

RASAM I ET II – A DATASET FOR THE
RECOGNITION AND ANALYSIS OF SCRIPTS
IN ARABIC MAGHREBI

RETOUR SUR DEUX EXPÉRIENCES DE
TRANSCRIPTION COLLECTIVE ET
DÉVELOPPEMENT DE DATASETS (2020-2022)

Noémie Lucas

14 Février 2024



THE UNIVERSITY
of EDINBURGH

Contexte



Moyen-Orient et
Mondes Musulmans
Groupement d'Intérêt Scientifique



CALFA

B U L A C
[תורה] [大学] [γλώσσες] [ལྷན་སྐྱེས་]

Bibliothèque universitaire
des langues et civilisations

VERS LA SCIENCE
OUVERTE?

La transition numérique
et la recherche
sur le Moyen-Orient
et les mondes musulmans
en France

État des lieux
et perspectives

Septembre 2020

OCR / HTR
et graphie
arabe

Les manuscrits arabes à l'heure
de la reconnaissance automatique
des écritures

Cahier du GIS N°3

Avril 2022

CONTEXTE

Objectifs de RASAM:

- Dimensions **technique** et **expérimentale**: HTR pour l'arabe manuscrit - preuve de concept
- Choix des écritures maghrébines - contexte de la **valorisation des collections patrimoniales** maghrébines de la **BULAC** et du travail en cours mené par le **GIS MOMM** sur le Maghreb et ses ressources
- **Formation** et Humanités numériques

INTELLIGENCE ARTIFICIELLE ET ÉCRITURE MAGHRÉBINE

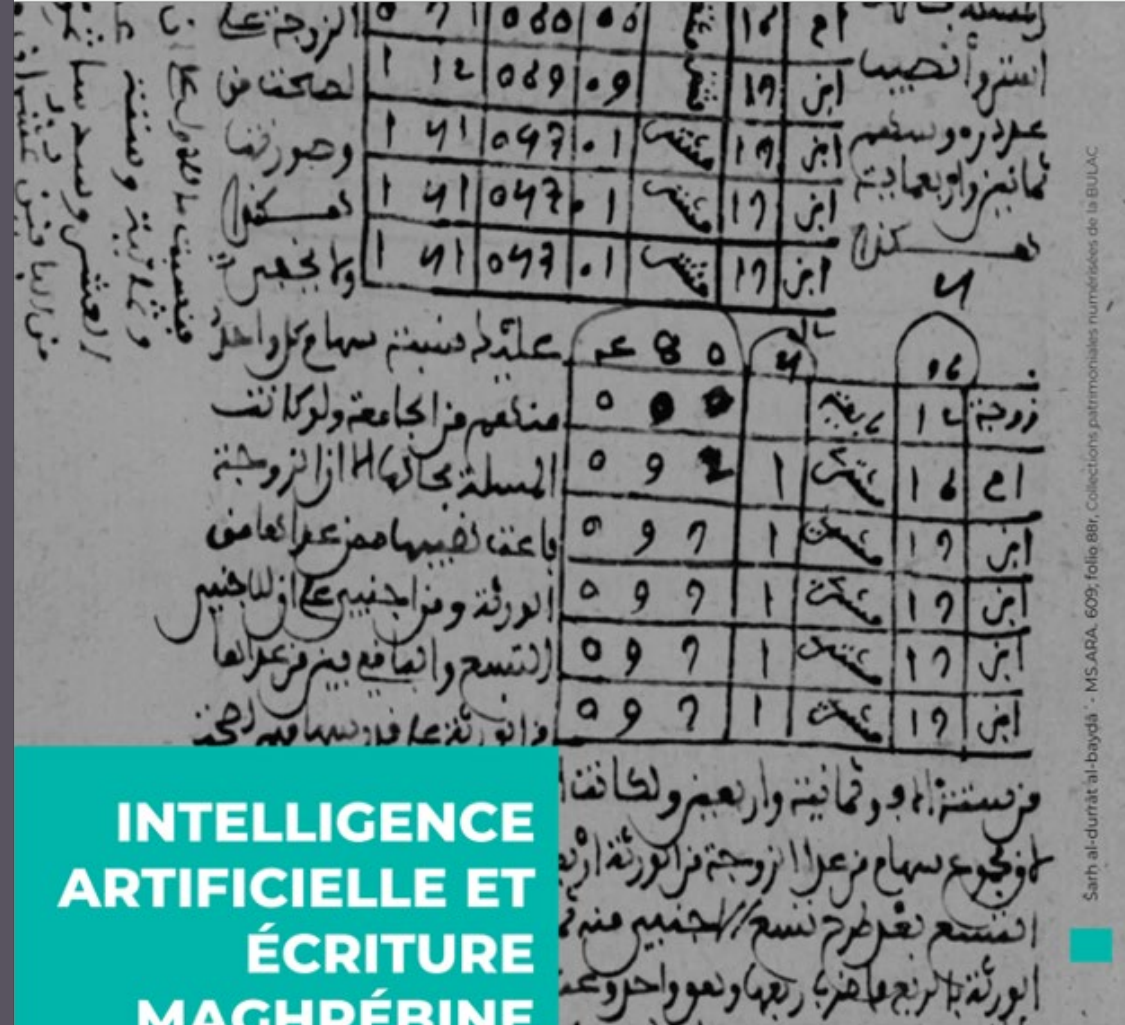
Hackathon pour l'OCR de l'arabe manuscrit

Janvier-Avril 2021



Plus d'informations

- [Site de la Bulac](#)
- [Blog de Calfa](#)
- [Carnet philaranum](#)

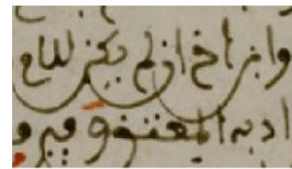


Sarh al-durrat al-bayda - MS.ARA. 609; folio 89r, Collections patrimoniales numérisées de la BULAC

Letters	Characteristics	Theoretical realization	Examples from mss ARA.1977, ARA.609 and ARA.417
<i>bā'</i> <i>tā'</i> <i>ṭā'</i> <i>fā'</i>	(i) Isolated position: concave form – (ii) Final position: closing denticle in the shape of an inverted comma	ب ت ط ف	مغلوب وركب لصاحب مغلوب
<i>dāl</i> <i>ḍāl</i>	Isolated, median and final positions: concave downstroke and final downward spur (<i>dāl kāfiyya</i>)	د ذ	الديسان دناير يزيد
<i>dāl</i> <i>ḍāl</i>	Final position: marked semicircular descender, resembling the letters <i>rā'</i> and <i>zā'</i>	ر ز	محمد فقد بعد
<i>sīm</i> <i>ṣīm</i> <i>ṣād</i> <i>ḍād</i> <i>qāf</i> <i>nū</i>	Final position: exaggerated semi-circular descenders, often described as 'swooping' or 'plunging', stretching below the following word	س ص ق ن	بن عبد كان من ماله
<i>ṣād</i> <i>ḍād</i> <i>ṭā'</i> <i>zā'</i>	Oval or semi-circular body and lack of denticle	ك س ط	قسنطينة اصطلاح سبط
<i>'ayn</i> <i>ḡayn</i>	Initial position: oversized curl	ع غ	عاشه (عاشة) عم
<i>kāf</i>	Initial and median positions: semicircle topped by a diagonal stroke	ك ق	كثير كذلك وكتب
<i>mīm</i>	Final and isolated positions: long curved tail in two variants (concave or convex)	م ح	اسلام تقدم ايام
<i>hā'</i> <i>tā'</i> <i>marbūṭa</i>	Isolated position: drawn in the shape of a '6', sometimes inverted	ه و	هاذه هاذه (هذه) يذكره

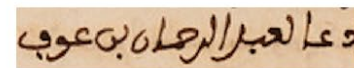
Graphie « maghrébine » / ḥaṭṭ maḡribī : spécificités et enjeux des écritures arabes

1. Realization of the letter *nūn* (and *sīn*, *shīn*, *ṣad*, *ḍād*) in final position :

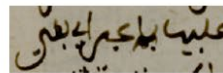


واين اخ ان لم يكن للام

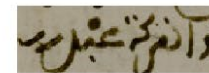
2. Confusion of some letters (*dāl / rā'* ; *dāl / zayn*)



دعا لعبد الرحمن بن عوف

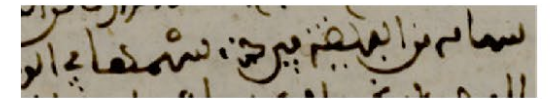


عليه بما عبر اي بقي



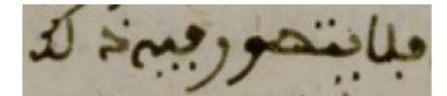
والتركة عبد

3. Same word made in different ways



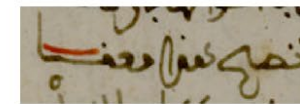
سهمه من الفريضة في جزء سهمها في

4. Spacing between words

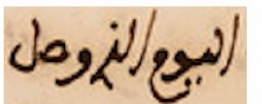


فلا يتصور فيه ذلك

5. Singular ligatures



نصح هذا معنا



اليوم الذي وصل

RASAM 1 (2020-2021)



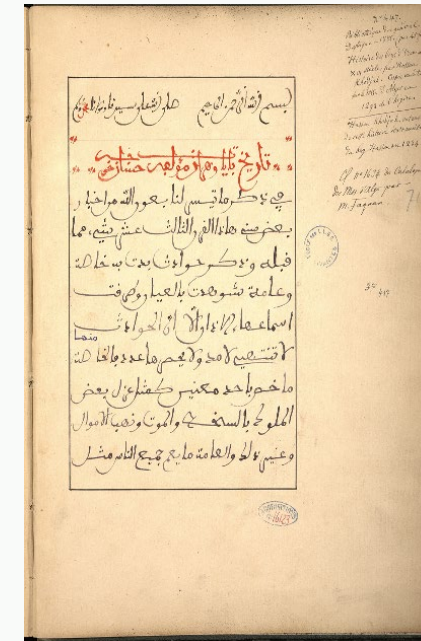
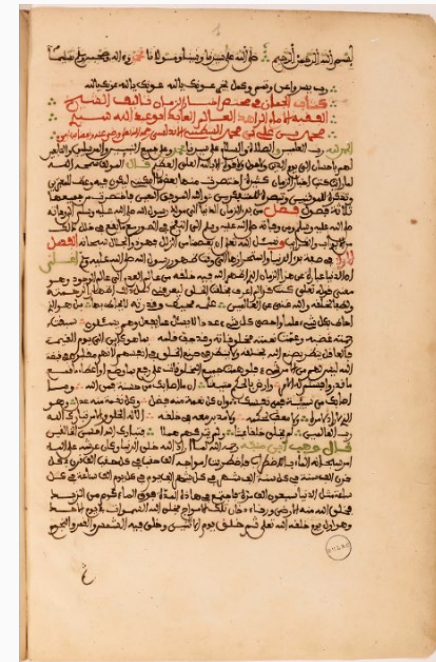
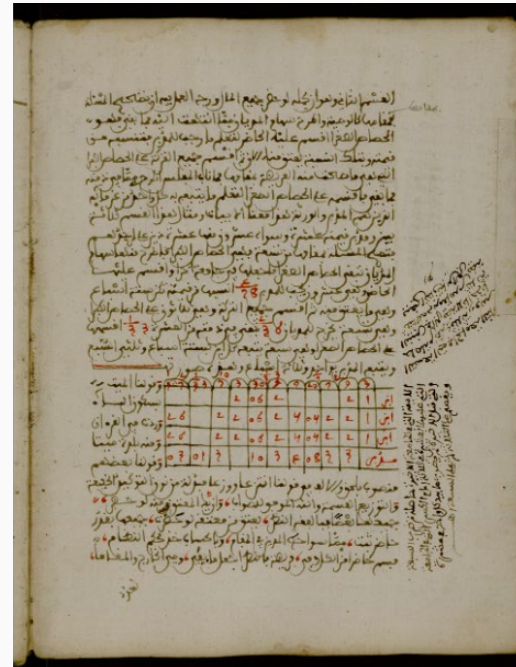
Séance de hackathon animée par Noémie Lucas et Chahan Vidal-Gorène (Maxime Ruscio / BULAC).

لصم يا احسان الذي يوم الدين والاهل والافول الابالته العلي العظيم **فصل** الموقد سمحه الله
لما رايت كتب اخبار الزمان كثيره اختصرت منها بعضه ما لم يكن ليكون فيه وعظ المعين
وتذكرة الموفين وتبصره للمعتبرين نواته الموقد المعين باختصرت من جميعها
ثلاثة **فصول** من در الزمان الدنيا التي مولد رسول الله صلى الله عليه وسلم والوفاته
صلى الله عليه وسلم ومن وجاته صلى الله عليه وسلم التي التبع في الصور مع ما يقع في ذلك ذلك
من الابواب والفرابي ونسئل الله تعالى ان يعصنا من الزلل جهود والجمال سبحانه **الفصل**
الاول في صفة برو الدنيا واستحرازها التي وقت حضور رسول الله صلى الله عليه وسلم **الفصل**
له الدنيا عبارة عن هذا الزمان الذي ارضهم الله فيه خلقه من عالم العدم التي عالم الوجود وهو
معنى قوله تعالى كنت كالم اعرف بخلقنا الخلق ليعرفني كل ذلك اظهر ان رحمة
وتكفنا خلقه والله عنى عن العالمين ع علمه بحسب وقدرته لا يحاكم بهما بن حوالته
اهاف بكل شيء وعلموا وحصى كل شيء عددا لا يسئل عما يجعل وهم يسئلون سبقت
رحمة غضبه وعنت نعمته مخلوقاته وقد جف قلبه بما هو كإبي التي يوم القيمة
والعاقبة يتضر بصنع الله بخلقها ولا يتغير في صنع الخلق في انفسهم لانهم مخلوقون فيفة
الله ليس لهم من امر قه و جلوه تمت جميع المخلوقات على رجع بها وضع اواعضاها ما يقع
ما قدر واجلس له الام وارغب بالخلق متيقنا ان ملاصابت من حسنة من الله وما
اصابت من سيئة جميع نفسك وان كل نعمة منه فضل وكل نعمة منه عذاب وهو
الذي افاض علينا من نعمه ما لا يحصى ولا يدرك ولا يدركه ولا يدركه ولا يدركه

Corpus

- ▶ Manuscripts issus des collections patrimoniales de la BULAC
- ▶ Deux (MS.ARA. 609 et MS.ARA.1977) étaient disponibles sur BINA <https://bina.bulac.fr>
- ▶ Le MS.ARA.417 a été numérisé à la demande

- Styles d'écriture différents et courants
- Variété des thèmes limitée (histoire et droit)
- Variété des mises en page





Search...



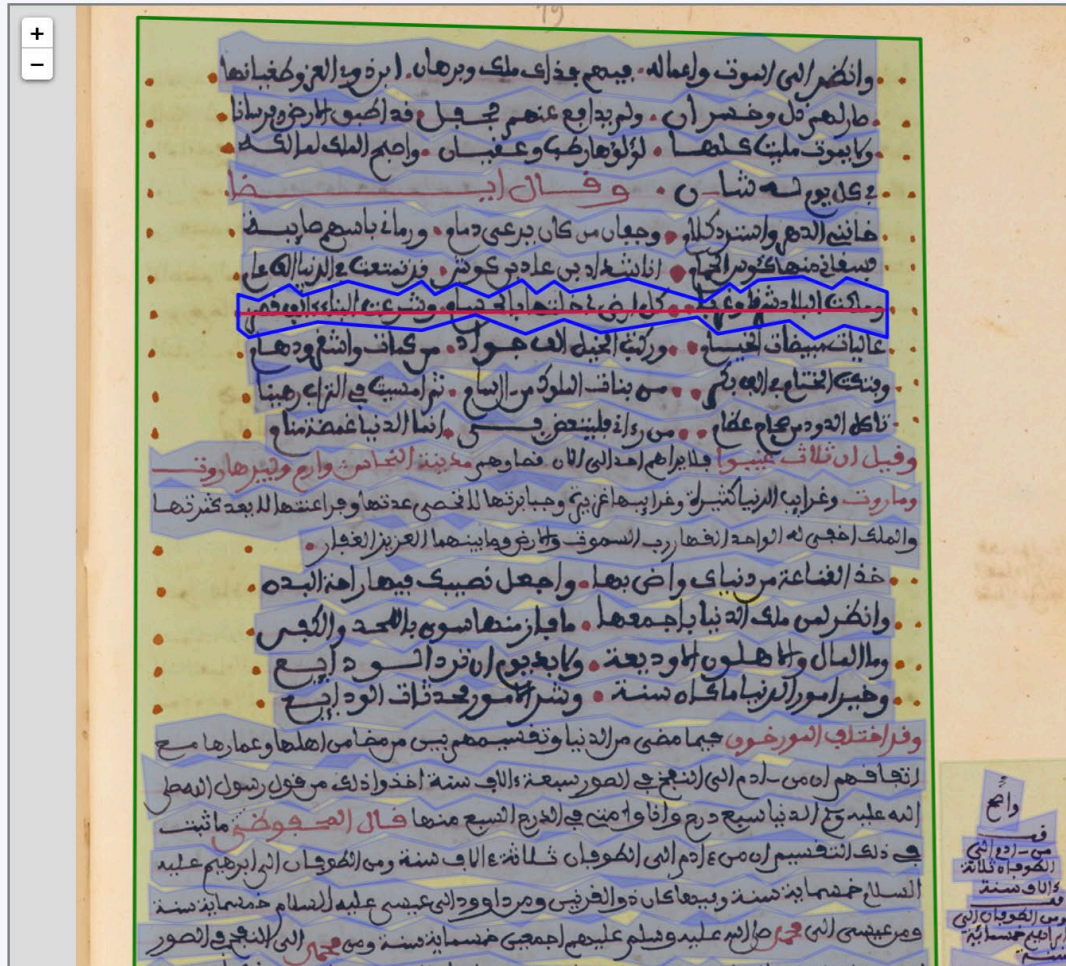
Chahan Vidal-Gorène

Projects / BULAC Hackathon - ms ARA 1977 / Images / BULAC_MS_ARA_1977_0028 / Labellize

← Back to BULAC_MS_ARA_1977_0028 1. Layout Analysis 2. Generate Lines 3. Text Recognition

BULAC_MS_ARA_1977_0028

Save All



- 5 خانى الدهر واسترد كلام وجفان من كان يرعى دمام ورماني باسهم صائبة
- 6 فسقانى منها كنوس الحمام انا شداد بن عاد بن كوش قد تمتعت في الدنيا الف عام
- 7 وملكت البلاد شرقا وغربا كل ارض دخلتها بالحسام وشرعت البناء الف قصر
- 8 عاليات مبيضات الخيام وركبت الخيل الف جواد من كمات واشقر ودهام
- 9 وفتحت الختام في الف بكر من بنات الملوك من السام ثم امسيت في التراب رهينا
- 10 تاكل الدود من صحاح عظام من رءانى فليتعض بى انما الدنيا غمضة منام
- 11 يل ان ثلاث غيبوا فلا يراهم احد الى الان فساو هي مدينة النحاس وارم وبيير هاروت
- 12 برائب الدنيا كثيرة وغرائبها غزيرة وجابرتها لا تحصى عدتها وفراعتها لا يعد كثرتها
- 13 والملك اخفى له الواحد القهار رب السموت والارض وما بينهما العزيز الغفار
- 14 خذ القناعة من دنياك واض بها واجعل نصيبك فيها راحة البده
- 15 وانظر لمن ملك الدنيا باجمعها ما فاز منها سوه باللحد والكفر

Dashboard

MAIN MENU

Projects

SECONDARY MENU

Settings

Guide

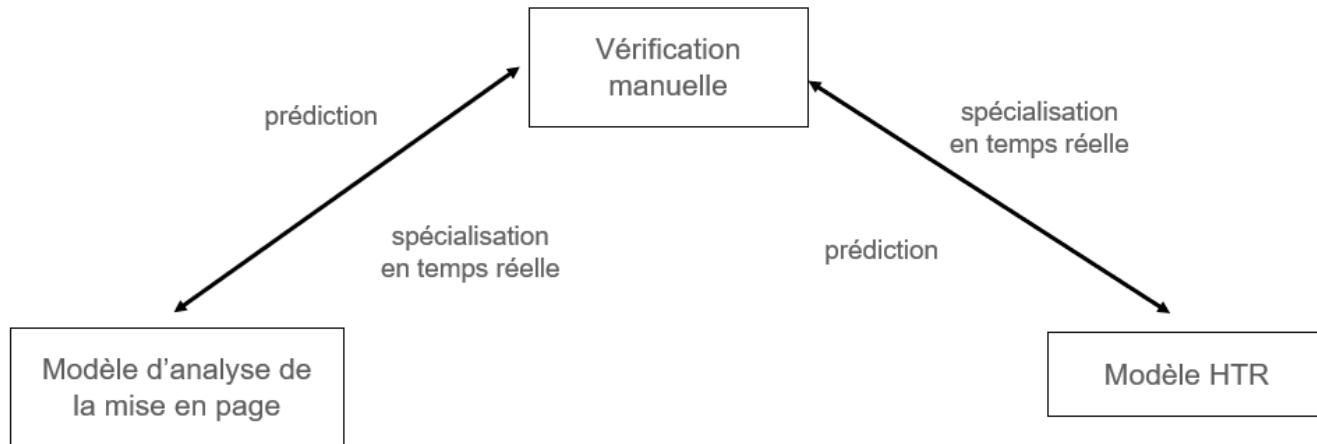
Terms of Service



واضح
فمن اراد ان
يكون له ثلاثة
اولاد فليصوم
يومين من الصوم
الاربعين

Déroulement du hackathon

**Objectif : créer un
dataset ouvert
pour les écritures
arabes maghribi
représentatif et
fonctionnel**

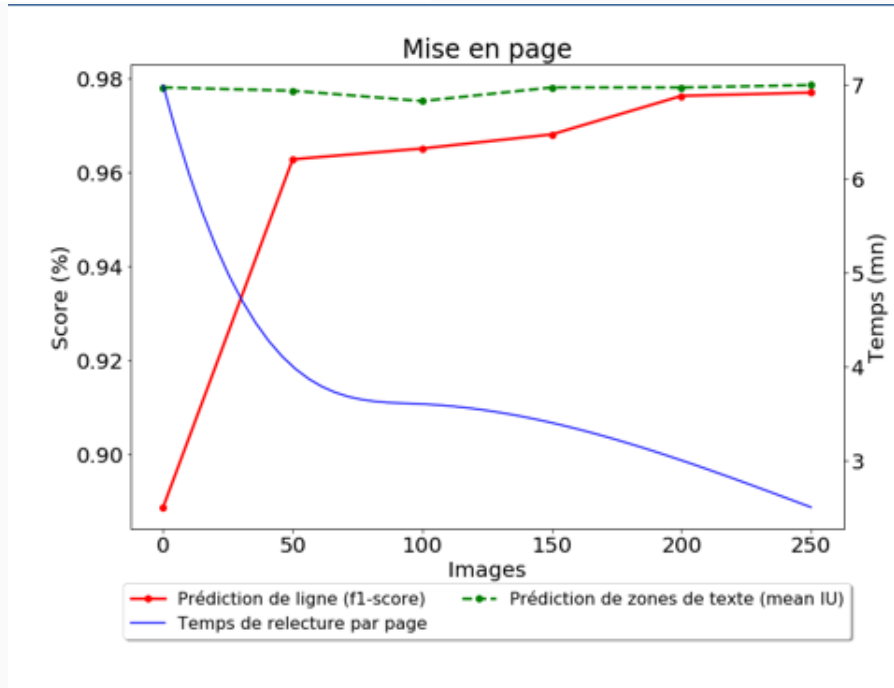




défaut



après 50 images relues



Déroulement

La spécialisation progressive des modèles a conduit à un gain de temps de **75%** dans la vérification de la mise en page.

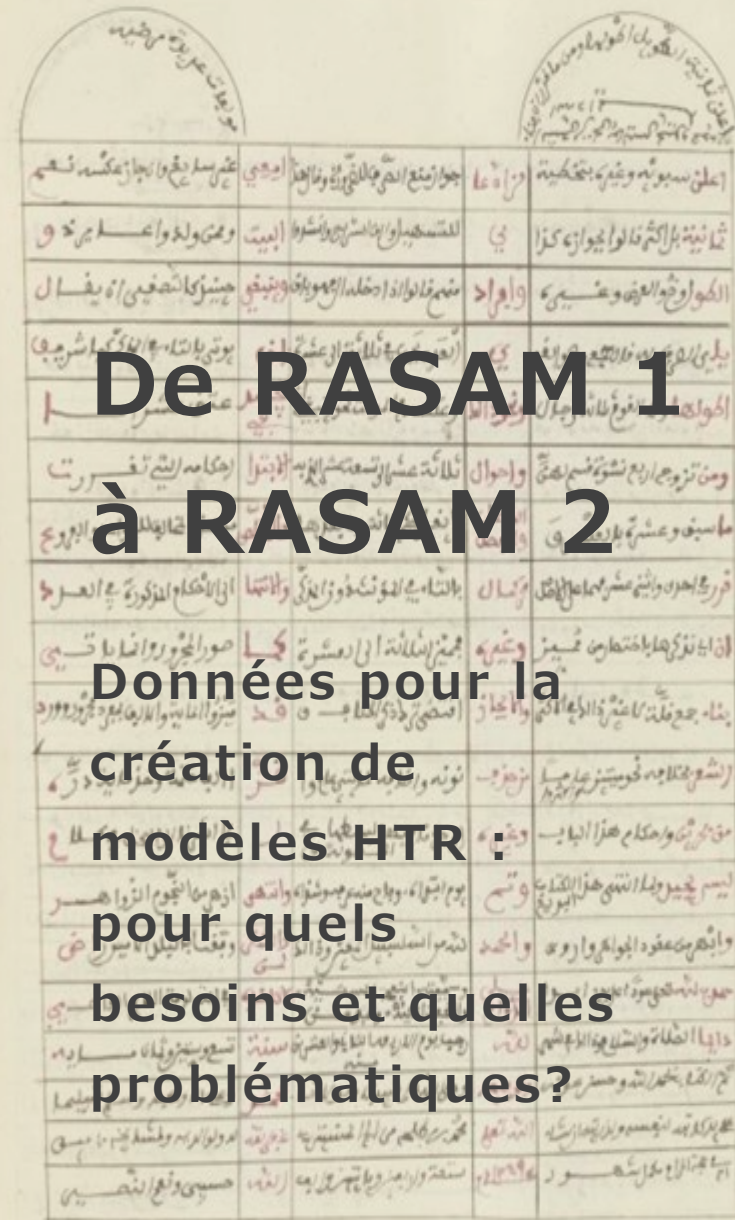
Rasam 1

- 3 ms, 300 pages
- Approche au mot
- Démonstration de faisabilité
- Enjeu des écritures arabes maghrébines
- Question mise en page

BASE SCIENTIFIQUE

Rasam 2

- 15 ms, 250 pages
- Renforcement de la polyvalence des modèles avec un hackathon (22 participants)
- Traitement HTR de ms de la Bulac
- Evaluation de la spécialisation: quels résultats? Quels coûts?



De RASAM 1 à RASAM 2

Données pour la
création de
modèles HTR :
pour quels
besoins et quelles
problématiques?

Corpus RASAM 2



BULAC.MS.ARA.6

BULAC.MS.ARA.9

BULAC.MS.ARA.23

BULAC.MS.ARA.24

BULAC.MS.ARA.45b

BULAC.MS.ARA.65

BULAC.MS.ARA.1926

BULAC.MS.ARA.1936

BULAC.MS.ARA.1943

BULAC.MS.ARA.1944

BULAC.MS.ARA.1946

BULAC.MS.ARA.1947

BULAC.MS.ARA.1960

BULAC.MS.ARA.1982

BULAC.MS.ARA.1983

- Diversité plus grande des genres =
nouveau vocabulaire

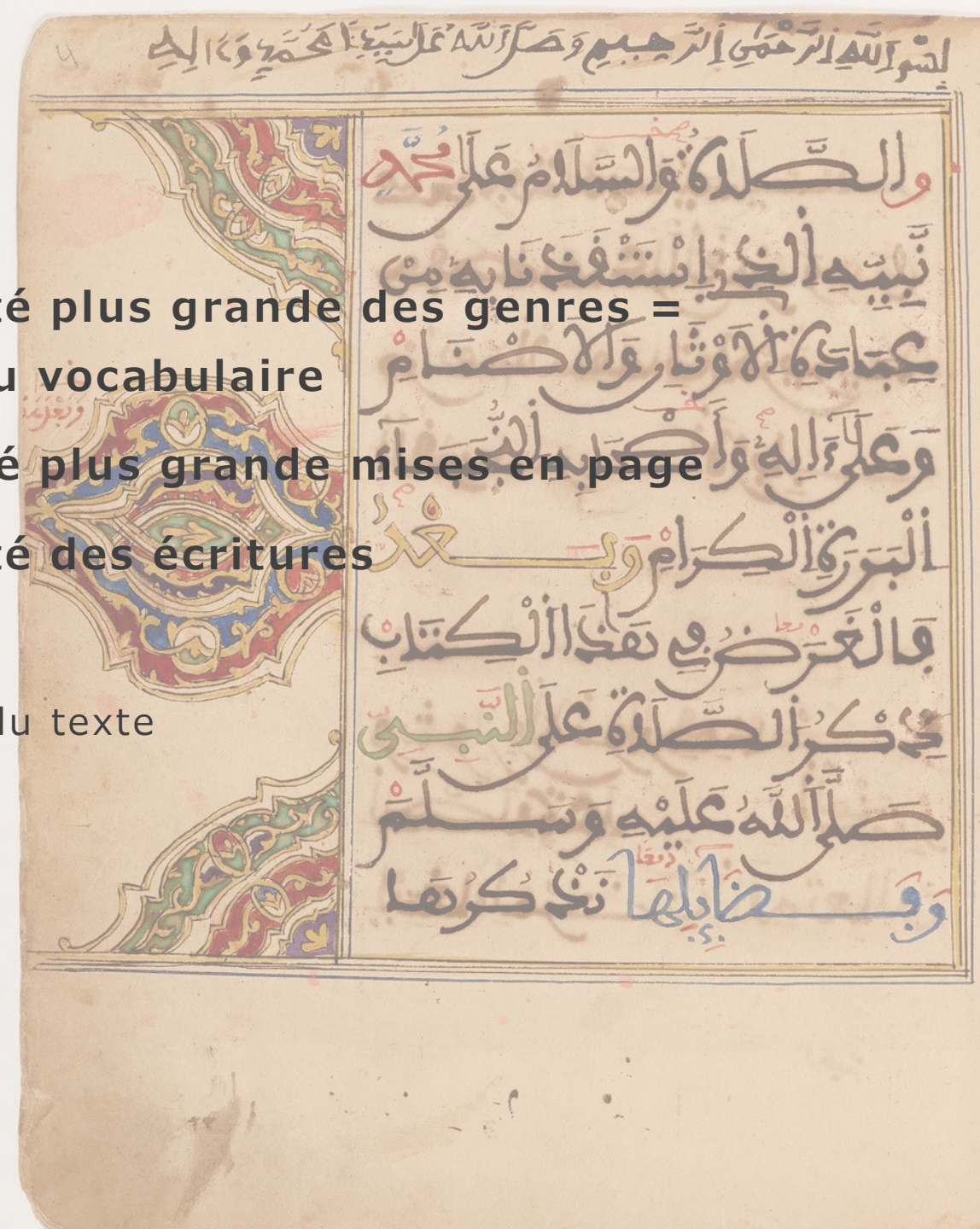
- Pluralité plus grande mises en page

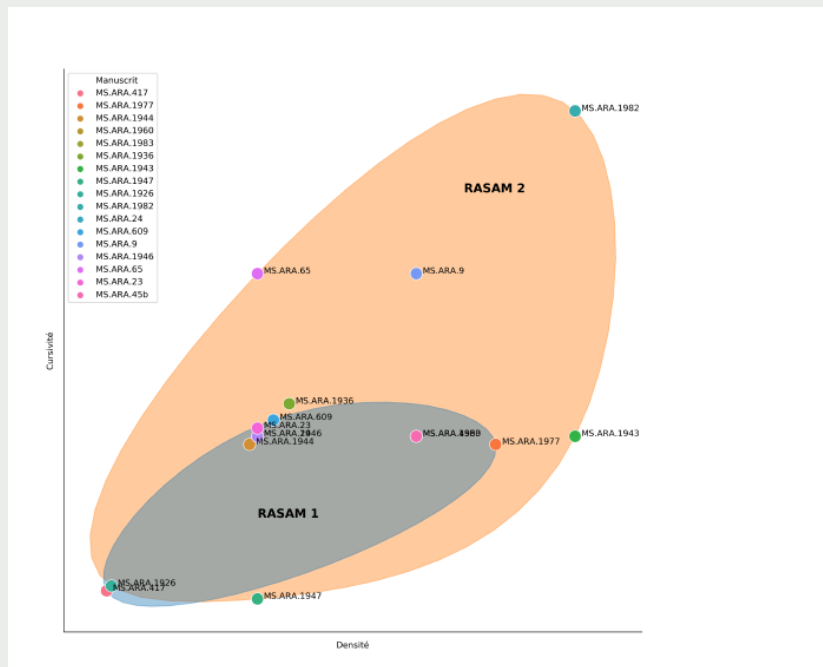
- Diversité des écritures

- Couleur

- Densité du texte

- Style





RASAM 2

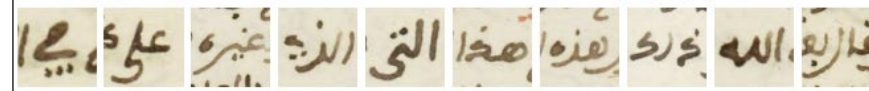
MS.ARA.6



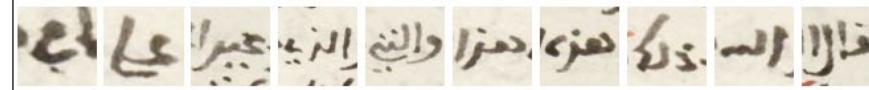
MS.ARA.9



MS.ARA.23



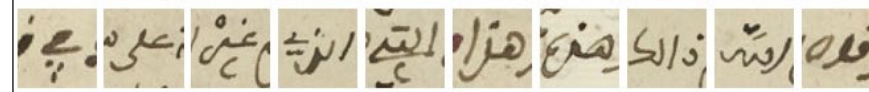
MS.ARA.24



MS.ARA.45b



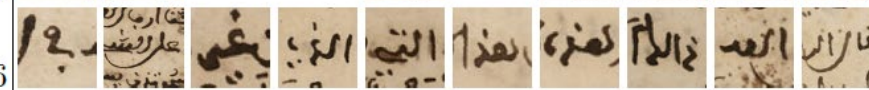
MS.ARA.65



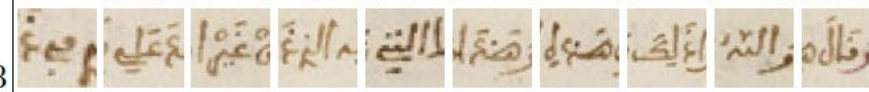
MS.ARA.1926



MS.ARA.1936

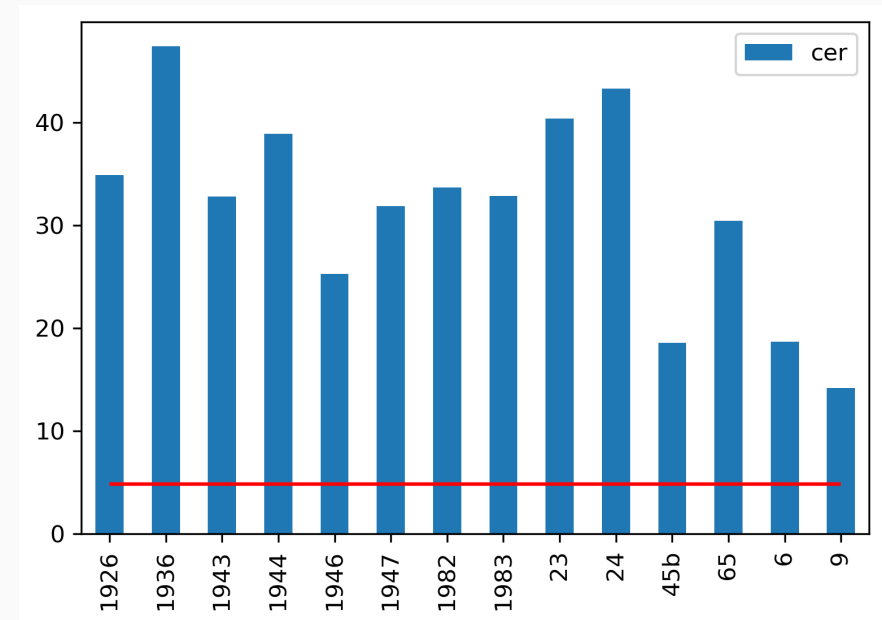


MS.ARA.1943



MS.ARA.6	18,68
MS.ARA.9	14,18
MS.ARA.23	40,34
MS.ARA.24	43,28
MS.ARA.45b	18,56
MS.ARA.65	0,44
MS.ARA.1926	34,86
MS.ARA.1936	47,39
MS.ARA.1943	32,78
MS.ARA.1944	38,86
MS.ARA.1946	25,28
MS.ARA.1947	31,83
MS.ARA.1982	33,67
MS.ARA.1983	32,84

Modèle RASAM1 appliqué aux manuscrits de RASAM2



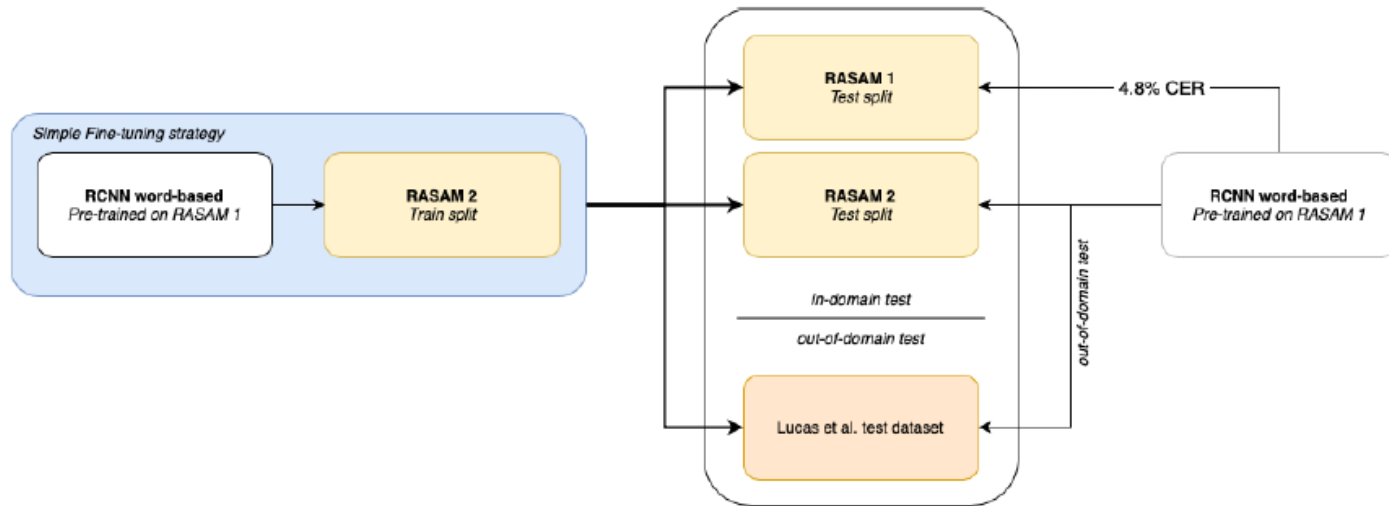


Fig. 3. Experiments conducted on the new dataset and comparison with the RASAM 1 and RASAM 2 models

=> Lucas, N., Salah, C., Vidal-Gorène, C.: *New Results for the Text Recognition of Arabic Maghribi Manuscripts - Managing an Under-resourced Script* (2022), <https://hal-enc.archives-ouvertes.fr/hal-03874725> , working paper or preprint

De RASAM 1 à RASAM 2: Méthodologie appliquée

Tendance actuelle: RASAM1 + spécialisation avec 10-20 images = ~ 6 - 12% de CER

Résultats RASAM 2

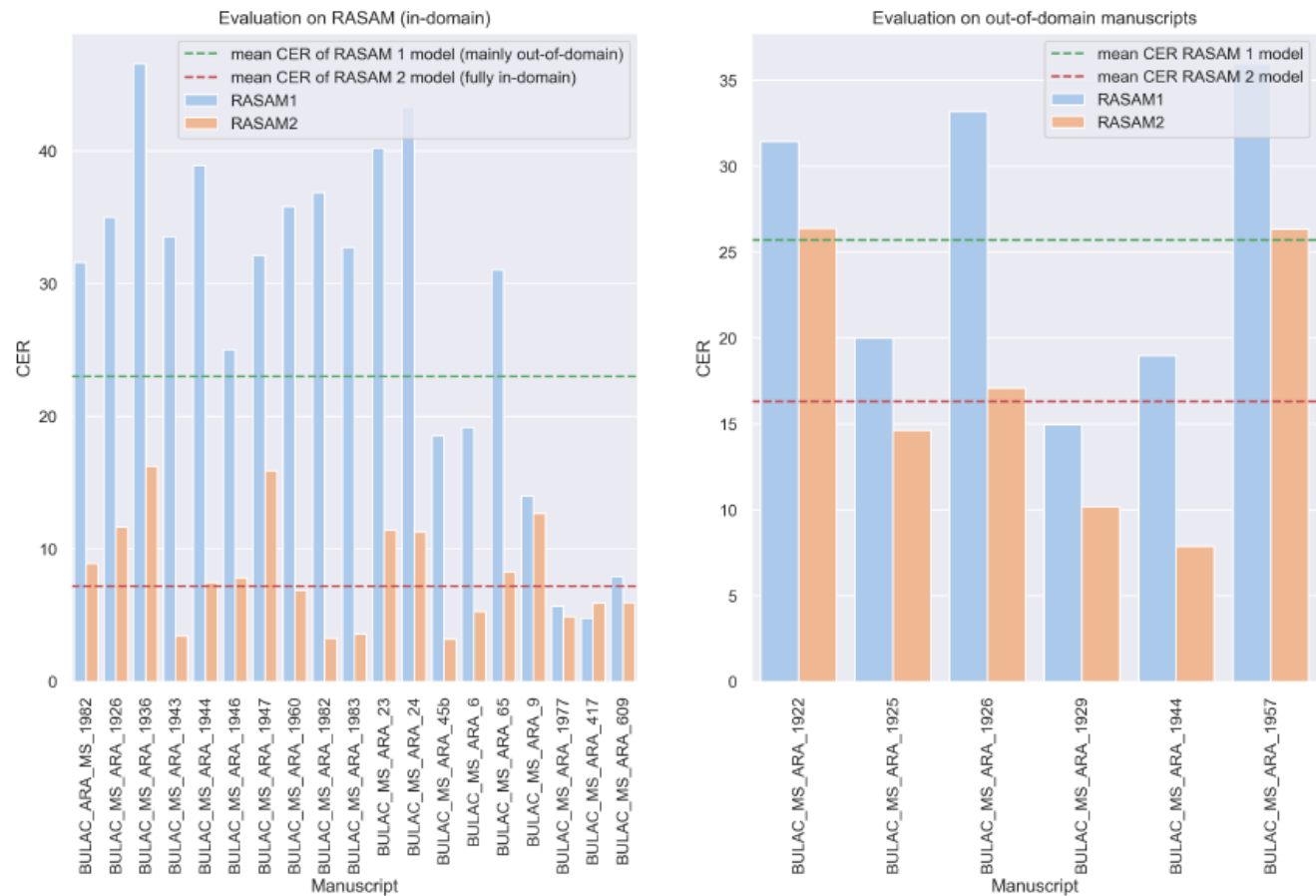


Fig. 5. Distribution of CERs obtained by RASAM 1 (in blue) and RASAM 2 (orange) for each in-domain and out-of-domain manuscript

Résultats

RASAM 2

- **Nouvelle tendance : RASAM2 + spécialisation avec 10 images = ~ 6 - 10% de CER. Gain de 50% dans l'annotation**
- Le jeu de données RASAM 2 dataset est représentatif de la variété des écritures arabes maghrébines manuscrites – Modèle plus polyvalent que RASAM1
- L'approche au mot démontre des performances in-domain et out-of-domains très satisfaisantes
- Le jeu de données constitue une base solide pour la spécialisation des modèles futures sur des manuscrits de ces graphies

مما لا شك فيه... ولا تشبهه... مالك ناصيتك... وتشتريه... انظر ان... او ينفذ... ملكك معشرك... انشجرت... فقلت... اعز اربا... لفضيبت... فيلحاه... انقضت...

Conclusion

- Deux datasets disponibles:
- RASAM 1: <https://github.com/califa-co/rasam-dataset>
- RASAM 2: En cours de mise en ligne
- Question de l'HTR des manuscrits arabes largement surmontées sous réserve d'un besoin identifié et d'un cahier des charges bien établi
- Pipeline désormais bien établi pour les écritures arabes



Facteurs de réussite

- 1 - Les projets et l'initiative plus large étaient soutenus et pilotés par le Groupement d'intérêt scientifique « Moyen-Orient et mondes musulmans »
- 2- Le projet a bénéficié d'une collaboration (par le biais d'un partenariat) entre trois acteurs complémentaires : le Groupement d'intérêt scientifique MOMM d'une part, la BULAC d'autre part et Calfa.
- 3- Les objectifs avaient été clairement énoncés et précisés au départ.