



**HAL**  
open science

## Experiments for automatic treatment of ancient Chinese sources

Marie Bizais-Lillig, Chahan Vidal-Gorène

► **To cite this version:**

Marie Bizais-Lillig, Chahan Vidal-Gorène. Experiments for automatic treatment of ancient Chinese sources. Documents anciens et reconnaissance automatique des écritures manuscrites, Ariane Pinche; Jean-Baptiste Camps, Jun 2022, Paris, France. hal-03792246

**HAL Id: hal-03792246**

**<https://hal.science/hal-03792246v1>**

Submitted on 16 Jan 2025

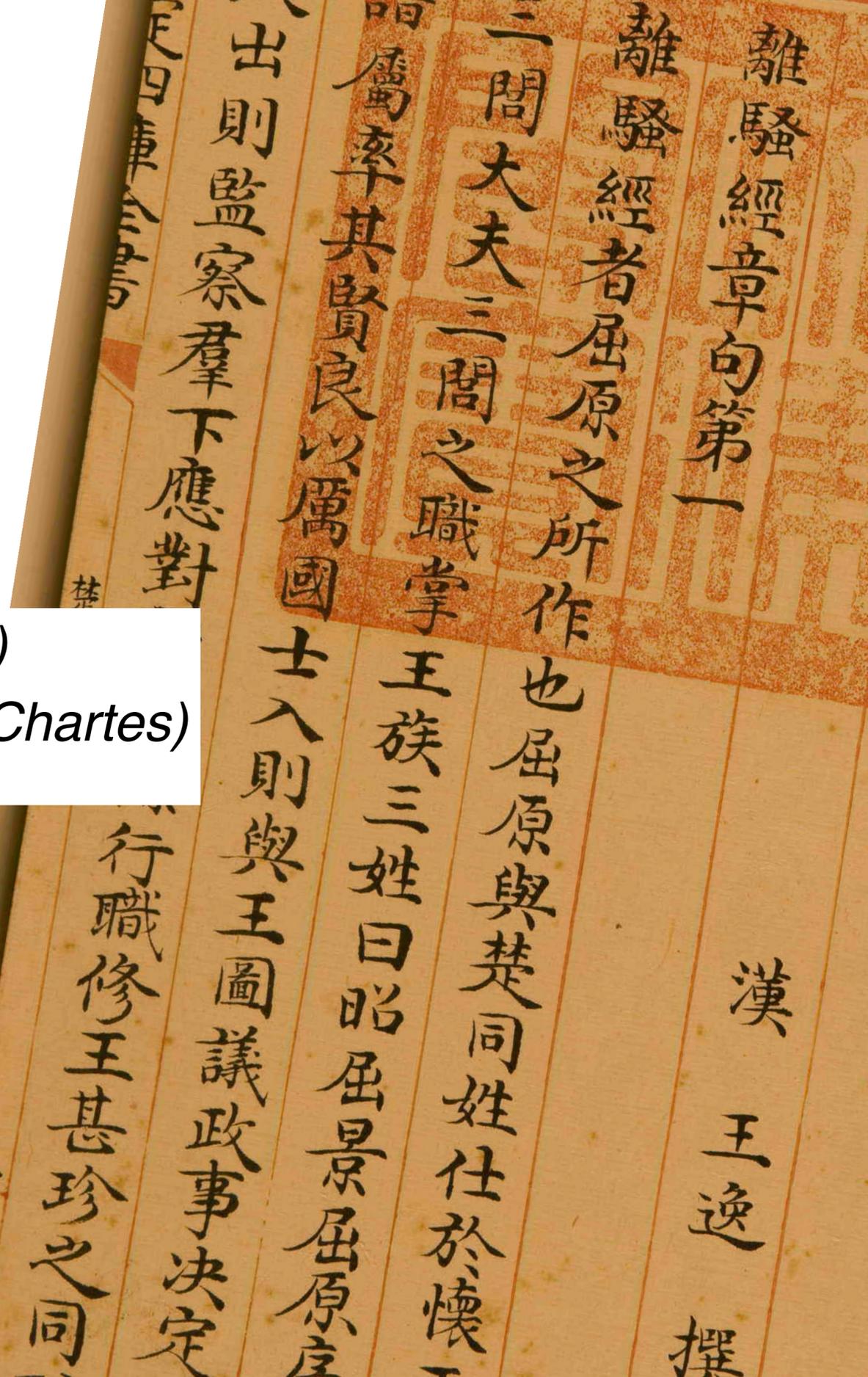
**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

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.

# Expérimentations pour l'analyse automatique de sources chinoises anciennes

Marie Bizais-Lillig (*Université de Strasbourg*)

Chahan Vidal-Gorène (*École nationale des Chartes*)

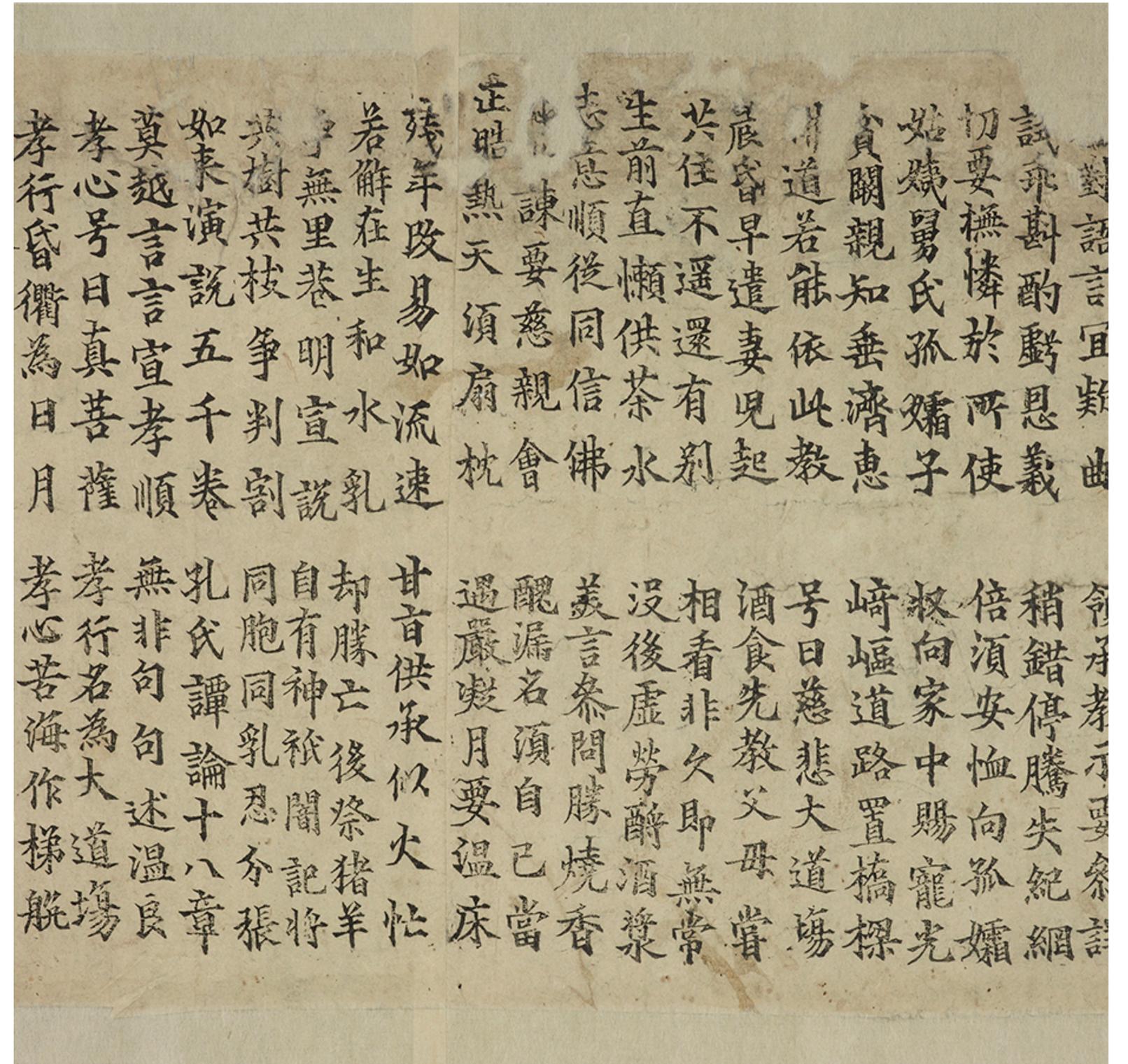


- Projet de fouille de textes chinois, pour l'essentiel du Moyen âge (IIe-Xe s.) : CHI-KNOW-PO
- Nécessite :
  - l'acquisition du corpus dans son intégralité et sans restriction de droits
  - la structuration du corpus (dimension éditoriale)
- Mai-juin 2021 : début d'une collaboration entre
  - des sinologues (de la Chine impériale)
  - l'équipe de Calfa pour l'HTR
- Nos craintes : les limites que poserait le nombre de classes

- Spécificités du corpus
- Les premières expérimentations
- La question des classes

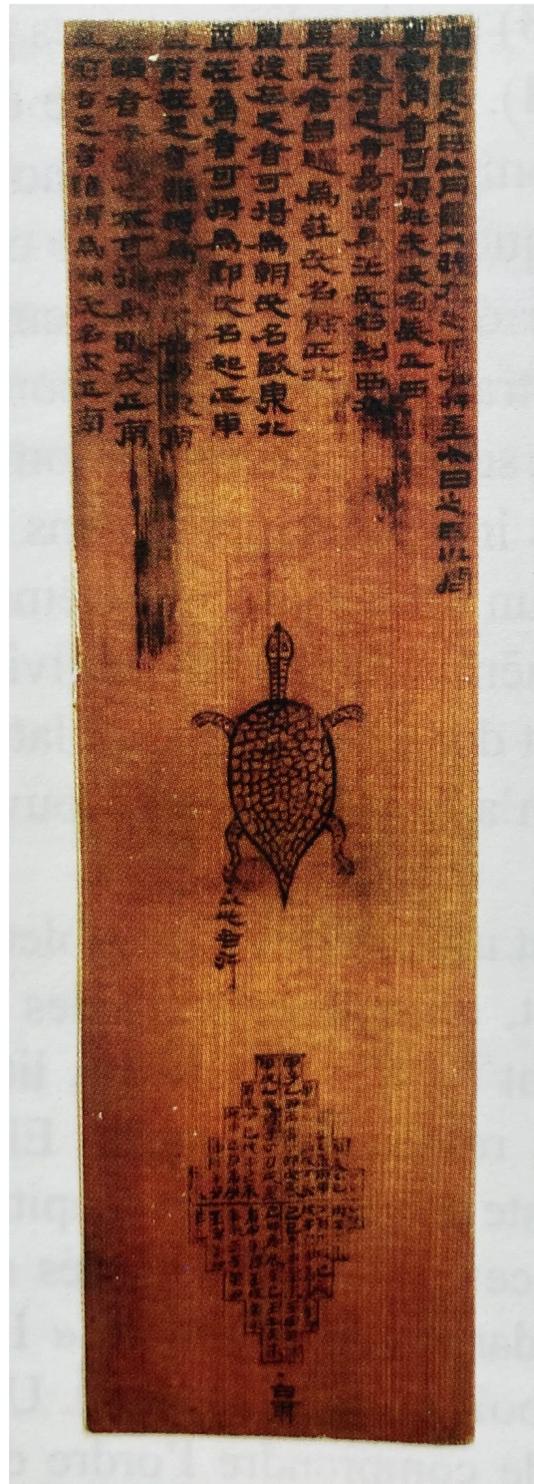
# Spécificités du corpus - des textes xylographiés

- Des textes établis entre 100 AE et 1000 NE
- Support : papier
- La xylographie au service de la conservation et de la dissémination
- Xylographie *versus* caractères mobiles



Impression xylographique (Dunhuang - British Library - OR8210)

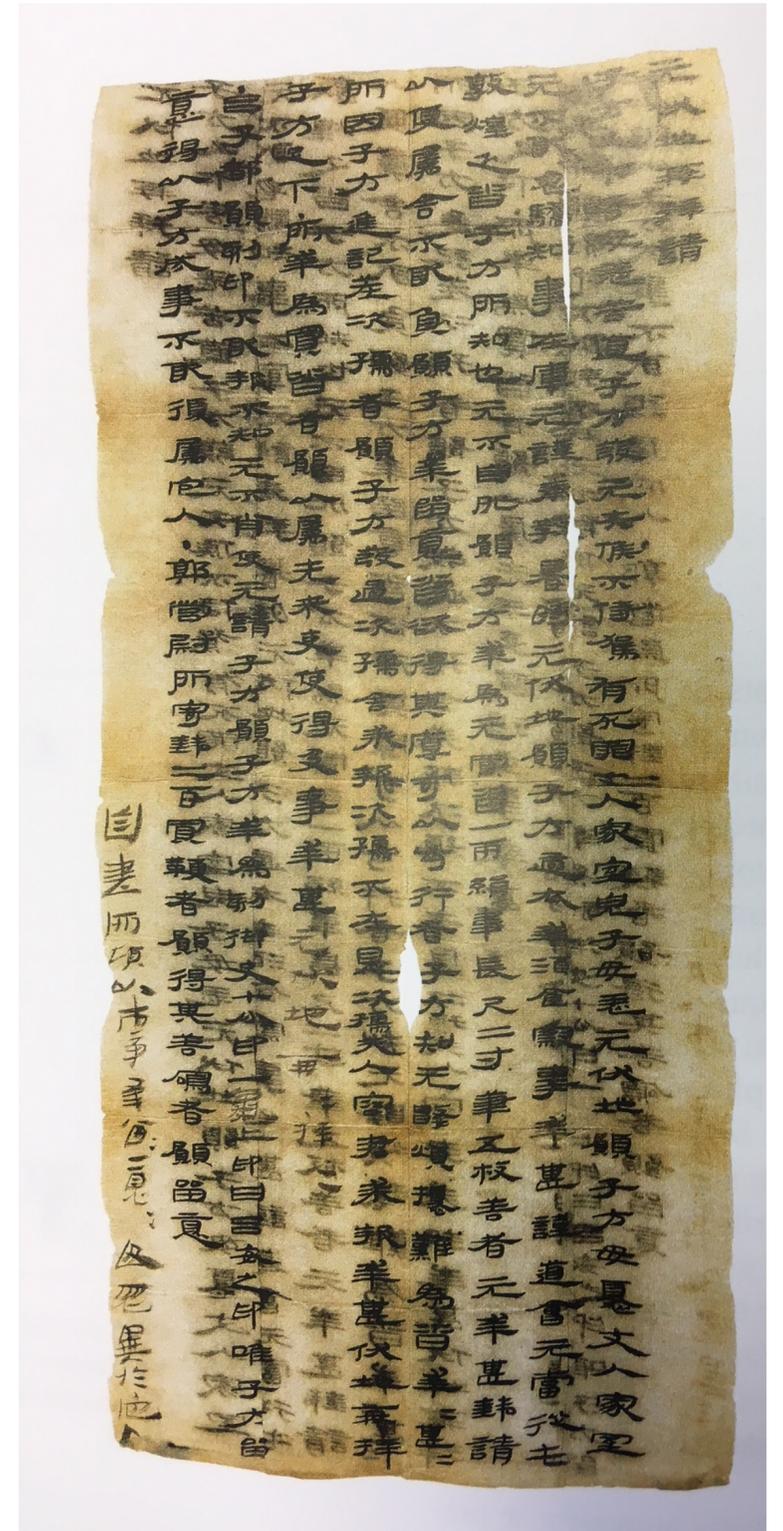
# Spécificités du corpus



plaquette de bois (ca. 10 NE - Musée de Lianyungang)  
source : *La Fabrique du lisible*



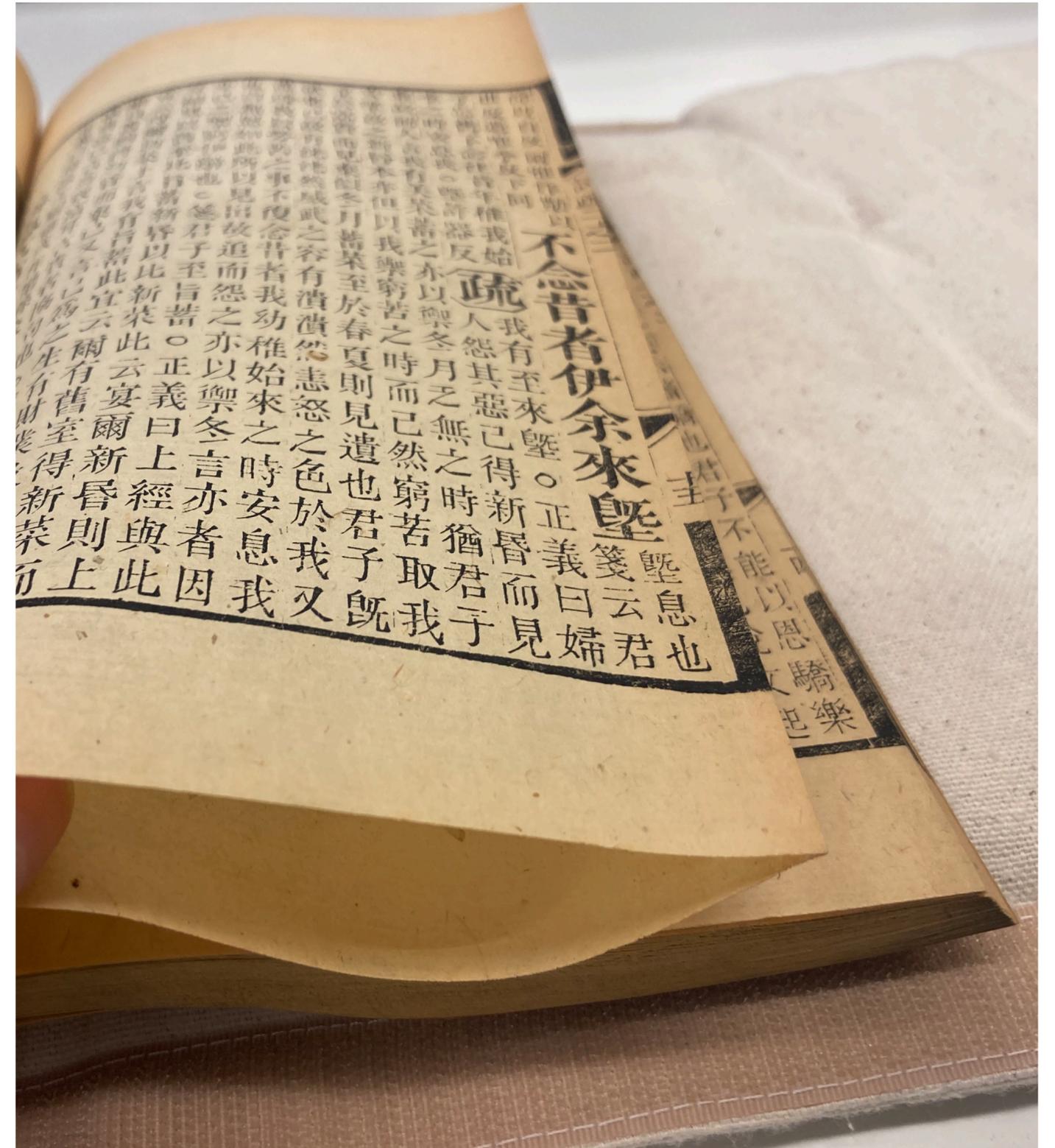
lamelles de bambou (*Shijing* 詩經 - ca. 250 AE - Université de l'Anhui)  
source : 安徽大学藏战国竹简 (一)



soie (lettre - 1er s. AE - Dunhuang)  
source : *La Fabrique du lisible*

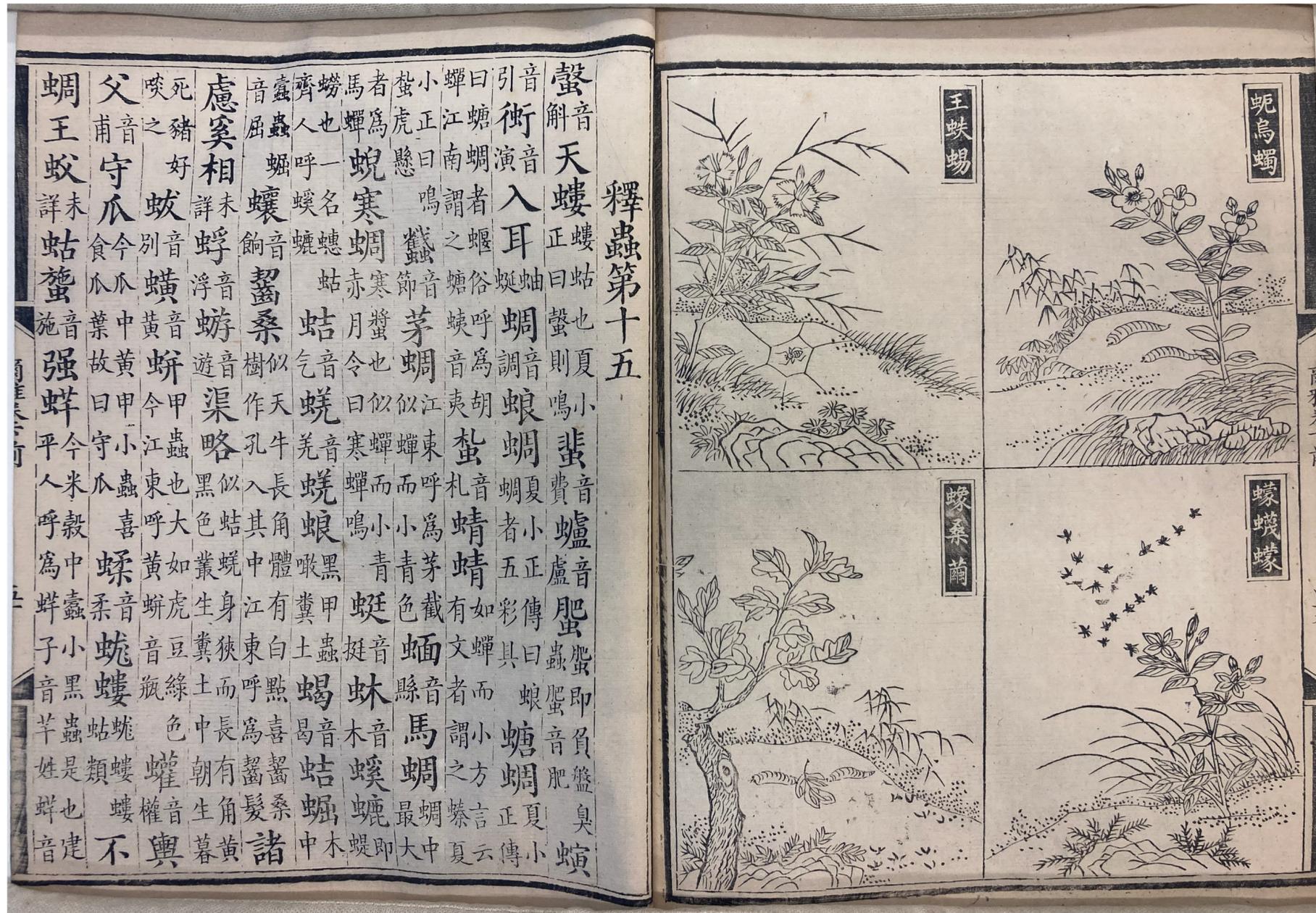
# Spécificités du corpus - mise en page

- Des textes transmis de pair avec un ou plusieurs commentaires
- Mise en page : entrelacement (à partir du IIe s. NE)
- Impression sur feuillets longs
  - pliés en accordéon
  - reproduits avec 2, 4 ou 6 double-pages par page (densité)

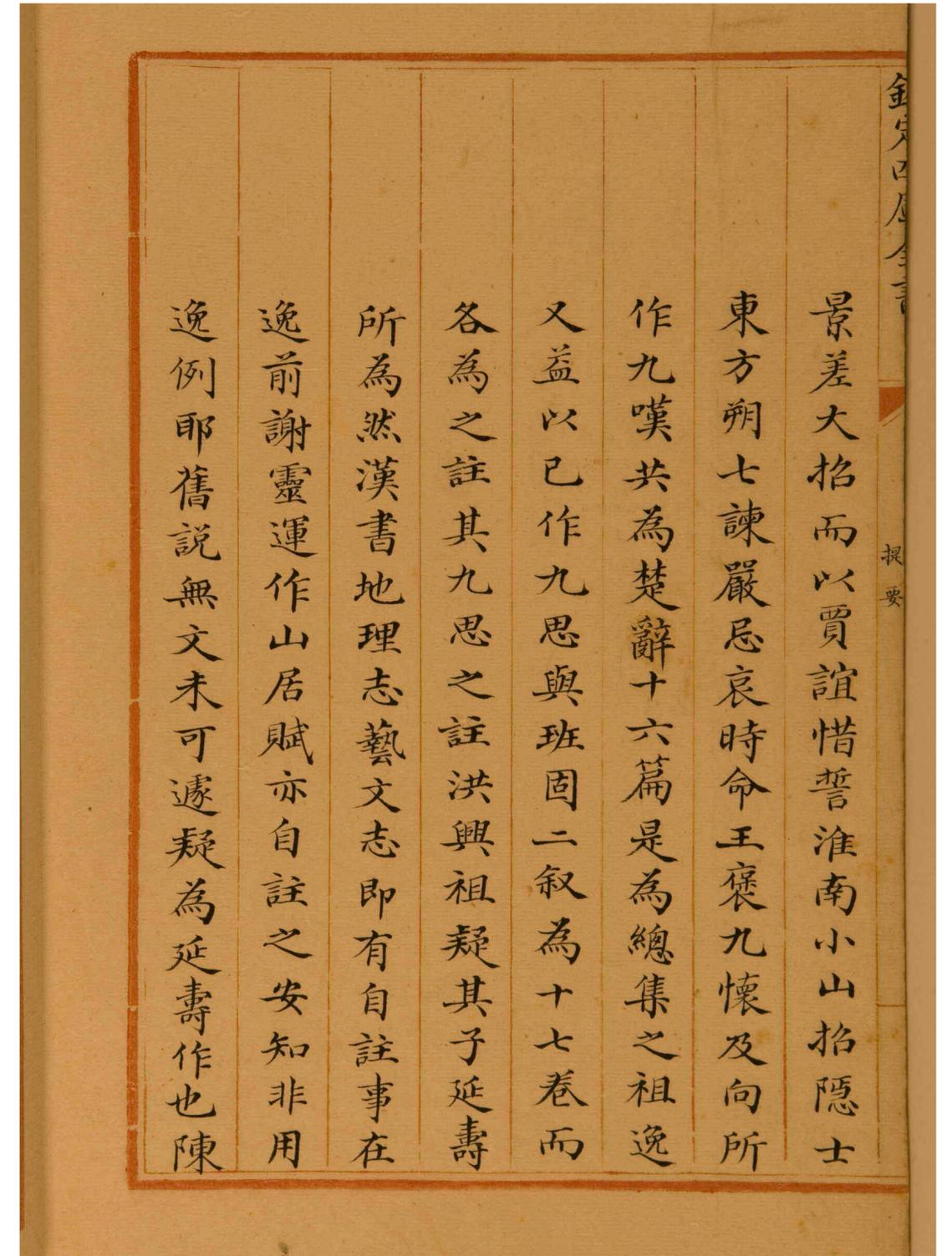


Reliure « accordéon » (BULAC)

# Spécificités du corpus



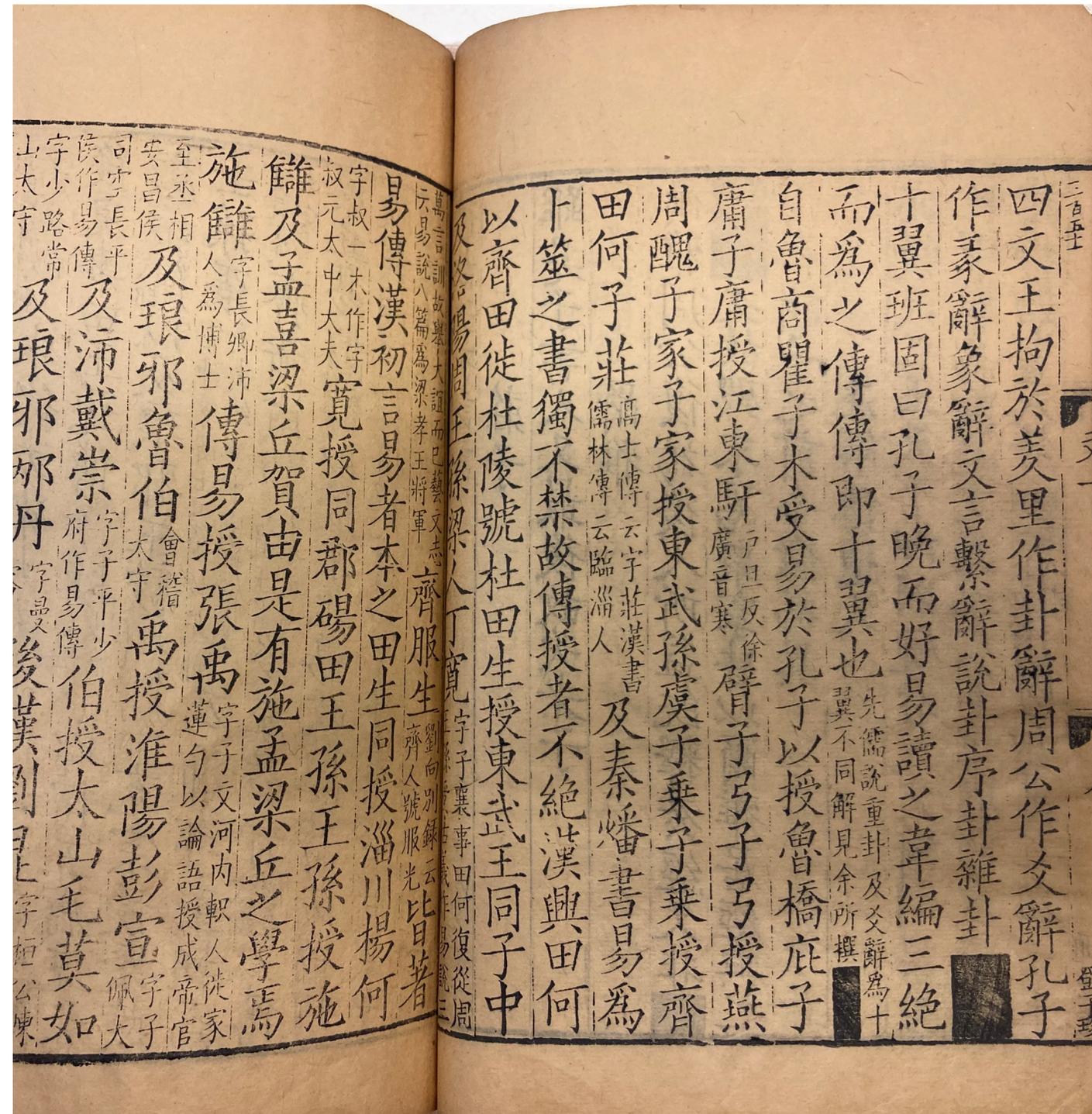
Entrelacement, découpage  
 爾雅 édition xylographiée illustrée (1801)  
 BULAC CHI1938(2)



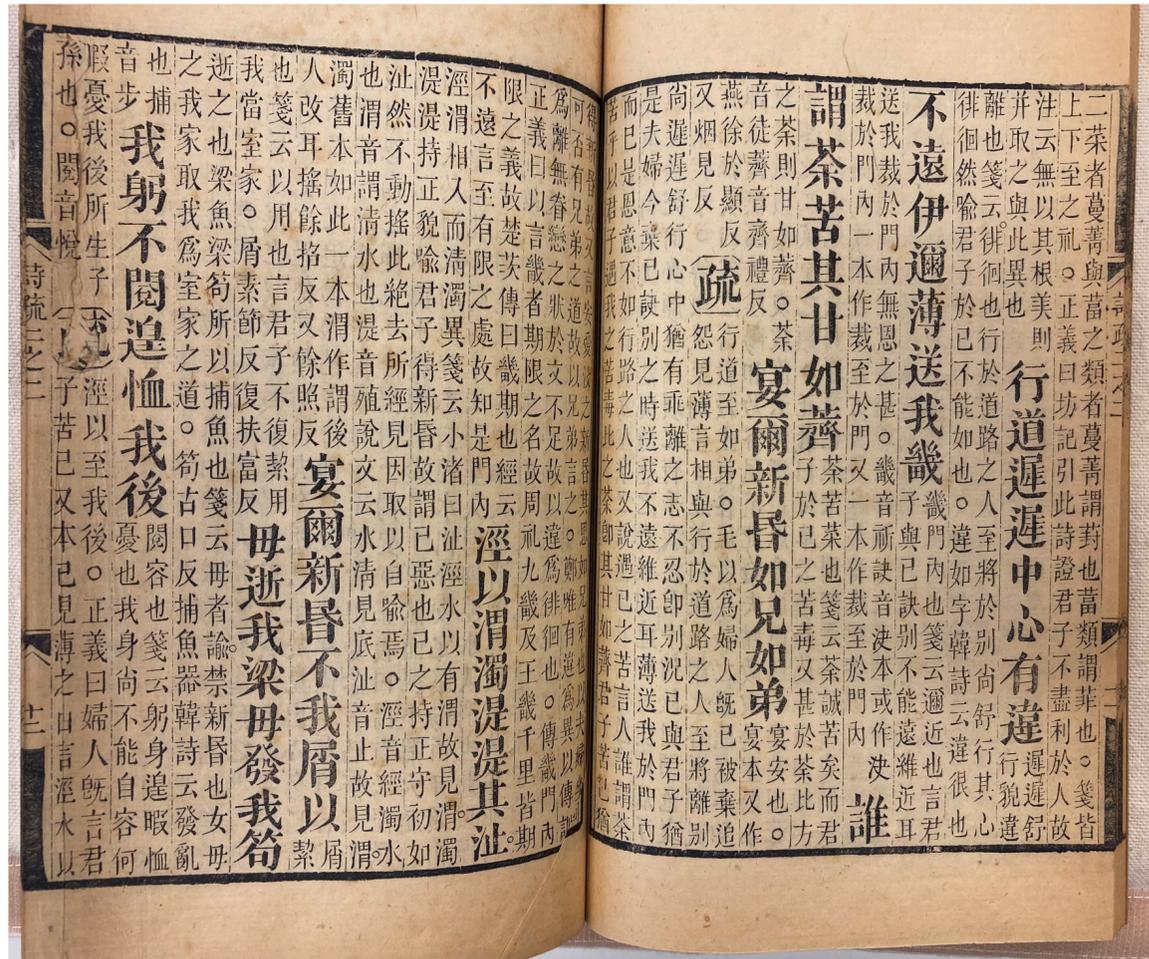
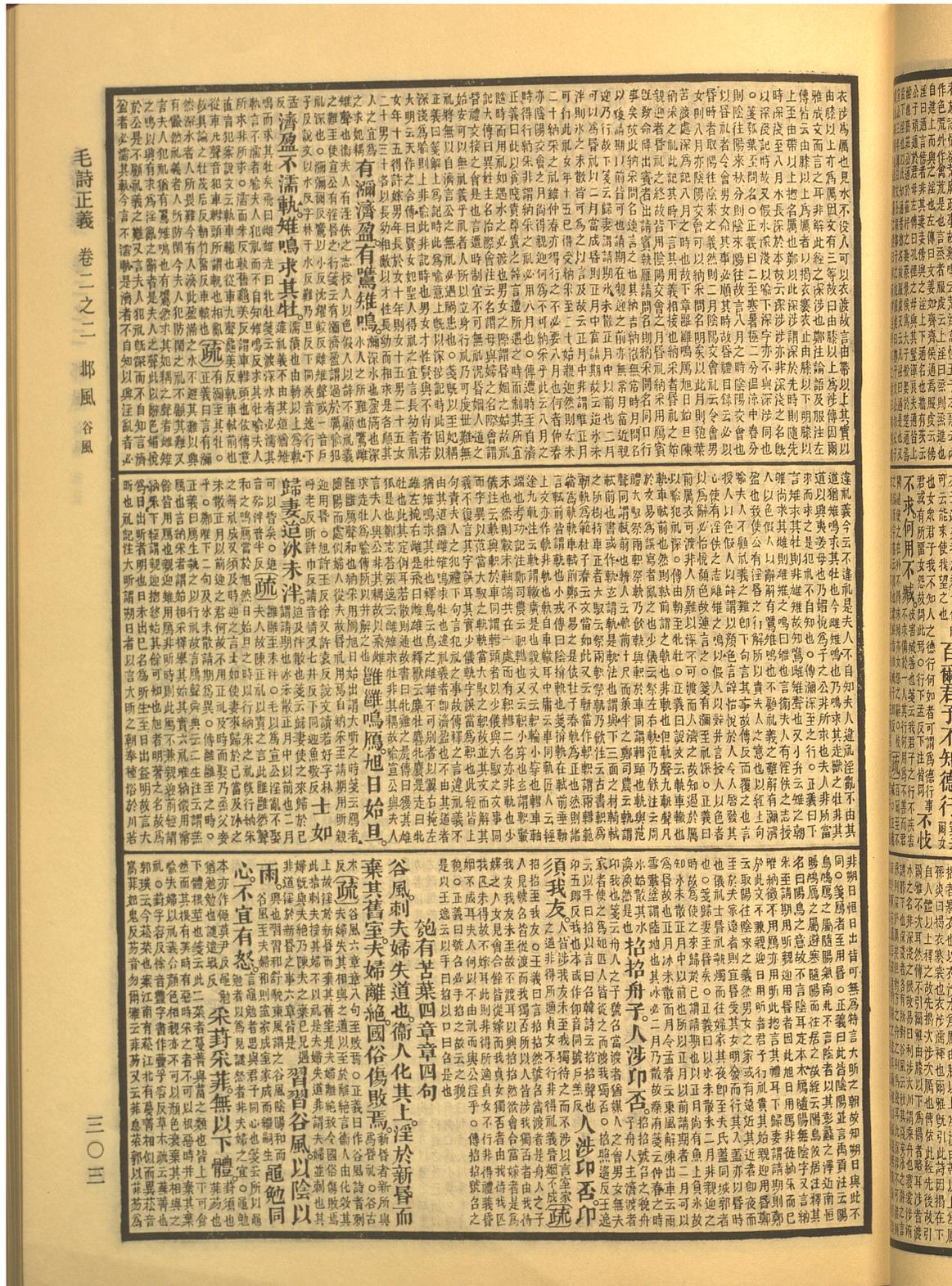
Édition impériale manuscrite  
 楚辭 dans le Siku

# Spécificités du corpus - l'écriture chinoise

- Un lexique s'appuyant sur plus de 54 000 caractères (glyphes) différents
- Dans les textes anciens, la variété des caractères est grande
- Corpus de départ : > 7000 glyphes distincts
- En outre : variétés de styles / graphies (avec des équivalents unicodes distincts)
- Des caractères tabous sous chaque empereur



# Spécificités du corpus - le Classique des Poèmes xylographié



Édition XIIe, impression 1815 (BULAC-CHI518(3))

Planche à xylographier (à partir des IXe-Xe s.)

毛詩正義 卷二之一 邶風 谷風

三〇三

# Les premières expérimentations - généralités

Corpus annoté de 50 images pour les expérimentations :

- 3 240 lignes en apprentissage
- 92 234 caractères

Approche par « baseline » initialement favorisée pour faciliter l'échange de ces données

Approche itérative sur la plateforme Calfa Vision : analyse automatique puis vérification manuelle

=> <https://vision.calfa.fr>

=> spécialisation rapide de la plateforme sur la détection de la mise en page (régions et lignes)

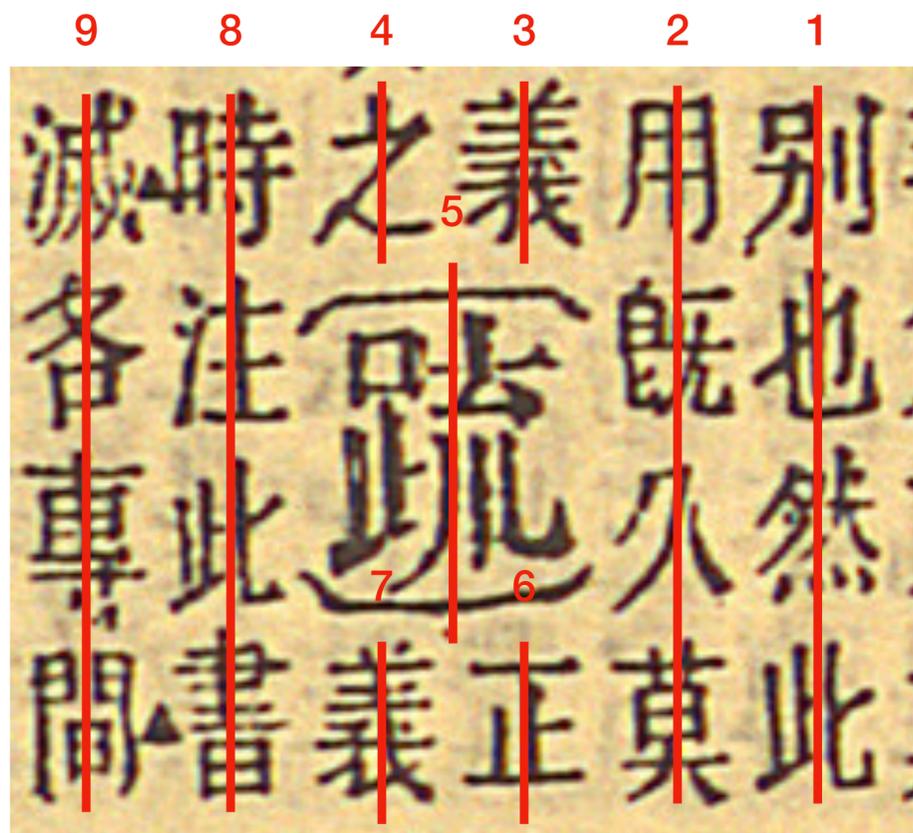
The screenshot displays the Calfa Vision web interface. On the left is a sidebar menu with options like Dashboard, Users, Projects, Settings, Guide, and Terms of Service. The main area shows a document with Chinese text, with a search bar at the top and navigation controls. Below the document is a table of detected lines, with columns for line number, navigation, and text. The table contains 13 rows of data, each representing a line of text from the document.

Line number	Navigation	Text
19	↑ ↓ U	之異而直言者非詩故更序詩必長歌之意情謂哀樂之情
20	↑ ↓ U	中謂中心言哀樂之情動於心志之中出口而形見於言初
21	↑ ↓ U	言之時直平言之耳平言之而意不足嫌其言未申志故咨
22	↑ ↓ U	嗟歎息以和續之嗟歎之猶嫌不足故長引聲而歌之長歌
23	↑ ↓ U	之猶嫌不足忽然不知手之舞之足之蹈之言身為心使不
24	↑ ↓ U	自覺知舉手而舞身動足而蹈地如是而後得舒心腹之憤
25	↑ ↓ U	故為詩必長歌也聖王以人情之如是故用詩於樂使人歌
26	↑ ↓ U	詠其聲象其吟詠之辭也舞動其容象其舞蹈之形也具象
27	↑ ↓ U	哀樂之形然後得盡其心術焉情動於中還是在心為志而
28	↑ ↓ U	形於言還是發言為詩上辨詩從志出此言為詩必歌故重
29	↑ ↓ U	其文也定本言之不足故嗟歎之俗本言之下有者字誤也
30	↑ ↓ U	定本永歌之不足下無故字有故字者亦誤也樂記云歌之
31	↑ ↓ U	為言也長言之也說之故言之言之不足故長言之長言之
32	↑ ↓ U	不足故嗟歎之嗟歎之不足故不知手之舞之足之蹈之其

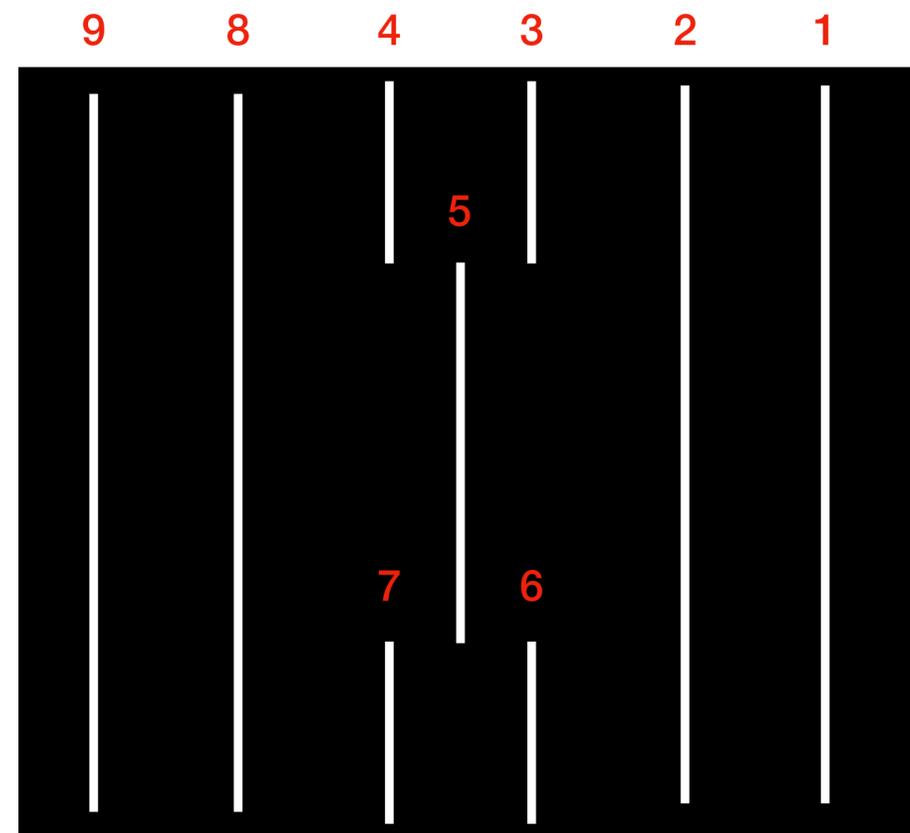
# Les premières expérimentations - l'ordre de lecture

Approche de base : tri des lignes avec le centroïde, solution de traitement d'image basique

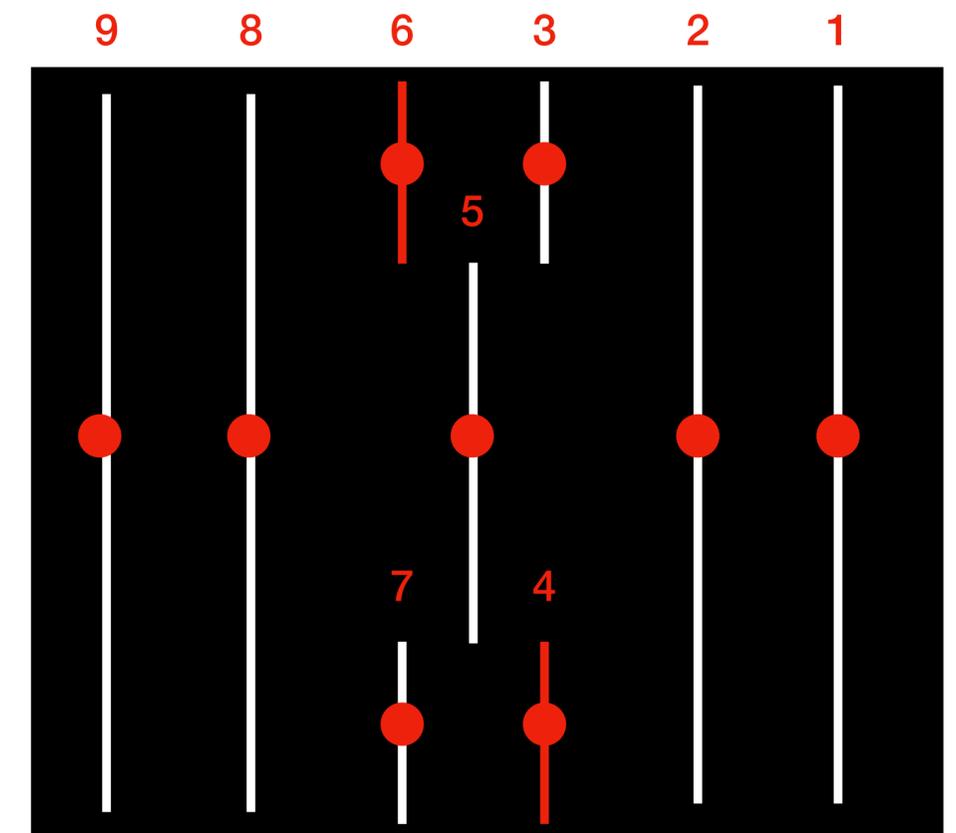
Efficace à **82,37%** sur le dataset, échoue y compris en situation simple



Détection des lignes



Sens de lecture GT

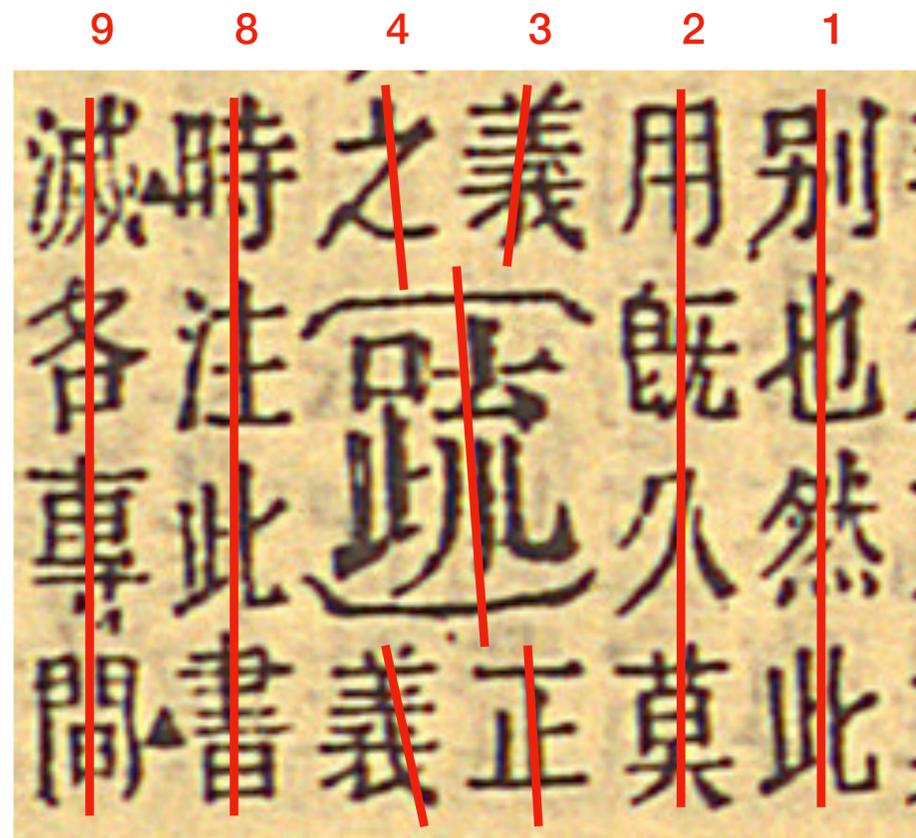


Tri des lignes par centroïdes

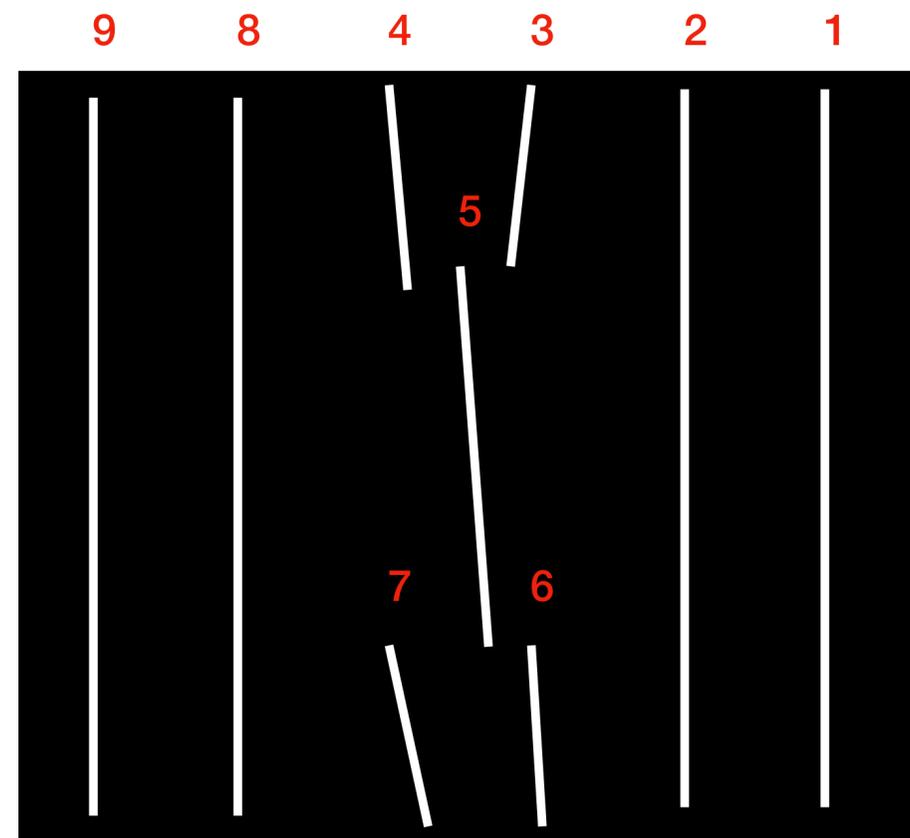
# Les premières expérimentations - l'ordre de lecture

Approche de base : tri des lignes avec le centroïde, solution de traitement d'image basique

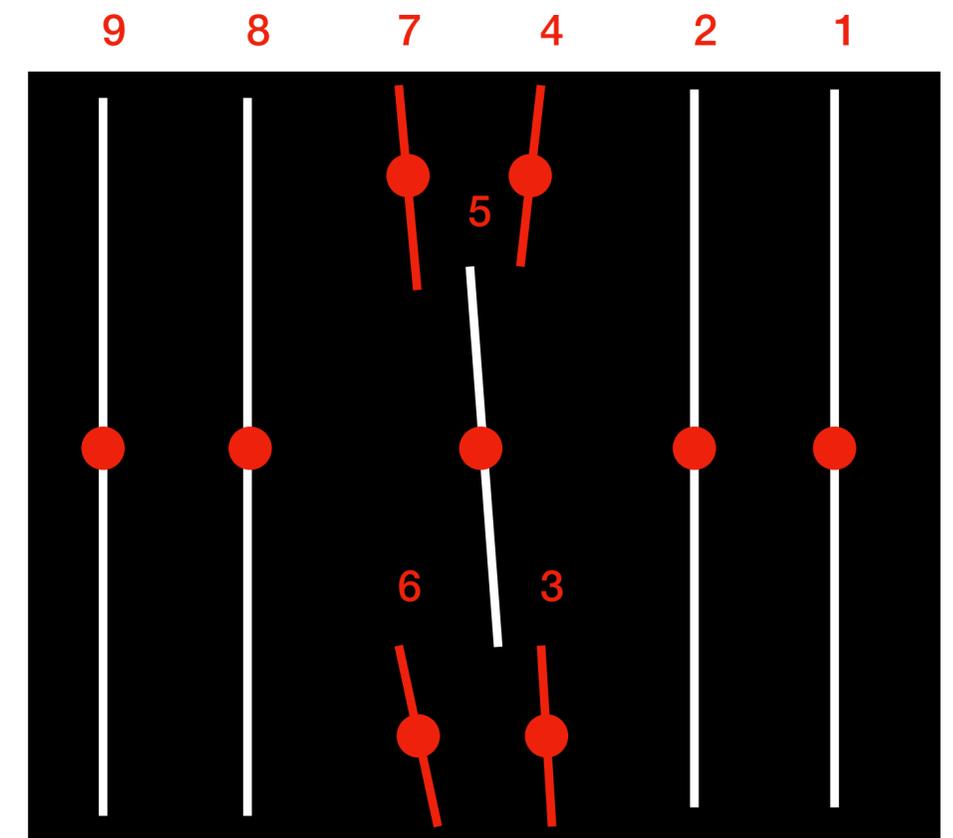
Efficace à **82,37%** sur le dataset



Détection des lignes



Sens de lecture GT



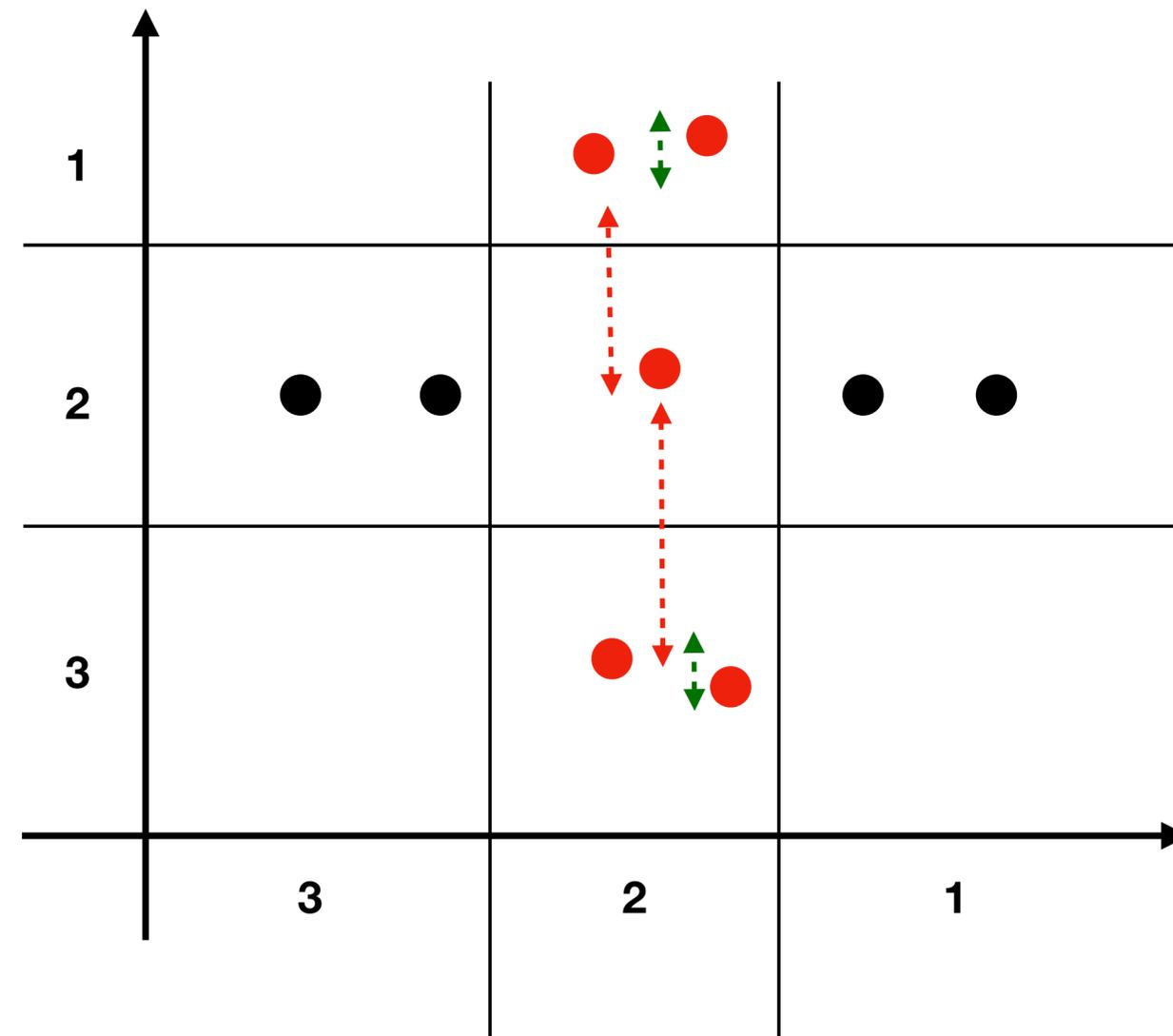
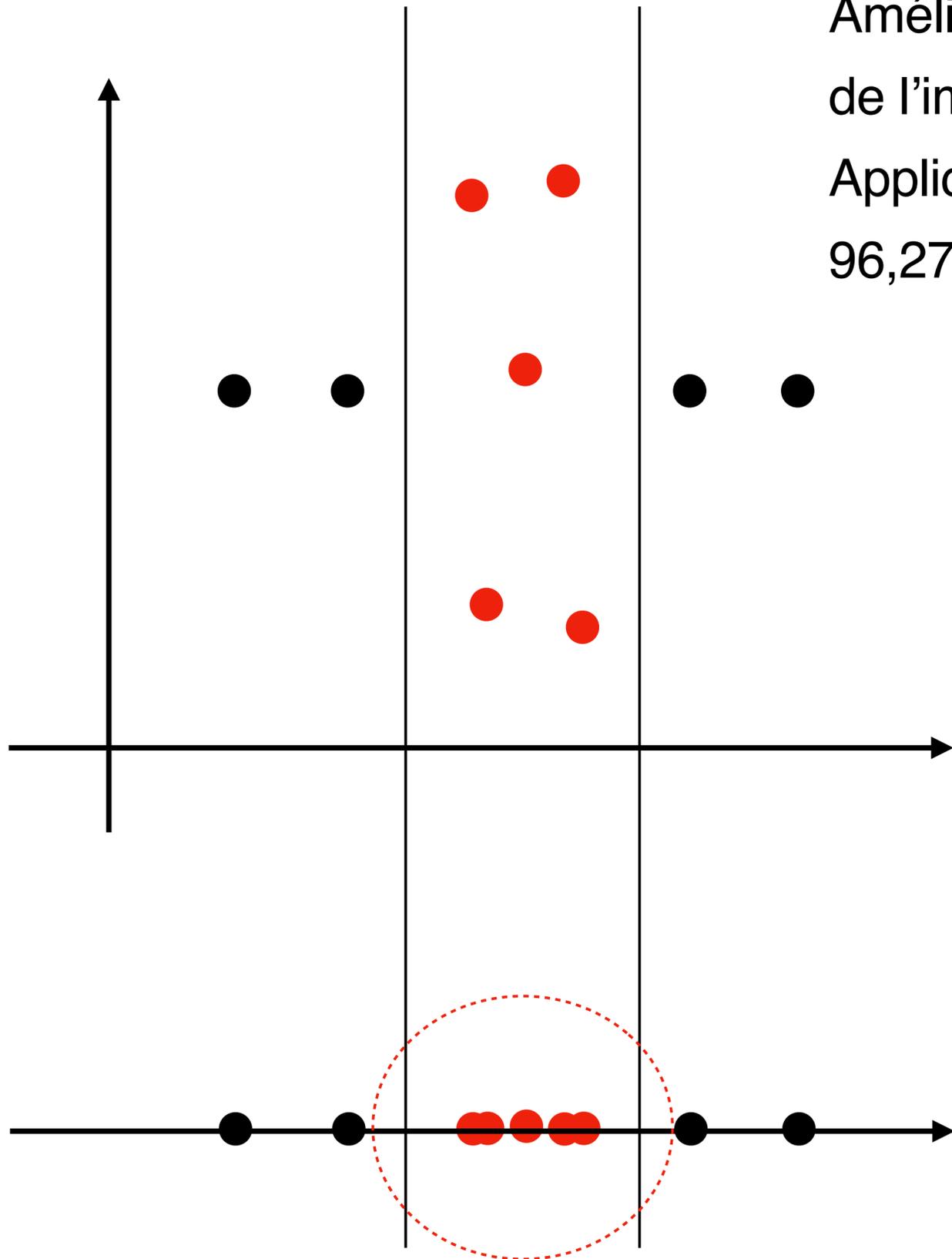
Tri des lignes par centroïdes

# Les premières expérimentations - l'ordre de lecture

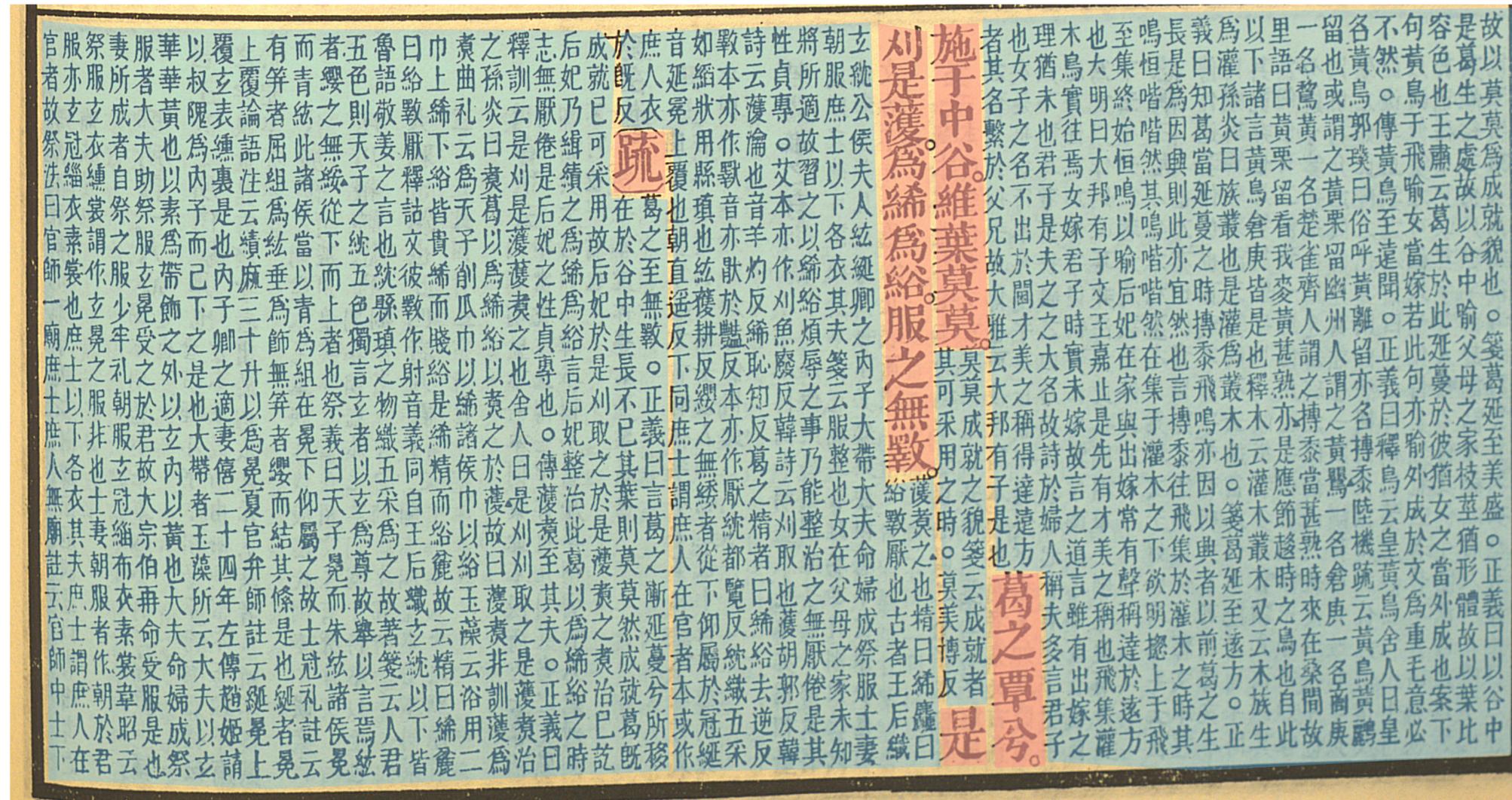
Amélioration du tri avec une solution de traitement

de l'image : projection des centroïdes et identification d'un cluster : 87,88%

Appliquée directement aux bounding-box vs baseline : de 91,31% à 96,27%



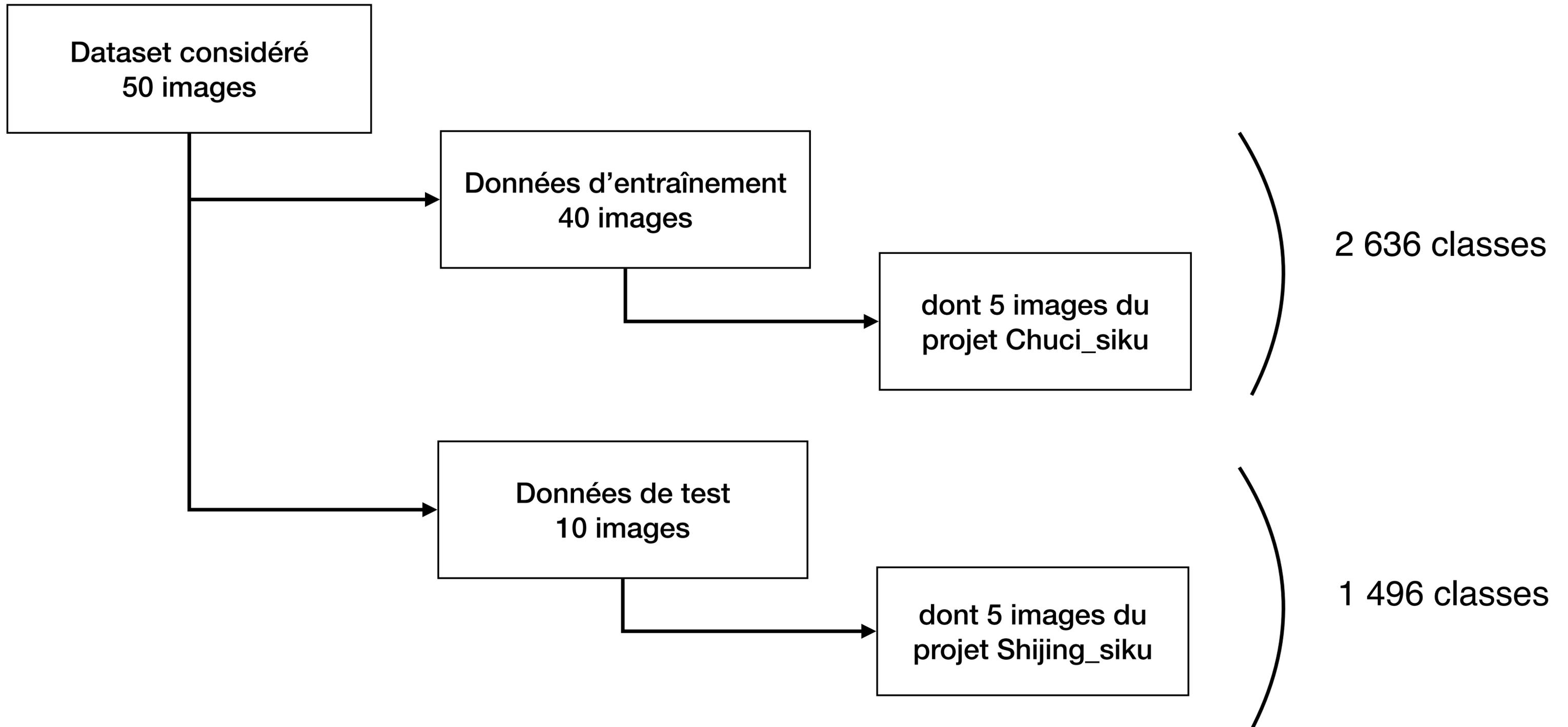
# Les premières expérimentations - l'ordre de lecture



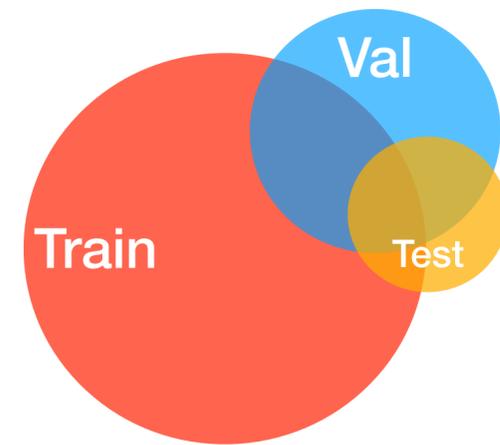
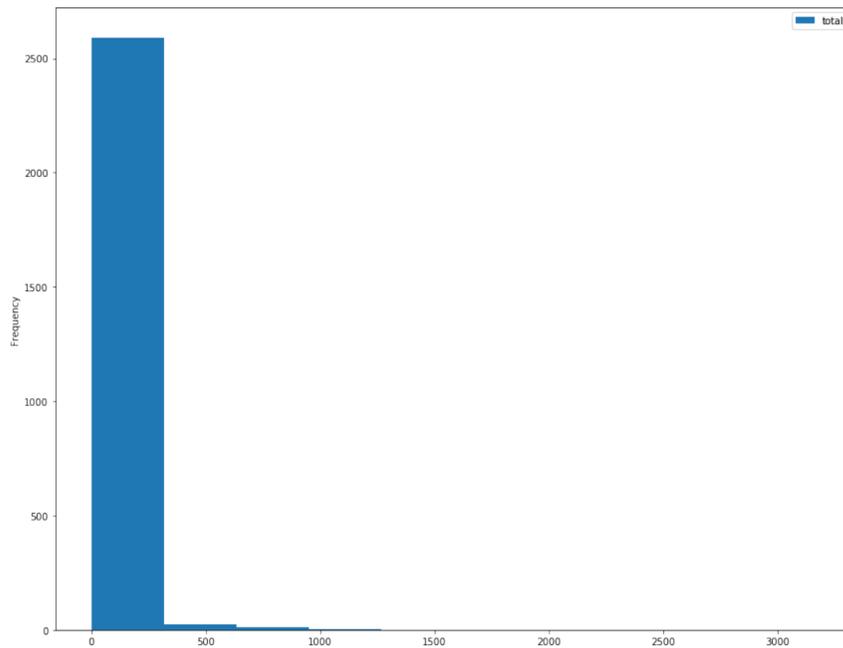
Solution en cours d'expérimentation : pré-classification sémantique au sein d'une région de texte puis tri indépendamment des lignes pour chaque zone identifiée. Abandon des baselines. Score obtenu sur un échantillon manuellement annoté, sans apprentissage : 99,67%

# Les premières expérimentations - la question des classes

Création d'un modèle de zéro, pas de données disponibles pour un *fine-tuning*



# Les premières expérimentations - la question des classes



## caractères les plus fréquents / les moins fréquents

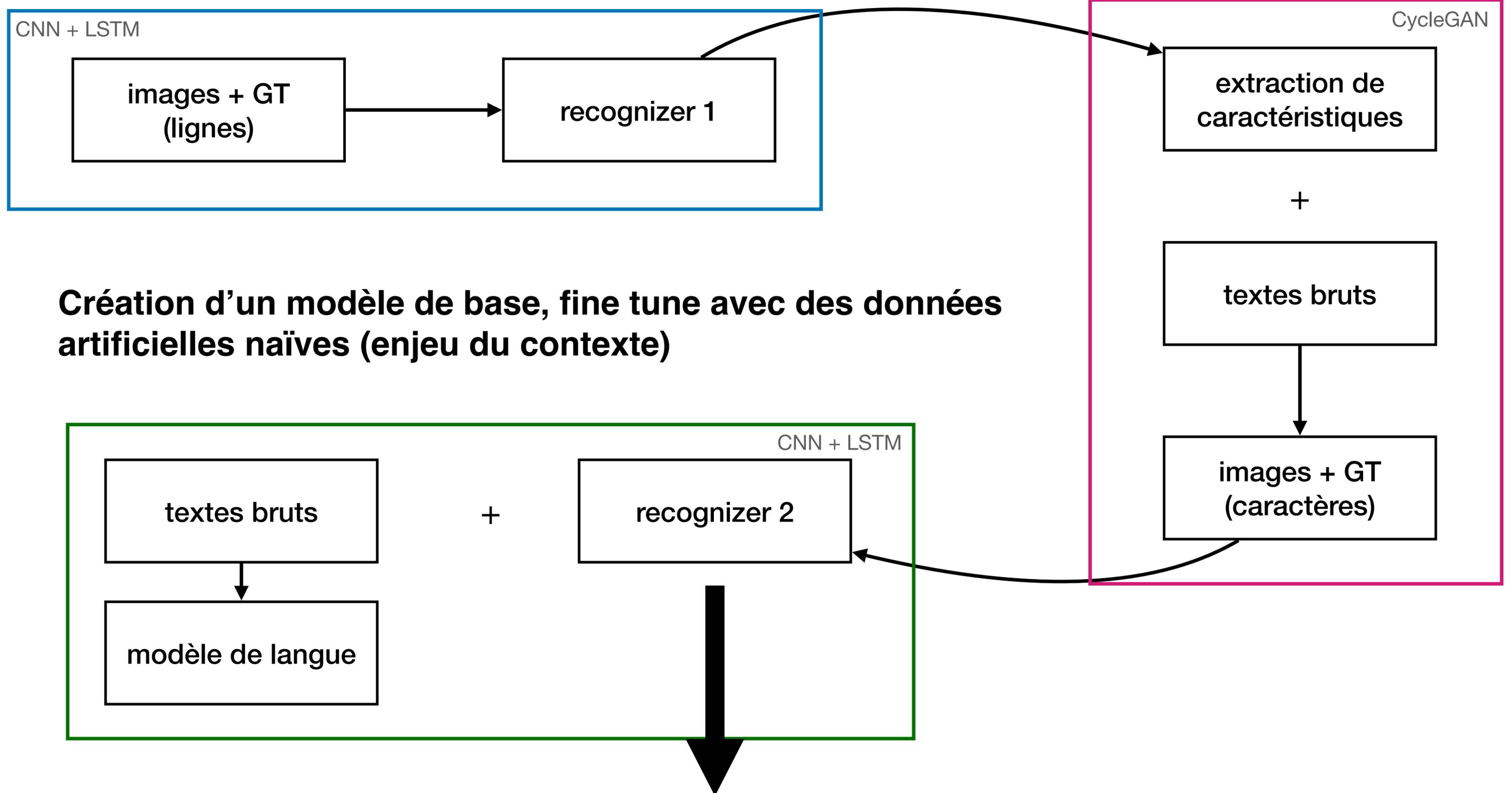
之	3170	芋	1
也	2032	恕	1
以	1463	刃	1
云	1110	訪	1
不	1074	隈	1
言	960	沔	1
者	944	斐	1
	...		

**2 095 classes sur 2636 avec un seul échantillon**

Seulement 1269 classes communes en moyenne entre les ensembles d'apprentissage et de test

**==> Situation de *one-shot learning* + classes non vues en apprentissage**

# Les premières expérimentations - la question des classes



# Les premières expérimentations - la question des classes

疏 國史至正。正義曰：上既言變詩之作，此又說作變之  
 由言國之史官皆博聞強識之士，明曉於人君得失善  
 惡之迹，禮義廢則人衛亂政，教失則法令酷，國史傷此人，儉  
 一之廢棄，哀此別政之青唐，哀傷之志，蘋積於內，乃除詠己二  
 情性以風刺其上，觀其改惡為善，所以作變詩也。國史者周  
 官大史、小史、外史、衛史之等皆是也。此承變風變雅之下，則  
 兼據天子諸侯之史矣。得失之述者，人君既往之所行也。明  
 時得失之述，哀傷而詠備性者，詩人也。非史官也。民勞常並  
 二公卿之作也。黃鳥碩人，國人之風，然則凡是臣民皆得風刺  
 二不必要其國史所為此文，特言國史者，鄭荅張逸云：國史采  
 二眾詩，時明其好惡，令誓媵歌之，其無作主皆國史主之，令可  
 二歌如此言，是由國史掌書，故託文史也。苟能制作文章，亦可  
 二謂之為史，不必要作史官。明云：史克作是頌，史官自有作蓋  
 一者矣。不盡是史官為之也。言明其好惡，令譬媵歌之，是國史  
 二還取善者始竹樂官也。言其無作主國史主之嫌，其作者無

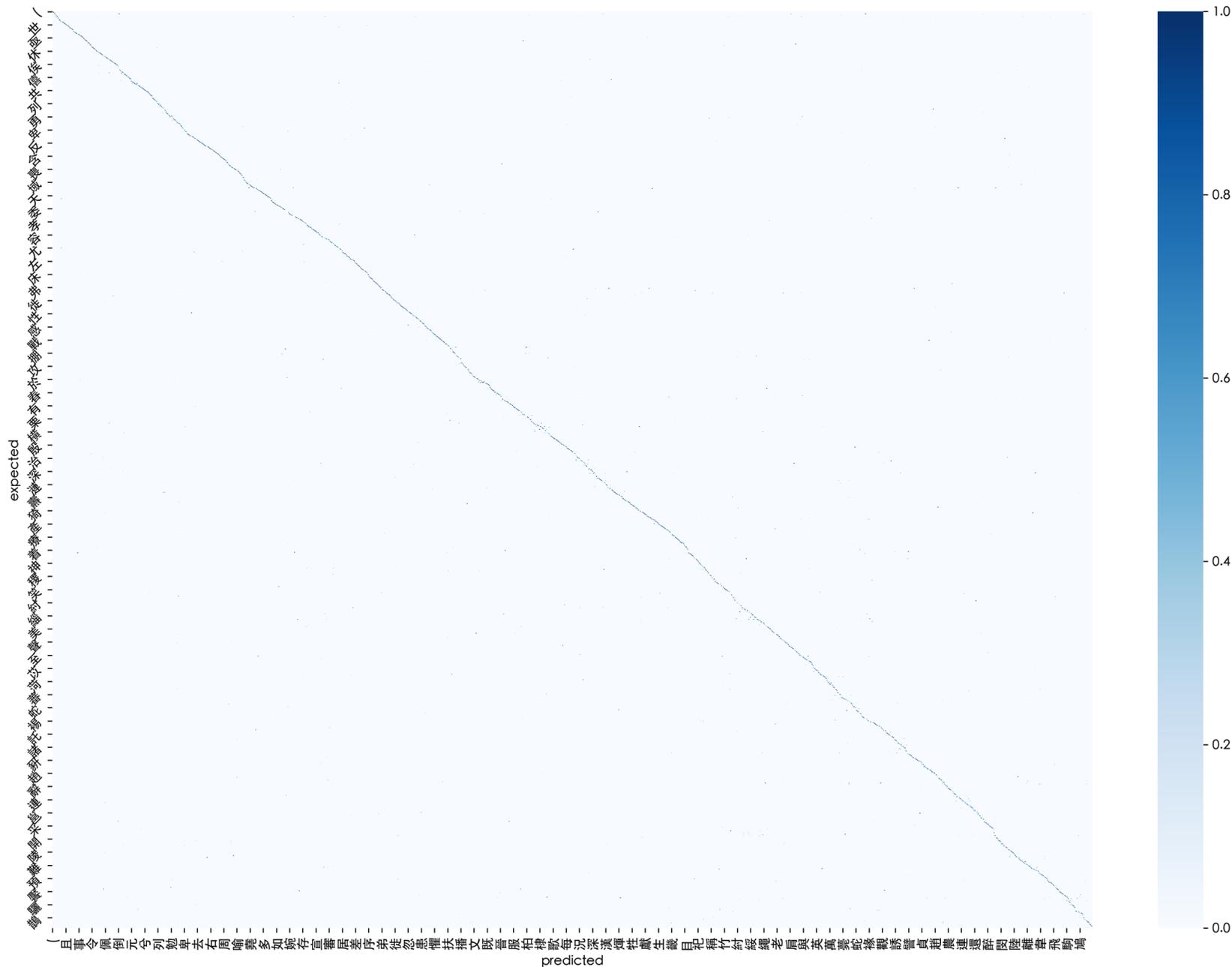
二 國史至\*正義曰上既言變詩之作此又說作變之  
 一 由言國之史官皆博聞強識之士明曉於人君得失善  
 二 惡之述禮義廢則人衛亂政教失則法令酷國史傷此人儉  
 二 一之廢棄哀此別政之青唐哀傷之志蘋積於內乃除詠己二  
 二 情性以風刺其上觀其改惡為善所以作變詩也國史者周  
 官大史小史外史衛史之等皆是也此承變風變雅之下則  
 二 兼據天子諸侯之史矣得失之述者人君既往之所行也明  
 時得失之述哀傷而詠備性者詩人也非史官也民勞常並  
 二 公卿之作也黃鳥碩人國人之風然則凡是臣民皆得風刺  
 二 不必要其國史所為此文特言國史者鄭荅張逸云國史采  
 二 眾詩時明其好惡令誓媵歌之其無作主皆國史主之令可  
 二 歌如此言是由國史掌書故託文史也苟能制作文章亦可  
 二 謂之為史不必要作史官明云史克作是頌史官自有作蓋  
 一 者矣不盡是史官為之也言明其好惡令譬媵歌之是國史  
 二 還取善者始竹樂官也言其無作主國史主之嫌其作者無

竹 : glyphe erroné

明 : glyphe vu une seule fois en apprentissage

制 : glyphe non vu en apprentissage

# Les premières expérimentations - la question des classes



Expérimentation	Accuracy
Sur les glyphes inconnus	<b>86,21 %</b>
Sur les glyphes connus (y compris <i>one shot learning</i> )	<b>94,49 %</b>
Données brutes	<b>91,5 %</b>
Données pondérées	<b>93,09 %</b>

Posent problème :

- les débuts de colonnes
- la ponctuation excentrée des colonnes
- certains caractères :

爲 / 為

荼 / 菜

苻 / 持

已 / 己

千 / 干

旋 / 施

## Considérations pratiques :

- ressources humaines (et financières) limitées
- densité du corpus, fatigabilité de l'œil et de l'esprit des annotateurs
- hétérogénéité des transcriptions

## Considérations philologiques :

- corpus constitué d'éditions de périodes différentes (caractères tabous différents)
- plusieurs mains, des simplifications qui ne sont pas uniformes au fil des pages

## Considérations liées à l'étape 1 du projet :

- l'objectif n'est pas d'étudier la réception au XVIIIe s. d'un texte antique (les « erreurs » seraient alors signifiantes et précieuses)
- pas d'approche diplomatique (édition choisie = la plus commune)
- structuration en vue d'une fouille de texte pour chercher des phénomènes d'emprunts d'un genre textuel à un autre

## Exemples :

爲 為

既 旣

即 卽

總 摠

**Layout** : jusqu'à 96% de bon tri des lignes mais résultats peu robustes - impact négatif de la « baseline »

**HTR** : 93,09 d'accuracy

- Consolider ces premiers essais dans le cadre du COLLEX-PERSÉE CHI-KNOW-PO-CORPUS
  - Corpus diversifié (mise en page, lexique)
  - Nouveaux développements en perspective pour améliorer les performances concernant l'ordre des colonnes
  - Enrichir la reconnaissance des types de textes (texte *versus* commentaires, titres, *etc.*)
- Partager ces modèles et données d'entraînement pour qu'ils servent aux chercheurs travaillant sur des textes anciens de l'Asie extrême-orientale (Chine, Corée, Japon, Vietnam)

## Outil utilisé :

- [vision.calfa.fr](http://vision.calfa.fr)

## Publication des données et modèles à venir :

- Nakala
- Github (et référencement dans HTR United)

## Remerciements :

- COLLEX-PERSÉE
- DISTAM (Digital Studies: Africa, Asia and the Middle East)
- USIAS