analyses visent à améliorer notre compréhension des variations linguistiques et des structures phonétiques à travers différentes langues, tout en fournissant des insights sur l'usage du français dans un contexte multilingue.

2 Constitution du corpus

Notre corpus se compose de cinq articles Wikipédia rédigés en cinq langues différentes, chacune utilisant un type de système d'écriture distinct :

Birman (système alphasyllabaire)

Hébreu (système consonantique, ou abjad)

Chinois (système logographique/syllabique)

Anglais (système alphabétique)

Grec (système alphabétique)

Pour constituer notre corpus, nous avons respecté plusieurs critères. Tout d'abord, parmi les cinq langues choisies, il était nécessaire qu'au moins deux nous soient totalement inconnues. De plus, tous les articles devaient porter sur le même thème, à savoir le changement climatique. Enfin, chaque fichier devait contenir au moins 2 000 caractères, espaces exclus. Pour récupérer le texte de ces pages, nous avons utilisé le script Python du professeur Ilaine Wang, puis nous l'avons adapté nous-mêmes.

```

import requests
from bs4 import BeautifulSoup
import re
# page request
response = requests.get(
    url="https://en.wikipedia.org/wiki/Climate_change"
)
# Parsing page content
soup = BeautifulSoup(response.content, 'html.parser')
# Get all paragraphs, unordered list and ordered list
allParagraphs = soup.find(id="bodyContent").find_all(["p", "ul", "ol"])
text = ""
for item in allParagraphs:
    # Remove possible numeric marks from labels
    clnItem = re.sub(r'[\d+]', "", item.get_text())
    text += clnItem + '\n' # 每个条目后加换行符
#print(text)
# 将The result is written to a text file
with open("wiki_changement_climatique_my.txt", "w", encoding='utf-8') as file:
    file.write(text)

```

Fig. 1 : Capture d'écran du script de scraping pour extraire du texte en ligne

Le texte de chaque langue est enregistré sous les noms suivants :

wiki_changement_climatique_my.txt (pour le birman)

wiki_changement_climatique_he.txt (pour l'hébreu)

wiki_changement_climatique_zh - Hans.txt (pour le chinois)

wiki_changement_climatique_en.txt (pour l'anglais)

wiki_changement_climatique_el.txt (pour le grec)

Ensuite, nous utilisons la commande suivante :

`grep -o '.*' wiki_changement_climatique_zh - Hans.txt`

Cette commande nous permet d'extraire chaque caractère du fichier, et les résultats sont enregistrés dans les documents suivants :

wiki_theme_bir3new

wiki_theme_chi3new

wiki_theme_en3new

wiki_theme_grec3new

wiki_theme_hebreux₃_new

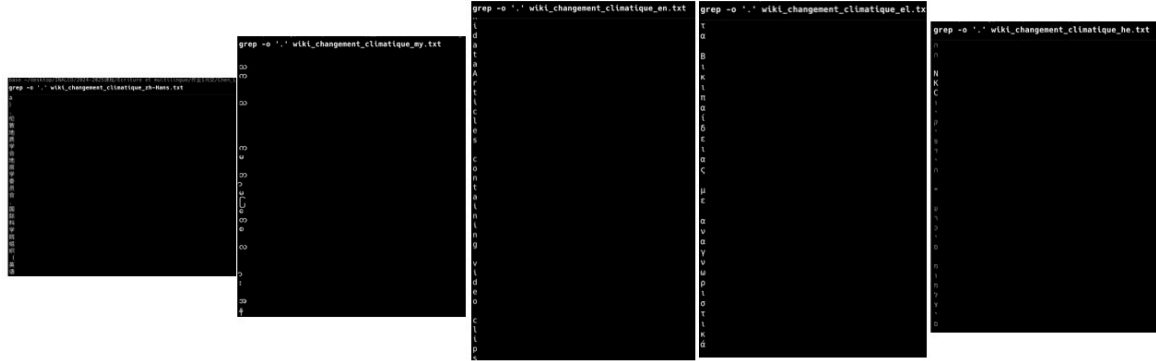


Fig. 2 : Capture d'écran des résultats de segmentation du texte dans ces cinq langues

Comme le montrent les images ci-dessus, en chinois, chaque caractère est séparé, y compris les signes de ponctuation, qui sont également considérés comme des caractères distincts. Il en va de même pour le birman, où chaque caractère est segmenté ligne par ligne. En revanche, en anglais, en grec et en hébreu, chaque lettre est traitée individuellement.

3 Analyse des caractères

3.1 Segmentation et analyse statistique

La première étape de cette partie consiste en la segmentation du fichier obtenu précédemment en caractères (non-mots compris). Pour ce faire, nous avons utilisé la commande (modifiée selon pour chaque langue) :

```
grep -o '.' wiki_changement_climatique_my.txt > wikithemebir3new
```

Les documents sont nommés suivants :

wiki_theme_bir₃_new

wiki_theme_chi₃_new

wiki_theme_en₃_new

wiki_theme_grec₃_new

wiki_thewme_hebreux₃_new

Pour pouvoir faire une analyse statistique des caractères, nous avons ensuite un tri :

1) Première étape : trier la liste de caractères

Dans cette étape, on utilise la commande `sort` pour organiser les caractères de manière ordonnée. La commande `sort wiki_theme_bir3_new.txt > sorted_wiki_theme_bir3_new.txt` permet de trier les caractères contenus dans le fichier `wiki_theme_bir3_new.txt` et d'enregistrer le résultat trié dans un nouveau fichier intitulé `sorted_wiki_theme_bir3_new.txt`. Comme le montrent les illustrations suivantes :

[illegible]

Fig. 3 : Captures d'écran de la liste des caractères pour chaque langue

2) Deuxième étape : compter les caractères \square

Après avoir trié la liste, nous utilisons la commande **uniq -c** pour compter l'occurrence de chaque caractère dans le fichier trié. La commande **uniq -c sorted_wiki_theme_bir3_new.txt > counted_sorted_wiki_theme_bir3_new.txt** (où "bir3" est remplacé selon la langue) génère un fichier qui répertorie le nombre d'occurrences de chaque caractère, en l'enregistrant sous le nom de **counted_sorted_wiki_theme_bir3_new.txt**.

Voici les résultats :

Birman	Chinois	English	Grec	Hébreux
644	3529	3529	1785	8414
8 "	10 "	10 "	4 "	1 !
1 %	14 %	14 %	3 %	116 "
7 (4 &	4 &	1 '	1 #
7)	25 '	25 '	17 (15 %
4 .	83 (83 (17)	5 &
1 /	83)	83)	92 ,	60 '
4 0	2 *	2 *	17 -	96 (
1 1	185 ,	185 ,	97 .	96)
2 2	732 -	732 -	9 /	843 ,
1 3	877 .	877 .	57 0	266 -
1 5	26 /	26 /	41 1	546 .
1 9	1259 0	1259 0	35 2	37 /
3 A	847 1	847 1	13 3	510 0
6 C	938 2	938 2	12 4	434 1
2 E	210 3	210 3	17 5	418 2
1 F	169 4	169 4	12 6	116 3
1 G	167 5	167 5	8 7	92 4
1 H	183 6	183 6	9 8	92 5
2 M	224 7	224 7	18 9	110 6
4 N	200 8	200 8	12 :	93 7
1 O	205 9	205 9	6 A	100 8
1 P	112 :	112 :	3 B	173 9
1 R	49 ;	49 ;	27 C	98 :
2 T	5 ?	5 ?	5 D	26 ;
3 W	143 A	143 A	3 E	1 =
25 a	88 B	88 B	4 F	11 ?
6 b	334 C	334 C	6 G	80 A
10 c	78 D	78 D	4 H	43 B
~ d	95 E	95 E	6 I	175 C

Fig. 4 : Illustration des occurrences de chaque caractère dans ces cinq langues

3) Troisième étape : trier par fréquence □

Dans cette dernière étape, on trie les caractères selon leur fréquence d'apparition. La commande `sort -rg counted_sorted_wiki_theme_bir3_new.txt > wiki_ecriture_bir3.txt` permet d'effectuer ce tri en ordre décroissant. Le résultat final est enregistré dans un fichier nommé `wiki_ecriture_bir3.txt`.

Birman	Chinois	English	Hébreux	Grec
1172 ၵ	142 [28246	8414	1785
644	138 温	17096 e	᾿ 3870	1000 α
625 ၵ	137 球	11626 a	י 3794	778 τ
	137 ^	10315 i	ן 2816	703 ι
611 ၵ	135 容	10242 t	נ 2199	690 o
	128 变	10129 n	ת 2164	668 ε
449 ၵ	126 地	9524 o	ז 2060	547 υ
444 c	124 P	9181 r	ך 1893	487 η
	120 度	7501 s	ס 1523	473 ρ
427 ၵ	118 大	6047 l	ב 1465	461 σ
	112 :	5975 .	ש 1363	442 ζ
413 ၵ	111 国	5479 c	י 1085	414 κ
	106 候	4579 h	ך 1048	384 μ
383 ၵ	105 全	4346 d	ד 1028	383 π
	105 中	3970 m	פ 919	328 λ
374 ၵ	99 是	3700 2	צ 859	271 υ
	99 R	3548 0	843 ,	224 ί
357 ၵ	95 E	3381 u	840 e	208 ó
292 ၵ	92 人	3272 l	ר 816	204 á
	92 F	3170 g	ת 805	198 é
283 ၵ	90 T	2641 p	כ 647	166 γ
	88 B	2422 ,	615 a	165 θ
277 ၵ	84 G	2367 C	ב 586	158 ή
	83)	2324 f	586 i	143 ω
259 ၵ	83 (2053 b	570 n	142 δ
	78 D	1910 y	546 .	97 .
250 ၵ	77 学	1591 v	λ 527	92 ,
		1456 S	510 0	90 χ
227 ၵ		1256 P	494 t	79 ú
		1236 A	ב 484	73 ó

Fig. 5 : Illustration du classement des caractères selon leur fréquence d'apparition dans ces cinq langues.

Nous analysons ces corpus en répondant les questions suivantes :

a) La segmentation en caractères pour ces langues est-elle aisée ?

Nous remarquons qu'en grec et en anglais, chaque lettre de chaque mot est séparée et les mots sont délimités par des espaces, parfois accompagnés de signes de ponctuation (comme les virgules et les points). En revanche, en chinois, les signes de ponctuation et les caractères sont également distincts. Comme déjà souligné, chaque caractère chinois est clairement séparé, y compris les signes de ponctuation, qui sont considérés comme des caractères autonomes. La segmentation des caractères chinois est relativement simple grâce à la nature du système d'écriture, où chaque caractère représente généralement un mot ou un concept unique. Ainsi, bien qu'il n'y ait pas d'espaces entre les mots, la segmentation reste facile tant que l'on dispose des caractères appropriés pour l'analyse.

En revanche, la situation semble moins claire pour le birman. En tant que personne n'ayant pas de connaissances en birman, je ne peux pas garantir que la segmentation soit correcte. Il m'est difficile de comprendre quand il faut insérer des espaces, et je ne suis pas sûr si un caractère est séparé ou non. D'après mes recherches sur Wikipédia (voir https://fr.wikipedia.org/wiki/%C3%89criture_birmane), bien que chaque caractère birman soit généralement une unité autonome, certains caractères composés peuvent représenter des combinaisons de sons, ce qui complique la segmentation.

Par ailleurs, la commande **grep -o '.'** a été utilisée pour extraire chaque caractère comme une ligne distincte, rendant ainsi cette tâche relativement simple. Cependant, cette méthode ne prend pas en compte les nuances grammaticales et syntaxiques propres à la langue birmane, qui peuvent poser des défis supplémentaires lors de l'analyse.

Dans l'ensemble, bien que la segmentation en caractères soit généralement simple dans ces cinq langues, des complexités peuvent surgir, en particulier avec le birman. Cela peut rendre la lecture du texte moins fluide et moins agréable.

b) Avez-vous rencontré des problèmes ? Si oui, les avez-vous résolus et comment ?

Des problèmes peuvent se manifester en raison de la diversité des caractères spéciaux et des glyphes présents dans les écritures birmane et chinoise. Par exemple, pour birman, certains caractères peuvent posséder des variantes de forme ou des diacritiques qui modifient leur prononciation. Pour surmonter ces défis, il est crucial de bien gérer les encodages en utilisant des outils appropriés afin d'assurer un traitement correct des données, notamment en vérifiant que les fichiers textes sont codés en UTF-8 un encodage largement utilisé pour le traitement de caractères multilingues.). Nous avons employé le script suivant :

```
# Lire le fichier
with open('/Users/lianchen/Desktop/INALCO/2024-2025 课程/Ecriture et multilingue/作业1月交/Chen Lian/
wiki_theme_bir3_new', 'r', encoding='utf-8') as file:
    content = file.read()
# Écrire dans le fichier
with open('/Users/lianchen/Desktop/INALCO/2024-2025 课程/Ecriture et multilingue/作业1月交/Chen Lian/
wiki_theme_bir3UTF_new.txt', 'w', encoding='utf-8') as file:
    file.write(content)
```

Fig. 6 : Script pour encoder des fichiers en UTF-8

Les problèmes potentiels peuvent inclure des caractères non reconnus ou mal encodés, surtout si le texte contient des caractères moins courants ou des symboles spécifiques. Pour résoudre de tels problèmes, il serait utile de s'assurer que le fichier est encodé avec UTF-8, d'utiliser des bibliothèques adaptées au traitement de caractères chinois, et de vérifier les résultats avec un affichage correct des caractères.

c) Sans compter l' espace, quels sont les 10 caractères les plus fréquents et les moins fréquents pour les 5 langues que vous avez choisies ? Dans quelle mesure ces statistiques sommaires vous informent-elles (ou pas) sur ces langues ou ces écritures ?

	Birman	Chinois	Grec	Anglais	Hébreux
Caractère le plus fréquents	1172 േ				
	644				
	625 ൈ	498 的			387 ם
	611 ൉	240 英	1000 α	17096 e	3794 ך
		229 语	778 τ	11626 a	2816 ן
	449 ൊ	208 气	703 ι	10315 i	2199 ם
		194 于	690 ο	10242 t	2164 ן
	444 ോ	192 存	668 ε	10129 n	2060 ך
		182 档	547 ν	9524 o	1893 ן
	427 ൌ	160 年	487 η	9181 r	1523 ם
Caractères les moins fréquents		160 化	473 ρ	7501 s	1465 ך
	413 ്	150 原	461 σ	6047 l	1363 ם
			442 ς	5975 .	1085 ך
	383 ൎ				1048 ם
	374 ൏				
	1 O	1 乏	1 —	2 ±	1 ,
	1 H	1 乌	1 Z	2	1 _
	1 G	1 举	1 H	2 >	1 ~
	1 F	1 丿	1 ±	2 =	1 é
	1 9	1 冫	1 z	1 ●	1 á
	1 5	1 ’	1 _	1 →	1 _
	1 3	1 ‘	1 V	1	1 Q
	1 1	1 ú	1 U	1 é	1 =
	1 /	1 é	1 K	1 ã	1 #
	1 %	1	1 ’	1 <	1 !

Fig. 7 : les 10 caractères les plus fréquents et les moins fréquents pour les 5 langues

Pour les moins fréquents, cela peut inclure des caractères comme les chiffres et les ponctuations.

d) Les signes de ponctuation sont-ils les mêmes qu’ en français ?

Dans notre analyse, il semble que les signes de ponctuation en anglais et en grec soient similaires à ceux du français. Concernant l’hébreu, cette langue utilise également des signes de ponctuation similaires, tels que les virgules et les points, mais leur emploi et leur position peuvent varier. Dans les textes hébreux, certains signes de ponctuation peuvent apparaître à des emplacements inversés par rapport à leur usage en français, ce qui est dû à la direction de l’écriture.

Par ailleurs, les signes de ponctuation utilisés en chinois diffèrent de ceux employés en français. Par exemple, la virgule (,) et le point (。) sont utilisés de manière distincte en termes de position et de fréquence par rapport à l'usage français. De plus, les guillemets (« » ou “ ”) présentent des différences tant dans leur forme que dans leur utilisation.

3.2 Relation graphème/phonème

Le travail effectué à l' étape 1 nous a aidé à identifier les caractères les plus fréquents et les moins fréquents dans un système d' écriture, basé sur un article Wikipédia. Cependant, cela ne montre pas comment ces caractères sont utilisés et leur valeur phonémique dans la langue. Nous avons donc tester le plateforme Wikicolor <https://wikicolor.alem-app.fr>. Cet outil utilise un code couleur pour mettre en évidence les graphèmes d' un mot et leur correspondance phonétique. Cela aide à visualiser comment les lettres se prononcent dans différents mots.

a) Testez et commentez le travail de l' outil WikiColor

Nous avons choisi la langue française et chinoise pour testet.
Nous avons d'abord tester cinq mots "Wikicolorisés" : WikiColor montre la relation entre l' orthographe et la prononciation. Par exemple :

La plateforme Wikicolor utilise un code couleur pour représenter visuellement la correspondance entre l' orthographe des mots et leur prononciation phonétique, en mettant en évidence les graphèmes. Les graphèmes sont les unités de son ou de graphe qui constituent un mot, et leur couleur permet de mieux comprendre les variations de prononciation. Voici un aperçu de la relation couleur-graphèmes, ainsi que des exemples que vous avez testés :

Exemples de mots wikicolorisés :

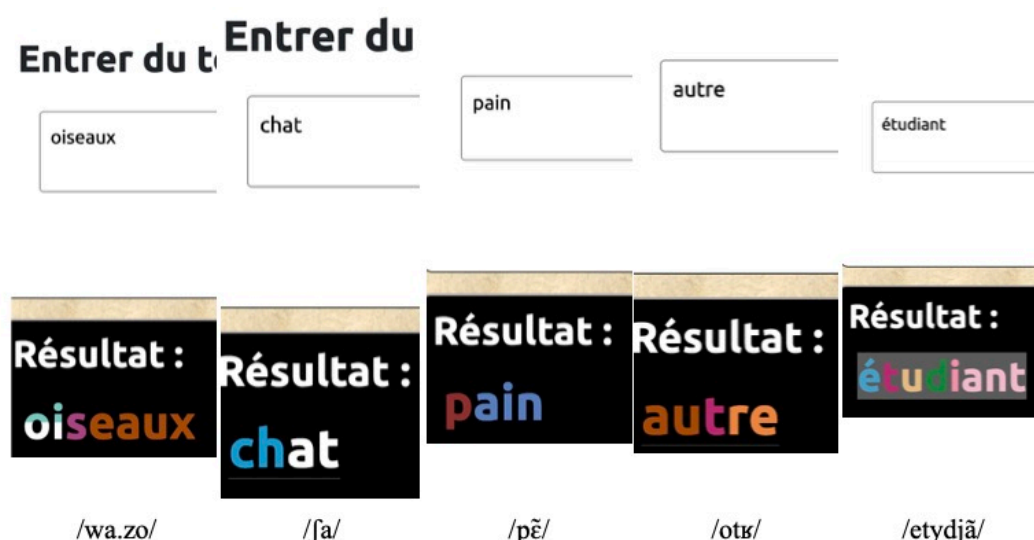


Fig. 8 : Exemples de mots wikicolorisés

Wikicolor facilite la compréhension de la phonétique en associant chaque graphème à une couleur spécifique, permettant ainsi aux utilisateurs de visualiser comment les lettres et les sons interagissent dans le français. Ce système est particulièrement utile pour l'apprentissage de la langue et l'enseignement de la prononciation correcte des mots.

Ensuite, nous avons testé le mandarin.

Nous avons d'abord sélectionné le chinois parmi les trois langues disponibles en haut à droite : chinois, anglais et français. Ensuite, nous avons testé des phrases simples telles que 你好 (Nǐ hǎo) pour «bonjour», 老师 (lǎoshī) pour «enseignant», et 今天 (jīntiān) pour «aujourd'hui». Nous avons constaté que les caractères chinois sont colorés, cependant, nous n'avons pas réussi à identifier de lien entre ces couleurs et le rythme des mots, car aucune explication n'était fournie.

D'après les résultats, nous avons noté que les caractères 老 (lǎo, «vieux») et 好 (hǎo, «bon») contiennent chacun trois sons, ce qui explique pourquoi ils sont associés à trois couleurs différentes : bleu, blanc et vert. En revanche, pour 今天 (jīntiān, «aujourd'hui»), bien qu'il se compose également de deux premiers sons, les couleurs qui lui sont attribuées semblent différentes.

Nous avons donc réessayé d'analyser leur pinyin et avons remarqué que tous les éléments étaient affichés en blanc. Cela indique qu'il y a une marge d'amélioration pour le système Wikicolor pour le chinois, afin d'assurer une meilleure correspondance entre les couleurs et les

phonèmes.



Fig. 9 : Exemples de mots en chinois wikicolorisés

L'interface de Wikicolor est très confortable et claire en un coup d'œil. WikiColor est un outil intéressant qui vise à illustrer la relation entre les graphèmes (les unités écrites, comme les lettres) et les phonèmes (les unités de son, comme les sons) dans différentes langues.

Nǐ hǎo

lǎoshī

jīntiān

结果:

Nǐ hǎo

lǎoshī

jīntiān

Fig. 10 : Exemples de pinyin chinois wikicolorisés

WikiColor utilise des couleurs pour représenter les différentes parties d'un mot, ce qui permet une visualisation immédiate des relations entre l'écriture et la prononciation. Cela peut être particulièrement utile pour les apprenants de langues, car cela facilite la compréhension des correspondances entre les lettres et les sons. En intégrant la transcription phonétique (API), WikiColor aide les utilisateurs à saisir non seulement comment un mot s'écrit, mais aussi comment il se prononce, ce qui est essentiel pour maîtriser correctement une langue, surtout celles avec des nuances de prononciation comme le français, l'anglais ou le mandarin. Cet outil peut potentiellement être utilisé pour diverses langues (français, anglais, chinois), ce qui le rend polyvalent.

Néanmoins, Pour des langues comme le mandarin, où les tonalités jouent un rôle crucial dans la signification, WikiColor peut rencontrer des défis lorsqu'il s'agit de représenter ces nuances. On a vu que dans les tests, l'outil a eu des difficultés à traiter correctement les tonalités, ce qui pourrait fausser l'apprentissage. Comme nous l'avons mentionné, des erreurs peuvent survenir, telles que des éléments affichés en rouge ou des symboles incorrects comme un "x". Ces problèmes techniques peuvent nuire à l'expérience utilisateur et nécessiter une attention continue pour garantir la fonctionnalité de l'outil. Aussi, La compréhension des graphèmes et phonèmes dépend aussi du contexte dans lequel ces mots sont utilisés. L'outil doit donc être accompagné d'explications sur les variations possibles en fonction du contexte pour être totalement efficace.

b) Un tel outil serait-il utile pour la langue inconnue que vous avez choisie pour le système alphabétique ?

Un tel outil serait indéniablement utile pour l'apprentissage de la langue grecque, notamment en raison de ses particularités orthographiques et phonétiques. Par exemple, en grec, chaque lettre correspond à un son spécifique, mais certaines lettres peuvent avoir des prononciations différentes selon leur position dans un mot. La lettre «γ» (gamma) est généralement prononcée comme un «g» dur, mais lorsqu'elle est suivie des lettres «ε» ou «ι», elle se prononce comme un «y» doux (ex : «γέρος» /□□en.os/) (voir https://fr.wikipedia.org/wiki/Alphabet_grec). Un outil qui utilise des couleurs pour représenter ces variations de prononciation permettrait aux apprenants de visualiser facilement ces différences et d'associer chaque graphème à son phonème approprié.

Les accents jouent un rôle crucial dans la prononciation et la signification des mots. Par exemple, le mot «πού» (pou, où) a un accent tonique, tandis que «που» (pu, que), sans accent, peut changer de signification selon le contexte. Un outil qui met en évidence les accents à l'aide

de couleurs pourrait aider les apprenants à comprendre l'importance de la prosodie et à prononcer correctement les mots en fonction de leur accentuation.

Le grec contient également de nombreuses diphtongues, telles que « αἰ » (ai) et « εἰ » (ei), qui peuvent être difficiles à maîtriser. En associant des couleurs spécifiques à ces diphtongues, un outil comme WikiColor pourrait aider les apprenants à les reconnaître et à les prononcer correctement. Par exemple, dans le mot « καί » (kai, et), le son « αἰ » pourrait être mis en évidence d'une manière qui aide l'apprenant à visualiser la transition entre les sons.

Ainsi, un outil comme WikiColor serait extrêmement bénéfique pour l'apprentissage du grec, car il offrirait une approche visuelle et interactive qui facilite la compréhension des relations entre l'écriture et la prononciation, tout en prenant en compte les spécificités phonétiques et orthographiques de la langue.

4 Translittération vers l' alphabet latin

Selon l'Organisation internationale de normalisation (ISO), la translittération est un « procédé qui consiste à représenter les caractères d'un système d'écriture alphabétique ou syllabique à l'aide des caractères d'un alphabet de conversion ». La dernière phase de ce projet vise à étudier la question de la translittération pour l'un des systèmes d'écriture que nous avons sélectionnés, parmi lesquels figurent : le système alphabétique inconnu, le système consonantique et l'alphasyllabaire. Nous avons opté pour la langue grecque et, par conséquent, Nous utiliserons le script suivant pour le transmettre et l'enregistrer sous le nom : **wiki_ecriture_grec3_translit.txt**.

```

translittération.py > ...
# Mapping pour la translittération
transliteration_map = {
    'Α': 'A', 'α': 'a',
    'Β': 'V', 'β': 'v',
    'Γ': 'G', 'γ': 'g',
    'Δ': 'D', 'δ': 'd',
    'Ε': 'E', 'ε': 'e',
    'Ζ': 'Z', 'ζ': 'z',
    'Η': 'I', 'η': 'i',
    'Θ': 'Th', 'θ': 'th',
    'Ι': 'I', 'ι': 'i',
    'Κ': 'K', 'κ': 'k',
    'Λ': 'L', 'λ': 'l',
    'Μ': 'M', 'μ': 'm',
    'Ν': 'N', 'ν': 'n',
    'Ξ': 'X', 'ξ': 'x',
    'Ο': 'O', 'ο': 'o',
    'Π': 'P', 'π': 'p',
    'Ρ': 'R', 'ρ': 'r',
    'Σ': 'S', 'σ': 's', 'ς': 's',
    'Τ': 'T', 'τ': 't',
    'Υ': 'Y', 'υ': 'y',
    'Φ': 'F', 'φ': 'f',
    'Χ': 'Ch', 'χ': 'ch',
    'Ψ': 'Ps', 'ψ': 'ps',
    'Ω': 'O', 'ω': 'o'
}

```

```

def transliterate(text):
    return ''.join(transliteration_map.get(char, char) for char in text)

# Chemin vers le fichier d'entrée
input_file_path = '/Users/lianchen/Desktop/INALCO/2024-2025课程/Ecriture et multilingue/作业1月交/2步/wi

# Lire le texte depuis le fichier
with open(input_file_path, 'r', encoding='utf-8') as input_file:
    text = input_file.read()
    transliterated_text = transliterate(text)

# Chemin vers le fichier de sortie
output_file_path = '/Users/lianchen/Desktop/INALCO/2024-2025课程/Ecriture et multilingue/作业1月交/2步/w

# Écrire le texte translittéré dans le nouveau fichier
with open(output_file_path, 'w', encoding='utf-8') as output_file:
    output_file.write(transliterated_text)

```

Fig. 11 : Script pour la translittération de la langue grecque.

Voici les résultats :

Avant	Après
1785	1785
1000 α	1000 a
778 τ	778 t
703 ι	703 i
690 ο	690 o
668 ε	668 e
547 ν	547 n
487 η	487 ī
473 ρ	473 r
461 σ	461 s
442 ς	442 s
414 κ	414 k
384 μ	384 m
383 π	383 p
328 λ	328 l
271 υ	271 y
224 ί	224 í
208 ό	208 ó
204 ά	204 á
198 έ	198 é
166 γ	166 g
165 θ	165 th
158 ή	158 ḥ
143 ω	143 ō
142 δ	142 d
97 .	97 .
92 ,	92 ,
90 χ	90 ch
79 ύ	79 ú
73 ώ	73 ō
71 β	71 v
69 α	69 a
67 φ	67 f
60 ε	60 e
57 θ	57 θ
55 ο	55 o
49 τ	49 t
49 ι	49 i
46 ξ	46 x
44 ν	44 n
43 ρ	43 r
41 λ	41 l
37 Ε	37 E
36 Α	36 A
35 2	35 2
33 c	33 c
30 λ	30 l
27 s	27 s

Fig. 12 : Résultats de la translittération de la langue grecque

5 Conclusion

Ce projet de traitement de corpus multilingue a permis d'explorer en profondeur les spécificités des systèmes d'écriture de cinq langues distinctes, à savoir le birman, l'hébreu, le chinois, l'anglais et le grec, à travers le prisme du changement climatique. En mobilisant des outils informatiques pour constituer un corpus, segmenter les textes, et analyser les caractères, nous avons pu recueillir des données significatives et effectuer une analyse comparative des

structures linguistiques.

L'analyse des caractères a révélé non seulement la fréquence d'utilisation de certains graphèmes, mais également des différences notables dans l'usage de la ponctuation et des règles de segmentation entre ces langues. Cela souligne la richesse et la diversité linguistique du corpus étudié. Les défis rencontrés, notamment liés à l'encodage et à la gestion des caractères spécifiques, ont également mis en lumière l'importance d'utiliser des outils adaptés pour de telles analyses.

De plus, l'utilisation de l'outil WikiColor a montré comment la visualisation des relations entre orthographe et prononciation peut enrichir l'apprentissage des langues. Cependant, il reste des obstacles à surmonter, notamment dans le traitement des langues tonales comme le mandarin, ainsi que l'amélioration continue de l'outil pour maximiser son efficacité.

La phase de translittération vers l'alphabet latin a également été une étape cruciale, facilitant ainsi la comparaison et l'interprétation des résultats de l'analyse. À travers ce projet, nous avons non seulement acquis une meilleure compréhension des dynamiques linguistiques, mais nous avons également mis en évidence l'importance de l'approche multilingue dans la recherche linguistique.

Références

- [1] Wikicolor. <https://wikicolor.alem-app.fr/fr>.
- [2] *Apprentissage des langues et multimodalité (ALeM)*. <https://alem.hypotheses.org/category/tutoriels/tutoriels-wikicolor>.
- [3] *Alphabetic grec*. https://fr.wikipedia.org/wiki/Alphabet_grec.
- [4] Keyman. *écriture birmane*. https://fr.wikipedia.org/wiki/%C3%89criture_birmane.