

(Département d'études arabes)))



Le projet HUNAI (Humanités Numériques pour l'histoire de l'Arabie Islamique)

Eric Vallet et Clarck Junior Membourou Moimecheme

L'HTR des langues peu dotées dans les programmes de recherche et dans les établissements de conservation français – 14 février 2024



arabes

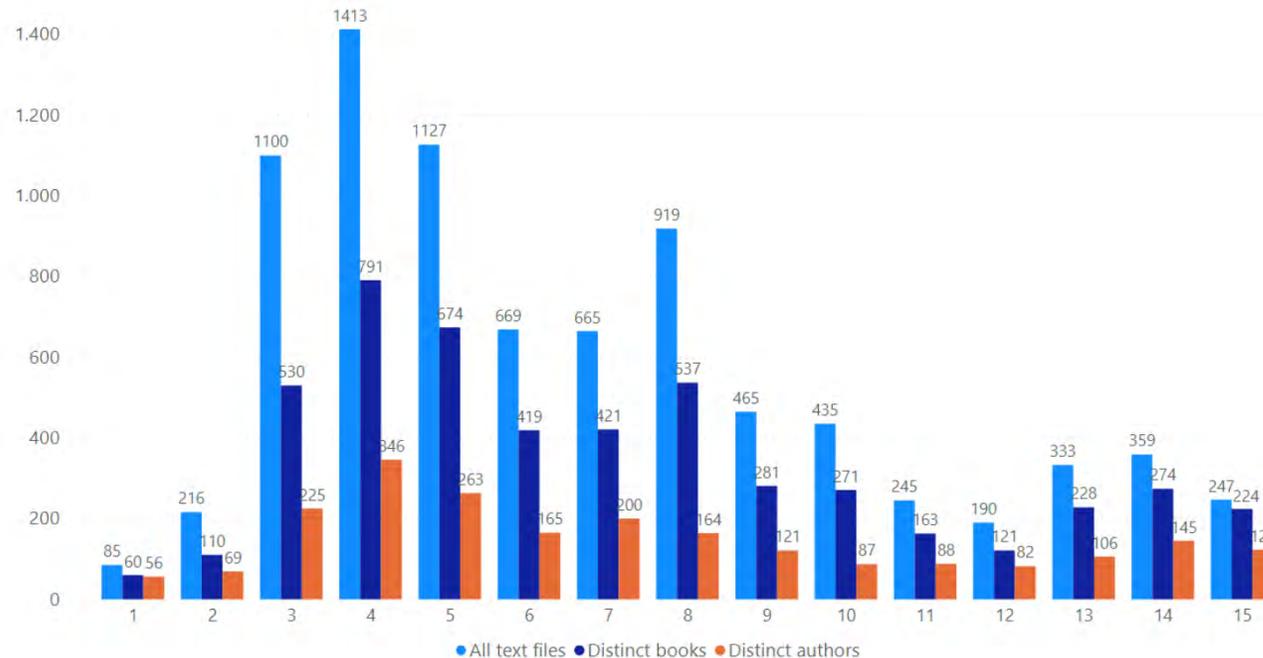
- 6 millions de manuscrits islamiques dans le monde ?

7280 manuscrits arabes à la BnF

- La “Constellation” OpenITI

KITAB Project > Open Islamicate Text Initiative OpenITI (Maxim Romanov, Sarah Savant and Matthew Miller)

10 202 fichiers textes ; 6 236 “oeuvres” de 2 582 auteurs (dernière livraison sur Zenodo).



Text files in OpenITI, per hijri century (light blue: number of text files; dark blue: number of distinct books; red: number of distinct authors). NB: pre-Islamic texts are in the first-century column, texts by authors still alive in the 15th-C column. (data: Dec. 2020)

Source :

<https://openiti.org/documentation/#faq>

arabes

Benjamin Kiessling, *Avancées en Reconnaissance Optique des Caractères pour les Documents Arabes Historiques*, thèse EPHE 2021.

Open Islamicate Texts Initiative Arabic-script OCR Catalyst Project (OpenITI AOCP) 1 et 2 (2019-2021)
escriptorium

Automatic Collation for Diversifying Corpora (ACDC) (Matthew Miller) 2024- – Financement National Endowment for the Humanities 2024 = Automatisation de l'entraînement de modèles HTR sur manuscrits arabes et persans

- Combinaison de Kraken, de l'outil Passim (text reuse) et des textes numérisés dans Open ITI (alignement automatique des textes issus de l'HTR à partir des textes déjà numérisés dans Open ITI)

<https://github.com/OpenITI/ACDC>

<https://www.youtube.com/watch?v=kNx4GyH5HSo>



Littératures populaires arabes

Hackaton Alexandre (2023)

<https://gitlab.huma-num.fr/lipa/iskandar>

2. HUNAI : genèse et objectifs

➤ Genèse

- Une vision étriquée de l'Arabie, terre sans histoire
- Un patrimoine manuscrit partiellement sauvegardé et valorisé par les institutions locales ou internationales
- Une production textuelle peu considérée, sous-exploitée et sous représentée

➤ Objectifs

- Réintégrer un corpus cohérent d'ouvrages majeurs **dans l'ensemble des ressources numériques arabes disponibles en *open access*** et annotées selon les standards internationaux
- Développer des **études et des analyses historiques et linguistiques reposant sur la fouille de données textuelles** appliquée à des corpus cohérents issus de cette production.



[Home](#) / [Yemeni Manuscripts Digitization Initiative](#)

Yemeni Manuscripts Digitization Initiative

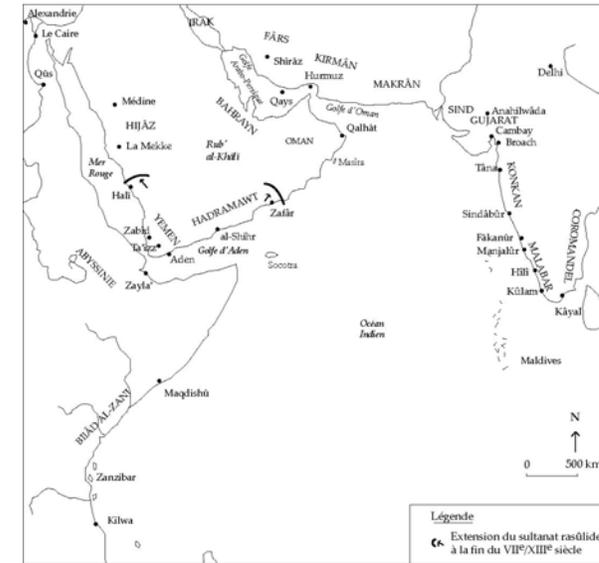
The Yemeni Manuscript Digitization Initiative (YMDI) is a collective of leading scholars of classical Islam, Middle Eastern history, and Arabic Literature from North America, Europe, and the Middle East whose mission is to preserve the Arabic manuscripts in the private libraries of Yemen.

Led by Director Dr. David Hollenberg (University of Oregon), in 2010, the initiative received a \$330,000 NEH/DFG Enriching Digital Collections grant shared between Princeton University Library, and the Freie Universität, Berlin to digitize 236 Arabic manuscripts in the fields of Islamic theology and law.

YMDI's mission is the preservation and dissemination of the Arabic manuscripts in the private libraries of Yemen. Working closely with a Yemeni non-profit organization which has endeavored to save Yemeni manuscripts for the past decade, YMDI digitally preserved three private libraries in the capital city of Sana'a, a total of 236 manuscripts. These digitized sources have been virtually conjoined to twelve manuscripts in the rare book collections of the Staatsbibliothek zu Berlin and the Princeton University Library, and made freely accessible in Princeton's Online Catalog and in DPUL's Digital Library of Manuscripts of the Islamic World.

Corpus

- Composé de données textuelles correspondant aux ouvrages historiographiques sunnites composés au Yémen à l'époque rasūlide (626-858/1229-1454)
- 6 manuscrits numérisés de la Bnf, publiés sur Gallica
 - al-Janādī, *Sulūk fī ṭabaqāt al-'ulamā' wa-l-mulūk*, BnF Arabe 2127 (copié en 820/1417)
 - al-Ashraf Ismā'īl/al-Khazraġī, *Al-Kifāya wa al-i'lām fī man waliya al-Yaman wa sakanahā fī al-islām*, Paris BnF 5831
- 9 mss (Yale ; Institut des manuscrits arabes du Caire)



3. L'annotation : genèse, méthodologie et adaptations

3.1. Genèse de l'annotation



1	 	لسعلمف الله الرحمن الرحيم ف و الحمد لله الملك العظم الاول الاحر العدم مما عب نسه	1 modèle arabe générique plus récent	 	 
2	 	محمدها عرلا الحليل والتحرير وداعنا الى الصراط المستقيم قبلع صلى الله عليه وسلم الرساله على النعمهم وهم لمن تعه		 	 
3	 	الدين القورهم مو عن على الكافه بصديقه و العام يواحب الصلوة لئه والتسليم فله الحمد حمد الحرج عن درك الاخصاه		 	 
4	 	ويدغم به طاييف من حمد وعسره وصلى الله على سة المنعوث صلوه استاسرها الحن ويحشه القنور واتميربها في الدارس		 	 
1	 	سمليقم اننه الحمن الركم و الحمديته المك العظم الاول الاحر العدمك ما عب نسه	2 modèle arabe générique ancien	 	 
2	 	محمد مههامه با الحليل والنحرير وداعنا الى الصراط المستقيم قبلع صلى الله عليه وسلم الرساله على النعمهم ونهم ان تعه		 	 
3	 	الدين القوريم فوع على الكافه يصديقه والقام يواجب الصلوة طيه والتسليم ه فله الحمد حمد الخرج عن دركالاحصاه		 	 
4	 	ويرغم به طرايف من حمد وعسي وصلى الله على مسة المبعون صلوا استاس هالحين او يحته القبور كواتمرتها في الدارس		 	 
1	 	لسمايفه الخمر الرجي والهية الك الحطم الولا الجر اليهع ما عب سه	3 modèle Hamawi Strasbourg	 	 
2	 	محبها لمرلا الحليل والنحرم ودا عبا الى الطراا المستعيم معبلع ملئ ابه عليه وسلم الرباله علي النغميم موتهم ل تعه		 	 
3	 	الدين العوه حمو ملن الكافه بضويقه والعيام يواج الطلواطيه والسشلم فله احمد ا بخرج عن در ا احصاه		 	 
4	 	ويرغميه انف من حد وعسي وطئ اه لي سة المبعوث صلوه استاس ماحتن ويحشه القور توا تميربا في الرارر		 	 

Annotations du f° 62 r. du BnF arabe 5832 dans un fichier .txt



فرحفت العساكر المنصورة على الحصن ثلاثة ايام متوالية وكتب الأمير علي ابن عبد الله الى كافة الأشراف كتبنا
 متتابعه يطلب منهم النصرة وهم يغالطونه ويعتذرون ثم حصل خطاب ومراسلات علي يعني الصلح
 واستقر الحال على أن الأمير علي بن عبد الله يواجه الصاحب موفق الدين فوصل اليه واتفق حضور الملك
 المظفر فاجتمعوا جميعا وساروا بأجمعهم إلى المقام الشريف السلطاني فلما علم السلطان بوصولهم ركب
 من مخيمه وقد صاروا بالقرب منه وأكرمه وانصفه وانفقد الصلح بينهم وأخذ للأشراف دمه سبعة اشهر
 وسلم لأجلها حصن ديفان لان السلطان امتنع من الدمه عليهم فلما استقر في المحطة طلب من السلطان دخول
 الإعلام الشريفه اظهرا للطاعة والتسليم فنصبت في أعلا الحصن وكذلك العظيمة فحقت ذواتها في أعالي
 الحصنين ولقد أحسن الحسن بن هاني يقول من كان بالسمر العوالي خاطبا ... حليت له بيض الحصون عرائسا
 وقال العفيف عبد الله ابن جعفر بمدح السلطان الملك المؤيد ويذكر أخذه العظيمة واللميعاف

إرت الخلافة في يدك مشاع ... وعرار سيفك شاهد قطاع
 منع النصيب من العد انصب القان ... والجرم القراع من اليوف قراع
 شمس رأت غلب الملوك شعاعها ... فقلوبها منها تطير شعاع
 تبع التتابع في عناصر حمير ... والي المناقب هم له أتباع
 عمرو وعمرو ذو الجناح ومنذر ... والأيمان وفايش وكلاع
 ماء السماء سقى منابت اصله ... ربا فأورق عرفه النزاع
 فلقد أعاض بيوسف يقطان فلا ... نكل ولا كل ولا مجراع
 أشرا إلى الشرق القصي يشرب ... خطواتها نحو المراع سراع
 والشمس من لمع الحديد كليله ... والجو من سمر البراع براع
 وفيالق سالت هوداي خيلها ... سبل الآتي تداولته تلاع
 تسري فمن زرق الأسنة فوقها ... نار ومن أسل الوشيخ شعاع
 غسلت مياه سيوفها ماء الدجي ... فتشابه الإصباح والأهراع
 ينحو بها مبدأ النجوم طوالعا ... ملك مطيع للإله مطاع
 ليس العظيمة بالعظيمة عند من ... بسيوفه ميفاعها ميفاع
 لم ينشق وافدهم إليه اتي وهل ... يشقى أمره وجليسه القعقاع
 فغمغت أدعية بأفواه لهم ... فيهن من ثدي البتول رضاع
 وحفظت حقا للنبي محمد ... فيهم ولست بما حفظت تضاع
 أيام الإسلام سيف وضع ... الوسر مجك من الشاسطاع
 أمؤيد الإسلام داود الذي ... للعالمين بفضله إجماع
 ما يلتقي شرق البلاد وغربها ... إلا إذا ما امتد منك الباع
 أهويت بالسيف العداة كما هوى ... ود بسيف محمد وسواع
 الله أعطاك السعادة كلها ... ماذا بضر وربك النعاع

وهي أطول مما ذكرت وهذه عيونها و أقبل السلطان على الأمير جمال الدين علي بن عبد الله بالبحر
 وأزال ما في خاطره وحدد له رفق الطليخانة وحمل معها من الكساوي والاموال شيئا كثيرا ولما
 كان أول يوم من شهر ربيع الأول ارتفع السلطان من محطة قاضي صنعاء
 أمام الكتيبة تزهى به ... مكان السنان من العامل

قال الشريف إدريس وسرت في خدمته مع والدي الي الوعدت من هناك وقد كنت خرجت
 اليه في محطة الميقات فقامت في حاك ...



البرقي محمدم

Annotations du f° 62 du BnF arabe 5832 importées dans le logiciel Calfa

The screenshot displays the Calfa software interface for managing manuscript annotations. The interface is divided into several sections:

- Search:** A search bar at the top left with a magnifying glass icon.
- Navigation:** Buttons for "Prev Image", "Next Image", "Full screen", and "Image info" at the top right.
- Project Path:** A breadcrumb trail showing "Projects / Hunai - BnF5832 / Images / ARA-BnF-5832_btv1b10030898x_62 / Labellize".
- Process Steps:** Three numbered steps: "1. Layout Analysis", "2. Generate Polygons", and "3. Text Recognition".
- Text Region List:** A table listing 9 text regions with their corresponding image thumbnails and control icons (lock, delete, settings).
- Text Region Info:** A panel on the right that provides details for the selected text region.
- Manuscript Preview:** A large image of the manuscript page (f° 62) with red and blue annotations overlaid on the text.

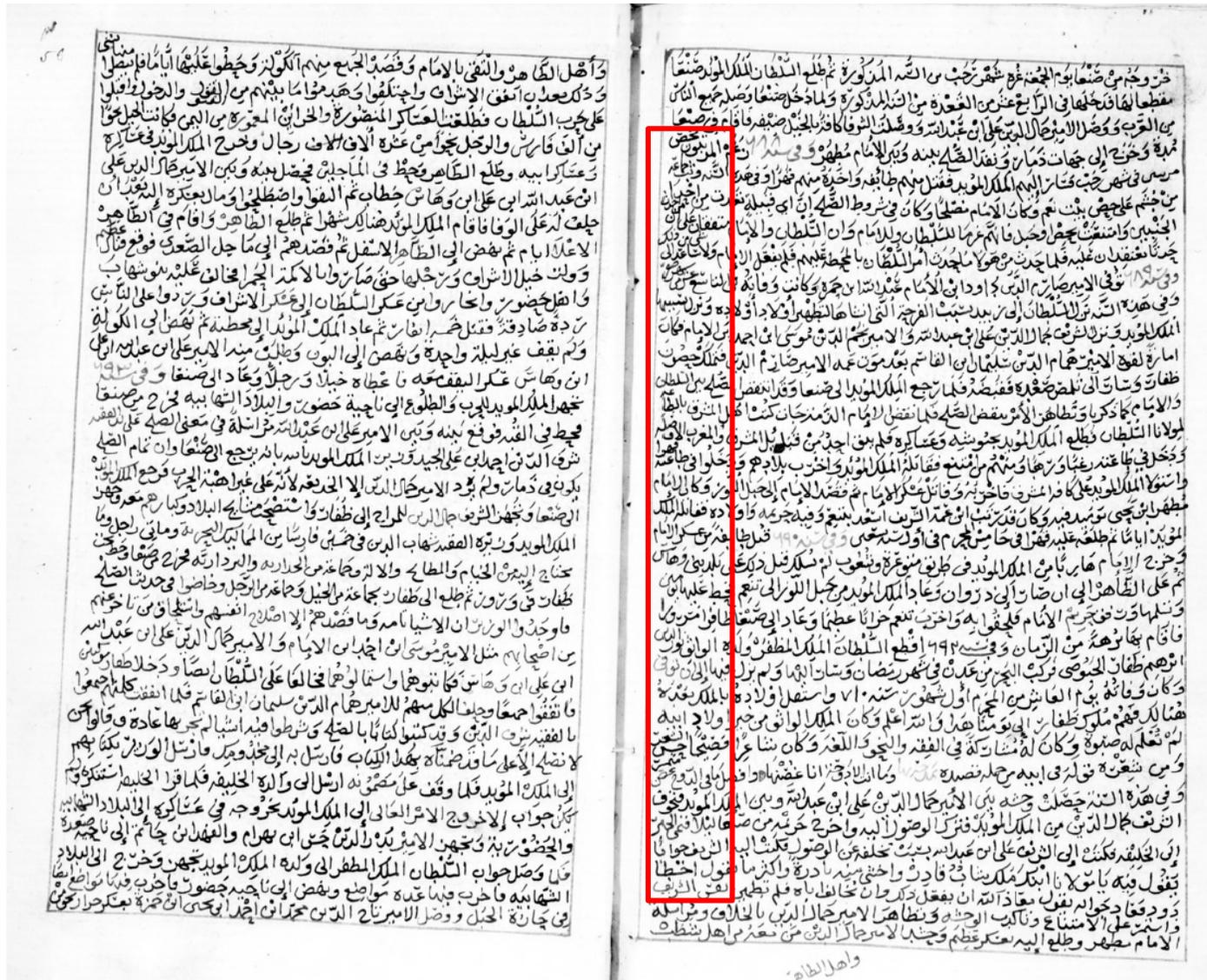
The text region list contains the following entries:

Region ID	Text	Image	Lock	Delete	Settings
1	فزحفت العساكر المنصورة على الحصن ثلاثة أيام				
2	متابعة يطلب منهم النصره وهم يغالطونه				
3	واستقر الحال على أن الأمير علي بن عبد الله يواجه				
4	المظفر فاجتمعوا جميعا وساروا بأجمعهم إلى				
5	من مخيمه وقد صاروا بالقرب منه وأكرمه وانصفه				
6	وسلم لأجلها حصن ذيفان لان السلطان امتنع من				
7	الإعلام الشريفه اظهرها للطاعة والتسليم فنصبت				
8	الحصنين ولقد أحسن الحسن بن هاني يقول من				
9	و قال العفيف عبد الله ابن جعفر يمدح السلطان				

3.3. Adaptations

Réclames, notes marginales et entre les lignes

L'ajout régulier des termes entre les lignes rend la lecture et l'annotation difficiles

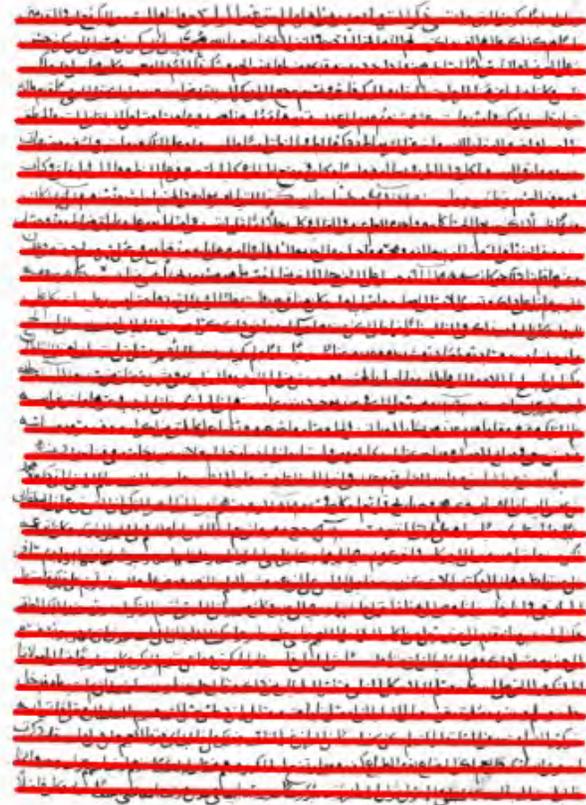


وأهل الظاهر

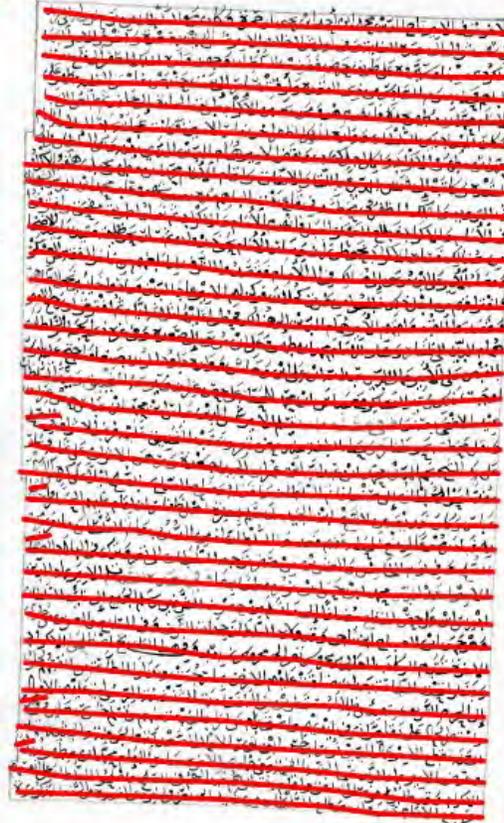
Les lignes de base appliquées sur les mss BnF arabe 2127 et 5832



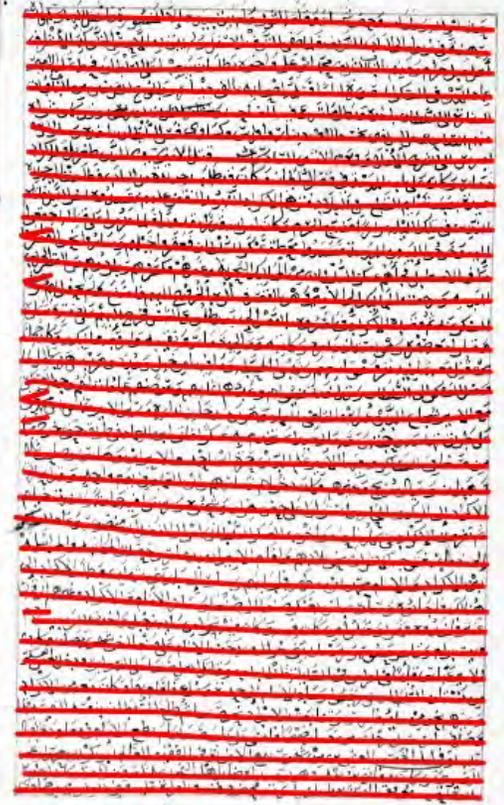
Page 1 of a manuscript showing Arabic text with red horizontal lines applied as a baseline. The text is written in a cursive script, and the lines are evenly spaced across the page.



Page 2 of a manuscript showing Arabic text with red horizontal lines applied as a baseline. The text is written in a cursive script, and the lines are evenly spaced across the page.



Page 3 of a manuscript showing Arabic text with red horizontal lines applied as a baseline. The text is written in a cursive script, and the lines are evenly spaced across the page.



Page 4 of a manuscript showing Arabic text with red horizontal lines applied as a baseline. The text is written in a cursive script, and the lines are evenly spaced across the page.

Les chiffres et les lettres

Utilisation irrégulière des points souscrits et suscrits obligatoires, ce qui génère des confusions entre les consonnes ب ت ث ج ح خ ص ض ي

Utilisation de ا à la place ي

Absence de la hamza ء en fin de mot

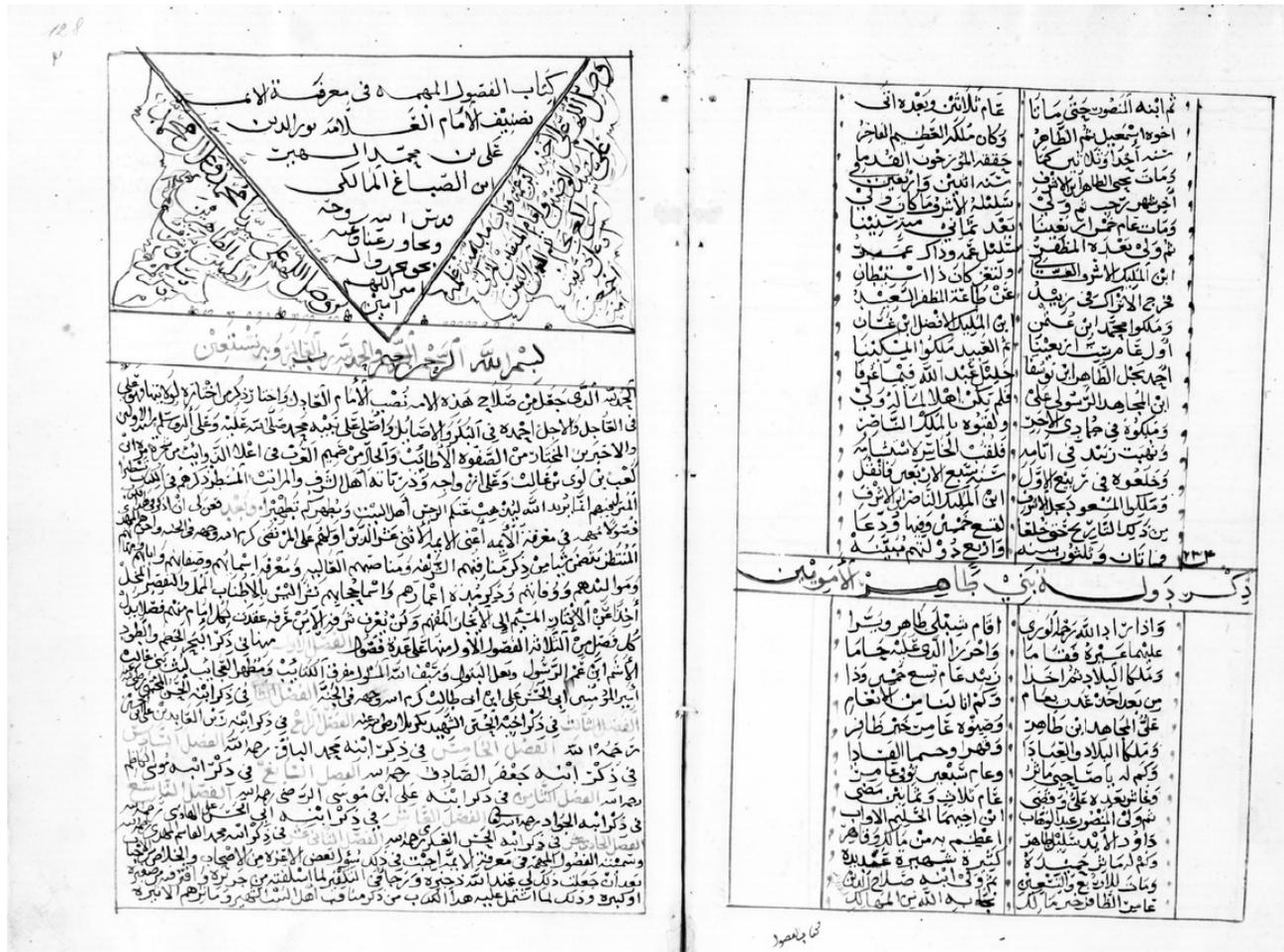
Les mots tels qu'ils sont écrits aujourd'hui	Les mots dans les manuscrits d'Hunai
علماء	علما
صنعاء	صنعا
شيء	شي

4. Un modèle adapté aux manuscrits du Yémen ?

Les annotations permettront à l'équipe de Calfa d'entraîner un modèle qui sera testé sur l'ensemble des manuscrits du projet HUNAI

Vérifications seront faites afin de déterminer le taux d'erreur, dans un premier temps sur les mss BnF arabe 2127, 5832 et peut-être 4609.

5. Des mises en page complexes à examiner avec Calfa



Comment organiser la sortie de texte ?

- Titres, sous-titres ...
- La poésie ...
- les schémas ...

