

## TYPOLOGIE D'UN ENSEMBLE DE TEXTES ARABES D'APRÈS LES OCCURRENCES DE FORMES DE MOTS ET DE LOCUTIONS

[TEXTES ARABES]

J.-P. BENZÉCRI  
A. CHABIR\*

### 1 Composition du corpus

L'article [MOTS ARABES] (in *CAD*, Vol XIX, n°1) rend compte d'une expérience stylométrique portant sur 13 fragments de textes longs, chacun, d'environ 3000 caractères (3k); et tous compris dans le domaine de la philosophie, même s'ils diffèrent quant au niveau de l'exposé.

Le présent travail, porte sur 91 fragments, longs également de quelque 3k, mais d'une plus grande diversité de genre. La cueillette, sans prétendre aucunement embrasser l'infinie diversité de ce qu'on peut lire en arabe, sort de la philosophie pour glaner dans la prose et la poésie classique; et ne fait pas fi de ce qu'ont produit les modernes; non seulement en littérature proprement dite, mais aussi dans la presse quotidienne.

En décrivant le corpus ainsi étendu, nous anticiperons, parfois, sur les résultats des analyses stylométriques.

#### 1.1 Poésie classique et poésie moderne

Du classique, outre les 13 fragments déjà analysés dans [MOTS ARABES], on a pris, d'abord, un recueil de vers et de prose choisi pour exercer à l'épreuve de récitation les élèves de 3-ème année de la Section d'Arabe Littéral de l'Institut National des Langues et Civilisations Orientales: l'INALCO, connu de tous pour être le domaine des Langues'O.

Le Recueil comprend 35 textes, dus à 28 auteurs différents, pour chacun desquels est donnée une brève notice. L'ensemble peut être distribué en trois parties, que nous marquons par les signes { $\pi$ , ¶, §}.

---

(\*) Maître de langue à l'Institut National des Langues et Civilisations Orientales.

### Indice des poètes : فحوريس الشعراء

$\pi_{sc1}$	الفرزدق	عمر بن أبي ربيعة
$\pi_{sc2}$	ابن الرومي	أبو العتاهية
$\pi_{sc3}$	ابو فراس الحمداني	أبو الطيب المتنبي
$\pi_{sc4}$	أبو العلاء المعري	أبو الحسن الانباري
$\pi_{sc5}$	ابن الفارض	ابن زيدون

Dans  $\{\pi\}$  sont 12 fragments de poésie ancienne. L'humaniste épris des lettres arabes reconnaîtra ici de très grands noms qui tous méritent l'attention: mais comme la plupart des fragments sont brefs, on s'est résigné à les cumuler en cinq blocs; lesquels, sous les sigles  $\{\pi_{sc1}, \pi_{sc2}, \pi_{sc3}, \pi_{sc4}, \pi_{sc5}\}$  répondent, à peu près, aux cinq premiers siècles après l'Hégire (la poésie anté-islamique étant absente).

Dans  $\{\mathfrak{f}\}$  sont 15 fragments de poésie postérieurs à 1800. La série s'ouvre par  $\mathfrak{f}Lac$ : le poème si bien connu de LAMARTINE, dans une surprenante traduction dont l'analyse stylométrique confirme le strict classicisme.

Suit  $\mathfrak{f}shu$ : quatre fragment de celui qu'on appela le *Prince des poètes*,  $\mathfrak{a}hmad \mathfrak{s}a\mathfrak{u}q\mathfrak{i}$ , أحمد شوقي, (1868-1932); qui, avec l'harmonie de ses vers pesés sur la balance classique, n'a pas le vocabulaire des poètes anciens.

Sort également de ce vocabulaire, malgré son thème, la  $qa\mathfrak{s}\mathfrak{i}d\mathfrak{d}\mathfrak{a}$ ,  $\mathfrak{f}Arb$ , où  $\mathfrak{h}\mathfrak{a}f\mathfrak{i}z\ \mathfrak{i}b\mathfrak{r}\mathfrak{a}h\mathfrak{i}m$ , حافظ ابراهيم, chante la langue du  $\mathfrak{d}\mathfrak{a}d$ .

Gibran,  $\mathfrak{g}\mathfrak{i}b\mathfrak{r}\mathfrak{a}n\ \mathfrak{k}\mathfrak{a}l\mathfrak{i}l\ \mathfrak{g}\mathfrak{i}b\mathfrak{r}\mathfrak{a}n$ , جبران خليل جبران, libanais de la diaspora, est connu de ceux-mêmes qui n'ont rien lu d'arabe: deux fragments de lui sont cumulés dans  $\mathfrak{f}gbr$ .

Aux sept autres fragments, la stylométrie n'a pas trouvé de proche parent: les sigles  $\{\mathfrak{f}occ, \mathfrak{f}dpl\}$  rappellent que les notices du Recueil attribuent aux auteurs soit une affinité occidentale, soit l'état de diplomate.

En passant, disons que ces brèves notices, cumulées suivant les trois parties,  $\{\pi, \mathfrak{f}, \mathfrak{s}\}$ , du Recueil, ont donné:  $\{\&Id\mathfrak{s}, \&Id\pi, \&Id\mathfrak{f}\}$ .

#### 1.2 Prose classique

$\{\mathfrak{s}KwD, \mathfrak{s}Hdt, \mathfrak{s}jhz, \mathfrak{s}rbh, \mathfrak{s}Mqm, \mathfrak{s}sfa, \mathfrak{s}twh, \mathfrak{s}tah\}$

Reste les huit proses de  $\{\mathfrak{s}\}$ , courant de l'aube du classicisme à notre XX-ème siècle. Nous renonçons à en donner ici plus qu'un aperçu:  $\mathfrak{s}Hdt$ , Hadith, ou traditions remontant au fondateur de l'Islam;  $\mathfrak{s}rbh$ : éloge des fortes sciences et des belles lettres;  $\mathfrak{s}twh$ , marque de la sollicitude d'un calife envers la plèbe de Bagdad;  $\mathfrak{s}sfa$ , anecdote extraite de l'encyclopédie des "Frères de la pureté et de la sincérité" (اخوان الصفاء)...

La diversité de thème paraît telle que, malgré la brièveté de certains fragments, nous n'en avons pas tenté de cumul.

Et pourtant, de §KwD, *kalīlah wa dimnah*, *كليلة و دمنة*, où *ibn-l-muqaffa'*, *ابن المقفع*, suivant la tradition de l'Inde et de l'Iran, met en arabe des fables dont héritera LA FONTAINE; à la scène familière, §tah, où, enfant, se découvre aveugle notre futur Docteur de la Sorbonne, *taha ḥuṣayn*, *طه حسين*; la stylométrie trouve une même prose dont ne s'éloignent ni le hadith ni le récit historique.

En sort seulement, pour suivre la poésie, §Mqm, cette scène précieuse - les arabe disent *maqāmā*, *مقامة*, séance - où un énigmatique jeune homme, d'abord attentif à un débat de lettrés, répond ensuite à ceux-ci en distillant, de chaque poète, quelques gouttes du parfum de son génie (*āl-maqāmah al-qariḍiyah*).

بديع الزمان الهمذاني : مقامات

Des mêmes *maqāmāt* de *badī' al-zamān al-ḥamdānī*, on a saisi, d'autre part, deux pièces: *mqmB*: un obséquieux filou se fait régaler par un paysan, ébloui dans Bagdad (*āl-maqāmah al-baḡdādiyyah*); et {*mqm1*, *mqm2*}: un nouveau riche vante sa demeure et tout ce dont il lui a plu de l'orner (*āl-maqāmah al-maḍiriyyah*).

À ce genre des séances, la stylométrie associe encore §jhz: où l'illustre *ḡaḥiḏ*, *الجاحظ*, met en scène un savant courtois fermant un jour son école pour porter le deuil d'une femme qu'il n'a jamais vue que de loin.

ابن المقفع : الأدب الكبير ؛ رسالة في الصحابة

Sans quitter encore le domaine de {§}, nous avons d'*ibn-l-Muqaffa'*, l'auteur déjà cité, de *kalīlah wa dimnah*, des propos auxquels ne suffiraient pas les langues des bêtes et des oiseaux: d'une part {1Adb, ..., 4Adb}, une grande partie de *al-'adab al-kabīr*, et d'autre part {1Rsl, ..., 4Rsl}, les premiers chapitres de la *risālah fī al-ṣaḥābah*: exposés de la norme morale et politique de la société islamique, dans ses débuts.

### 1.3 Prose moderne

Enfin, la prose moderne a fourni trois blocs homogènes considérables: {Grb, Sgn, Prs}.

Du libanais *Mīka'il Nu'aymah*, *ميخائل نعيمة*, mort quasi centenaire en 1988, on a, subdivisés en {1Grb, ..., 9Grb, aGrb, ..., kGrb}, les trois quarts du livre *al-ḡirbāl*, *الغربال*, - le tamis - : théorie des belles lettres dont le titre reprend cette formule:

فمهمة الناقد إذن هي غربلة الآثار الأدبية .

“la mission du critique est dépasser au tamis les œuvres littéraires”.

De ce même auteur, une poésie, النهر المتجمد , “le fleuve pris en glace”, est comptée dans la tranche §occ du Recueil.

Dans sign al-‘umr, سجن العمر : la prison des jours, (littéralement: de l'âge) dont on n'a saisi que le début {1Sgn, ..., 9Sgn, aSgn,..., eSgn}, tawfiq al-ḥakīm, توفيق الحكيم, un des maîtres du roman égyptien au XX-ème siècle, déroule son autobiographie philosophique, entre le décor et les personnages de la ville et des champs.

Enfin, {1Prs, ..., 9Prs, aPrs, bPrs} est une chrestomathie de la presse égyptienne contemporaine; que nous avons découpée en veillant, dans la mesure du possible, à distinguer les sujets: affaires nationales ou étrangères; et les genres: éditoriaux ou nouvelles brèves.

## 2 Élaboration des textes et création de tableaux de correspondance

Les données analysées ont exactement le même format que dans [MOTS ARABES] (cité désormais ici:[MA]): il suffit donc de rappeler brièvement les principes adoptés dans ce premier article; en tenant compte de ce que l'on a aujourd'hui, sur Macintosh, (notamment avec Wintext, Teachtext, ou l'éditeur Nasher,) une représentation du texte arabe en octets différente de celle suivie par le premier logiciel 'Alkaatib'.

### 2.1 *Scriptio minima* et découpage en formes et locutions

Il est dit, dans [MA:§1.2], que, pour l'analyse des textes arabes, une élaboration totalement automatique, ne requérant lors de la saisie aucune interprétation grammaticale, doit être fondée, non sur la *scriptio plena* (comportant tous les signes diacritiques flottant de part et d'autre de la ligne; et, notamment, les voyelles brèves); mais, au contraire, sur une *scriptio minima*, d'où sont éliminées, avec les signes diacritiques, des marques de distinctions essentielles (telles que celles affectant les variantes de y sans deux points); si ces marques manquent constamment dans maintes éditions classiques; pour ne rien dire de la presse contemporaine.

De plus, notre *scriptio minima* conserve le découpage usuel de la phrase arabe: ainsi, dans le titre du présent article, l'expression “formes de mots et locutions” renvoie aux segments minimaux compris entre deux blancs: les conjonctions de coordination, و et ف (wa et fa), se lient au mot suivant à quelque catégorie que celui-ci appartienne; l'article défini se lie au nom; etc... Les locutions ainsi distinguées, ne le cèdent en rien aux formes isolées, prises exclusivement, dans l'étude d'autres langues, pour caractériser le style.

{ ʔ ä a b p t t ġ ċ ħ k d d r z ž s š š d t z c ġ f  
 á ġ k g l m n h w y }

À la *scriptio minima* romanisée sert une police ‘diacr’; laquelle, à des variantes mineures près (issues, pour la plupart, du dictionnaire de Hans WEHR), suit les normes de transcription phonétique internationale. (Les signes {ğ, g, p, ċ, ž} étant réservés pour des mots non arabes: berbères, persans, urdus...).

Voici, par exemple, dans cette transcription, la phrase citée plus haut du “*Tamis*”:

fṡhnä alnaqd aḡn, hy ġrblä alaṡar aladbyä

## 2.2 Tri alphabétique des occurrences

Dans [MA], par un programme ‘triarab’, on passe directement du fichier d’octets, créé par ‘Alkaatib’, à la suite ordonnée des formes, romanisée et étiquetée par versets (i.e. phrases ou paragraphes) et chapitres.

Pour les textes saisis par d’autres logiciels, et analysés ci-après, on a bénéficié d’une option ‘Texte seul’, qui, éliminant les variantes de tracé et les ligatures si diverses de la calligraphie arabe, (variantes et ligatures dont une partie subsiste dans le texte imprimé,) donne des fichiers dont la structure est très proche de la *scriptio minima*; et peut y être exactement réduite par un programme approprié, ‘dearab’, avec le code des lettres en octets propre à la police ‘diacr’.

D’après le fichier créé par ‘dearab’, un programme ‘triarap’, crée finalement une liste de formes qui se présente exactement comme celle que crée ‘triarab’ d’après un fichier d’Alkaatib.

Un programme ‘fus’ permet de fondre en une seule les listes afférentes à plusieurs textes; créant ainsi une liste ordonnée qui ne diffère en rien de celle qu’on obtiendrait à partir des textes mis bout à bout (; à supposer que ceux-ci soient saisis dans le même format).

Ainsi, on peut considérer que le corpus décrit au §1, (ou une partie de ce corpus), assimilé à un texte unique, est transformé en la suite, ordonnée alphabétiquement, de toutes ses occurrences de formes; répétées, chacune, autant de fois qu’on la trouve, avec les références, par chapitre et verset, se succédant dans l’ordre du texte.

Ces références ne sont, en fait, utilisées que suivant le découpage du corpus en un ensemble de 91 segments; eux-mêmes déjà décrits, plus haut, avec leurs sigles. On trouve dans [MA:§3.1] le fichier de commande qui définit le découpage du corpus des textes philosophiques en 13 fragments: on a procédé de même ici.

### 2.3 Dictionnaire et lexiques

Partant du fichier trié des occurrences, le programme 'qamus' crée un dictionnaire du texte, ou suite ordonnée alphabétiquement des formes, chacune accompagnée de sa fréquence. Le programme 'trimu§' range cette dernière suite dans l'ordre croissant des fréquences: c'est d'après le résultat de ce tri que l'on choisit, selon divers critères, un sous-ensemble  $\Delta$ , ou lexique, de formes destinées à être prises en compte dans l'analyse.

Enfin, le programme 'tridic' crée un tableau de correspondance croisant le lexique  $\Delta$  retenu; pris comme ensemble des lignes, avec l'ensemble J des segments:  $k(i, j)$  étant le nombre des occurrences de la forme  $i$  dans le segment  $j$ .

### 2.4 Enchaînement des analyses

Les textes et fragments dont nous disposons différant grandement quant à la taille et au genre, on a cru prudent de n'en approcher que par palliers l'analyse globale.

Dans une première analyse, rentrent les 13 fragments philosophiques de [MA]; et les 22 fragments découpés dans le Recueil, avec {mqmB, mqm1, mqm2} des deux  $m\bar{a}q\bar{a}m\bar{a}t$  (cf. *supra*, §1.2); soit un ensemble J1 de 38 fragments. Le lexique  $\Delta_1$ , de 120 formes, est choisi, avec un seuil de fréquence à 10, en veillant à éliminer tout ce qui, dans certains fragments, peut prendre valeur de mot plein, lié au contenu. L'axe 1 étant dominé par les poésies modernes de { $\mathbb{1}occ$ ,  $\mathbb{1}dpl$ }, avec  $CTR_1(\mathbb{1}Occ)=504\%$  et  $CTR_1(\mathbb{1}dpl)=99\%$ , on reprend l'analyse du tableau  $\Delta_1 \times J_1$  en mettant en supplément les deux colonnes { $\mathbb{1}occ$ ,  $\mathbb{1}dpl$ }: d'où, pour J1 ainsi réduit à J1', une représentation clairement interprétable; tant à l'analyse factorielle qu'à la classification: cf. *infra*, §3.1.

Les 45 fragments distingués dans l'ensemble des trois grands blocs de prose moderne, {Grb, Sgn, Prs} = GSP, ont fait l'objet d'une analyse,  $\Delta_2 \times GSP$ , fondée sur un lexique  $\Delta_2$  de 162 formes, choisies, également, avec 10 pour seuil de fréquence. Analyse factorielle et CAH distinguent nettement les trois blocs: cf. §3.2.

En adjoignant à J1', textes philosophiques et Recueil, les 8 fragments de {Rsl, Adb},  $r\bar{i}s\bar{a}l\bar{a}$  et  $a\bar{d}a\bar{b}$ , d' $i\bar{b}n-l-muq\bar{a}ff\bar{a}'$ , le corpus classique est au complet, pour un tableau  $\Delta_3 \times Cls$ , avec un lexique de 139 mots; le seuil restant à 10. Les trois genres de la philosophie, des belles-lettres (adab) et de la poésie sont bien vus d'après ce tableau: cf. §3.3.

Il semblait difficile de traiter ensemble le classique et le moderne, (i.e. Cls et GSP, ce qu'on notera C&M): car, d'une part, en descendant au seuil de 10, on est conduit à retenir un lexique  $\Delta_5$  comprenant 232 formes - ou agrégats - de mots outil; d'autre part, les trois blocs {Grb, Sgn, Prs}, chacun

homogène, relativement à la diversité du reste du corpus, pouvaient dominer l'analyse, imposant les deux premiers axes factoriels, et partageant la totalité des fragments entre trois pôles.

C'est pourquoi, concurremment à  $\Delta 5$ , on a considéré un lexique plus restreint,  $\Delta 4$ , comptant 141 formes ou locutions, toutes de fréquence  $\geq 20$ ; et, dans certaines analyses, on a réduit le poids des colonnes des blocs afférents à {Grb, Sgn, Prs}, en les multipliant par des coefficients; que, sans prétendre suivre une règle certaine, on a pris à (1/4) pour Grb, (1/2) pour Sgn et (3/4) pour Prs.

Au §4, on rend compte, avec plus ou moins de détails des résultats issus des tableaux  $\Delta 4 \times C\&M$  et  $\Delta 5 \times C\&M$ , repondérés ou non. Il nous suffit d'annoncer ici, en bref, que l'analyse et la CAH reconnaissent les principaux genres: philosophie, prose artistique ou poésie, roman et presse; mais avec, entre les variantes, des nuances qui nous font préférer  $\Delta 5 \times C\&M$  repondéré.

### 3 Analyses partielles

#### 3.1 Textes philosophiques et recueil de lecture

```
mots de  $\Delta 1 \times$  fragments de  $J1$  : textes philosophiques et recueil, avec { $\mathbb{f}occ$ ,  $\mathbb{f}dpl$ }
trace : 1.917e+0
rang : 1 2 3 4 5 6 7 8 9 10
lambda : 2498 1757 1459 1128 1011 906 823 782 687 654 e-4
taux : 1303 917 761 589 528 473 430 408 359 341 e-4
cumul : 1303 2220 2981 3570 4097 4570 5000 5407 5766 6107 e-4
```

Nous nous bornerons à rappeler que, comme on l'a dit au §2.4, l'analyse du tableau  $\Delta 1 \times J1$ ,  $120 \times 38$ , distingue les deux fragments { $\mathbb{f}occ$ ,  $\mathbb{f}dpl$ }, (plus exactement, deux cumuls de fragments) qui créent le 1-er axe et contribuent encore fortement au suivant. Avec { $\mathbb{f}occ$ ,  $\mathbb{f}dpl$ } vont les formes { $lst$ ,  $am$ ,  $knt$ ,  $any$ } qui évoquent le dialogue du poète avec un objet, animé ou non.

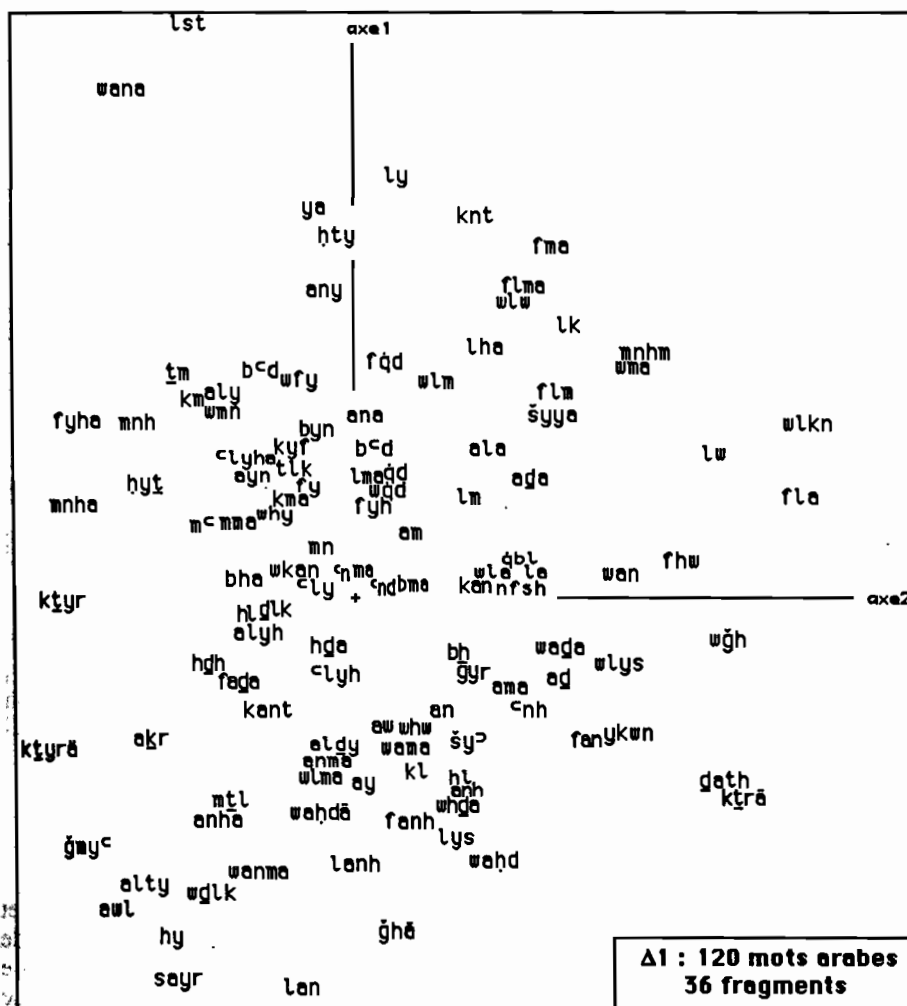
Comme il s'agit de poésies modernes, sans parenté étroite avec le reste du corpus, l'analyse est reprise sans celles-ci: soit un tableau  $\Delta 1 \times J1'$ ,  $120 \times 36$ .

```
mots de  $\Delta 1 \times$  fragments de  $J1'$  : textes philosophiques et recueil, sans { $\mathbb{f}occ$ ,  $\mathbb{f}dpl$ }
trace : 1.726e+0
rang : 1 2 3 4 5 6 7 8 9 10
lambda : 2094 1504 1094 1001 930 887 781 729 723 589 e-4
taux : 1213 871 634 580 539 514 452 422 419 341 e-4
cumul : 1213 2084 2718 3298 3837 4351 4803 5225 5644 5986 e-4
```

Sur la suite des valeurs propres et taux, les deux premiers axes se détachent nettement, la décroissance étant, ensuite, lente et régulière. [Relativement à l'analyse de  $\Delta 1 \times J1$ , l'inertie totale et les taux afférentes aux trois premiers axes diminuent; mais l'intervalle entre les axes 2 et 3 s'accroît]. C'est pourquoi on publie d'une part, les deux ensembles des mots et des fragments dans le plan (1, 2); et d'autre part la classification de chacun des deux ensembles en correspondance (effectuée dans l'espace des profils rapporté à la totalité des 35 axes extraits par l'analyse).







Passons dans le demi-plan ( $F1 > 0$ ). Se signalent les fragments de maqâmâ (séances...) vers l'extrémité de l'axe 1. Les trois cumuls, par parties, des notices du Recueil, {&Idπ, &Id¶, &Id§} sont dans le quadrant ( $F1 > 0, F2 < 0$ ). Sans être étroitement groupés, les cumuls de poésies par siècle, {πsc1, ..., πsc5}, sont dans le quadrant ( $F1 > 0, F2 > 0$ ), à l'exception de πsc4 ( $F2 \approx 0$ ). La prose du recueil, {§}, n'occupe pas de position excentrique (§Mqm excepté).

Mais parce que le plan (1, 2) ne rend compte que de 21% de l'inertie, il convient de chercher dans la classification une vue systématique du corpus fondée sur l'ensemble des associations entre fragments et formes.

c	Partition de J1' en 13 classes : Sigles des fragments de la classe c
1	isf†
59	¶jbr rZt3
13	rZt4
49	rZt2 rZt1
31	§sfa
52	§rbh ¶fGh ¶shu ¶Arb §Hdt §KwD Stah Stwh
48	¶fHy Rmq2 Rmq1 Rmq3
56	ghz2 ghz1
4	ghz3
17	πsc1
57	πsc5 πsc2 ¶Lac πsc4 §Mqm πsc3
44	&Id§ &Idπ &Id¶
58	§jhz mqm2 mqm1 mqmB

isf†	F1--	F4--	67	F1< F2< F3>	70
59	F1-	F2-	F3+	64	
rZt4	F1---	F2---	F3++	62	
49	F1-	F2-			
§sfa	F1+	F2+	F4---	65	66
52	F1+			60	
48	F1-	F2+			
56	F1-	F2+++		63	
ghz3	F1-	F2++++	F3---		
πsc1	F1+++	F2+++	F3+++	F4--	61
57	F1+++	F2+	F3++		F2> 69
57	F1+++	F2+	F3++		F1>
44	F1+	F2---	F3--		68
44	F1+	F2---	F3--		F2<
58	F1+++	F2-	F4+++		

ci-dessous: étiquetage sommaire d'après les classes de formes

64: 198++++ 229++ ; 1: 228++++ ; 65: 188++ 229++ ; 63: 185++++ 230++ ;  
61: 223++++ 226+++ ; 44: 225++++ ; 58: 217++++ ;

Au sommet de la hiérarchie, le corpus se sépare en deux branches: j70 et j69. L'étiquetage en terme de facteurs montre que j69 est caractérisée par de fortes valeurs de  $F1 > 0$ . On a, dans j61, du côté ( $F2 > 0$ ) les 5 siècles de poésie classique, avec la traduction du ¶Lac; et §Mqm, la maqâmâ sur les poètes. Le reste, j69, constitue la classe j68, ( $F2 < 0$ ); se partageant, à un haut niveau, en, d'une part, j58: deux autres maqâmâ et §jhz, الجاحظ, le deuil du savant courtois; et, d'autre part, j44: les notices du recueil, {&Idπ, &id¶, &Id§}.

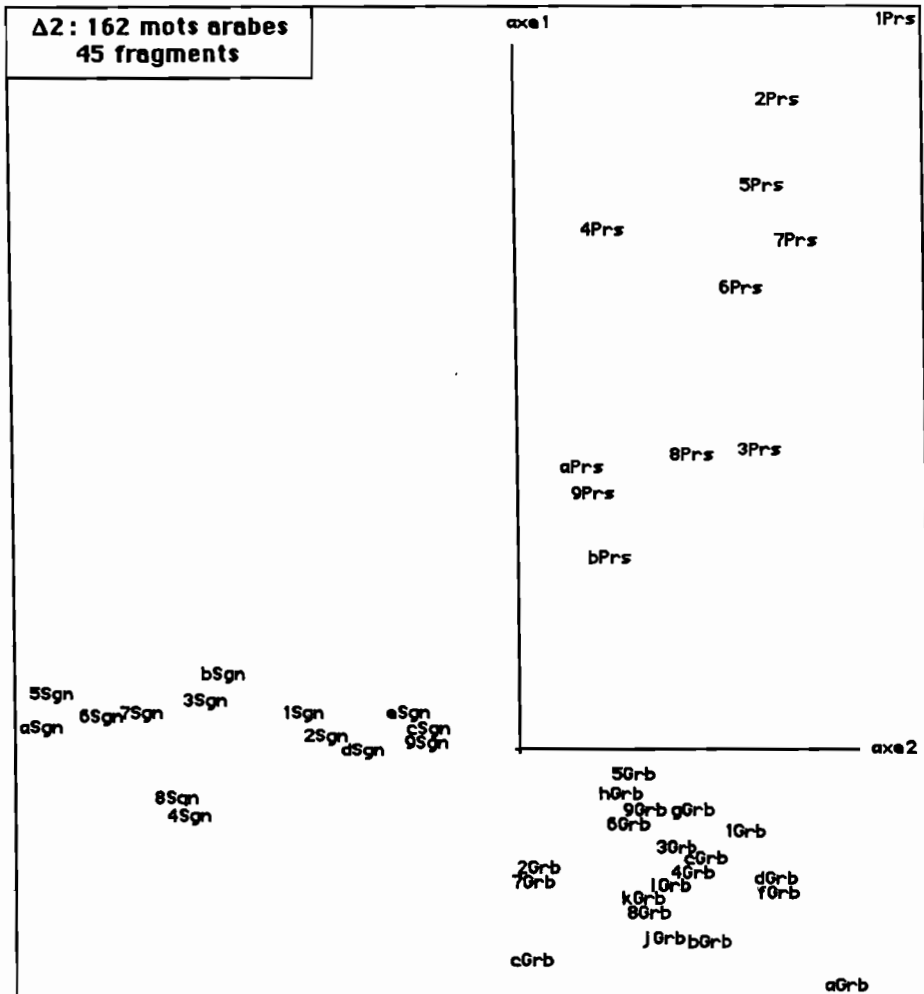
La branche j70 comprend j67 et j66. Dans j66, on a, d'une part, j63, Ghazâlî; et, d'autre part, j65: avec j48, philosophie sans appareil technique, et j52, la prose artistique (maqâmâ exceptée) avec des œuvres poétiques modernes: ¶shu, le Prince des poètes; et ¶Arb, à la gloire de la langue arabe.

Enfin j67 comprend les dissertations formelles de philosophie, rZt, ou de théologie, isf† (texte chrétien expliqué dans [MA]). À ces fragments, dont

c	Partition de Δ1 en 16 classes : Sigles des formes de la classe c	
221	akr	aw kant alyh clyh cly ma fanh c'n
202	why	hgh hga mtl sayr anha wdlk lan
222	ay	wlma hw aldy bl whw wama anh ad whda lys an nfsh lanh   wkan kan bh bma c'nd
-----		
198	awl	ktyrā ġmy <sup>c</sup> alty wanma hy
-----		
211	šy <sup>p</sup>	ğyr fan kl wahdä
201	ktyr	m <sup>c</sup> ġhā wahd
-----		
188	mna	flm šyya
-----		
185	wada	fhw dath ktrā c'nh
-----		
93	hl	
224	wlys	lh ykwn b <sup>c</sup> d tlk wan lm qd fyh ala la
-----		
223	wğh	wlkn lw ama fla wla km wlv wlm fma ada
-----		
193	am	knt faga any
218	mnhm	flma qbl lk lma ly lha wma
-----		
217	clyha	im lst wana ya ana ayn hty
-----		
215	anma	byn b <sup>c</sup> d kma fy kyf wmn fqd mn
220	bha	wqd dlk mnh hyt wfy mnha fyha aly
-----		
221	F1-	≈CdG 229 232235
202	F1--	F2--
222	F1--	F2+
198	F1---	F2---- F3++
211	F1---	F2+ 228 231
201	F1---	F2--
188	F1++	F2+ F4----
185	F1-	F2++++ F3---- 234 237 238
93	F1-	F2-- F3+++ 230
224	F1+	F2++
223	F1++	F2+++ F3+++ 233
193	F1++	226 F3+++
218	F1+++	F2++
217	F1++++	F2-- 236
215	F1+	F2- 225
220	F1++	F2---

l'association a déjà été signalée dans le demi-plan (F1<0), s'agrège fjjbr, poésies de Gibran: (F1>0), mais proximité avec rZt3 sur d'autres axes.

Quant à la CAH du lexique Δ1 des formes, nous nous bornerons, en attendant d'autres éléments de comparaison, à la soumettre au lecteur arabisant; en suggérant de l'examiner avec le plan (1, 2) afin de reconnaître les affinités entre textes et locutions; sans négliger l'étiquetage qui signale, entre autres, un lien entre j63, Ghazâlî, et i185.



### 3.2 Trois blocs de prose moderne

mots de Δ2 × fragments de J2: textes modernes en prose: GSP = {Grb, Sgn, Prs}  
 trace : 1.487e+0  
 rang : 1 2 3 4 5 6 7 8 9 10  
 lambda : 1888 1255 828 715 632 580 530 483 441 413 e-4  
 taux : 1269 844 557 481 425 390 356 324 296 278 e-4  
 cumul : 1269 2113 2670 3151 3576 3966 4322 4646 4943 5220 e-4

Deux œuvres, Grb et Sgn, cf. *supra* §1.3, et une chrestomathie de la presse, Prs: les trois blocs se voient bien dans le plan (1, 2); et la Classification Ascendante Hiérarchique les distingue sans aucune confusion.

c	Partition de GRB en 13 classes : Sigles des fragments de la classe c
73	1Grb gGrb fGrb kGrb 3Grb iGrb 2Grb
72	eGrb 5Grb 6Grb jGrb cGrb
4	4Grb
-----	
9	9Grb
76	bGrb hGrb 7Grb 8Grb dGrb
-----	
10	aGrb
-----	
22	2Sgn
75	cSgn bSgn 6Sgn 4Sgn eSgn dSgn 9Sgn 3Sgn
66	1Sgn aSgn 8Sgn 7Sgn 5Sgn
-----	
45	bPrs
50	8Prs 9Prs aPrs
37	3Prs
-----	
77	4Prs 6Prs 7Prs 5Prs 2Prs 1Prs
-----	
73	79 81 86 F1<0 F2>0 88
72	F3=0   86: F3≤0
4	
9	83 85   85: F3>0
76	
10	
22	80 F1=0 F2<0
75	78
66	
45	84 F3<0 87 F1>0 F2>0
50	82
37	
77	F3>0

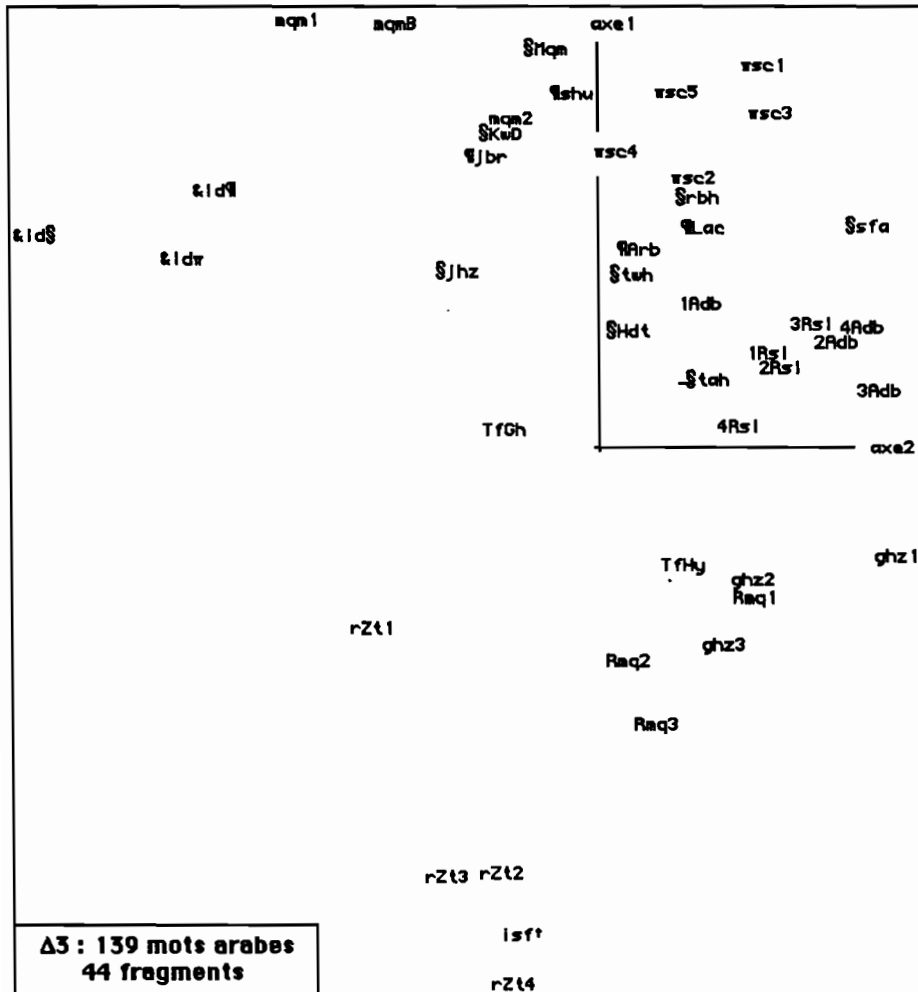
Il apparaît ainsi que la différence de style entre les deux œuvres, Grb et Sgn, se conserve au fil des chapitres; tandis que, malgré la diversité des thèmes abordés, les articles des journaux n'approchent guères de la littérature proprement dite.

Le bloc de la presse, j87, est partagé en deux; division qu'on retrouvera, (voire un peu modifiée,) dans toutes les analyses globales du §4. La classe j77, la plus écartée vers (F1>0) dans le quadrant (F1>0, F2>0) du plan (1, 2), comprend quasi exclusivement des nouvelles relatives à l'activité du gouvernement égyptien. La classe j84, moins écartée que j77 sur (F1>0) et qui s'oppose, de plus, à celle-ci, sur l'axe 3, comprend, au contraire, des nouvelles sur l'étranger: USA, 3Prs; ou Liban, Palestine, Israel.

À propos du lexique Δ2, composé de 162 formes, nous citerons seulement un détail: la CAH distingue une classe de 4 formes:

{ yqwm cqb m<sup>c</sup> klal } { يقوم، عقب، مع، خلال }

{ fait (ou: effectue), après, avec, durant }, fortement associée à j77; et caractérisant cette classe relativement à j84.



### 3.3 Philosophie, recueil de lecture et prose classique

mots de  $\Delta 3 \times 44$  fragments de CLS: textes classiques: J1' et {Rsl, Adb}  
 trace : 1.774e+0  
 rang : 1 2 3 4 5 6 7 8 9 10  
 lambda : 1871 1516 1113 906 861 818 708 678 606 548 e-4  
 taux : 1054 854 627 510 485 461 399 382 342 309 e-4  
 cumul : 1054 1909 2536 3047 3532 3993 4392 4774 5116 5425 e-4

Adjoindre à J1', (textes philosophiques et Recueil) les 8 fragments de {Rsl, Adb}, risālā et adab, d'ibn-l-muqaffa', modifie peu les agrégats que montrent le plan (1, 2) et la CAH; mais le schéma global de la littérature classique nous paraît meilleur.

c	Partition de Cls en 15 classes : Sigles des fragments de la classe c
1	isft
56	TfHy Rmq2 Rmq1 Rmq3
71	ghz3_ghz2 ghz1
13	rZt4
12	rZt3
59	rZt2 rZt1
31	§sfa
66	§tah §twh §KwD ¶Arb §Hdt 1Adb ¶shu §rbh TfGh
72	3Adb 4Adb 1Rsl 2Adb 4Rsl 2Rsl 3Rsl
68	&Id¶ &Idπ &Id§
70	§jhz mqm2
62	mqm1 mqmB
23	¶jbr
73	πsc5 ¶Lac πsc4 πsc2 §Mqm πsc3
17	πsc1

1	F5>0	82_F4<0	84	F1<0
56		79_F2>0		
71				
13		80_F2<0		
12		77		
59				
31	F5>>0	81	F3<0	86_F1>0
66		76		
72				
68		F2<<0	F4<0	F3<0
70		74_F6>0	83_F3>0	85
62				
23		78_F6<0		
73		75		
17				

Au sommet, dans j84, la philosophie, avec tous ses genres, se sépare du reste: j86. L'ensemble des belles-lettres en prose, confirmé par 8 fragments, forme une classe j81, agrégée à un bas niveau. Comme au §3.1, ¶shu, le Prince des poètes; et ¶Arb, à la gloire de la langue arabe, vont avec cette prose; où il est juste de trouver TfGh, Ghazâlî vu par Ibn Tofayl.

Du reste, j85, les notices {&Id} se séparent à un niveau élevé. Il n'y a plus, dans j78, que maqâmä et poésie classique; et ¶jbr, Gibran, et §jhz, الجاحظ, le deuil du savant courtois, sont maintenant mieux placés qu'au §3.1.

Enfin, sans entrer dans les détails, il vaut la peine de relever, dans l'étiquetage de l'arbre, que des facteurs de rang élevé contribuent à distinguer certaines classes. Notamment, maqâmä et poésie classique ne se séparent que suivant l'axe 6; et des fragments originaux, isft, §sfa, sortent sur l'axe 5.

#### 4 Analyses de l'ensemble des textes classiques et modernes

##### 4.1 Analyses fondées sur un lexique restreint de 141 formes

En prenant un lexique restreint,  $\Delta_4$ , ne comptant que 141 formes ou locutions, toutes de fréquence  $\geq 20$ , on pensait éviter de mettre l'accent sur les particularités peu fréquentes d'un texte. En pondérant les blocs {Grb, Sgn, Prs}, les coefficients étant (1/4) pour Grb, (1/2) pour Sgn et (3/4) pour Prs, on voulait réduire la force attractive de ces blocs pour les autres fragments.

Sans présenter de graphique plan, nous apprécierons, d'après la CAH des fragments, la vue globale du corpus offerte par l'analyse. [Dans les tableaux de partition, des subdivisions interprétables sont marquées entre parenthèses.]

```
141 mots de  $\Delta_4 \times 89$  fragments de C&M ; lexique au seuil 20 ; corpus non repondéré
trace : 1.667e+0
rang : 1 2 3 4 5 6 7 8 9 10
lambda : 1411 1127 881 808 583 534 479 458 437 408 e-4
taux : 846 676 528 485 350 320 288 275 262 245 e-4
cumul : 846 1522 2050 2535 2885 3205 3493 3767 4030 4274 e-4
```

c	Partition de C&M en 12 classes : Sigles des fragments de la classe c
158	isft ghz2
164	ghz1 ghz3 TfHy Rmq2 Rmq3 Rmq1
162	rZt2 rZt1 rZt3
13	rZt4
165	4Grb 1Grb gGrb fGrb eGrb 3Grb kGrb dGrb bGrb ¶jbr cGrb hGrb 7Grb ¶Lac   $\pi$ sc4 9Grb 8Grb aGrb 2Grb iGrb jGrb 5Grb 6Grb Grb + ...
153	SHdt SKwD 3Adb 1Rsl 4Adb 1Adb 2Adb 4Rsl 2Rsl 3Rsl
166	(bPrs 3Prs 8Prs 9Prs aPrs) (&Id $\pi$ &Id¶ &Id\$) notices &Id + Presse(international)
161	$\pi$ sc1 $\pi$ sc3 $\pi$ sc5 SMqm
151	mqm2 mqmB ¶shu mqm1
31	Ssfa
163	2Sgn eSgn Srbh TfGh ¶Arb 1Sgn 9Sgn \$tah 4Sgn dSgn \$twh cSgn $\pi$ sc2 \$jhz   5Sgn aSgn bSgn 6Sgn 7Sgn 3Sgn 8Sgn Sgn + .....
144	4Prs 6Prs 7Prs 5Prs 2Prs 1Prs
	Presse (nouvelles nationales)
158	_____170_____172_____toute la philosophie_____176_.
164	_____170_____172_____176_____
162	_____169_____
rZt4	_____169_____
165	_____175_____
153	_____174_____
166	notices et Presse int _____173_____
161	_____168_____171_____
151	_____168_____171_____
Ssfa	_____167_____
163	_____167_____
144	_____167_____





#### 4.2 Analyses fondées sur un lexique étendu de 232 formes

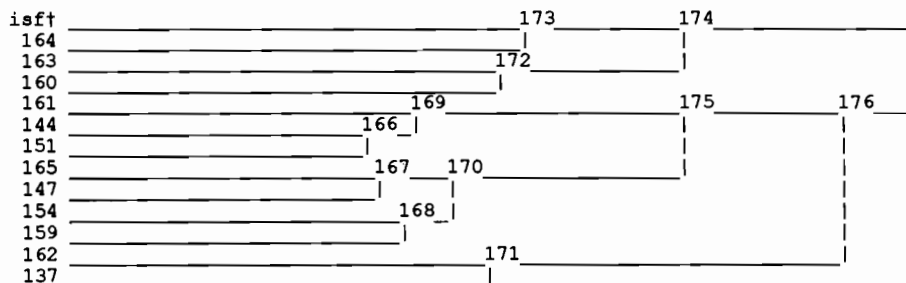
Afin de faciliter la comparaison avec l'autre lexique, nous considérons d'abord les classifications de l'ensemble C&M des fragments, classiques et modernes, fondées sur les tableaux brut ou repondéré.

La structure d'ensemble de la presse est bien vue dans les deux cas (à ceci près qu'avant pondération, manque le fragment 3Prs); les notices du Recueil allant avec les nouvelles internationales, plus concises. Après pondération, la presse s'oppose même à tout le reste, en formant une branche.

Nous savons qu'il y a, en philosophie deux niveaux: le plus rigoureux étant dans rZt; l'autre, dans {Rmq ghz}. Avant pondération, le tout forme une branche 174, où se trouvent également deux traités en belle forme, {Adb Rsl}.

232 mots de  $\Delta 5 \times 89$  fragments de C&M ; lexique au seuil 10 ; corpus non repondéré  
 trace : 2.329e+0  
 rang : 1 2 3 4 5 6 7 8 9 10  
 lambda : 1529 1314 1045 944 707 640 622 580 567 541 e-4  
 taux : 657 564 449 405 304 275 267 249 243 232 e-4  
 cumul : 657 1221 1669 2075 2378 2653 2920 3169 3412 3645 e-4

c	Partition en 13 classes : Sigles des fragments de la classe c
1	isft
164	rZt4 rZt3 rZt1 rZt2      commentaire d'Averrhoès sur le livre Z de la Métaphysique
163	Rmq1 Rmq3 Rmq2 ghz3 ghz1 ghz2      Traité décisif et Effondrement des Philosophes
160	1Adb 3Adb 3Rsl 1Rsl 2Rsl 4Rsl 4Adb 2Adb      adab el kabîr et risâlâ
161	fGrb 9Grb 5Grb 6Grb 3Prs kGrb 3Grb iGrb gGrb hGrb jGrb lGrb eGrb 4Grb
144	fjbr aGrb
151	fArb 7Grb 8Grb dGrb cGrb bGrb
165	\$sfa 2Sgn TfHy 2Grb \$Hdt \$tah 9Sgn fshu \$twh \$KwD dSgn \$rbh TfGh cSgn 4Sgn 1Sgn eSgn 3Sgn
147	\$jhz 6Sgn 8Sgn bSgn aSgn 5Sgn 7Sgn
154	$\pi$ sc1 fLac $\pi$ sc3 $\pi$ sc4 $\pi$ sc2 $\pi$ sc5 \$Mqm
159	mqmB mqm1 mqm2
162	&Id $\pi$ &Idf &Id\$ bPrs 8Prs 9Prs aPrs      manque 3Prs!
137	6Prs 4Prs 5Prs 7Prs 2Prs 1Prs



Après pondération, {isft rZt} s'isolent dans j174, au sein de la branche j176, complémentaire de celle de la presse; tandis que, dans j173, {Rmq ghz} se retrouvent, associés à {Adb Rsl}; auxquels s'agrègent, opportunément, dans j162, divers fragments de prose classique.

La poésie classique forme une subdivision, comprenant le ¶Lac et §Mqm, le débat de lettrés; les autres maqâmâ en étant proches.

Enfin les blocs de prose moderne, {Grb Sgn}, sont reconnus; mais, avant pondération, s'agrègent à Sgn de multiples fragments qui trouvent ensuite une place sans doute meilleure au sein de l'adab (belles lettres classiques).

232 mots de  $\Delta 5 \times 89$  fragments de C&M ; lexique au seuil 10 ; corpus repondéré

```

trace : 2.514e+0
rang : 1 2 3 4 5 6 7 8 9 10
lambda : 1789 1538 1136 934 849 769 724 692 650 604 e-4
taux : 712 612 452 372 338 306 288 275 258 240 e-4
cumul : 712 1323 1775 2147 2484 2790 3078 3353 3612 3852 e-4

```

c	Partition de C&M en 13 classes : Sigles des fragments de la classe c		
1	isft	443+++	352++++
13	rZt4	443++	448++++
159	rZt3 rZt1 rZt2		448++++
147	Rmq1 Rmq2 Rmq3		451++++
156	ghz3 ghz1 ghz2	443++++	451++
31	\$sfa		412+++++
162	1Adb \$tah TfHy \$KwD \$twh ¶shu 2Rsl ¶Arb \$rbh TfGh 3Adb 3Rsl 1Rsl 4Adb 4Rsl 2Adb		451+ 450++
163	¶sc1 ¶Lac ¶sc3 ¶sc4 ¶sc2 ¶sc5 §Mqm		447+++++
165	mqm2 mqmB \$Hdt mqm1		445+++++
161	¶jbr aGrb 7Grb 8Grb dGrb cGrb bGrb fGrb 9Grb 5Grb 6Grb jGrb 1Grb 2Grb 3Grb gGrb iGrb kGrb hGrb eGrb 4Grb		449+++++
157	\$jhz 4Sgn cSgn 1Sgn 9Sgn eSgn 3Sgn dSgn 2Sgn 6Sgn 8Sgn bSgn aSgn 5Sgn 7Sgn		420+++++
164	(bPrs 3Prs 8Prs 9Prs aPrs) (&Id\$ &Idπ &Id¶)		441++++
136	6Prs 4Prs 5Prs 7Prs 2Prs 1Prs	439+++++	441++++

isft	174	F1<<_F2<<_F3<<	176
rZt4	166	livre Z	
159			
147	traité décisif	F2<	167
156	tahafot		173
\$sfa		prose	F1<_F3>_175
162	adab etc...	classique	F2>0
163	poésie classique ¶Lac §Mqm		(sauf 147)
165	deux maqâmâ et hadith		170
161	Grb avec Gibran		169
157	Sgn avec jahiz		
164	Presse étr et notices du recueil		172
136	Presse (national)	F1>>_F2<<	

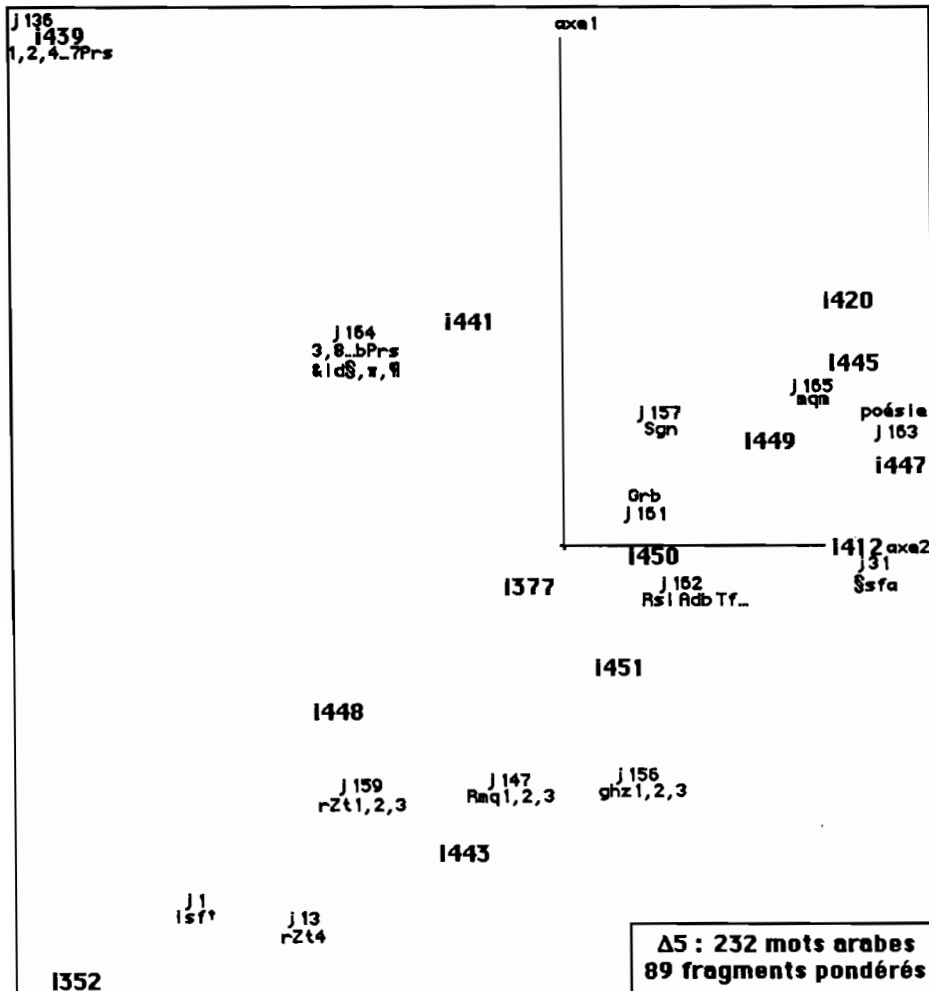
c	Partition de $\Delta 5$ en 13 classes : formes de la classe c
449	abda m <sup>c</sup> h fhm ant <sup>c</sup> ny wnhn lna fwa <sup>c</sup> whkda kla hwla lkm am anna   <sup>c</sup> ndna anfsna swy nhn lkn dak wbyn lyst <sup>c</sup> ndkm whm hl
420	agn <sup>c</sup> ndma wknt alyna kdik tht knt mra ly mnd fada by nfsy ana
447	bynna lqd wkn lw fla ega b <sup>c</sup> d wma lma flma <sup>c</sup> lyk lk km <sup>c</sup> ndy   wlv wfy wlm fma fqd
445	alan kyf kna kda fh <sup>c</sup> l wkyf ya lst ldlk <sup>c</sup> lyha tm ma <sup>c</sup> wana bkl   twl ana ayn hty
412	šyya bšy <sup>c</sup> flm mma lhm kanwa <sup>c</sup> lyhm
450	dlk mnh <sup>c</sup> ndh ga bma wkl ffy lkl fmn <sup>c</sup> nd wan fyh ala mnhm fyma fan   dayma alkbyr amam wkant why lha kan kant tik ykn lm dwn qd <sup>c</sup> ly mn   wla ma awla aqr aw <sup>c</sup> lyna mna tkn alağ fanh <sup>c</sup> nha alyh w <sup>c</sup> ly <sup>c</sup> lyh <sup>c</sup> n
377	kbyrā aldyn bynhm gmy <sup>c</sup> ahd nfsh
451	afdl fkan fhw alyha <sup>c</sup> nh waga la wlys lh nfs bd ykwn bhgh bhm   lmn bgyr hyn lanh hkda wama fama kha bdlk bh swah hm anma qbl   whgh anhm bha wahda lys gy <sup>c</sup> r whw aldy klh an
352	fhy kašā ghā wahdā wmnhm
443	wqt bl whga bnfsh eđ an <sup>c</sup> h fana sayr alšy <sup>c</sup> wgh šy <sup>c</sup> alaw wahd ama kl
448	wkan alty wanh bhga hga hgh wdlk fqt lan anha wgyr ml aktr awl   wanm hy wma ay ktyrā hw
441	ktyr mnha ktyra fyha aly nsbā hyt wmn wqd b <sup>c</sup> d fy
439	wkašā klal ban hnak byn kma hwl m <sup>c</sup>

449	453	455	459	461	462
420					
447					
445					
412		457			
450		456			
377	454				
451					
352			458		
443	452				
448					
441			460		
439					

Nous publions ci-dessus le tableau du lexique  $\Delta 5$ , soumis à la CAH d'après l'analyse du tableau pondéré,  $\Delta 5 \times C\&M$ .

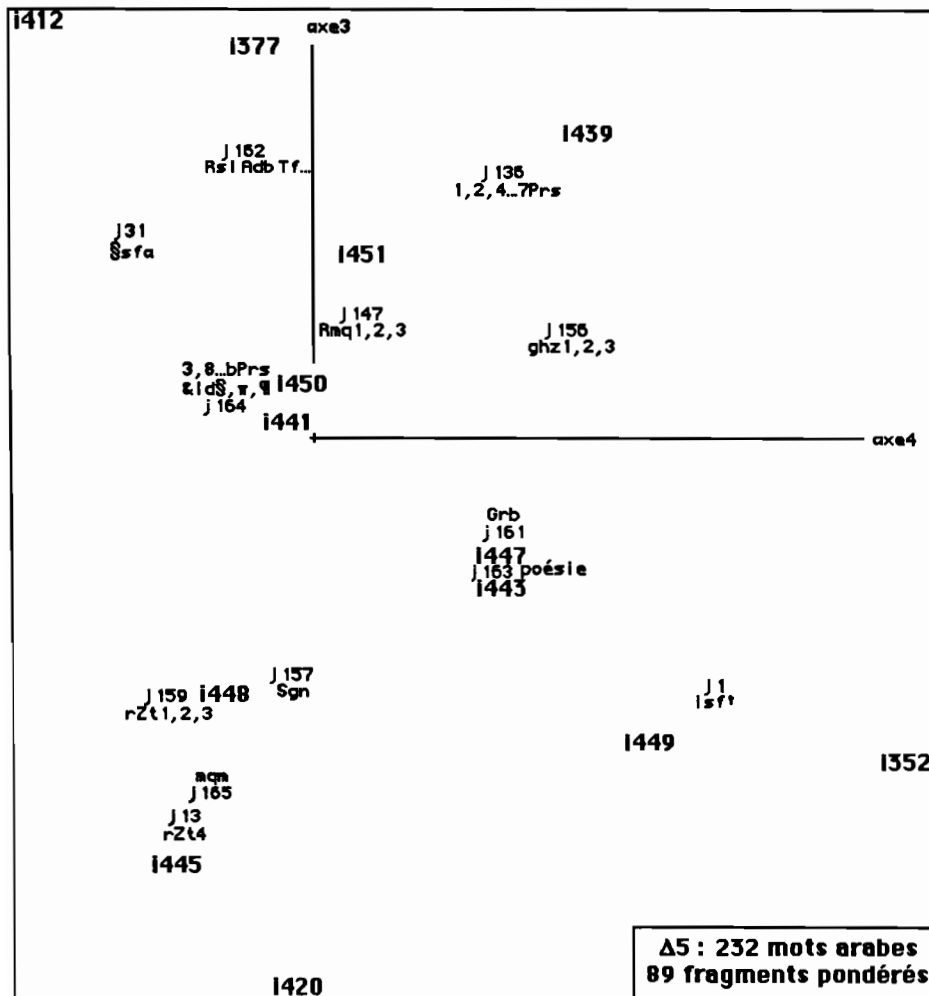
Les 13 classes de la partition retenue ont servi à étiqueter la CAH des 89 fragments; et leurs centres sont marqués, dans les plans (1, 2) et (3, 4), avec ceux des classes de fragments. Le lecteur arabisant appréciera par lui-même les associations ainsi manifestées entre formes ou locutions, d'une part; et genres littéraires, d'autre part.



Dans le plan (1, 2), le nuage des 89 fragments présente un amas, à peu près compris dans le quadrant ( $F1 > 0$ ;  $F2 > 0$ ) mais voisin du demi-axe ( $F2 > 0$ ), et deux pointes; dont l'une, dessinée par la philosophie, tend vers ( $F1 < 0$ ;  $F2 < 0$ ), isf† étant le fragment le plus écarté; et l'autre, située dans le quadrant ( $F1 > 0$ ;  $F2 < 0$ ), comprend la Presse et les notices brèves du Recueil; la position extrême étant celle des articles concernant l'activité des ministres égyptiens.

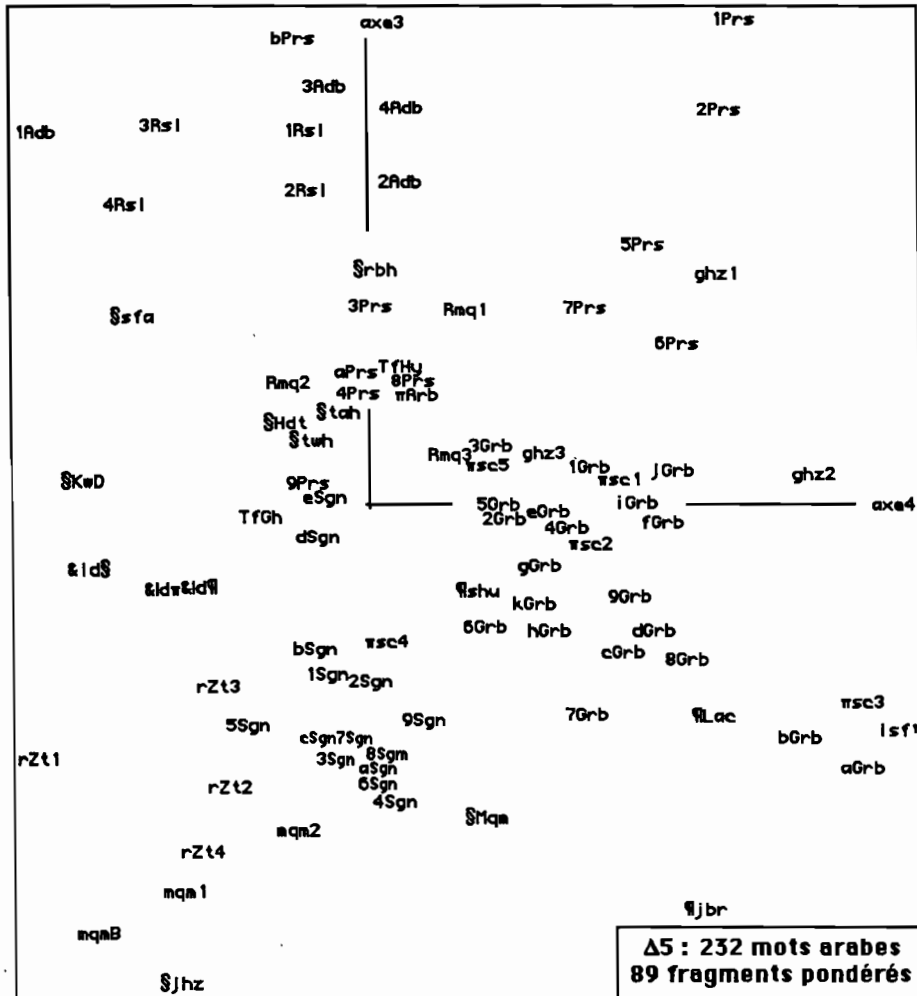
Au sein de l'amas, le facteur  $F2$  est maximum pour la poésie classique (j169); tandis que l'adab (j162) et les œuvres modernes (Grb;j161, Sgn;j157) se projettent sur le demi-axe ( $F2 > 0$ ) au voisinage de l'origine.

Soit un triangle {philosophie poésie presse}, dont le côté ( $F1 < 0$ ) est vide.



Dans le plan (3, 4), se disperse ce qui formait l'amas du plan (1, 2): j162 (adab) et j167 (Sgn) s'opposent nettement suivant l'axe 3, dans le demi-plan ( $F4 < 0$ ); Grb (j161) et poésie (j163) sont dans le quadrant ( $F3 < 0, F4 > 0$ ). Plus précisément, sur le nuage des fragments individuels, on voit que j163 se distribue dans toutes les directions autour de son centre; tandis que les sections qu'on a distinguées au sein de Grb dessinent plutôt une pointe, vers ( $F3 < 0, F4 > 0$ ).

Reste le cas des poésies modernes { $\{occ, \{dpl\}$  que, depuis le §3.1, on a mises en supplément comme perturbant l'analyse: par le programme 'discr', d'analyse discriminante, on affecté chacun ces points d'une part au centre de



classe le plus proche, d'autre part au fragment individuel le plus proche. Soit:

( $\text{\textcircled{occ}} \rightarrow \text{mqmB}$ )( $\text{\textcircled{dpl}} \rightarrow 8\text{Sgn}$ ) ; ( $\text{\textcircled{occ}} \rightarrow \text{j157}$ )( $\text{\textcircled{dpl}} \rightarrow \text{j157}$ ) ;

la distance d'affectation étant moins forte pour  $\text{\textcircled{dpl}}$  que pour  $\text{\textcircled{occ}}$  (très excentrique). [On se souviendra que  $\text{j157}$  s'identifie à  $\text{Sgn}$ .]

Nous appellerons enfin l'attention sur les tableaux de valeurs propres qui illustrent le §4. Lorsqu'on passe du lexique  $\Delta 4$  à  $\Delta 5$ , la trace et les valeurs propres augmentent nettement: ce qui atteste un contraste plus prononcé. La pondération augmente également l'inertie, mais moins fortement.

#### 4 Conclusions et perspectives

La présente étude, fondée sur un corpus plus étendu et plus varié que celui considéré dans [MOTS ARABES], confirme qu'une stylométrie de l'arabe peut être fondée sur le dénombrement automatique des formes ou locutions (définies comme suites de caractères délimitées par des blancs ou des signes de ponctuation).

Nous espérons que l'occasion nous sera offerte de traiter un corpus beaucoup plus étendu; comprenant un grand nombre d'œuvres saisies dans leur intégralité; et non seulement de brèves anthologies ou des textes reçus sans avoir été systématiquement choisis.

#### Références bibliographiques

[MOTS ARABES]: "Sur l'étude des textes arabes d'après les occurrences des formes de mots", in *CAD*, Vol. XIX, n°1, pp. 65-84; (1994).

*Textes et Poèmes choisis pour l'épreuve de récitation*; Vocalisation, présentation et commentaire de Ayadi CHABIR et Jean TARDY; recueil distribué par la Section d'Arabe Littéral de l'Institut National des Langues et Civilisations Orientales (INALCO) ; Centre d'Asnières; (1995).

*Documents pour l'étude et l'histoire de la langue arabe* ; 1, prose; par Ayadi CHABIR , INALCO, Section d'Arabe; (1995).