

Le projet CollEx CHI-KNOW-PO CORPUS : acquérir un large corpus pour analyser la circulation des savoirs dans la Chine médiévale

Marie Bizais-Lillig

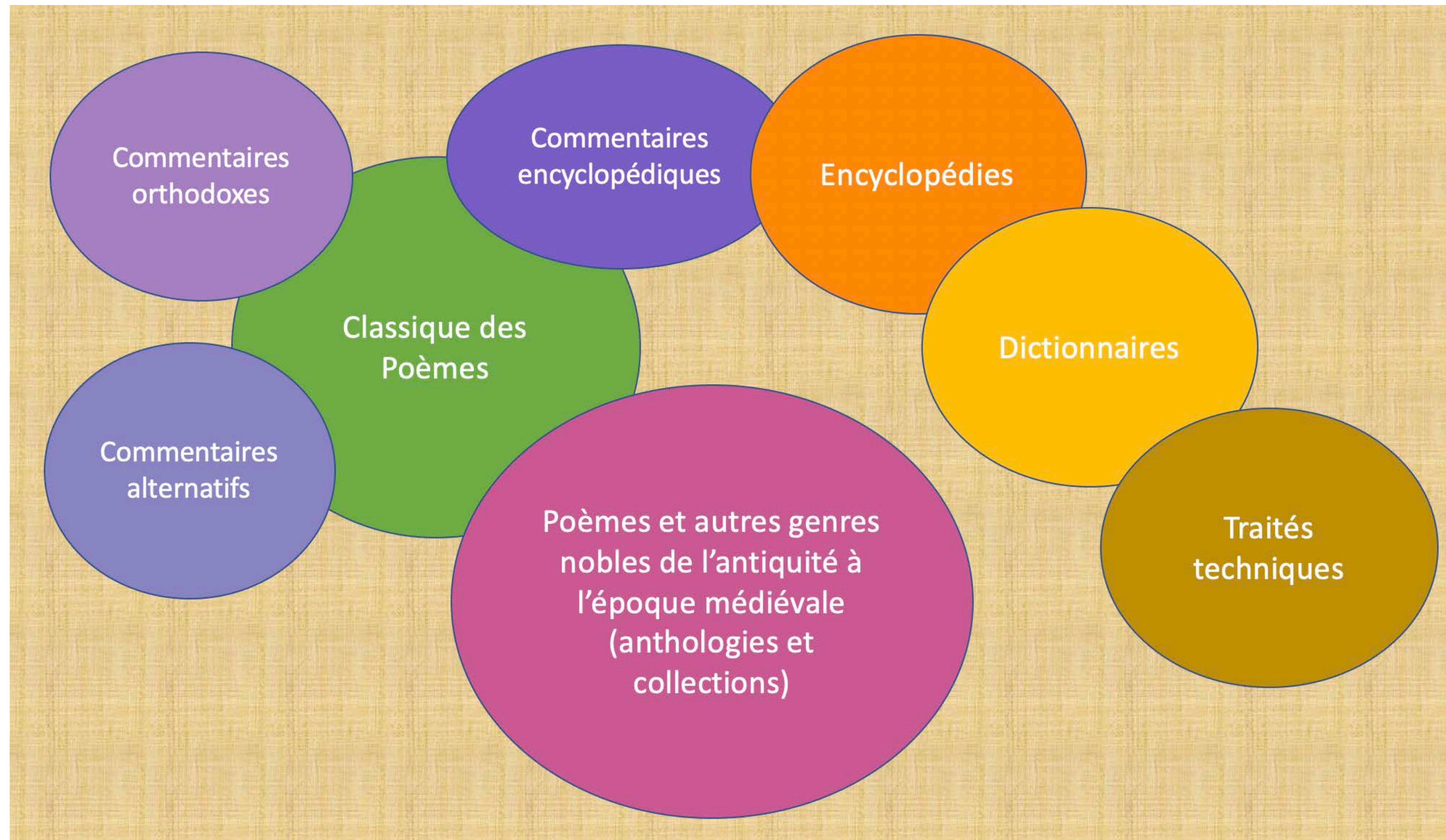
Université de Strasbourg - USIAS - consortium Distam (Huma-Num, CNRS) - porteuse du projet COLLEX CHI-KNOW-PO

D'un mot

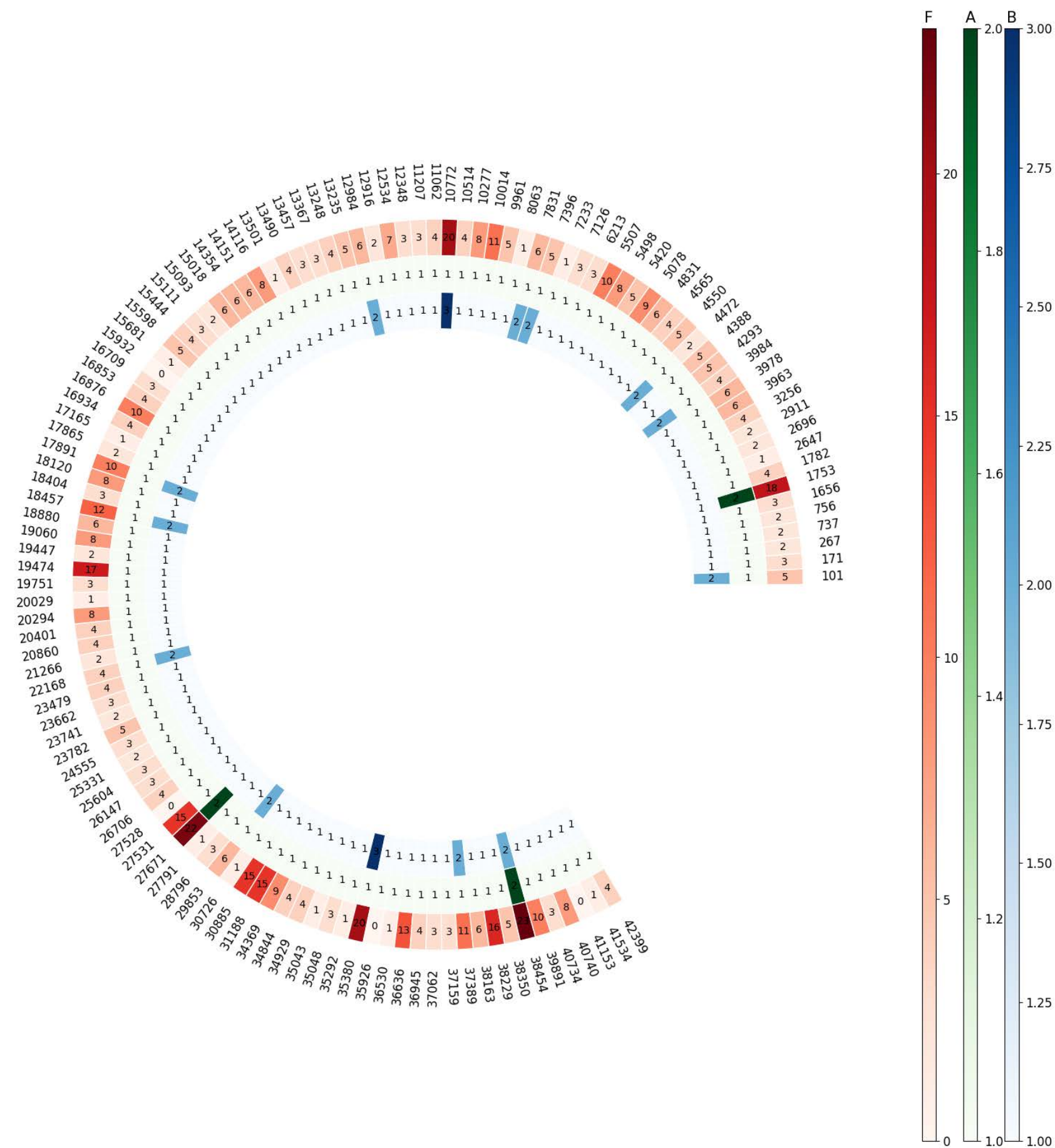
- Dates : juillet 2022 - juin 2024
- Objectifs : numérisation, HTR, édition, archivage, exposition
- Partenaire : Calfa pour l'HTR
- Co-porteuses : Marie Bizais-Lillig (Université de Strasbourg) et Soline Lau-Suchet (BULAC)
- Collections concernées : BULAC (fonds chinois ancien), BIHEC (Collège de France), BNU (fonds JP et Colette Diény)
- 1 ouvrage de la Bibliothèque d'études chinoises de l'Université de Strasbourg (don du Collège de France) ne sera finalement pas pris en compte.

1. Pourquoi ? Les besoins de la recherche

Cadre : Le projet CHI-KNOW-PO



Explorations : cooccurrences



孟浩然 《夜泊廬江，聞故人在東寺，以詩寄之》

poem 5 of 81, no. 6703 (40 words)

江路經廬阜，松門入虎溪。聞君尋寂樂，清夜宿招提。
石鏡山精怯，禪枝佈鶴棲。一燈如悟道，為照客心迷。

孟浩然 《李少府與楊九再來》

poem 6 of 81, no. 6785 (39 words)

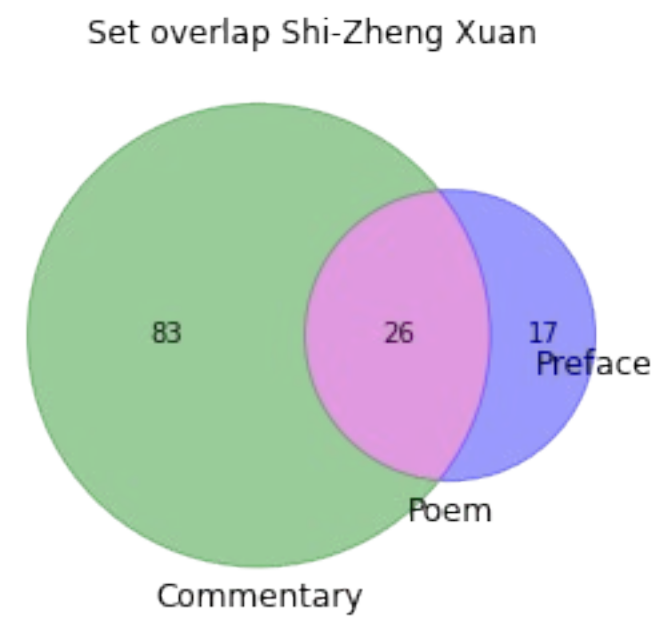
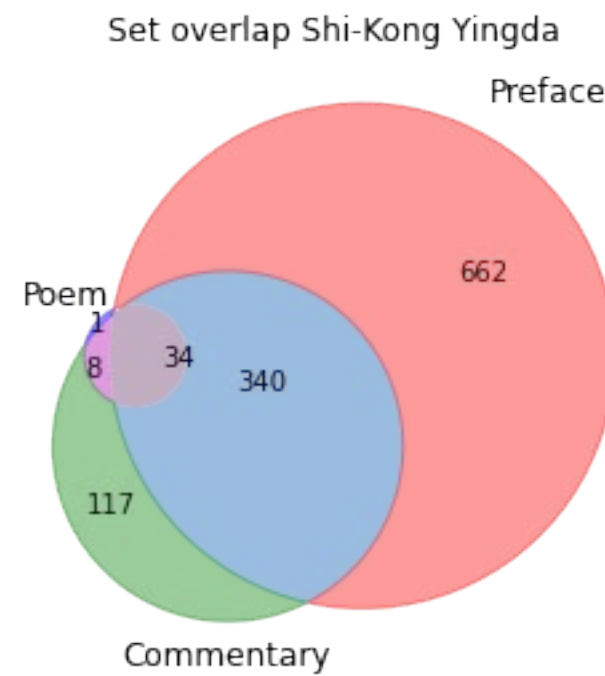
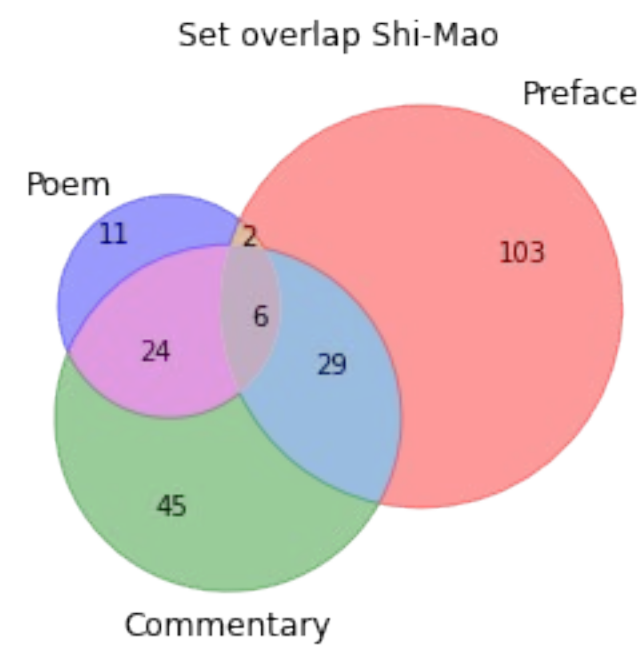
弱歲早登龍，今來喜再逢。如何春月柳，猶憶歲寒松。
煙火臨寒食，笙歌達曙鐘。喧喧鬥雞道，行樂羨朋從。

李白 《永王東巡歌十一首》

poem 7 of 81, no. 7037 (302 words)

永王正月東出師，天子遙分龍虎旗。樓船一舉風波靜，江漢翻為雁鷺池。
三川北虜亂如麻，四海南奔似永嘉。但用東山謝安石，為君談笑靜胡沙。
雷鼓嘈嘈喧武昌，雲旗獵獵過尋陽。秋毫不犯三吳悅，春日遙看五色光。
龍蟠虎踞帝王州，帝子金陵訪古丘。春風試暖昭陽殿，明月還過鳩鵲樓。
二帝巡遊俱未回，五陵松柏使人哀。諸侯不救河南地，更喜賢王遠道來。
丹陽北固是吳關，畫出樓臺雲水間。千岩烽火連滄海，兩岸旌旗繞碧山。
王出三山按五湖，樓船跨海次陪都。戰艦森森羅虎士，征帆一一引龍駒。
長風掛席勢難回，海動山傾古月摧。君看帝子浮江日，何似龍驤出峽來。

Explorations : reprises



POEM 3 卷耳 MAO

卷耳后妃之志也又當輔佐君子求賢審官知臣下之勤勞內有進賢之志而無險詖私謁之心朝夕思念至於憂勤也憂者之興也采采事采之也卷耳苓耳也頃筐畚屬易盈之器也懷思真置行列也思君子官賢人置周之列位陟升也崔嵬土山之戴石者虺隤病也姑且也人君黃金疊永長也山脊曰岡玄馬病則黃兕觥角爵也傷思也石山戴土曰砢瘡病也痛亦病也籲憂也

+++++

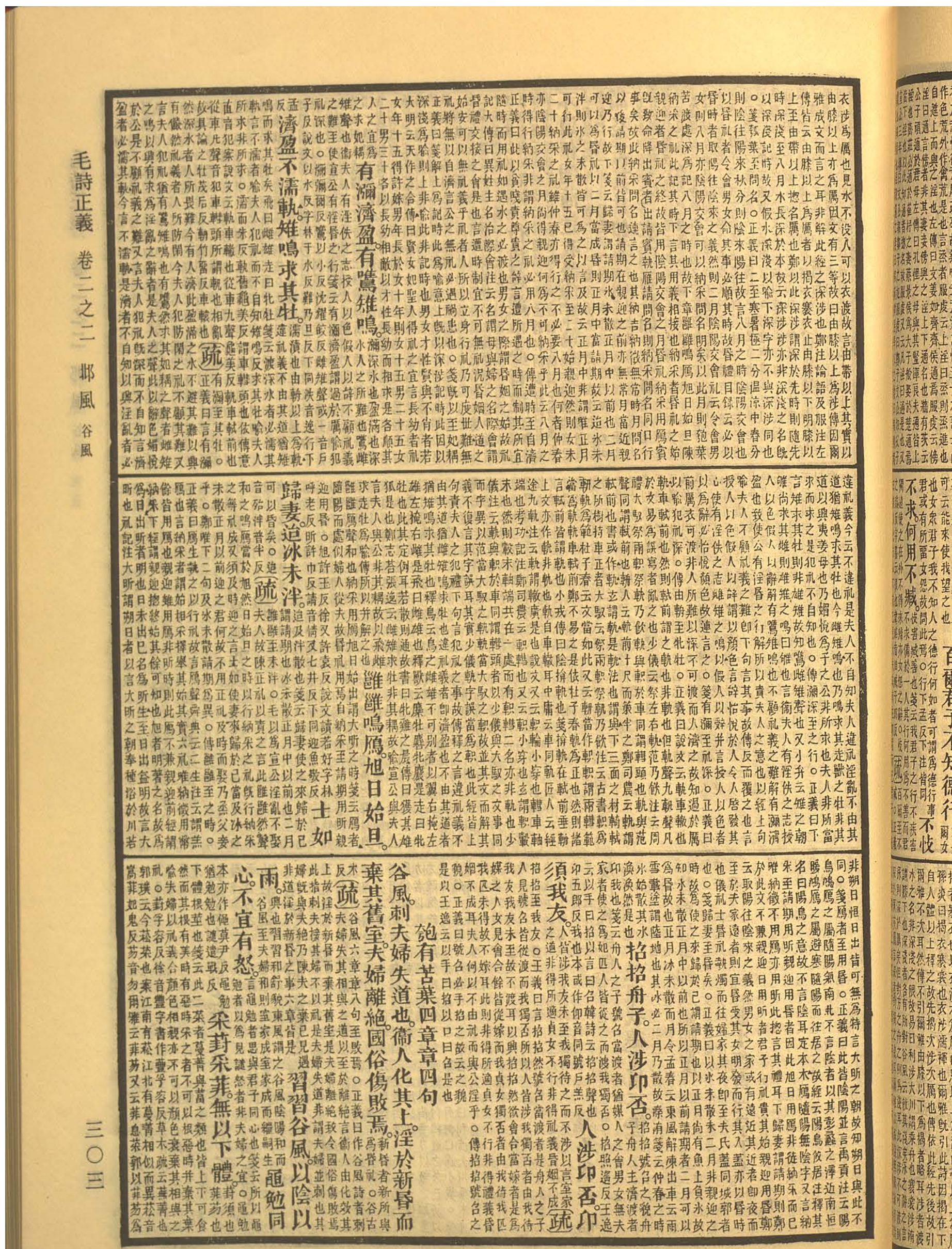
KYD

作卷耳詩者言后妃之志也后妃非直憂在進賢躬率婦道又當輔佐君子其志欲令君子求賢德之人審置於官位復知臣下出使之勤勞欲令君子賞勞之內有進賢人之志唯有德是用而無險詖不正私請用其親戚之心又朝夕思此欲此君子官賢人乃至於憂思而成勤此是后妃之志也言又者繫前之辭雖則異篇而同是一人之事故言又為亞次也輔佐君子總辭也求賢審官至於憂勤皆是輔佐君子之事君子所專后妃志意如然故云后妃之志也險詖者情實不正譽惡為善之辭也私謁者婦人有寵多私薦親戚故厲王以豔妻方嬖七子在朝成湯謝過婦謁盛與險詖私謁是婦人之常態聖人猶恐不免后妃能無此心故美之也至於憂勤勤為勞心憂深不已至於勞勤后妃之篤志也至於憂勤即首章上二句是也求賢審官即首章下二句是也經敘倒者敘見后妃求賢而憂勤故先言求賢經主美后妃之志能為此憂勤故先言其憂也言有人事采此卷耳之菜不能滿此頃筐頃筐易盈之器而不能滿者由此人志有所念憂思不在於此故也此采菜之人憂念之深矣以興后妃志在輔佐君子欲其官賢賞勞朝夕思念至於憂勤其憂思深遠亦如采菜之人也此亦后妃之憂為何事言后妃嗟呼而嘆我思君子官賢人欲令君子置此賢人於彼周之列位以為朝廷臣也我者后妃自我也(ZX S2 E1+E2)下箋云我我使臣我我君此不解者以詩主美后妃故不特言也言彼者后妃主求賢人為故以周行為彼也不云興也而云憂者之興明有異於餘興也餘興言采菜即取采菜喻言生長即以生長喻此言采菜而取憂為興故特言憂者之興言興取其憂而已不取其采菜也言事采之者言勤事采此菜也此與芣苢俱言采采彼傳云非一辭與此不同者此取憂為興言勤事采菜尚不盈筐言其憂之極故云事采之彼以婦人樂有子明其采者眾故云非一辭其實采采之義同故鄭志答張逸

Fondements : HTR 2021



- Calfa
- Maoshi Zhengyi 毛詩正義
- Corpus annoté de 50 images pour les expérimentations :
 - 3 240 lignes en apprentissage
 - 92 234 caractères
- Intégration des variantes dans la transcription
- 4 annotateurs : Shueh-Ying Liao, Weihang Wu, Hugo Dubois-Mouro, Marie Bizais-Lillig

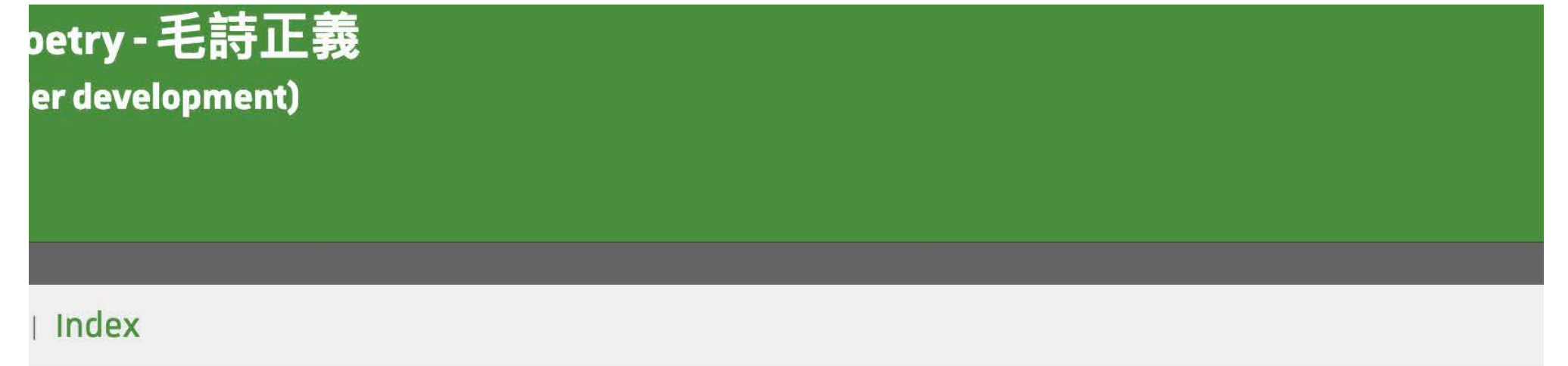


毛詩正義 卷二之二 鄘風 谷風

三〇三

Publier les textes, en relation avec des éditions

- Un site Web (Plateforme ESTRADES, avec l'aide de Pikselkraft)
- Donner à voir les images (manifestes IIIF) et les textes structurés (XML dans BaseX)
- Offrir des possibilités de circulation à travers le corpus grâce à des index liés à notre base de données
- Renvoyer vers un espace de correction TACT
- Donner accès aux fichiers sources



· 真彼周行 ·

箋：
器之易盈而不盈者，志在輔佐君子，憂思深也。
懷，思。真，置。行，列也。思君子官賢人。

箋：
周之列位，謂朝廷臣也。
箋「周之」至「延臣」。正義曰：知者，以其言周行是周之列位，非朝廷，故知官人是朝廷臣也。襄十三年：「《詩》曰『嗟我懷人，真彼周行也。王及公、侯、伯、子、男、采、衛，各居其列，所謂周行也』。彼非周行者，傳證楚能官人，引《詩》與此同。

[疏] 「采采」至「周行」。正義曰：言有人事采此卷耳之菜，不能滿此頃筐。頃筐，易盈之器，而不能滿者，由此人志有所念，憂思不采菜之人憂念之深矣，以興后妃志在輔佐君子，欲其官賢賞勞，朝夕思念，至於憂勤。其憂思深也。此后妃之憂為何事，言后妃嗟呼而嘆，我思君子官賢人，欲令君子置此賢人於彼周之列位，以者，后妃自我也。下箋「我，我使臣」，「我，我君」。此不解者，以詩主美后妃，故不特言也。主求賢人為此，故以周行為彼也。

Structurer et éditer les textes

- Conservation d'une partie des données du XML-Alto pour l'alignement texte/image
- Structuration des textes issus de l'HTR sur la base des types de régions et colonnes selon un schéma TEI préalablement défini
- Enrichissement des fichiers XML : balisage des noms de personnes et titres (cf. base de données)

```
270 <div resp="#KYD" type="subcommentary">
271 <div source="#sj034 #sj034_com_Mao">
272 <head>疏「《匏有苦葉》四章，章四句」至「淫亂」。正義曰：</head>
273 <ab>並為淫亂，亦應刺夫人，獨言宣公者，以詩者主為規諫君，故舉君言之，其實亦刺夫人也。故經首章、三章責公不依
274 之。</ab>
274 </div>
275 <div source="#sj034_com_ZX">
276 <head>箋「夫人謂夷姜」。正義曰：</head>
277 <ab>知非宣姜者，以宣姜本適伋子，但為公所要，故有魚網離鴻之刺。此責夫人，云「雉鳴求其牡」，非宣姜之所為，明
278 </div>
279 </div>
280 </div>
281 <div type="poem" xml:id="sj034">
282 <lg type="stanza" xml:id="sj034s01">
283 <lg type="couplet" xml:id="sj034s01e1">
284 <l xml:id="sj034s0111">匏有苦葉，</l>
285 <l xml:id="sj034s0112">濟有深涉。</l>
286 </lg>
287 <lg type="couplet" xml:id="sj034s01e2">
288 <l xml:id="sj034s0113">深則厲，</l>
289 <l xml:id="sj034s0114">淺則揭。</l>
290 </lg>
291 </lg>
```


2. La préparation du corpus

Définition du corpus

Belles Lettres

- *Wenxuan* 文選 (commenté par Li Shan (1809), commentée par les 6 ministres (1923)) — BNU
- *Yutai xinyong* 玉臺新詠 (1879) — BIHEC (V XIV 69 (1-8))
- *Quan Tang shi* 全唐詩 (1707) — BIHEC (SB 4002 (1-12) (1-120))
- *Yuefu shiji* 樂府詩集 / Wikisource

Textes techniques

- *Gujin shiwen leiju* 古今事文類聚: (1604) — BIHEC(SB 3705 (1-26))
- *Qimin yaoshu* 齊民要術 (1896) — BIHEC (V I 22 (1-4))
- *Xinzhai shizhong* 心齋十種 (1785/1788) — BIHEC (V I 53 (1))
- *Shennong bencao jing jizhu* 神農本草經集注 / Wikisource

Définition du corpus

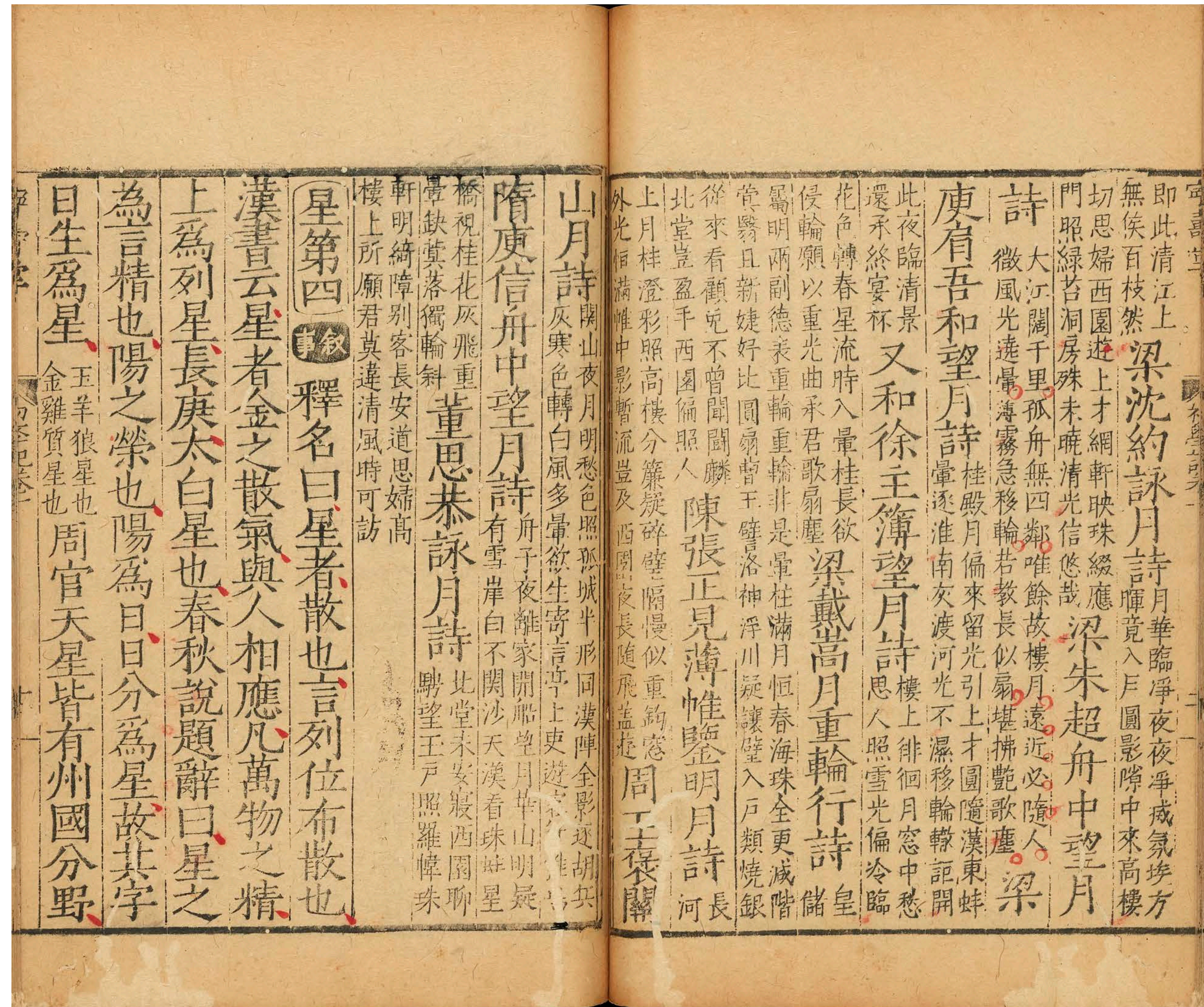
Textes de savoirs

- *Beitang shuchao* 北堂書鈔 (1888) — BULAC (BIULO CHI.1087)
- *Bowu zhi* 博物志 (1875) — BULAC (BIULO CHI.1140)
Erya yintu 影宋鈔繪圖爾雅 (1801) — BULAC (BIULO CHI.1938(1)-(3))
- *Mao Shi caomu niaoshou chongyu shu* 毛詩草木鳥獸蟲魚疏 (1857) — BIHEC Paris (V I 111 (1) 5)

Textes de savoirs

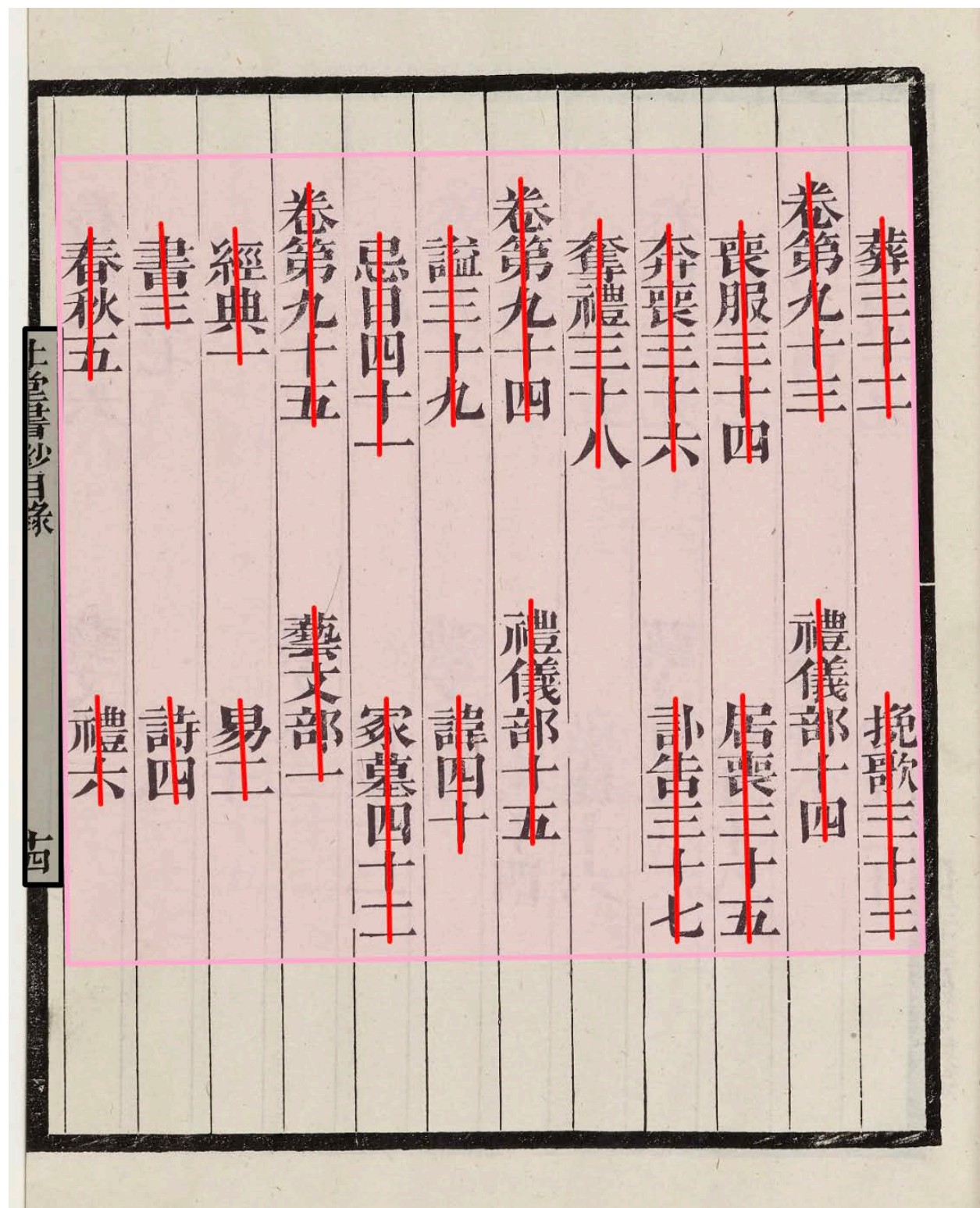
- *Yiwen leiju* 藝文類聚 (1879) — BIHEC (CIII 5-7 (1-8))
- *Chuxue ji* 初學記 (?) — BIHEC in Paris (SB 3701 (1-2))
- *Zhi bu zu zhai cong shu* 知不足齋叢書 (1921) — BIHEC (F X 2 (1-15) 1-120)
- *Mao Shi zhengyi* 毛詩正義 / Wikisource
- *Shuowen jiezi* 說文解字 / Wikisource

Numérisation du corpus

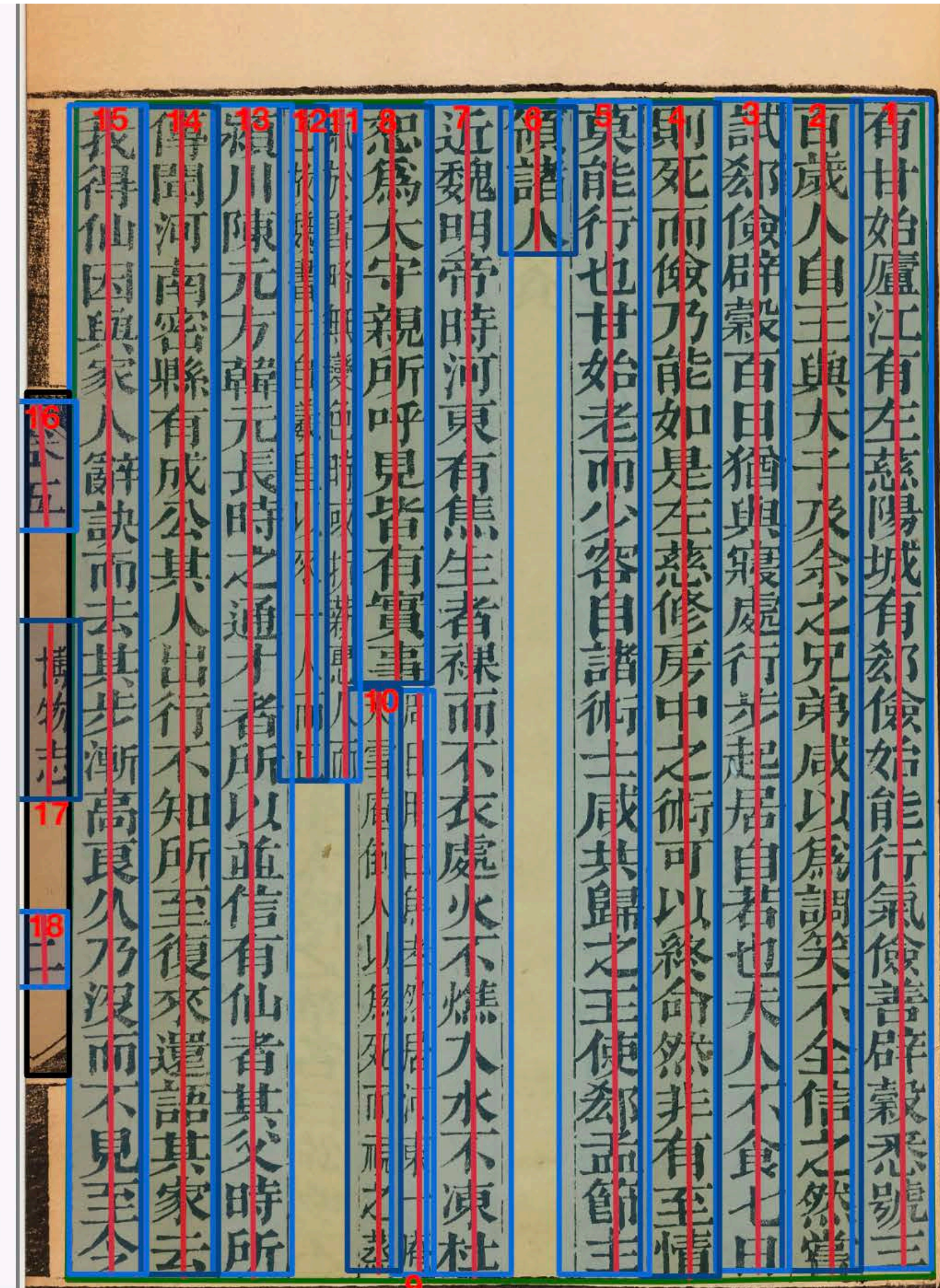


- Environ 56 000 pages
- Janvier-Juin 2023
- Modalités différentes selon les institutions de conservation (simple/double page, format de destination)
- Nommage uniformisé
- Une version conservée et exposée par les bibliothèques
- Une version versée dans l'entrepôt Nakala (modalités et métadonnées variées)

Annoter et transcrire

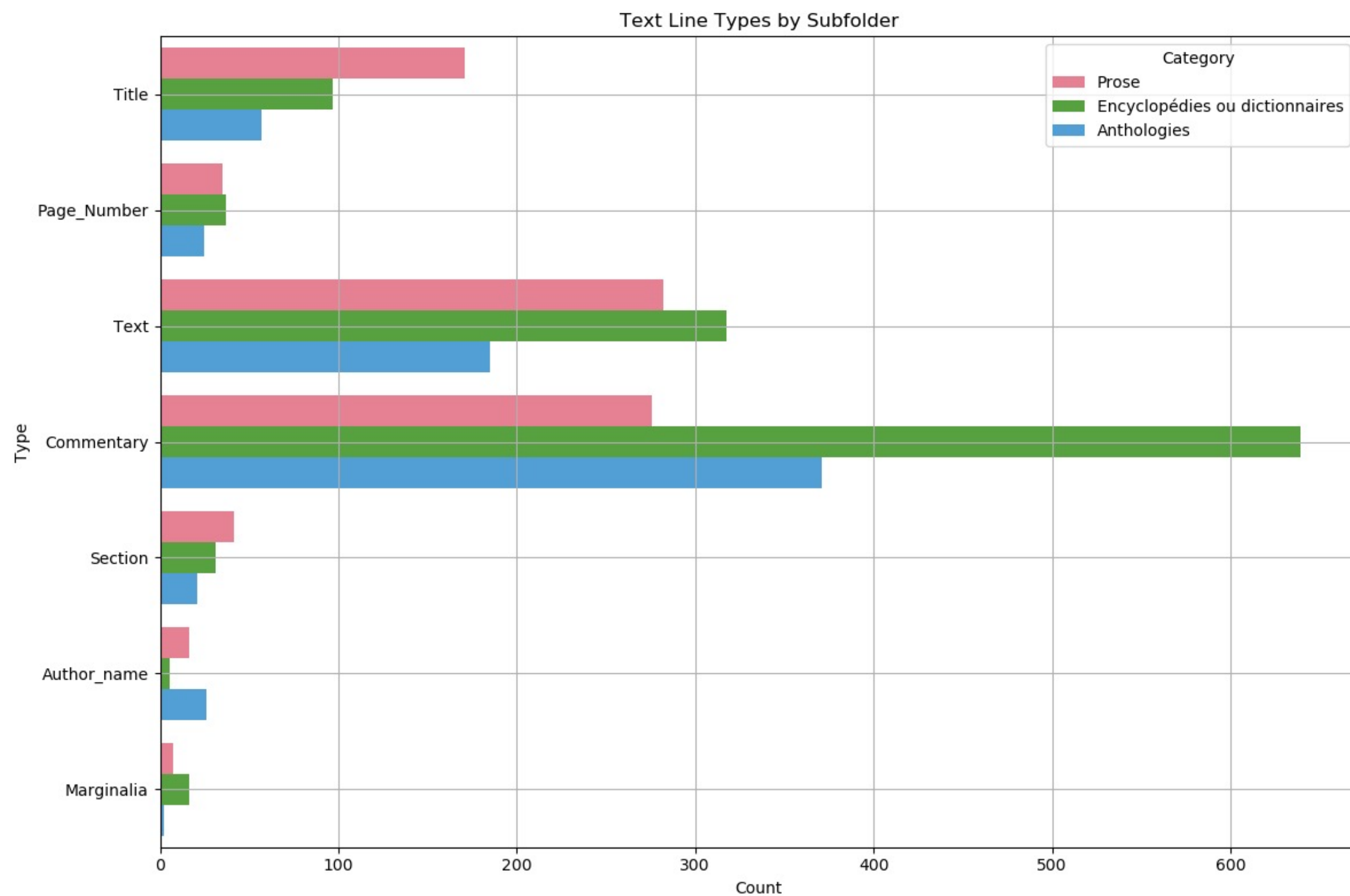
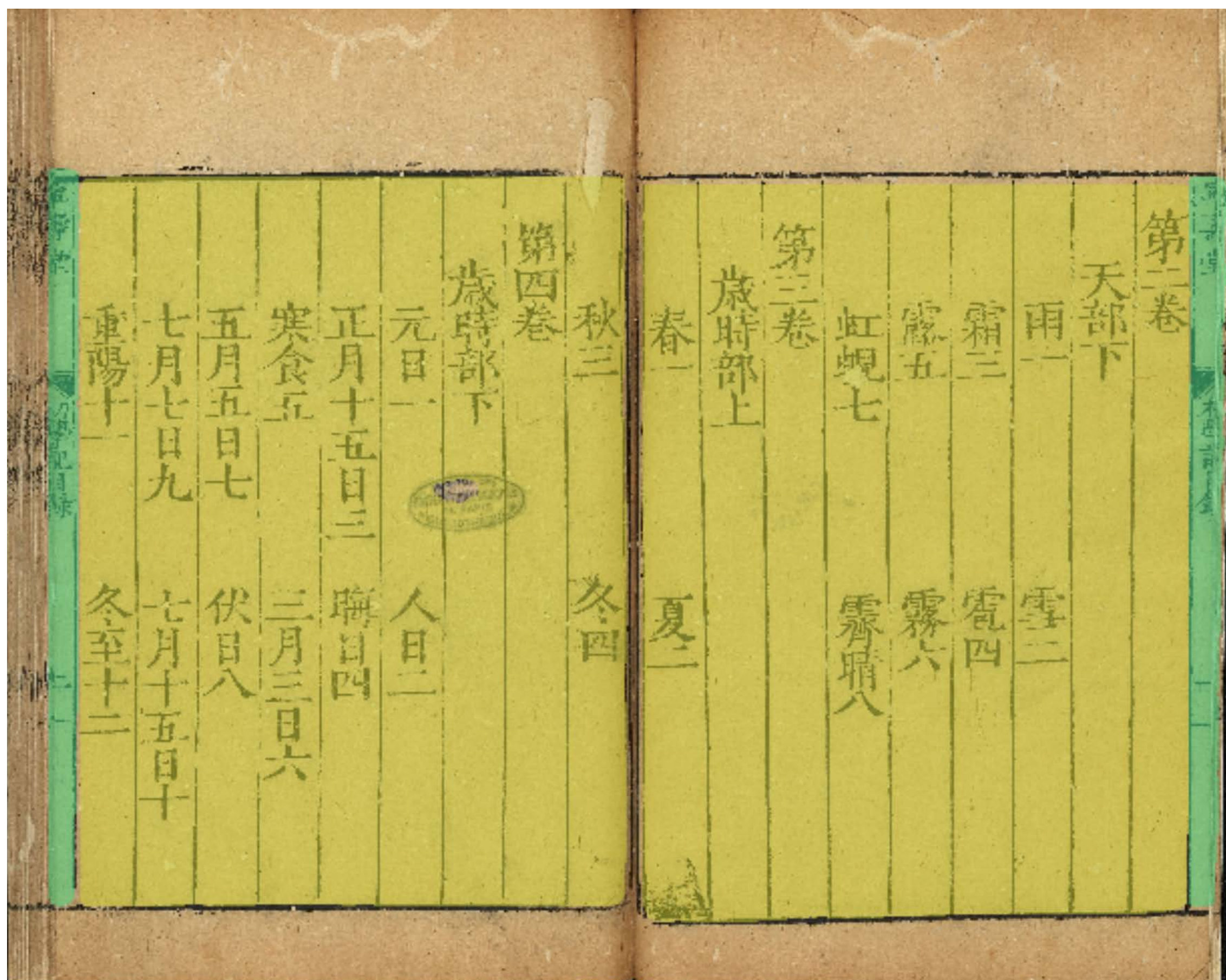


TextRegion 1 ^	MainText			
1	<u>有甘始廬江有左慈陽城有郗儉始能行氣儉善辟</u>			
2	<u>百歲人自王與太子及余之兄弟咸以為調笑不全</u>			
3	<u>試郗儉辟穀百日猶與寢處行步起居自若也夫人</u>			
4	<u>則死而儉乃能如是左慈修房中之術可以終命然</u>			
5	<u>莫能行也甘始老而少容自諸術士咸共歸之王使</u>			
6	<u>領諸人</u>			
7	<u>近魏明帝時河東有焦生者裸而不衣處火不焦入</u>			
8	<u>怨為太守親所呼見皆有實事</u>			
9	<u>周日用曰焦孝然居河東一庵</u>			
10	<u>大雪庵倒人以為死而視之蒸</u>			
11	<u>氣於雪略無變色時或析薪惠人而</u>			
12	<u>已故魏書云自羲皇以來一人而已</u>			
13	<u>潁川陳元方韓元長時之通才者所以並信有仙者</u>			

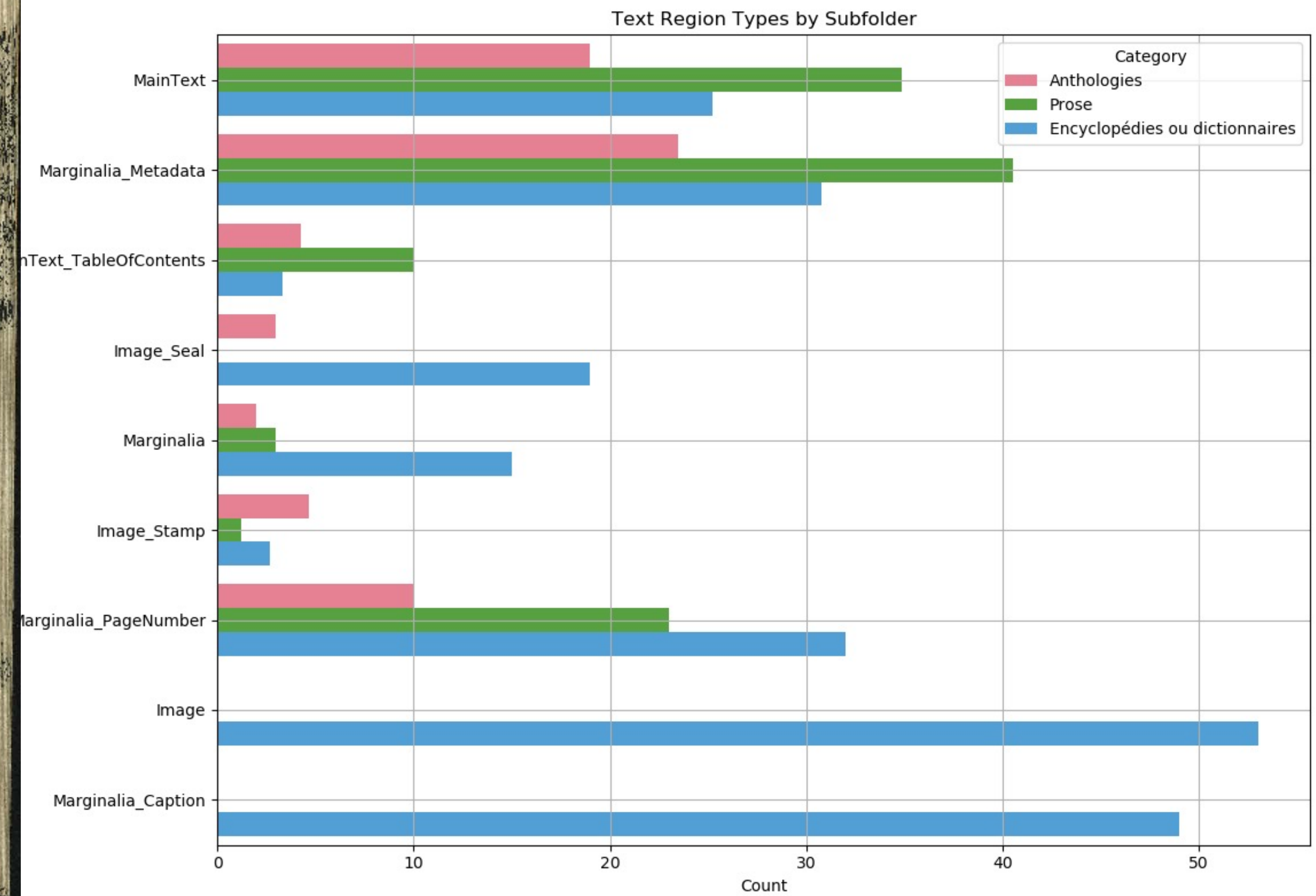
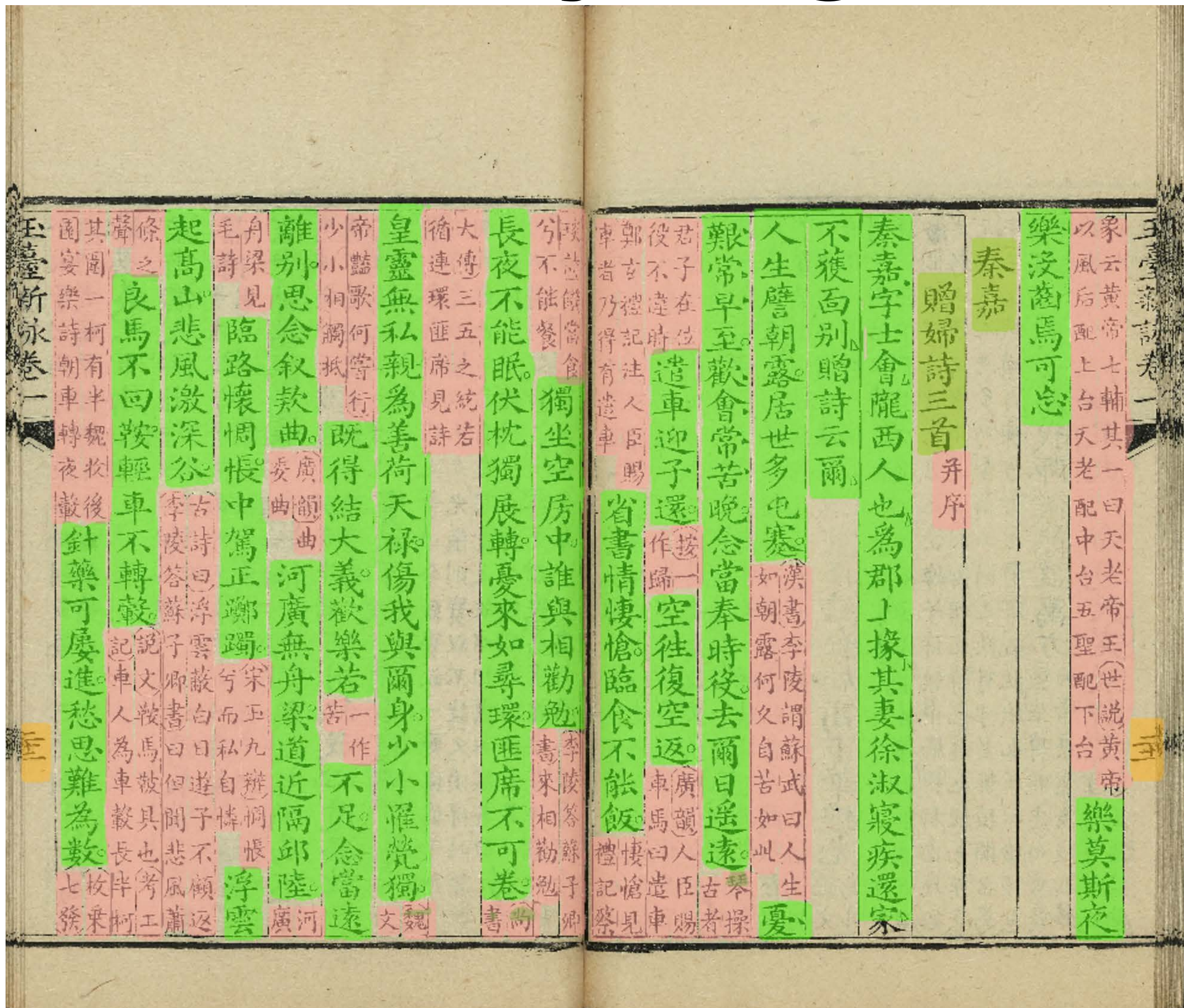


Annotateurs 2023 : Xinmin Hu,
Elsa Cuillé, Marie Bizais-Lillig,
Ani Tanelian, Anahide Kasparian,
Chahan Vidal-Gorène

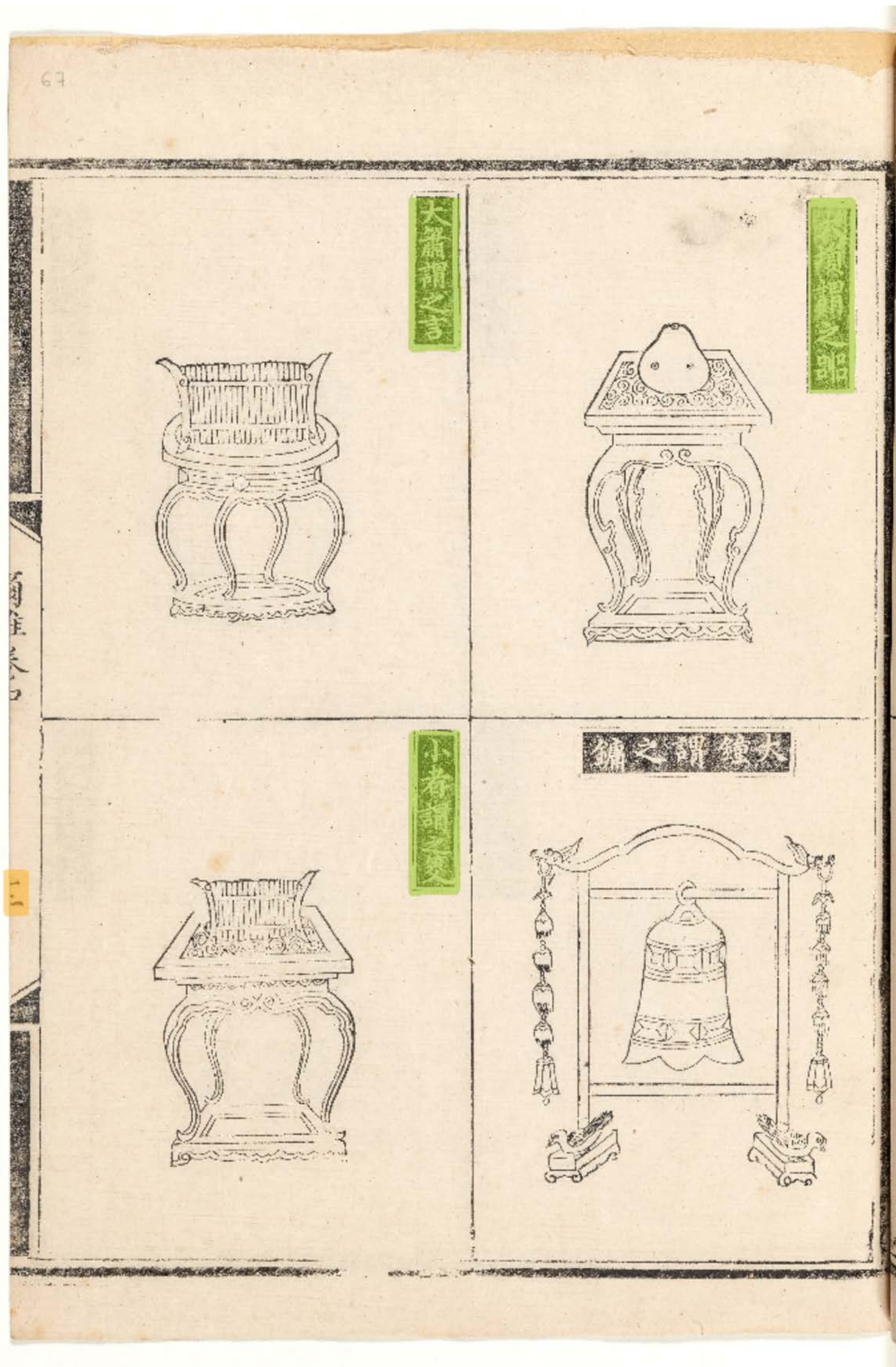
Typage des régions



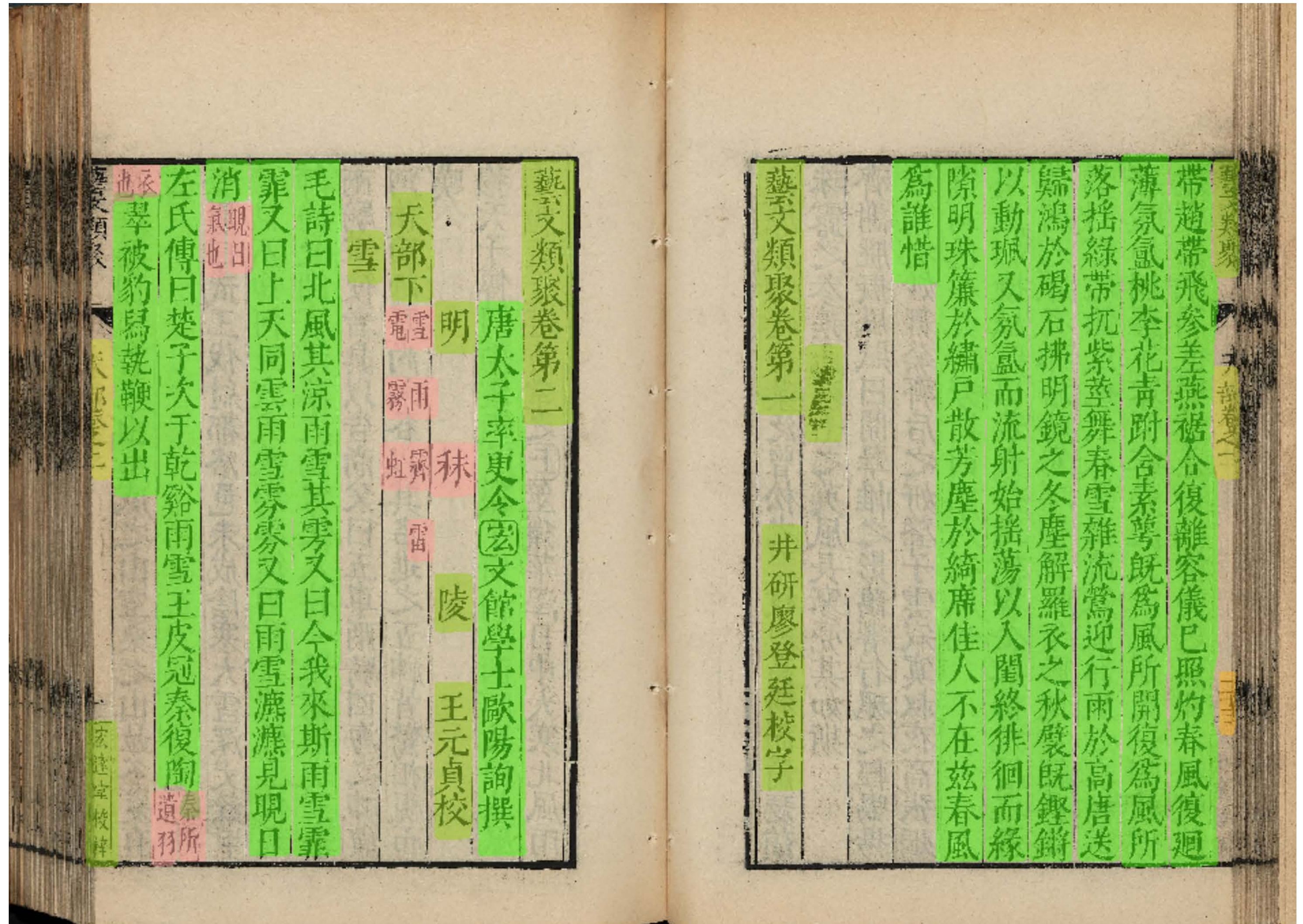
Typage des baselines



Typages



BULAC-BIULO-CHI-1938



BIHEC-C_III_5-7

Transcribe

- Les caractères tabous
- Les graphies historiques
- Les variantes

國家教育研究院 | NATIONAL ACADEMY for EDUCATIONAL RESEARCH | 最新消息 常見問題 編輯說明 字典附錄 顯示模式 網站導覽 請輸入查詢

教育部 異體字字典

國家教育研究院 維護

部首查詢 | 筆畫查詢 | 單字查詢 | 複合查詢 | 注音查詢 | 漢語拼音查詢 | 倉頡碼查詢 | 四角號碼查詢 | 附收字查詢

部 首 查 詢 【部首字音讀表】

部首筆畫數：9 部首：頁 部首外筆畫數：10

查詢結果：96字(共3頁) 每頁 40 ▾ 筆

顛	獺	額	嬾	頤	嬾	頤	頤	額	額
		正	附	2		正		1	2
額	額	額	額	額	鬢				
正		正		1	2				

上一頁 1 2 3 下一頁

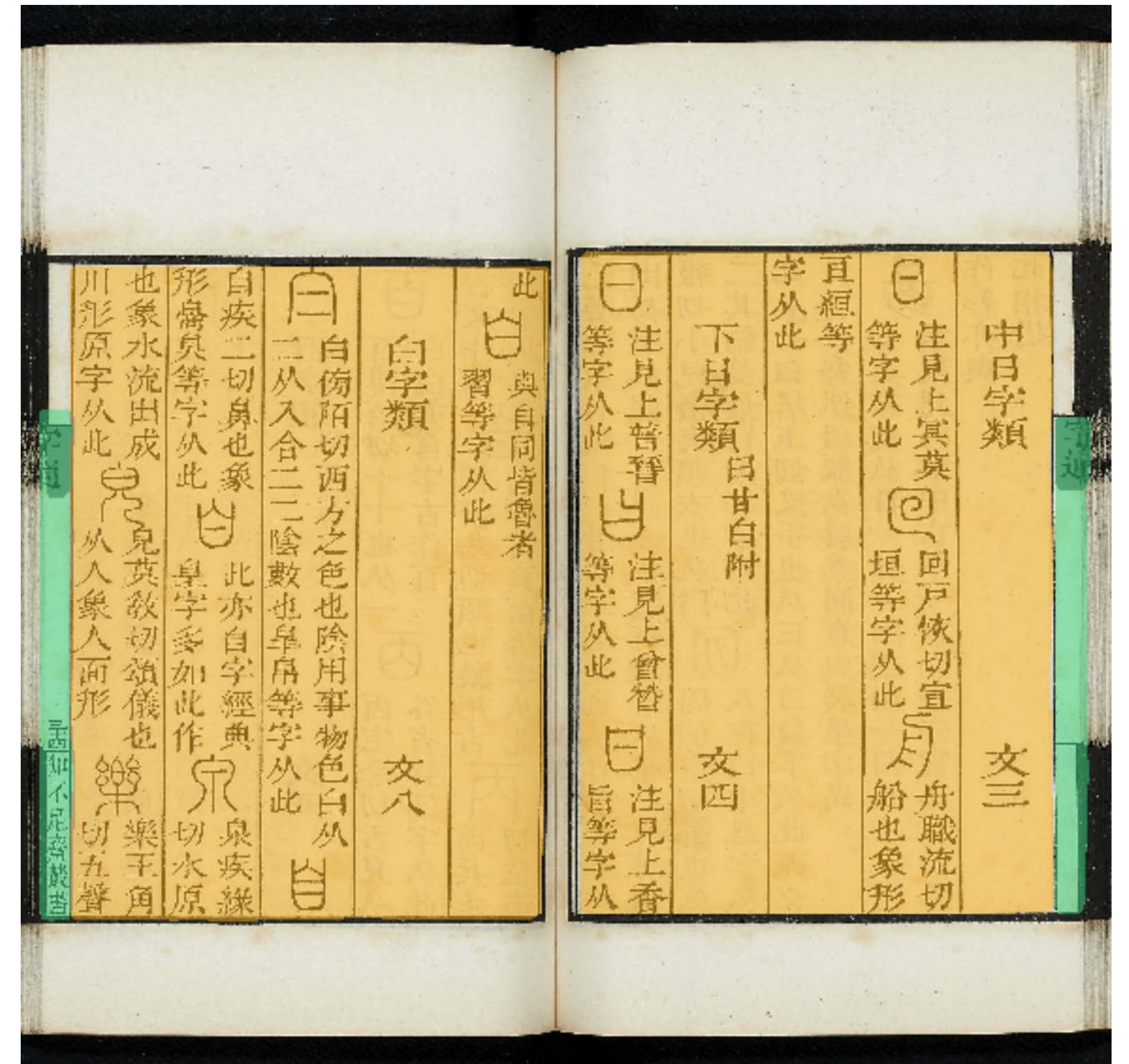
(異體字字典) 正式六版 中華民國教育部 版權所有 ©2017 Ministry of Education, R.O.C. All rights reserved.

國家教育研究院
NATIONAL ACADEMY for
EDUCATIONAL RESEARCH

個資法及隱私聲明 辭典公眾授權網 意見交流 網網相連

三峽總院區地址：新北市三峽區三樹路2號
臺北院區地址：臺北市大安區和平東路一段179號
中部辦公室：臺中市豐原區師範街67號
電話總機：(02)7740-7890、傳真：(02)7740-7064、TANet VoIP：9009-7890

線上人數：9
瀏覽總人次：55300929



BIHEC-FX2-27-214

Transcrire

- Images annotées : 306
- Régions annotées : 1 175
- Lignes annotées : 12 198
- Glyphes : 97 523
- Glyphes différents : 5 334

10 glyphes les plus fréquents

之	2 115	occurrences
日	1 464	occurrences
也	1 343	occurrences
不	885	occurrences
一	851	occurrences
以	841	occurrences
十	720	occurrences
有	720	occurrences
而	690	occurrences

1421 caractères n'existent qu'une seule fois dans la dataset.

Soit 26,6% de classes uniques

10 glyphes les plus rares

佯	
?	
膀	
窈	
?	
榧	
獲	
帕	
伉	
懾	

Premiers résultats

Reconnaissance

Pour les régions

Precision = 0.955

Rappel = 0.96

mAP = 0.968 (= mean average precision)

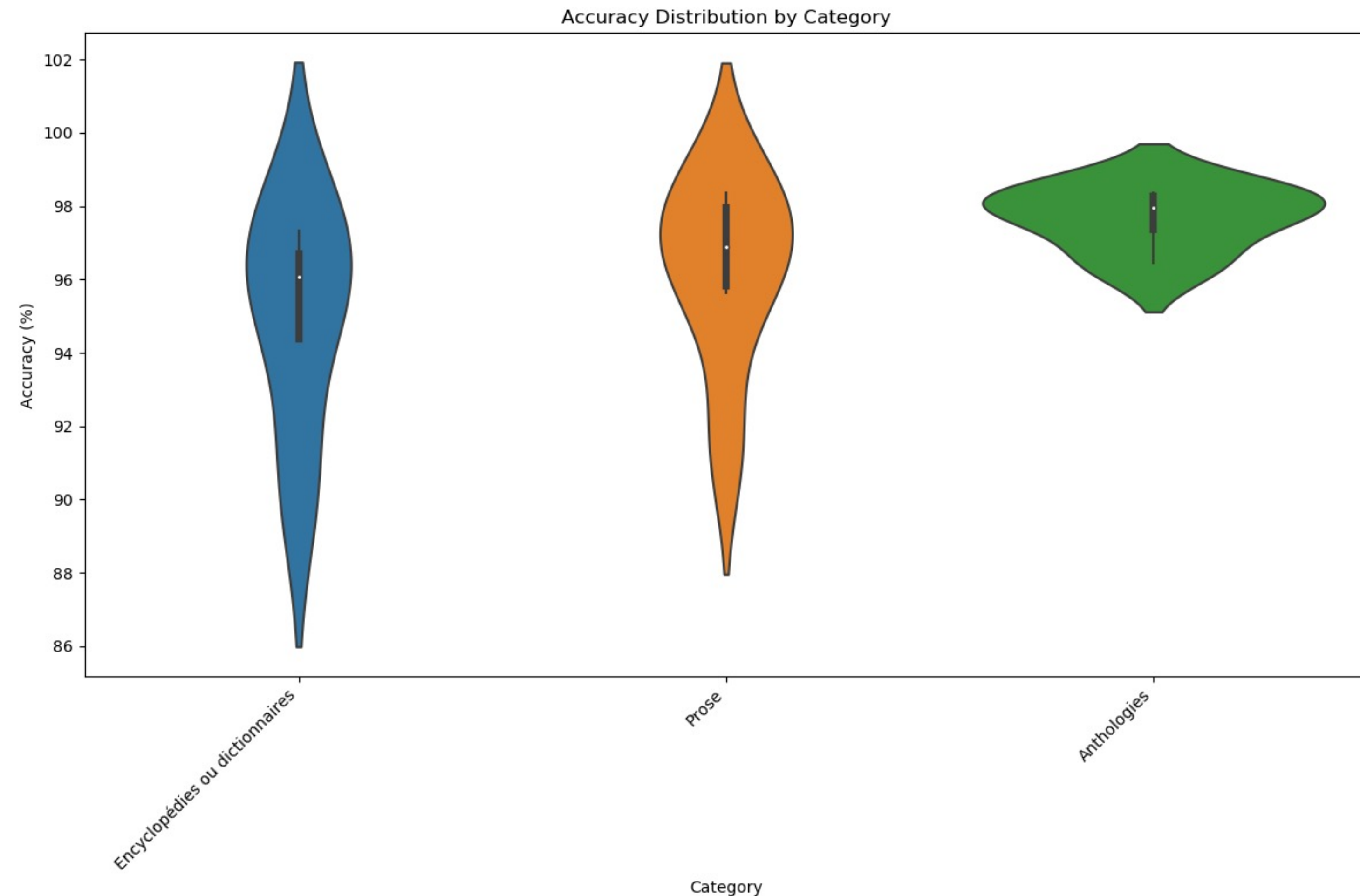
Pour les caractères

CER moyen 2%

Entre 86% et 100% selon les lignes

Grande variété selon genres

Pour les caractères



Illustration

鈞 埃塵言輕千鈞喻重也列子楊朱曰貴非所貴賤非所賤齊貴齊賤漢書曰十六兩為一斤三十斤為一鈞
 主父宦不達骨肉還相薄 史記或說主父偃曰太橫主一
 父偃曰臣結髮游學四十年 身不得遂親不以為子昆弟不收杜預左氏傳注曰官(宦)
 仕也呂氏春秋曰父母之於子也子之於父母也此之謂骨肉之親薄輕鄙之也史記曰君薄淮陽邪
 買臣困采樵伉儀(儷)不安宅 漢書曰朱買臣家貧常刈薪樵賣以給食擔束薪行且誦書妻亦負戴相隨數止買臣無謳歌道中買臣愈益疾歌妻養(羞)之求去買臣笑曰我年五十當富貴也今已四十餘矣汝苦曰(日)久待我富貴報汝功力妻恚怒曰如公等終餘(餓)死溝中耳能何富貴買臣不能留即聽去左氏傳曰施氏之婦怒施氏曰己不能庇其伉懼(儷)杜預曰曜(儷)偶也位(伉)敢(敵)也
 陳平無產業歸來翳負郭 漢書曰陳平家貧好讀書負郭窮巷以席為門然門外多長者車轍方言曰翳憂(憂)也郭璞曰謂蔽變(憂)也音愛鄭玄禮記注曰負之言背也長卿還成都壁立何寥廓 史記曰卓文君奔司馬相如相與馳歸成都居徒四壁立郭璞曰貧窮也楚辭文二二(十一)六

鈞

埃塵言輕千鈞喻重也列子楊朱曰貴非所貴賤非所賤齊貴齊賤漢書曰十六兩為一斤三十斤為一鈞

主父宦不達骨肉還相薄

史記或說主父偃曰太橫主一

父偃曰臣結髮游學四十年

身不得遂親不以為子昆弟不收杜預左氏傳注曰官(宦)

仕也呂氏春秋曰父母之於子也子之於父母也此之

謂骨肉之親薄輕鄙之

也史記曰君薄淮陽邪

買臣困采樵伉儀(儷)不安宅

漢書

曰朱

買臣家貧常刈薪樵賣以給食擔束薪行且誦書妻亦

負戴相隨數止買臣無謳歌道中買臣愈益疾歌妻養(羞)

之求去買臣笑曰我年五十當富貴也今已四十餘矣

汝苦曰(日)久待我富貴報汝功力妻恚怒曰如公等終餘(餓)

死溝中耳能何富貴買臣不能留即聽去左氏傳曰施

氏之婦怒施氏曰己不能庇其伉懼(儷)杜預曰曜(儷)偶也位(伉)

敢(敵)

也

陳平無產業歸來翳負郭

漢書曰陳平家貧好讀書

負郭窮巷以席為門

然門外多長者車轍方言曰翳憂(憂)也郭璞曰

謂蔽變(憂)也音愛鄭玄禮記注曰負之言背也

長卿還成

都壁立何寥廓

史記曰卓文君奔司馬相如相與馳歸

成都居徒四壁立郭璞曰貧窮也楚辭

文二二(十一)

六

Illustration

都壁立何寥廓
 成都居徒四壁立郭璞曰貧窮也楚辭
 謂蔽夢也音愛鄭玄禮記注曰翳夢也郭璞曰
 然門外多長者車轍方言曰翳夢也郭璞曰
陳平無產業歸來翳負郭
 漢書曰陳平家貧好讀
 也敵氏死汝之負買也謂仕身不主
 陳平無產業歸來翳負郭
 漢書曰陳平家貧好讀
 也敵氏死汝之負買也謂仕身不主
主父宦不達骨肉還相薄
 史記或說主父偃曰太橫主
 父偃曰臣結髮游學四十年
 身不得遂親不以為子昆弟不收杜預左氏傳注曰官(宦)
 仕也呂氏春秋曰父母之於子也子之於父母也此之
 謂骨肉之親薄輕鄙之
 也史記曰君薄淮陽邪
買臣困采樵伉儷不安宅
 漢書
 曰朱
 買臣家貧常刈薪樵賣以給食擔束薪行且誦書妻亦
 負戴相隨數止買臣無謳歌道中買臣愈益疾歌妻養(羞)
 之求去買臣笑曰我年五十當富貴也今已四十餘矣
 汝苦曰(日)久待我富貴報汝功力妻恚怒曰如公等終餘(餓)
 死溝中耳能何富貴買臣不能留即聽去左氏傳曰施
 氏之婦怒施氏曰己不能庇其伉懼(儷)杜預曰曜(儷)偶也位(伉)
敢(敵)
 也
 陳平無產業歸來翳負郭
 漢書曰陳平家貧好讀
 書負郭窮巷以席為門
 然門外多長者車轍方言曰翳憂(憂)也郭璞曰
 謂蔽變(憂)也音愛鄭玄禮記注曰負之言背也
 長卿還成
 都壁立何寥廓
 史記曰卓文君奔司馬相如相與馳歸
 成都居徒四壁立郭璞曰貧窮也楚辭
 文二二(十一)
 六

鈞

埃塵言輕千鈞喻重也列子楊朱曰貴非所貴賤非所賤齊貴齊賤漢書曰十六兩為一斤三十斤為一鈞

主父宦不達骨肉還相薄

史記或說主父偃曰太橫主

父偃曰臣結髮游學四十年

身不得遂親不以為子昆弟不收杜預左氏傳注曰官(宦)

仕也呂氏春秋曰父母之於子也子之於父母也此之

謂骨肉之親薄輕鄙之

也史記曰君薄淮陽邪

買臣困采樵伉儷不安宅

漢書

曰朱

買臣家貧常刈薪樵賣以給食擔束薪行且誦書妻亦

負戴相隨數止買臣無謳歌道中買臣愈益疾歌妻養(羞)

之求去買臣笑曰我年五十當富貴也今已四十餘矣

汝苦曰(日)久待我富貴報汝功力妻恚怒曰如公等終餘(餓)

死溝中耳能何富貴買臣不能留即聽去左氏傳曰施

氏之婦怒施氏曰己不能庇其伉懼(儷)杜預曰曜(儷)偶也位(伉)

敢(敵)

也

陳平無產業歸來翳負郭

漢書曰陳平家貧好讀

書負郭窮巷以席為門

然門外多長者車轍方言曰翳憂(憂)也郭璞曰

謂蔽變(憂)也音愛鄭玄禮記注曰負之言背也

長卿還成

都壁立何寥廓

史記曰卓文君奔司馬相如相與馳歸

成都居徒四壁立郭璞曰貧窮也楚辭

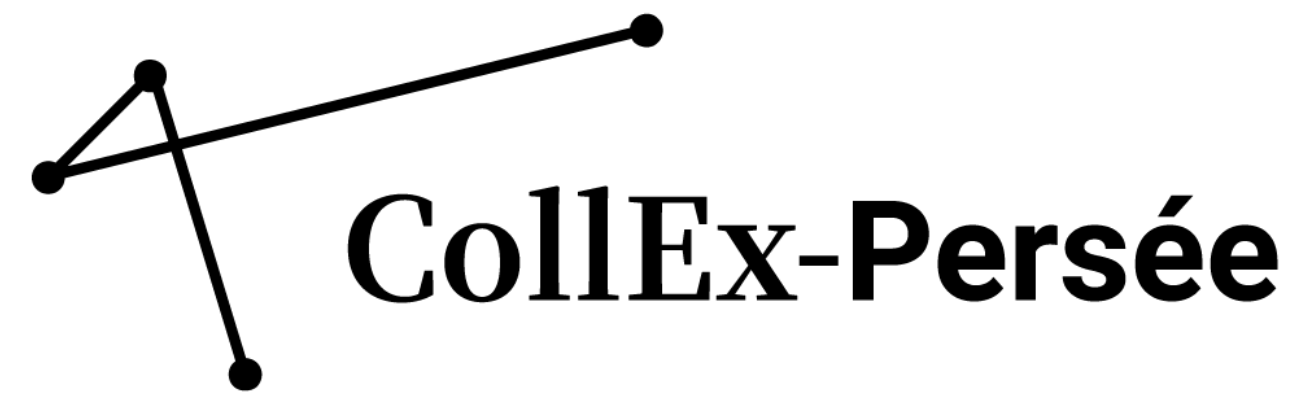
文二二(十一)

六

En cours et à venir

Partager / Mettre à disposition

- Les modèles et données d'entraînement de l'HTR (Github de Calfa + Zenodo) avec signalement (HTR United)
- Les images (entrepôt Nakala)
- Les fichiers en XML-TEI (entrepôt Nakala)
- La base de données (bio-bibliographique et lexicale)
- Sous licence CC BY-NC 4.0



Merci.



bizais@unistra.fr

<https://gitlab.huma-num.fr/chi-know-po>
<https://github.com/calfa-co/chi-know-po>

Les soutiens

Travail d'exploration fondé sur

- Une base de données :
 - lexique
 - base temporelle
 - titres
 - personnes
- Des outils d'exploration associés
- Contributeurs :
 - Tilman Schalmey
 - Shueh-Ying Liao
 - Marie Bizais-Lillig
- Des scripts (Python) d'exploration :
 - cooccurrences
 - reprises ou échos
- Des scripts associés de visualisation
- Contributeurs :
 - Mariana Zorkina
 - Tilman Schalmey
 - Xinmin Hu
 - (Marie Bizais-Lillig)
- Un corpus :
 - 4 schémas modèles XML-TEI
 - odd + RelaxNG
 - fichiers image + manifestes IIF
- Contributeurs :
 - Ilaine Wang
 - Shueh-Ying Liao
 - Xinmin Hu
 - Marie Bizais-Lillig

Inviter à la correction des textes

TACT, PLATEFORME DE TRANSCRIPTION ET D'ANNOTATION DE CORPUS TEXTUELS

TACT est une plateforme de transcription et d'annotation collaborative qui accueille des projets de recherche et permet à chacun d'apporter sa contribution ou de solliciter la communauté de la plateforme.

Vous voulez participer à un projet ? Vous pouvez transcrire et annoter des textes originaux, relire et vérifier des transcriptions faites par d'autres internautes, consultez ce [manuel du contributeur](#) et allez [explorer les projets](#) !

Vous souhaitez déposer un projet sur TACT ? N'hésitez pas à nous en faire la demande via [ce formulaire](#). Vous deviendrez gestionnaire de votre projet TACT, vous pourrez en définir les modalités et gérer les contributions. Nous vous invitons à lire ce [manuel du gestionnaire](#) qui vous donnera toutes les informations nécessaires ainsi que la marche à suivre. Bonne lecture !

Vous pouvez aussi nous contacter à cette adresse : tact@univ-grenoble-alpes.fr

Dépôt/liens sur TACT

Correction des erreurs de l'HTR

Correction de la structuration automatisée

Correction/enrichissement de la détection automatisée des entités nommées