



-HAT-

corpus 2020

by H. Sayoud / EDT

Description of the HAT Corpus

Our textual dataset is composed of 100 groups of Arabic texts that are extracted from 100 different Arabic books. The books were written by 100 different authors and each group contains 3 different texts that are written by the same author, which means that each group belongs to only one author. This set of 300 text documents has been collected in 2019 from “*Alwaraq digital library*” (www.alwaraq.net).

We called this corpus "HAT" (*i.e. Hundred of Arabic Travelers*). Furthermore, the HAT corpus could represent a reference dataset for author style analysis in Arabic, which could be used by NLP researchers for a purpose of comparative evaluation. More details on HAT can be found at www.univ-edt.com.

For concreteness, here is a piece of text belonging to Author #92 (*figure 1*).

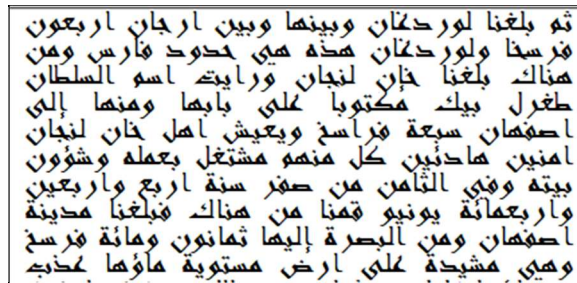


Fig. 1. Example of Arabic text belonging to Author #92 (*N. Khasru*)



Fig. 2. Ancient portrait and sheet of paper containing a text of Author #92 (*N. Khasru*).

List of the authors and their corresponding topics

The different authors are listed in table 1.

Table 1. List of Authors.

Reference Number	Author Name	Author name in Arabic	Reference Number	Author Name	Author name in Arabic
1	Ibrahim Basha	إبراهيم رفيع باشا	51	Sara Sirayt	ماراسيرايت
2	Ibrahim Almazini	إبراهيم محمد القادر المازني	52	Said ben Ali	سعيد بن علي المغربي
3	Ibn Almujafer	ابن المجاور	53	Saloum Addahdah	سلوم الحداد
4	Ibn Batuta	ابن بطوطة	54	Suleiman bnu Sayam	سليمان بن سيام
5	Ibn Djubair	ابن جبير	55	Sadeq Bacha	ساذق باشا المؤيد العظم
6	Ibn Fathlane	ابن فضلان	56	Taleb Muchtraq	طالب مشتاق
7	Ahmed Atelemsani	أحمد بن صلال التلمساني	57	AbdAlHalim Almisri	محمد الحليم المصري
8	Ahmed Albaqri	أحمد ناصر محمود البكري	58	AbdAlGhani Alnabulsi	محمد الغني بن اسماعيل النابلسي
9	Ahmed Hassaneen	أحمد محمد حسنين	59	AbdAllatif Albaghdadi	محمد اللطيف البغدادي
10	Ikhraj Arradhi	إخراج أبقراط عمر طاهر الراعي	60	AbdAllah Afandi	محمد الله أفندي إلياس رعد
11	Ismael Mustapha	إسماعيل محمد مصطفى	61	AbdAllah Alkandary	محمد الله بن محمد بن سالم باختيار الكندي
12	Alab Malon	الأمير المسمى مالون	62	AbdAllah Heshima	محمد الله حشيمه
13	Alab Rene	الأمير رينه مورتد اليسوعي	63	AbdAlMuhsein Albarkaty	محمد المحسن البركاتي
14	Alab Luis	الأمير لويس جلابره اليسوعي	64	Abdalmasseeh Antaky	محمد المسيح أنطاقي
15	Alab Michel	الأمير ميشال جوليان اليسوعي	65	Azam AbdAlWahab	عزاه محمد الوهاب
16	Alab Henry	الأمير هنري لامنس اليسوعي	66	Ali Ibrahim	علي إبراهيم عردي
17	Ameer Mohamed Ali	الأمير محمد علي واه	67	Ali Aldjerjawy	علي أحمد الجرجاوي
18	Albalwee	البليوي	68	Fouad Ghosn	فؤاد غسن
19	Alhadj Ahmed Kamal	الحاج أحمد جمال	69	Cartsten Nebur	كارستن نيبور
20	Alkhouri Batras	الخوري بطرس العنقاري	70	Clodis Jems	كلوديس جيمس ريج
21	Almudhfir Almusili	الرحالة العدي المظفر الموسلي	71	Alexander Kinglake	كينغليك ، نقلنا إلى العروبة محمود العليدي
22	Seddeek Hassan Khan	السيد حديق حسن خان	72	Abu AbdAllah Alabdary	أبي محمد الله محمد بن محمد العبدري الجدي
23	Muhsen Abutabeebh	السيد محسن أبو طيخ	73	Lessan Eddin bnu Alkhatib	لسان الدين ابن الخطيب
24	Muhsen Alameen	السيد محسن الأمين	74	Alab Luis Ranz	للأمير لويس رنز
25	Muhamed ibn Assayed	السيد محمد ابن السيد أحمد الحسيني	75	Leonhard Rauwolf	للرحالة الهولندي الدكتور ليونهارد راوولف
26	Sheikh Muhamed Alqayati	الشيخ محمد عبد الجواد القاياتي	76	Mar Athanasius	مار اثاناسيوس انطاقيوس ثوري
27	Alqas Ishaq	القاس اسحق أرمله	77	Muheb Eddine Alhamwy	محمد الدين الحموي
28	Alqas Antonio	القاس أنطونيو شيني اللبناني	78	Muhamed Alhajwy	محمد الحجوي
29	Alqas Suleiman	القاس سليمان حايغ السلطاني الموصلّي	79	Muhamed Alhur	محمد الحر
30	Othman Multadha	المطربة الأميرية بالقاهرة / عثمان مرتضى	80	Muhamed Alhashemy	محمد الهاشمي
31	Alyussy	اليوسي	81	Muhamed Amine	محمد أمين فخري بلن
32	Benyamin Alandalussy	بنامين يونا التطيلي النجاري الأندلسي	82	Muhamed Aumissawy	محمد بن محمد الله الحسيني الموسوي
33	Imad Raouf	تحقيق د. عماد محمد السلام رؤوف	83	Muhamed Thabet	محمد ثابت
34	Seyar Aldjameel	ترجمة د. سيار الجميل	84	Muhamed Rashed Reda	محمد رشيد رضا
35	Batras Haddad	ترجمة وتعليق الأمير د. بطرس حداد	85	Muhamed Saoud	محمد سعود العوري
36	Mustapha Juda	تعريب مصطفى محمد جودة	86	Muhamed Kurd	محمد عردي علي
37	Djerji Serkis	جرجي بن يوسف اليان سركيس	87	Muhamed Labeeb	محمد لبيب البتوني
38	Djerjee Zidan	جرجي زيدان	88	Mahieddine Reda	مهي الدين رضا
39	John Luis	جون لويس بورخماره	89	Mustapha Muhamed	مصطفى محمد
40	HarethnYussuf	حارث يوسف خديمة	90	Naji Jawed	ناجي جواد
41	Habeeb Afendy	حبيب أفندي فريده البطلاني	91	Nasser Djarjees	ناصر جرجيس
42	Hassan Muhamed	حسن محمد جوهر	92	Nasser Khasru	ناصر خسرو
43	Hassan Fawzi	حسن فوزي	93	Nazih Muayed	نزيه مؤيد العظم
44	Alkhouri Ibfahim	حضرة الخوري إبراهيم عرفوش	94	Newab Hamid	نواب محمد بار جونك بصادر
45	Alkhouri Luis	حضرة الخوري لويس ملحة	95	Nawal Assaadawy	نوال السعداوي
46	Siraz Dubler	دوقما سيزار دوبلر	96	Henry Tristram	هـ. بـ. تريسترام
47	Hamad Aljasser	حمد الجاسر	97	Hor Diver	هور ديفر
48	Sherif Hetata	د. شريف حنّانة	98	Yahia bnu Muthassen	يحيى بن أبي السفا بن أحمد المعروفه وابن مناسن
49	Refaa Attahtawi	رفاعة الططاوي	99	Yussuf Aboud	يوسف عبود
50	Arrabie ben Suleyman	رواية تلميذه الربيع بن سليمان الجيزي	100	Youssof Kamal	يوسف جمال

The HAT text documents have a relatively short/medium size: the average text length is about 1100 words per document and there are 3 documents per author, which corresponds to 300 documents in the entire corpus. This situation involves severe experimental conditions, since it has been shown, in previous research works (Eder 2010;

Signoriello 2005), that the minimum number of words per text should be at least 2500 words to get good attribution performances. In this investigation, the use of medium-short texts is interesting in order to evaluate the different classifiers with small documents in Arabic. In fact, when short texts are used, the AA performances decrease and it becomes difficult to make an efficient identification.