ministère de l'éducation nationale, de l'enseignement supérieur et de la recherche



en ligne en ligne

# AnIsl 29 (1995), p. 283-311

#### Christian Gaubert

Analyse morphologique d'un texte arabe par ordinateur: méthode d'évaluation, résultats.

#### Conditions d'utilisation

L'utilisation du contenu de ce site est limitée à un usage personnel et non commercial. Toute autre utilisation du site et de son contenu est soumise à une autorisation préalable de l'éditeur (contact AT ifao.egnet.net). Le copyright est conservé par l'éditeur (Ifao).

#### Conditions of Use

You may use content in this website only for your personal, noncommercial use. Any further use of this website and its content is forbidden, unless you have obtained prior permission from the publisher (contact AT ifao.egnet.net). The copyright is retained by the publisher (Ifao).

#### **Dernières publications**

9782724710922	Athribis X	Sandra Lippert
9782724710939	Bagawat	Gérard Roquet, Victor Ghica
9782724710960	Le décret de Saïs	Anne-Sophie von Bomhard
9782724710915	Tebtynis VII	Nikos Litinas
9782724711257	Médecine et environnement dans l'Alexandrie	Jean-Charles Ducène
médiévale		
9782724711295	Guide de l'Égypte prédynastique	Béatrix Midant-Reynes, Yann Tristant
9782724711363	Bulletin archéologique des Écoles françaises à	
l'étranger (BAEFE	")	
9782724710885	Musiciens, fêtes et piété populaire	Christophe Vendries

© Institut français d'archéologie orientale - Le Caire

# ANALYSE MORPHOLOGIQUE D'UN TEXTE ARABE PAR ORDINATEUR : MÉTHODE D'ÉVALUATION, RÉSULTATS

Les possibilités d'analyse morphologique automatique de l'arabe et d'extraction de racines mises en évidence par Claude Audebert et André Jaccarini <sup>1</sup> dans leurs travaux récents ouvrent de nombreuses perspectives dans le domaine du traitement automatique de textes arabes. La force de la démarche adoptée réside autant dans le respect des ambiguïtés naturelles et graphiques de la morphologie que dans l'efficacité de l'algorithmique employée. Une grammaire morphologique peut être considérée comme un objet susceptible de modifications et donc de mises au point et d'affinements : ce sont les variations de grammaire <sup>2</sup>.

La morphologie du nom trilitère retient l'attention par la position prédominante de cette catégorie de mot par rapport aux autres, que ce soit les verbes, les tokens, au sens défini par Audebert et Jaccarini <sup>3</sup>, ou les autres noms, quadrilitères, bilitères ou propres. Une première grammaire <sup>4</sup> non déterministe <sup>5</sup> évoluée, baptisée G2, a permis de reconnaître la plupart des obstacles créés par la prise en compte de l'ambiguïté, celle du modèle programmé comme celle due à la représentation graphique de la morphologie arabe. Nous proposons ici une méthode d'évaluation et de développement des grammaires morphologiques fondée sur l'analyse de textes. Cette méthode a pu être mise au point grâce à l'élaboration préalable d'un logiciel d'analyse aisément transformable en un module linguistique indépendant inclus dans un programme plus général.

- 1. Cette étude s'inscrit dans la suite logique des travaux publiés depuis 1987 par Claude-F. Audebert et André Jaccarini. Elle en exploite la terminologie et les résultats principaux. Nous nous contenterons donc de rappels généraux sur le fonctionnement des grammaires formelles sans en détailler les fondements théoriques.
- 2. Ce concept a été exposé par Audebert et Jaccarini dans « Méthode de variations de grammaire et
- algorithme morphologique, vers un extracteur de racine en arabe » in *BEO* XLVI, 1994, p. 77-91.
- 3. Cf. *idd*. « De la reconnaissance des mots outils et des token » in *AnIsl* XXIV, 1988, p. 269 n. 2.
  - 4. Cf. BEO XLVI, p. 82-86.
- 5. Est non déterministe une grammaire qui peut présenter au cours de l'analyse du mot plusieurs possibilités qui devront être explorées systématiquement.

#### 1. LA GRAMMAIRE FORMELLE Gn3: UNE DESCRIPTION SOUPLE DES NOMS À RADICAUX TRILITÈRES.

Une grammaire représentée par un automate peut être décrite comme un ensemble d'états dont un état initial  $E_i$  et un état final  $E_f$ . Les variables d'entrée sont des mots ou suites de caractères pris dans un ensemble désigné par *alphabet*. Chaque état constitue une étape de l'analyse du mot, l'étape suivante étant atteinte par la lecture d'un caractère selon une catégorie linguistique précise : ce sont là les transitions ou règles de réécriture, aisément représentables par arcs orientés dans un diagramme. Chaque règle ou arc correspond à une catégorie linguistique précise à l'intérieur du mot. Des transitions particulières appelées *epsilon-transitions* permettent le passage direct d'un état à un autre lorsque la ou les catégories intermédiaires ne figurent pas dans le mot.

Si le dernier caractère du mot analysé emprunte un arc qui peut mener directement à l'état final, le mot est accepté et donc conforme à la grammaire. Notons qu'il peut l'être pour plusieurs raisons, en vertu de l'indéterminisme du modèle. Si au contraire le dernier caractère ne peut être relié à l'état final, le mot est refusé et n'est donc pas formé suivant les règles de cette grammaire.

La grammaire non déterministe Gn3 est une version légèrement modifiée <sup>6</sup> de G2. Elle procède dans ce cadre à une analyse fine mais sans souci d'exhaustivité des préfixes et suffixes éventuels du nom arabe, dans le cas réduit du trilitère sain. Les objets soumis à son verdict sont des mots formés d'une suite de caractères arabes, à l'exclusion des signes de vocalisation externe *fatha*, *damma*, *kasra*, *sukūn*, *waṣla*, *tanwīn et šadda*; la *hamza* est en revanche tolérée sous toutes ses formes. Les catégories du nom arabe sont ainsi clairement mises en évidence et, une fois détectées, doivent permettre de cerner les possibilités de radicaux du nom et d'émettre, dans certains cas, des hypothèses sur la fonction du mot dans la phrase.

Le modèle Gn3(c,c,c) <sup>7</sup>, reproduit dans la figure 1 sous la forme de graphe de transitions, constitue un des états du développement de cette grammaire ; il sera pris comme référence pour la suite de cet article.

La partie radicale de cette grammaire doit retenir toute notre attention. L'influence néfaste de la présence d'éléments agglutinés conduit à distinguer deux variétés de consonnes.  $\Sigma_r$  est l'ensemble des consonnes nommées « solides » car elles ne peuvent qu'appartenir à la racine, à la rare exception de certains infixes produits par la dérivation de la forme VIII  $^8$ .  $\Sigma_f$  est son complémentaire dans l'ensemble des consonnes « saines » : il s'agit

renthèses font référence aux élément radicaux reconnus respectivement en R1, R2 et R3, soit toutes les consonnes saines dans le cas présent.

8. C'est le cas de « د » dans « إزدهار » et de « ط » dans « ط ».

<sup>6.</sup> Seules les modifications apportées à cette grammaire seront détaillées au cours de cette étude.

<sup>7.</sup> Nous utiliserons lorsqu'elle s'avère nécessaire cette notation fonctionelle. « n3 » est mis pour « nominal trilitère » et les caractères entre pa-

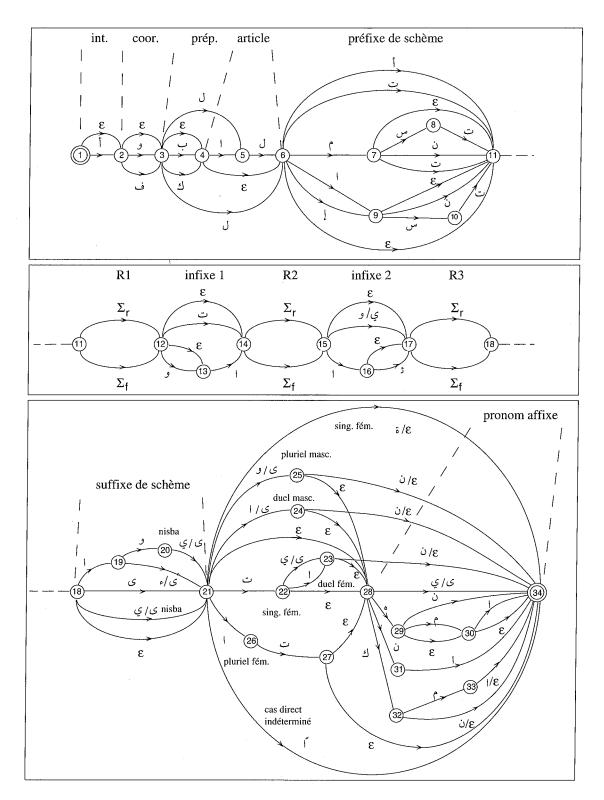


Fig. 1.

	catégorie	abréviation		exemple
	interrogatif	int	f	أبقلمه ؟
pré-concaténation	coordonnant	coor	و	وغريزة
	préposition	prep	ب	وغريزة بزيارة
	article	art	ال	الشيخ
	préfixe de schème	pref	ت	التقليد
	première radicale	R1	ق	التقليد
racine et schème	infixe de schème entre R1 et R2	inf1	1	القاهرة
	seconde radicale	R2	J	التقليد
	infixe de schème entre R2 et R3	inf2	ي	التقليد
	troisième radicale	R3	د	التقليد
	suffixe de schème	suff	اء	والغرباء
	nisba	nisb	ي	القرويين
	singulier féminin	sgfm	ت	القرويين بقبلاته
	pluriel masculin	plma	ین	الداخلين
post-concaténation	pluriel féminin	plfm	ات	كرامات
	duel masculin	duma	ان	الزائران
	duel féminin	dufm	تان	الساحرتان
	cas direct non déterminé	obnd	١	ناقها
	postfixe	post	لما	ونسائها

Fig. 2. – Catégories de Gn3.

donc des consonnes susceptibles d'être agglomérées au nom, au titre d'article, préposition, postfixe, désinence de genre ou de nombre mais aussi d'élément du schème de dérivation.

$$\begin{split} & \sum_{r} = \{ \underline{t}, \, \underline{\check{g}}, \, \underline{h}, \, \underline{h}, \, \underline{d}, \, \underline{q}, \, \underline{r}, \, \underline{z}, \, \underline{\check{s}}, \, \underline{s}, \, \underline{d}, \, \underline{t}, \, \underline{z}, \, \underline{\check{r}}, \, \underline{\dot{g}}, \, \underline{q} \} \\ & \sum_{f} = \{ \underline{b}, \, \underline{t}, \, \underline{s}, \, \underline{f}, \, \underline{k}, \, \underline{l}, \, \underline{m}, \, \underline{n}, \, \underline{h} \} \end{split}$$

Les arcs  $\Sigma_d$  et  $\Sigma_f$  en R1, R2 et R3 représentent donc les possibilités de détection par Gn3s de lettres radicales saines.

Les éléments ajoutés au mot ne faisant partie ni du schème ni de la racine seront désignés par éléments de pré-concaténation ou post-concaténation suivant leur position. L'ensemble  $\Sigma_f$  peut être considéré comme la réunion de deux sous-ensembles de consonne

« à risques », l'un  $\Sigma_{f1}$  au regard de la post-concaténation, l'autre  $\Sigma_{f3}$  pour la préconcaténation.

$$\begin{split} & \Sigma_{f1} = \{ b, \, t, \, s, \, f, \, k, \, l, \, m, \, n \} \\ & \Sigma_{f3} = \{ t, \, k, \, m, \, n, \, h \} \end{split}$$

De même, une nouvelle catégorie intitulée « suffixe de schème » permet la reconnaissance de schèmes tel « فعلاء » et leur comportement après l'ajout d'une nisba.

Le tableau de la (fig. 2) récapitule les catégories introduites avec un exemple d'occurrence.

Remarquons dès maintenant que la grammaire Gn3(c,c,c) peut être facilement étendue à un modèle qui permettrait la reconnaissance dans certains cas de radicaux incluant les semi-consonnes w ou y et les lettres contenant la hamza. C'est dans cette direction que nous effectuerons des variations de grammaire. Il est clair que si la racine est parfaitement saine, aucune de ces lettres ne saurait être tolérée en R1, R2 et R3.

#### 2. MISE EN ŒUVRE INFORMATIQUE.

La programmation d'un analyseur de textes s'appuyant sur des grammaires non déterministes tel Gn3 a déjà été effectuée dans le cadre d'un langage de programmation interprété, le Lisp <sup>9</sup>. Cet analyseur tirait parti, dans une certaine mesure, de la structure naturelle de liste propre au Lisp. La réalisation d'un programme indépendant d'analyse de l'arabe pouvant être intégré à des logiciels tels les gestionnaires de bases de données ou les traitements de textes imposait une programmation au moyen d'un langage algorithmique courant aux possibilités étendues. Notre choix s'est porté sur le langage C, aujourd'hui très répandu ; il est dès lors envisageable de faire fonctionner nos programmes sur différents types de matériel sans en changer le moteur que constitue l'analyseur <sup>10</sup>.

#### L'ANALYSEUR.

La difficulté essentielle de l'analyse à l'aide d'une grammaire telle que Gn3 vient de son indéterminisme : il faut pouvoir détecter et enregistrer les choix multiples qui apparaissent au cours du traitement afin de pouvoir exhiber les interprétations <sup>11</sup> du mot

9. Cf n. 2.

10. Soulignons toutefois qu'une adaptation de l'analyseur programmé en Lisp pourra s'avérer pertinente lors de la phase d'optimisation des pro-

grammes morphologiques.

11. Le terme « interprétation » sera précisé plus has

selon la grammaire utilisée. C'est une structure de donnée en « arbre » que nous avons adoptée pour développer les interprétations licites de chaque mot <sup>12</sup>.

L'exemple d'analyse proposé ci-dessous dans la figure 3 montre chacune des ramifications de l'arbre solution dans deux cas donnant lieu respectivement à deux et trois interprétations. Elles sont rangées dans un ordre propre à l'investigation qui n'a aucune signification linguistique ou même alphabétique. Chaque caractère du mot est accompagné d'une catégorie abrégée.

```
فراشه
1/ن:(coor/):R1 /l:inf1/ش:R2 /ه:R3
2/ن:(R1 /post):R3 /s:post الترك:
1/ن:(R1 -pref/ب:R1 / R1:(ك:(R3 - Pref/ب:R1 / R3):(ك:(R3 - Pref/ب:R1 / R3):(ك-(R3 - Pref/ب:R1 / R2):(ك-(R3 - Pref/ب:R1 / R2):(-(R3 - Pref/ب:R1 / R3):(-(R3 - Pref/ب:R1 / R3):(-(R3 - Pref/ R3):
```

Fig. 3. - Exemple de réponse produite par l'analyseur.

#### L'EXPLOITATION DE L'ANALYSEUR.

La mise au point des grammaires nominales et verbales comme des analyseurs est difficilement imaginable sans le recours au traitement systématique de textes variés et la validation des résultats obtenus. Après une première phase d'élaboration de l'analyseur à l'aide de quelques mots isolés, nous avons réalisé un logiciel expérimental baptisé « Sarfeyya » qui permet d'analyser un texte entier de quelques pages.

Le principe de variations de grammaire se traduit par la possibilité donnée à l'opérateur de composer la partie radicale de la grammaire en autorisant ou interdisant les consonnes, les lettres hamzées et les semi-consonnes.

L'analyse proprement dite est affichée ou enregistrée dans des fichiers au format neutre où sont reportées les interprétations des mots du texte, les racines et les schèmes détectés.

Il est important de remarquer que les progrès récents de la micro-informatique ont rendu fort acceptables les temps d'exécution de tels traitements. En effet, l'ordre de grandeur du temps d'analyse d'un texte de la taille de cet article est la seconde <sup>13</sup>.

12. Un exposé détaillé de l'algorithme employé n'est pas ici notre propos. L'analyseur est donc considéré dans cette étude comme une « boîte noire » dont les entrées sont des textes et la grammaire employée, et dont la sortie est le détail de l'analyse de chaque mot détecté.

13. Ces performances sont fortement dépendantes du matériel utilisé; leur détail relève de l'optimisation de l'analyseur et des techniques de programmation.

#### 3. L'ANALYSE D'UN TEXTE : COMMENT ÉVALUER UNE GRAMMAIRE ?

#### CLASSIFIER LES RÉPONSES.

Rappelons ici les hypothèses faites sur le mot graphique soumis à la grammaire Gn3:

- 1) appartenance à la langue arabe standard, sans archaïsme ;
- 2) non vocalisé : si toutefois il l'est totalement ou partiellement, cette vocalisation est ignorée dans l'analyse ;
- 3) la hamza est correctement orthographiée, soit entre autres « أحد » et non « إحد », « أحارة » et non « إحارة » ;
- 4) le mot ne contient pas de fautes d'orthographe ou de frappe, notre propos n'étant pas ici la correction orthographique.

Ces hypothèses sont assez restrictives dans la mesure où bien des textes modernes maltraitent les difficultés orthographiques en les simplifiant, en particulier dans la presse.

Il est clair qu'une grammaire du type de Gn3 fonctionne comme un filtre vis-à-vis de toute suite de caractères qui lui est proposée : il peut soit refuser un mot, soit l'accepter. Un refus signifie simplement que le mot n'est pas conforme aux seules règles de la grammaire programmée ; de même, une acceptation peut être multiple mais ne pas contenir l'analyse correcte du mot. Or cette grammaire opère à ce jour sans aucun lexique. Il n'y a donc aucune confirmation de l'existence de la racine extraite. Le schème de dérivation du mot obtenu, reconstitué à partir des catégories préfixes, infixe 1, infixe 2 et suffixe de schème ne subit pas non plus de contrôle d'existence : un préfixe détecté peut être incompatible avec un second infixe et l'interprétation se révéler erronée.

L'opérateur est seul juge de la pertinence des solutions. Il apparaît dès lors indispensable d'adopter une terminologie permettant de classifier les réponses en fonction de leur validité.

# NOTIONS DE SILENCE ET DE BRUIT APPLIQUÉES À LA MORPHOLOGIE ARABE.

#### a. Terminologie.

Adoptons quelques définitions :

- $-\Gamma$  est un ensemble de productions écrites ou mots graphiques de l'arabe répondant à un ou plusieurs critères linguistiques : l'ensemble des noms trilitères sains, l'ensemble des verbes contenant une *hamza* en radicale, etc. Il est désigné par *champ* d'application par opposition au domaine des mots graphiques.
- G est une grammaire formelle morphologique conçue dans le but d'accepter les éléments de  $\Gamma$  en en détaillant les ambiguïtés et de refuser ceux qui n'en font pas partie.

Le tableau de la figure 4 permet de préciser dans une situation concrète les définitions courantes : le champ  $\Gamma$  est l'ensemble des noms formés à partir d'une racine trilitère saine et G est la grammaire Gn3(c,c,c).

20

			ļ							
ex.	Mot	Catégorie Nr	ż	Réponse	Racine	Racine Schème	Classe	Cause du bruit	Conséquences	Remarques
-	اصطدام	nom sain	0	vide			silence			imposs.en inf1
3	رجال	nom sain	0	1/2.R1 /2:R2 / \tinf2/ L:R3	<u>ئ</u>	فعال	solution			
4	أشفال	nom sain	2	1/ i.pref/.å:R1 /z.R2 /:inf2/.j:R3	شفل	أفعال	solution			
				2/ \tint/\tint /2:R2 / \tinf2/\tink3	عا	فعال	bruit	ambiguïté interne	B. radicalement exact	
				R3 / الا: هـــ/ 3R2 (13: الا:ن/duma) الا: الا: هـــ/ 18: هـــ/ 18: الا: الا: الا: الا: الا: الا: الا: ال	:बै	فعل	bruit	post-concat. : الله	B. résiduel : 🚣 🗟 non attesté	ध en R1 est pris pour préposition
5	كنفها	nom sain	3	2/s.prep/.:R1 (4:8.82 / 18:0.8nd	:बं	فعل	bruit	post-concat.	B. résiduel	
				3/4:R1/.:R2/:±/83:-;post/:post	i.id Zi	فعل	solution bruitée			
				1/نprep/رart /ع:د/ A:بد/Pref/عاتد/ العنل/ا	يرك	تفعل	solution bruitée			
9	للتبرك	nom sain	ε.	2/يiprep/ اR1 ( تا:ب/1311). (R3 /كا:ب/2). اR3 الماني /2	73;	فتعل	bruit	pré-concat. : ਹੀ	B. résiduel : ما لبر B. résiduel	bruit de schème erroné
				3/ اrep/رiart (ج. 13.1). (31.1/3). 3/ العنل/غانون].	<b>'</b> 1;	فعل	bruit	ambiguïté interne	ambiguïté interne B. radicalement attesté	incorrection : مtfixe دران affixe
2	جدي	nom sourd	0	vide			rejet			imposs. en R3
7	البيت	nom concave	- 1	1/ ا:pref/J:R1 / ب/ R2 الــــــــــــــــــــــــــــــــــــ	7	فعيل	bruit sans solution	pré-concat. : ປ	B. résiduel : ت ا non attesté	
∞ .	يخرج	verbe	0	vide			rejet			ي imposs. en préf. Sch.
6	استقر	verbe	-	1/ !:pref/.w:R1 / ::nf1/.R3	ıď	افتعل	bruit sans solution	ambiguïté interne	B. radicalement attesté	bruit de schème erroné
10	تفطن	verbe	-	1/ R3:ن/ R3:د// 17:ح/ 18:د/1	نوط فط	تفعل	bruit sans solution		ambiguïté interne B. radicalement exact	
11	فعندئذ	token	0	vide			rejet	A control		نے imposs. en suff. Sch.
12	<u>"J</u>	token	1	/3:R1 /4:R2 /J:R3	ين ا	فعل	bruit	ambiguïté interne	ambiguïté interne B. radicalement exact	

Fig. 4. - Exemples de classement des réponses de la grammaire Gn3(c,c,c) appliquée aux noms trilitères sains et aux autres mots.

La réponse R de G à un mot arabe M quelconque est constituée de Nr éléments, Nr positif ou nul. Plusieurs cas sont alors à distinguer :

- Si Nr = 0, R est un *silence* si M appartient à  $\Gamma$  (exemple 1, fig. 4); un *rejet* si M n'appartient pas à  $\Gamma$  (ex. 2);
- Si Nr > 0, il reçoit le nom d'ambiguïté de la réponse : R est alors composée de Nr interprétations qui sont qualifiées de :
  - solution si M appartient à  $\Gamma$  et l'interprétation est une analyse juste de M (première interprétation de l'ex. 3, soit 3.1);
  - bruit si M appartient à G et l'interprétation est une analyse erronée (ex. 4.2) ou si M n'appartient pas à  $\Gamma$  (ex. 7.1).

La solution, si elle existe, est bien entendu unique : il n'est pas tenu compte des cas très rares où l'auteur jouerait sur l'écriture non vocalisée d'un mot pour lui affecter plusieurs sens.

On peut donc parler de *solution bruitée* (ex. 5.3) si R contient une solution parmi du bruit ; de même, une réponse constituée uniquement de bruit sera qualifiée de *bruit sans solution* (ex. 7.1).

#### b. Les causes naturelles du bruit.

L'étude des causes du bruit est capitale pour la compréhension des phénomènes d'ambiguïté engendrés par le système d'écriture graphique de l'arabe.

Un bruit, analyse erronée du mot, peut avoir plusieurs origines parfois mêlées entre elles. Pour les distinguer, nous procéderons ainsi :

- si le mot à analyser, une fois débarrassé de tout élément de post ou préconcaténation, produit la même réponse erronée dépourvue elle aussi de ces éléments, le bruit sera considéré comme *bruit d'ambiguïté interne* (ex. 4.2, 6.3, 9.1);
- si le bruit persiste lorsque l'on ôte du mot à analyser les seuls éléments de postconcaténation, nous parlerons de *bruit de pré-concaténation*. C'est un phénomène de décalage de la racine : une lettre placée avant R1 est prise pour R1 et R3 elle-même est interprétée comme élément de concaténation car elle appartient à l'ensemble à risque  $\Sigma_{f3}$ (ex. 6.2, 7.1) ;
- réciproquement, le *bruit de post-concaténation* désigne les interprétations restées erronées malgré l'élimination des éléments de pré-concaténation <sup>14</sup> : les rôles sont inversés et c'est R1 qui appartient à l'ensemble à risque  $\Sigma_{f1}$  (ex. 5.1, 5.2).

14. Cf. l'exemple « كنفها » de la fig. 4.

Certaines causes contenues dans l'une des trois raisons majeures décrites ci-dessus peuvent faire l'objet d'un décompte particulier : c'est le cas ici du bruit engendré par la présence de l'article « ال » et de la préposition-article « , 以 ».

#### c. Les bruit propres au modèle.

Le fait que le schème de dérivation ne subisse pas de contrôle d'existence rend parfois possible des interprétations au schème erroné mais toléré par la grammaire. Ce phénomène n'est jamais une cause directe de bruit mais s'apparente plutôt à un défaut de la grammaire formelle qui pourra être éliminé par l'introduction d'une liste de schèmes. Ce bruit sera qualifié de *bruit de schème erroné* (ex. 6.2, 9.1).

Un second défaut important de ce modèle est sa tolérance vis-à-vis de la présence simultanée d'un article et d'un pronom affixe, cas impossible car il conduirait à une double détermination du nom. Ce bruit, ainsi que tout autre bruit se développant à cause d'un défaut de nature linguistique, sera noté bruit d'incorrection morphologique (ex. 6.3). Seule une modification des règles de réécriture ou une augmentation <sup>15</sup> de la grammaire pourra l'éliminer.

#### d. Les conséquences du bruit.

Si l'objectif de l'analyse est l'extraction de la racine, le bruit doit être catégorisé suivant son incidence dans cette extraction.

Il arrive en effet que parmi le bruit apparaissent des interprétations dont le radical est bien celui de la solution, que celle-ci soit par ailleurs détectée ou non : nous qualifierons ce bruit de *bruit radicalement exact* (ex. 4.2).

Au sein du bruit non radicalement exact peuvent se présenter des interprétations dont la racine est attestée, cas particulièrement embarrassant. Ces interprétations désignées par bruit radicalement attesté correspondent à l'ambiguïté même du système graphique de l'arabe non vocalisé, celle qui peut faire hésiter parfois un lecteur même averti (ex. 6.3, 9.1). L'ambiguïté radicale totalise le nombre de racines attestées différentes dans R.

Le *bruit résiduel* correspond aux interprétations dont la racine n'est pas attestée : le seul recours à une liste des racines attestées suffit à l'éliminer (ex. 5.1, 5.2, 6.2, 7.1).

La terminologie adoptée prend tout son sens lorsque l'on étudie la réponse des verbes et des tokens, qui ne sont donc pas éléments de ce champ  $\Gamma$ . Certains d'entre eux sont heureusement rejetés, mais d'autres sont acceptés et la réponse est alors clairement un bruit. Ce dernier est radicalement exact si le schème non vocalisé du mot se confond avec un schème non vocalisé nominal prévu par la grammaire. La cause de ce bruit est classée parmi les ambiguïtés internes (ex. 10.1).

15. Une grammaire est dite « augmentée » lorsque l'on sort du cadre des langages réguliers pour introduire des lexiques ou des arcs ne pouvant

être empruntés que sous certaines conditions afin de modéliser des exceptions.

Ainsi, tout dépouillement de texte à l'aide d'une grammaire orientée vers l'analyse des mots et l'extraction de leur racine devra comporter la comptabilité précise de ces bruits et silences. Il est clair par ailleurs qu'une telle étude devra porter sur tous les mots du texte, sauf les noms propres et emprunts étrangers, et ceci afin de mesurer en quelque sorte le pouvoir de rejet ou « pouvoir séparateur » de la grammaire vis-à-vis des autres catégories de mots que celles pour laquelle elle est destinée.

#### UNE PROPRIÉTÉ IMPORTANTE DES GRAMMAIRES FORMELLES.

Une certaine catégorie de variations des grammaires s'accompagne d'une variation *li-néaire* de la réponse. En effet, supposons qu'une grammaire G accepte un nombre T de transitions en deux quelconque de ses états E et F. Sa réponse R à un mot se décompose en deux réponses : celle R<sub>1</sub> contenant les interprétations qui empruntent les T arcs permis et celle R<sub>2</sub> des interprétations utilisant une autre voie (un arc epsilon-transition entre deux états situés de part et d'autre de E et F par exemple). Si l'on ajoute une transition possible supplémentaire entre E et F, la réponse de l'automate obtenu G' au même mot se décompose de la même façon entre une R'<sub>1</sub> et R'<sub>2</sub>. Il est clair que R'<sub>1</sub> est R<sub>1</sub> et que R'<sub>2</sub> est R<sub>2</sub> augmentée des interprétations utilisant la possibilité nouvelle de transition entre E et F.

La portée de cette propriété est importante dans notre cas : elle signifie que pour une seule variation concernant l'ajout d'un arc entre deux états, il suffit d'étudier la réponse de l'automate auxiliaire pour lequel les transitions entre ces états sont limitées au seul arc introduit. La réponse de l'automate complet est alors la somme des réponses de l'automate original et de l'auxiliaire. Ce procédé peut être généralisé à plusieurs arcs à condition qu'ils relient toujours les deux mêmes états. Il peut aussi par extension s'appliquer à l'introduction de sous-automates entre deux états <sup>16</sup>.

#### MÉTHODE DE DÉPOUILLEMENT D'UN TEXTE.

L'opérateur désireux d'évaluer la portée d'une grammaire pourra suivre naturellement ces étapes :

1) une analyse préalable consistant en l'extraction « manuelle » des racines et leur report dans une base de donnée. Celle-ci devra comporter au moins le mot, sa ca-tégorie (nom, verbe, token et autres), sa racine et le type de cette dernière (saine, concave, etc.) ;

16. Une étude complète de ces propriétés dépasse la portée de cet article. Remarquons toutefois que le recours à la formulation algébrique (Cf. *BEO* XLVI,

annexe I) permet de démontrer rapidement ce cas de linéarité.

20A

- 2) l'utilisation du logiciel pour produire l'analyse « automatique » du texte ;
- 3) la validation du résultat de 2 par la vérification de chaque interprétation et l'enrichissement de la base des mots conformément à la terminologie exposée cidessus ;
- 4) L'étude et la synthèse des statistiques issues de 3, au moyen par exemple de taux mesurant la réussite des extractions, la fréquence de telle cause de bruit, etc.

Les variations de grammaire ont pour but l'élimination d'une source de bruit ou la réduction des silences observés à l'issue de l'étape 4. Les tentatives de modification du champ d'étude des grammaires s'effectuent également par variations.

Ces variations impliquent la reprise des phases 2 à 4 et la comparaison des résultats obtenus en 4. Si l'une d'elles entre dans le cadre de la propriété de linéarité, il suffit d'étudier l'analyse du texte produite par la seule présence du ou des arcs nouveaux entre les états souhaités.

L'étude des origines du bruit peut amener à définir de nouvelles variables de dépouillement afin de pouvoir repérer les causes majeures et leur éventuelle interpénétration. Un dialogue s'amorce entre l'étape 3 et la suivante jusqu'à l'obtention d'une description satisfaisante des phénomènes survenant lors de l'analyse.

On peut mesurer au moyen de taux adéquats la réussite globale de la grammaire dans sa mission d'analyse morphologique. Ces taux peuvent ensuite servir d'éléments de comparaison de deux grammaires à l'épreuve d'un même texte, ou d'une même grammaire devant deux textes de nature différente tel un article de presse ou une poésie.

C'est la démarche que nous avons adoptée ici. Le dépouillement a été facilité par le fait que l'analyse d'un texte dans sa forme de fichier informatique standard peut être relue par un simple gestionnaire de base de donnée. Les fonctions de recherche et de tri permettent ensuite d'établir les statistiques souhaitées.

#### 4. ANALYSE D'UN TEXTE LITTÉRAIRE CONTEMPORAIN

ÉTAPE 1: LE DÉPOUILLEMENT « MANUEL ».

C'est sur un texte littéraire contemporain que nous avons procédé à une première analyse nominale. Il s'agit des deux premiers chapitres de la nouvelle de Yaḥyā Ḥaqqī intitulée « Qindīl Umm Hāšem » <sup>17</sup> dont l'intégralité est reproduite en annexe. Ce texte répond aux hypothèses exposées ci-dessus mais comporte naturellement de nombreux noms propres et emprunts de dialecte cairote. Les dialogues, qui représentent une part négligeable de la masse des mots, ont été coupés. Cependant une comptabilité précise

17. Mu'allafāt Yaḥyā Ḥaqqī, al-qiṣaṣ 1, al-Hay'a al-miṣriyya al-ʿāmma li-l-kitāb, Le Caire, 1990, p. 59-69.

s'imposait pour mesurer l'ampleur réelle de ce qui sera toujours l'ennemi du traitement automatique d'un texte.

Le tableau I (fig. 5) distingue pour chaque catégorie de mots le nombre des premières occurrences du nombre total d'apparitions dans le texte. Dans le premier cas comme dans le second *a fortiori* les suites de caractères telles « جدى » et « بلدى » sont comptées pour deux mots distincts. Il s'avère ainsi, et ceci ne doit pas nous surprendre, que pour un texte d'un total de 1300 mots environ près de la moitié est constituée de noms trilitères, 30 % de tokens et moins de 20 % de verbes trilitères. Les noms propres et toutes les autres catégories peuvent être pris pour quantité négligeable, leur total étant circonscrit autour de 5 %. La priorité donnée à la mise au point de grammaires nominales apparaît dès lors pleinement justifiée.

Le tableau II donne le résultat du dépouillement radical « manuel » de chaque nom trilitère. Les abréviations adoptées sont les suivantes :

- « c » désigne une consonne élément de  $\Sigma_d$  ou  $\Sigma_f$ ;
- « s » désigne une semi-consonne w ou y ;
- « h » désigne la hamza;
- « = » en troisième position signifie que la racine est redoublée : R3 = R2.

Les noms sourds à consonnes solides sont distingués suivant deux catégories : la forme dissimilée « مرور », notée cc=2, et la forme assimilée « مرق », notée cc=1.

On peut constater la forte proportion des racines saines ccc, qui totalisent près de la moitié des cas, et donc 22 % du texte total. Les autres cas les plus fréquents sont les radicaux concaves csc, puis les sourds cc=, défectueux ccs, assimilés scc et à première radicale hamzée hcc. Les cas présentant une double irrégularité ne totalisent que 3 % de la masse des racines nominales trilitères de ce texte.

L'étude similaire pour les verbes trilitères présentée dans le tableau III montre que là encore, les racines sont saines pour presque la moitié des cas, concaves pour le quart tandis que les autres irrégularités se présentent avec une fréquence comparable à celle du nom.

Ces faits sont à comparer aux statistiques portant sur les 4814 racines trilitères attestées <sup>18</sup> d'un dictionnaire arabe classique : la proportion de racines trilitères saines s'élève à 62 %. La suite des autres irrégularités rangées par fréquence décroissante respecte le même ordre que celui issu de notre texte.

L'enjeu de cette recherche commence ainsi à se dessiner : doit-on mettre au point autant de grammaires qu'il y a de cas d'irrégularités dans la morphologie – qu'elle soit

18. D'après une étude menée par 'Alī Hilmī Mūsā à partir du dictionnaire *Aṣ-ṣihāh* d'Al-Ğawharī. « Dirāsat 'iḥṣā'iiya liğudūr mufradāt al-lūga al- 'arabīyya », Université de Koweit, 1971. Cette étude

est complétée par les travaux de Hussein Habaili, « Phonologie et morphologie de l'arabe », thèse de doctorat de 3ème cycle, Université de Paris III.

catégorie	type racine	1ère occ	répétitions	total	part	cumul
nom	trilitère	540	60	600	46,7	
	bilitère	10	2	12	0,9	
٠	quadrilitère	12		12	0,9	48,5
	propre	20	14	34	2,6	
	dialecte	16		16	1,2	3,9
verbe	trilitère	221	13	234	18,2	
	quadrilitère	2		2	0,2	18,4
token	trilitère	37	40	77	6,0	
	autres	88	211	299	23,3	29,2
totaux		946	340	1286	100,0	

Tableau I. Répartition des catégories de mots dans le texte.

type de rac	ine trilitère	1ère occ	répétitions	total	part	cumul
sain	ccc	274	10	284	47,3	
sourd	cc=2 dissim.	26		26	4,3	
Soura	cc=1 assim.	25	7	32	5,3	
une	csc	89	19	108	18,0	
semi-	ccs	42	6	48	8,0	
consonne	scc	25	3	28	4,7	
	hcc	20	6	26	4,3	
une hamzée	cch	6		6	1,0	
	chc	4		4	0,7	93,7
	cs=	6	3	9	1,5	
]	ssc	5	1	6	1,0	
	csh	5		5	0,8	
deux	hc=	5	5	10	1,7	
irrégularités	css	3		3	0,5	
	hsc	2		2	0,3	
	scs	1		1	0,2	
	sc=	1		1	0,2	
	hcs	1		1	0,2	6,3
		540	60	600	100,0	

Tableau II.

Répartition des types
de racines des noms
trilitères.

type racine	1ère occ	répétition	total	part	cumul
ccc	102	2	104	44,4	
cc=2 dissim.					
cc=1 assim.	14		14	6,0	
esc	48	7	55	23,5	]
ccs	22	3	25	10,7	
scc	12		12	5,1	
hee					
cch	11		11	4,7	
chc	1		1	0,4	94,9
css	5		5	2,1	
chs	2	1	3	1,3	
csh	2		2	0,9	
hss	1		1	0,4	
hcs	1		1	0,4	5,1
	221	13	234	100,0	

Répartition des types de racines des verbes trilitères.

Tableau III.

Fig. 5. – Résultats de l'étape de dépouillement du texte étudié.

d'ailleurs nominale ou verbale – ou doit-on au contraire tenter de rassembler, au prix sans doute de compromis, les cas embarrassants en un seul modèle ? Quels sont les cas réellement « embarrassants » vis-à-vis d'un traitement automatique ? Nous tenterons de répondre en montrant jusqu'à quel point le concept de variations de grammaire peut s'adapter à des problèmes précis de traitement de l'arabe.

#### PREMIÈRE ANALYSE: LA GRAMMAIRE DU NOM TRILITÈRE SAIN Gn3(c,c,c).

Chaque analyse se présentant sous forme d'une liste d'interprétations pour l'intégralité des mots du texte, nous nous contenterons d'en reproduire en annexe un extrait significatif. Les réponses vides, silences ou refus, n'apparaissent pas dans les listes d'interprétations. Les tableaux A de résultats montrent la répartition des réponses en termes de silence, rejet, succès et bruit concernant les noms et verbes trilitères ainsi que les tokens : c'est donc l'aboutissement de l'étape 3.

#### a. Gn3(c,c,c) et les noms trilitères sains.

La solution est toujours trouvée, à l'exception d'un silence dû à l'ignorance du phénomène de transformation d'infixe de la forme VIII <sup>19</sup>. Cette valeur minimale du silence est le résultat de l'affinement du modèle d'origine G2 par la présente méthode.

Le bruit détecté provient en grande partie d'éléments de post-concaténation qui ne sont pas des pronoms affixes, tels « ا », « ين », « ين » après R3, ce dernier cas pouvant donner lieu à cinq interprétations (cf. l'analyse de « شحاذی » montrée dans l'extrait). L'influence de la présence du « ال » est presque nulle. La plupart des cas d'ambiguïté interne sont le fait du « أ » avant R1 dont la catégorie est toujours préfixe de schème dans ce texte mais qui peut s'interpréter, rarement il est vrai, comme une particule interrogative. Le bruit le plus embarrassant, le bruit radicalement attesté, s'avère ici très faible.

#### b. Gn3(c,c,c) et les autres mots.

Le cas des noms à radical sourd mais dissimilé est sans surprise car il constitue une sous catégorie des trilitères sains : il présente donc des résultats tout à fait comparables.

Tous les autres noms devraient être rejetés par Gn3(c,c,c). Ce n'est pourtant pas le cas, plus d'un tiers des analyses donnant lieu à du bruit. Celui-ci est essentiellement dû à la présence du « 🗸) » qui trouble fortement les réponses des *csc* et des *ccs*. Cette source de bruit est néanmoins contrôlable car elle produit des racines parasites à première radicale « 🗸 » rarement attestées. Les éléments de post-concaténation, et à leur tête les pronoms

19. Cf. n. 8. Ce phénomène peut être modélisé mais nécessite une augmentation de la grammaire pour ne pas engendrer de bruit.

Tableaux A. – Dépouillement des réponses de Gn3(c,c,c).

	т.																
						classe		typ	e bru	it	<u></u>		causes	i		évalı	ıation
NOMS/Gn3(c,c,c)	type racine	total occurrences	total réponses (Nr)	silence	rejet	solution	bruit	bruit rad. exact	bruit rad. attesté	bruit résiduel	schème inexact	pré-concat. (sauf article)	うぉヿ	post-concat.	ambiguïté interne	succès extraction	taux bruit attesté
sains	ccc	284	340	1		283	57	41	6	10	9		4	36	14	94,7	1,8
sourds	cc=2	26	31			26	5	5					.,	5		100,0	0
	cc=1	32	31		11		31		7	24	9	7	10	16	3	0	22,6
une radicale S	csc ccs scc	108 48 28	53 34 10		58 28 21		53 34 10		15 16 7	38 18 3	38 24 4	5	34 24 2	11 6 2	3	0 0 0	28,3 47,1
une radicale H	hec che cch	26 4 6	6 0 0		22 4 6		6		1	5	1			4	1	0 0 0	-
	ssc csh hc= css	6 5 10 3	0 0 0 1		6 5 10 2		1		. 1			1				0 0 0	- - -
reste	cs= hsc scs sc= hcs	9 2 1 1	0 0 0 0 2		9 2 1 1		2			2	2			1		0 0 0 0	-
totaux sauf ccc et cc=2	1	<b>600</b> 290	<b>508</b> 137	1 0	<b>186</b> 186	<b>309</b> 0	<b>199</b> 137	<b>46</b> 0	<b>53</b> 47	<b>100</b> 90	<b>87</b> 78	13 13	<b>74</b> 70	<b>81</b> 40	<b>25</b>		10,43 34,31

VERBES/Gn3(c,c,c)	type racine	total occurrences	total réponses (Nr)	silence	rejet	solution	bruit	bruit rad. exact	bruit rad. attesté	bruit résiduel	schème inexact	pré-concat. (sauf article)	うぉつ	post-concat.	ambiguïté interne	succès extraction	taux bruit attesté
sains	ccc	104	59		53		59	50	3	6	5			2	49	48,1	5,08
sourds	cc=2	0															
	cc=1	14	11		5		11		3	8	1			2	8	0	27,3
une radicale S	csc ccs scc	55 25 12	25 25 4		33 13 8		25 25 4		2	23 25 2	2 2 1	Ī			24 25 4	0 0 0	8 0 -
une radicale H	hec che cch	0 1 11	0	]	1				:							0	- - -
reste	css chs csh hss	5 3 2 1	0 0 0 0		5 3 2 1								•			0 0 0 0	-
totaux sauf ccc et cc=2		234 130	124 65	<b>0</b> 13	<b>36</b> 33	<b>0</b> 0	124 65	<b>50</b> 0	<b>10</b> 7	<b>64</b> 58	11 6	1 !	<b>0</b> 0	<b>4</b> 2	110 61		8,06 10,8

	l .				c	lasse		typ	e brui	it			causes	i		éval	uation
TOKENS/Gn3(e,c,c)	type racine	total occurrences	total réponses (Nr)	Silence	1001	solution	bruit	bruit rad. exact	bruit rad. attesté	bruit résiduel	schème inexact	pré-concat. (sauf article)	うぉつ	post-concat.	ambiguïté interne	succès extraction	taux bruit attesté
sains	ccc	17	11		6		11	[]			1			2	11	64,7	0
sourds	cc=2																
	cc=1	20	17	1	1		17		9	8		4		12	9	0	52,9
s	csc ccs scc	10 21 1	1 3 0		9 9 1		3		1	1 2		1		1	3	0 0 0	33,3
Н	hee ehe ech	4	4 0		2		4			4				4	4	0	-
reste trilitères	hsc	3	0		3											0	-
non trilitères		299	34	27	3		34		12	22		2		19	32	0	35,3
total ni sain ni sourd		<b>376</b> 359	<b>70</b> 59	<b>0 32</b> 0 31		<b>0</b> 0	<b>70</b> 59	11 0	<b>22</b> 22	<b>37</b> · 37	1 0	7 7	<b>0</b> 0	<b>38</b> 36	<b>60</b> 49		31,4 122

affixes tels « 。 », « ه », « ه », « ه », produisent un fort bruit radicalement attesté pour les noms à racine csc et les sourds assimilés. Les autres catégories de racines sont clairement rejetées, la taille réduite des occurrences et du bruit rendent d'ailleurs les statistiques peu significatives.

Le rejet attendu des verbes n'est pas flagrant : 40 % d'entre eux admettent des interprétations qui sont donc clairement des bruits. Il s'agit essentiellement de bruit radicalement exact dû à la coïncidence de schèmes de dérivation des racines saines décrite plus haut. Les schèmes les plus fréquemment pris pour des schèmes nominaux sont classés par ordre de fréquence dans le tableau ci-dessous.

Schèmes dont la forme non vocalisée est commune aux noms et aux verbes

Forme non vocalisée	فعل	تفعل	أفعل	فاعل
Exemple nominal	فَعْل	تَفَعُّل	أَفْعَل	فَاعِل
Exemple verbal	فَعَلَ	تَفْعَلُ	أَفْعَلَ	فَاعَلَ

Le rejet des verbes *csc*, *ccs* et *cc=1* n'est pas plus franc. Il s'explique par la présence fréquente du « ت » suffixé au verbe pour marquer l'accompli, lequel peut être pris pour une radicale R3. Il s'agit toutefois principalement de bruit résiduel.

Les tokens ne sont pas non plus totalement rejetés. En effet, nombre d'entre eux sont formés sur une racine trilitère saine et selon le schème ambigu « فعل ». Les racines sourdes de forme assimilée sont bien représentées par tokens issus de « کل » et subissent comme les noms et les verbes la forte influence de la post-concaténation. Les tokens non trilitères sont assez nettement rejetés car souvent constitués de deux lettres seulement. Ils se comportent comme les sourds assimilés dès qu'ils sont concaténés.

Parmi les bruits propres au modèle, le bruit de schème erroné est très important : il totalise 14 % des réponses. Le bruit d'incorrection morphologique s'avère en revanche très marginal : seuls trois cas se sont produits sur l'ensemble du texte.

#### c. Quelques mesures de la validité de l'analyse.

Une évaluation de la réussite de l'extraction de la racine peut s'effectuer au moyen d'un taux  $\tau\sqrt{}$  défini ainsi :

 $\tau \sqrt{}$  = (solutions accompagnées éventuellement de bruit radicalement exact)

/nombre de noms trilitères étudiés

L'influence du bruit le plus embarrassant, le bruit radicalement attesté, peut être mesurée par  $\tau_{bra}$ :

 $\tau_{bra}$  = bruit radicalement attesté/total des réponses

Le calcul de ces deux taux pour chacun des types de racine montre la grande capacité de Gn3(c,c,c) à extraire les racines saines de ce texte, mais aussi sa perplexité devant les autres cas qu'elle ne parvient pas vraiment à refuser, puisque  $\tau_{bra}$  varie alors de 20 à 50 %. Ces taux se réduisent à 15 % si l'on suppose un contrôle de l'existence du schème, qui devrait contrer efficacement le bruit dû à l'article.

#### d. Un exemple d'application : la recherche par racine.

Le fait que la racine saine soit détectée, même parmi du bruit, rend possible la réalisation d'une fonction de recherche des mots par leur racine, autrement dit une indexation restreinte aux cas sains d'un texte par les racines. Se donnant une racine saine attestée c<sub>1</sub>c<sub>2</sub>c<sub>3</sub>, l'opérateur désire trouver tous les mots du texte formés à partir de celle-ci. Il peut composer la grammaire Gn3(c<sub>1</sub>,c<sub>2</sub>,c<sub>3</sub>) pour laquelle les seuls arcs permis en R1, R2, R3 sont respectivement c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>. Les réponses non vides de cette analyse sont en effet toutes les interprétations fondées sur cette racine, et un calcul simple de probabilité montre que pour un texte de taille raisonnable il est très peu probable qu'un bruit radicalement attesté vienne troubler cette recherche.

#### EXTENSION DE Gn3 À LA MORPHOLOGIE NOMINALE NON SAINE.

Les cas de racines ne présentant qu'une seule irrégularité forment la majeure partie de la morphologie non saine. Les comportements de ces racines se manifestent de deux manières distinctes :

- soit sous une forme saine où la lettre défectueuse s'insère telle une consonne radicale dans le mot ;
- soit sous une forme irrégulière où la lettre revêt une autre forme ou disparaît par assimilation ou élision.

Remarquons que dans le premier cas la grammaire Gn3 serait capable de reconnaître le mot par simple adjonction d'un arc autorisant les semi-consonne w et y pour les racines cwc et cyc. Il est aussi possible de reconnaître la forme que prend éventuellement la lettre radicale dans le second cas. Le cas des racines concaves csc retiendra ici notre attention car il est le seul dans la morphologie irrégulière nominale qui ne soit pas sujet aux assimilations et élisions.

Une précision de vocabulaire s'impose : il est important de faire maintenant la distinction entre la racine telle qu'elle est apparue jusqu'ici et sa représentation dans un mot, qu'il soit nom ou verbe. Nous parlerons donc de représentants de R2 pour le cas des racines concaves. R2 étant alors soit w, soit y, son représentant noté R'2 peut être l'une des quatre lettres suivantes : « و », « ا », « و ». Les différentes possibilités de représentation sont données dans le tableau suivant <sup>20</sup>:

représentant R'2	R2	exemples : forme I	formes dérivées
، و	و	جوع	تكوين
ي	ي	أطياب جيل	مكيف
	و	أطياب جيل ثياب جيران	مكيف مستعيذ مقيم
1	و	باب	إعارة
	ي	عار	إمالة مهابة
5	و	سائق	
	ي	سائق بائع	

20. Dans le dernier cas, la lettre « بسائل » peut aussi représenter une hamza radicale R2 comme dans « سائل ».

Il est clair que les différents cas de représentation répondent à des conditions précises s'appuyant sur le type de dérivation. Si la simulation de ces conditions n'est pas possible actuellement sur Gn3, on peut en revanche par simple ajout des quatre représentants de R2 dans la catégorie R2 mesurer les conséquences d'un élargissement du champ d'application de la grammaire en terme de bruit vis-à-vis des mot à racines non concaves.

Considérons la nouvelle grammaire Gn3(c,cs,c) ainsi obtenue : cette variation de Gn3 répond aux conditions de la propriété de linéarité. Nous pouvons en conséquence nous contenter de restreindre dans une grammaire Gn3(c,s,c) la catégorie R2 aux seuls arcs nouvellement introduits pour analyser l'effet de l'introduction des représentants des semi-consonnes en seconde radicale  $^{21}$ .

#### SECONDE ANALYSE:

LA GRAMMAIRE DU NOM TRILITÈRE SAIN OU CONCAVE Gn3(c,cs,c).

Les tableaux B donnent les résultats de l'analyse du texte par la grammaire réduite Gn3(c,s,c). La classification des réponses a été appliquée comme si le champ d'application était bien réduit aux seuls noms concaves : c'est pour cette raison que les réponses aux noms ccc sont classées soit comme bruits, soit comme refus bien que que la partie saine de Gn3(c,cs,c) les accepte.

Si tous les noms *csc* fournissent au moins une réponse, les trois quarts d'entre eux sont tels que le représentant R'2 est R2. On accède donc dans ce cas directement à la racine de type *cwc* ou *cyc*. Les autres cas sont principalement ceux pour lesquels le représentant est « 1 » ou « 5 » mais aucun d'entre eux n'est véritablement ambigu car même si les deux possibilités *cwc* et *cyc* sont attestées, une seule des deux produit la forme détectée. Aussi avons-nous pris le parti dans ce dépouillement de ne plus nous restreindre aux seules attestations de cas de racines concaves, mais d'inclure aussi les formes I de ces racines. Il faut donc pour chaque interprétation d'un mot vérifier si elle peut provenir du *cwc* ou du *cyc* corrrespondant, et attribuer à chacun des deux cas le bruit adéquat.

Un mot tel « بيع » est ambigu car, bien que reconnu directement par la grammaire avec le bon radical R2 y, le recours à un lexique des formes I rappelle l'existence de « بيع » issu de « بوع ». Ce mot est compté comme solution mais aussi comme bruit radicalement attesté bien que la réponse ne comporte qu'une seule interprétation. À l'opposé, « ساق » est un bruit car la racine n'est reconnue qu'à travers son représentant alif, lequel n'est produit qu'avec la racine « سوق ». Cette dernière donne lieu à un bruit radicalement exact, et l'autre possibilité « سيق » n'étant pas attestée est comptée comme bruit résiduel.

21. La linéarité de cette variation s'exprime ici simplement par Gn3(c,cs,c) = Gn3(c,c,c) + Gn3(c,s,c). Cf. n. 16.

Tableaux B. – Dépouillement des réponses de Gn3(c,s,c).

					classe		typ	e brui	it			causes	i		évalı	ation
NOMS/Gn3(c,c,c)	type racine	total occurrences	total réponses (Nr)	silence rejet	solution	bruit	bruit rad. exact	bruit rad. attesté	bruit résiduel	schème inexact	pré-concat. (sauf article)	15 to 17	post-concat.	ambiguïté interne	succès extraction	taux bruit attesté
sains	ccc	284	18	267		18		3	15	4				17	0	16,7
sourds	cc=2	26	1	25		1			1					1	0	-
	cc=1	32	2	30		2			2				**************	2	0	-
une radicale S	csc ccs scc	108 48 28	125 7 11	44 19	80	45 7 11	52	3 2 2	56 5 9	1 1 4	1	4	10	30 6 6	74,1 0 0	2,4
une radicale H	hec chc cch	26 4 6	1 0 0	25 4 6		1			1					1	0 0 0	-
	ssc csh hc=	6 5 10 3	4 0 0	5 5 10 2		4			4	4		4	1		0 0 0	- - -
reste	cs= hsc scs sc=	9 2 1	2 0 1	7 2 1	***************************************	1 1			1		1 1			2	0 0 0 0	-
3 .	hcs	1	1	1		1			1				1		0	-
totaux sauf csc		<b>600</b> 492	175 50	<b>0 454</b> 0 454	<b>80</b> 0	<b>95</b> 50	52 0	10 7	<b>99</b> 43	14 13	3	<b>8</b> 8	13 3	<b>65</b> 35		5,71 14,00

VERBES/Gn3(c,s,c)	type racine	total occurrences	total réponses (Nr)	silence	rejet	solution	bruit	bruit rad. exact	bruit rad. attesté	bruit résiduel	schème inexact	pré-concat. (sauf article)	うぉつ	post-concat.	ambiguïté interne	succès extraction	taux bruit attesté
sains	ccc	104	2		102		2		1	3	1				1	0	-
sourds	cc=2	0														٠	
	cc=1	14	0		14				***********							0 -	0
une radicale S	esc ecs sec	55 25 12	30 0 0		29 25 12		30	27	3	30					30	49,1 0 0	10 - 0
une radicale H	hec ehe ech	0 1 11	0		1 11											0	- 0
reste	css chs csh hss hcs	5 3 2 1	1 0 0 0		4 3 1 1		1		-	3	1					0 0 0 0	
totaux sauf esc		<b>234</b> 179	<b>33</b> 3		<b>204</b> 175	<b>0</b> 0	<b>34</b> 4	<b>27</b> 0	<b>4</b> 1	<b>37</b> 7	<b>2</b> 2	<b>0</b> 0	0	0	31 1		-

#### **ANNEXE**

Texte traité : les deux premiers chapitres de Qindīl Umm Hāšem de Yaḥyā Ḥaqqī

١

كان جدى الشيخ رجب عبد الله إذا قدم القاهرة وهو صبى مع رجال الأسرة ونسائها للتبرك بزيارة أهل البيت ، دفعه أبوه إذا أشرفوا على مدخل مسجد السيدة زينب ، وغريزة التقليد تغنى عن الدفع - فيهوى معهم على عتبته الرخامية يرشقها بقبلاته ، وأقدام الداخلين والخارجين تكاد تصدم رأسه . وإذا شاهد فعلتهم أحد رجال الدين المتعللين أشاح بوجهه ناقها على الزمن ، مستعيداً بالله من البدع والشرك والجهالة ، أما أغلبية الشعب فتبسم لسذاجة هؤلاء القرويين - ورائحة اللبن والطين والحلبة تفوح من ثيابهم - وتفهم ما في قلوبهم من حرارة الشوق والتبجيل ، لا يجدون وسيلة للتعبير عن عواطفهم إلا ما يفعلونه : والأعمال بالنيات . وهاجر جدى - وهو شاب - إلى القاهرة سعياً للرزق . فلا عجب أن اختار لإقامته أقرب المساكن لجامعه المحبب . وهكذا استقر بمنزل للأوقاف قديم ، يواجه ميضأة المسجد الخلفية ، في الحارة التي كانت تسمى (حارة الميضة) . « كانت » لأن معول مصلحة التنظيم الهدام أتى عليه من معالم القاهرة . طاش المعول وسلمت للميدان روحه ، إنها يوفق في المحو والإفناء حين تكون ضحاياه من حجارة وطوب! ثم فتح جدى متجرا للغلال في الميدان أيضا . وهكذا عاشت الأسرة في ركاب « الست » وفي هماها : أعياد المسجد ساعتنا .

اتسع المتجر وبورك لجدى فيه - وهذا من كرامات أم هاشم - فها كاد يرى أبنه الأكبريتم دراسته في الكتاب حتى جذبه إلى تجارته ليستعين به ، وأما ابنه الثاني فقد دخل الأزهر ، واضطرب فيه سنوات وأخفق ، ثم عاد للدتنا ليكون فقيهها ومأذونها . بقى الابن الأصغر - عمى إسهاعيل آخر العنقود ، بهيئه القدر واتساع رزق أبيه لمستقبل أبهى وأعطر . لعله خشى في مبدأ الأمر ، عندما أجبره أبوه على حفظ القرآن أن يدفع به إلى الأزهر ، لأنه يرى صبية الميدان تلاحق الفتية المعممين بهذا الهتاف البذى :

[...]

ولكن الشيخ رجب سلمه ، بقلب مفعم بالآمال ، إلى المدارس الأميرية ، وعندئذ أعانته تربيته الدينية وأصله القروى فسرعان ما امتاز بالأدب والاتزان وتوقير معلميه ، مع حشمة وكبير صبر . إن حرم التأنق لم تفته

النظافة . وهو فوق ذلك أكثر رجولة و أقوم لساناً وأفصح نطقاً من زملائه (المدلعين) أولاد الأفندية المبتلين بالعجمة وعجز البيان ، فها لبث أن بذ الأقران وتلألأت على سيهائه نجابة لا تخطئها العين ، فتعلقت به آمال أسرته .

أصبح ، وهو لم يزل صبياً ، لا ينادى إلاب (سي إسهاعيل) أو إسهاعيل أفندى ، ولا يعامل إلا معاملة الرجال . له أطيب ما في الطعام والفاكهة .

إذا جلس للمذاكرة خفت صوت الأب ، وهو يتلو أوراده إلى همس يكاد يكون ذوب حنان مرتعش ، ومشت الأم على أطراف أصابعها ، حتى فاطمة النبوية – بنت عمه ، اليتيمة أبا و أما – تعلمت كيف تكف عن ثرثرتها وتسكن أمامه في جلستها صامتة كأنها أمة وهو سيدها . تعودت أن تسهر معه كأن الدرس درسها ، تتطلع إليه بعينيها المريضتين المحمرتي الأجفان ، وأصابعها تعمل في حركة متصلة لا تنقطع في بعض أشغال (التريكو) من ذا الذي يقول لإسهاعيل . تنبه إلى هاتين اليدين كيف دبت فيها خلسة حياة غريبة و حساسية يقظة ، ألا تفهم ألا تفطن إلى أن دليل اقتراب عاهة العمى في السليم هو أن تبدأ يده في الإبصار؟

[...]

بين حين وآخر تحيل دمعة مترقرقة شخصه إلى شبح مبهم فتمسحها بطرف كمها وتعود إلى تطلعها . الحكمة عندها تتمثل في كلامه إذا نطق .

يالله! كيف تحوى الكتب كل هذه الأسرار والألغاز؟ وكيف يقوى اللسان على الرطانة بلغة الأعاجم؟ وكلما كبر في نظرها انكمشت أمامه وتضاءلت. قد يعلق بصره بضفيرتيها فيتريث ويبتسم. هؤلاء الفتيات! لو يعلمن كم هى فارغة رؤوسهن!! إذا أوى إلى فراشه فعندئذ، وعندئذ حسب، تشعر الأسرة أن يومها قد انقضى، وتبدأ تفكر فيها يلزمه في الغد. كل حياتها وحركاتها وقف على توفير راحته. جيل يفنى نفسه لينشأ فرد واحد من ذريته. محبة وصلت من قوتها إلى عنفوان الغريزة الحيوانية. الدجاجة القلقة ذات النظرة المتجسسة الحذرة ترقد على بيضها مشلولة الحركة ذليلة العين، كأنها راهبة تصلى . . . هل هى هبات من فيض كرم؟ أم جزية جبار مستبد، إرادته حديد، له فى كل عنق طوق، وفى كل ساق قيد؟ تعلق هذه الأسرة بولدها تعلق مسلوب الحرية والإرادة! فأين بربك جماله؟ جواب هذا السؤال عند قلبى . فها من مرة تمثلت فيها هذه الأيام البعيدة إلا وجدته يخفق يذكراها، ويبدولي وجه جدى الشيخ رجب وحواليه هالة من وضاءة ونور . أما جدتى - الست عديلة ، بسذاجتها و طيبتها ، فمن السخف أن يقال إنها من البشر، وإلا فكيف إذا تكون الملائكة! ما أبشع الدنيا وأبغضها لو خلت من مثل تسليمها وإيهانها.

۲

سنة بعد سنة وإساعيل يفوز بالأولوية فإذا أعلنت النتيجة دارت أكواب الشربات على الجيران ، بل ربا شاركتهم المارة أيضا ، وزغردت ( ماشا الله ) بائعة الطعمية والبصارة وفاز الأسطى حسن - الحلاق ودكتور الحى - بحلوانه المعلوم وأطلقت الست عديلة بخورها وقامت بوفاء نذرها لأم هاشم. فهذه الأرغفة تعد وتملأ بالفول النابت وتخرج بها أم محمد تحملها في مقطف على رأسها : ما تهل في الميدان حتى تختطف الأرغفة ، ويختفى المقطف و تطير ملاءتها ، وترجع خجلة تتعثر في أذيالها غاضبة ضاحكة من جشع شحاذى السيدة وتصير حادثتها فكاهة الأسرة بضعة أيام يتندرون بها . وكذلك نشأ إسهاعيل في حراسة الله ثم أم هاشم . عياته لا تخرج عن الحي والميدان ، أقصى نزهته أن يخرج إلى المنيل ليسير بجانب النهر أو يقف على الكوبرى . إذا أقبل المساء وزالت حدة الشمس وانقلبت الخطوط والانعكاسات إلى انحناءات وأوهام ، أفاق الميدان إلى نفسه وتخلص من الزوار والغرباء إذا أصخت السمع وكنت نقى الضمير فطنت إلى تنفس خفى عميق يجوب الميدان لعله سيدى العتريس بواب الست - أليس اسمه من أسهاء الخدم ؟ - لعله في مقصورته ينفض يديه وثيابه من عمل النهار ، ويجلس يتنفس الصعداء . فلو قيض لك أن تسمع هذا الشهيق والزفير فانظر عندئذ الميا المقبد، لألاء من نور يطوف بها ، يضعف ويقوى كومضات مصباح يلاعبه الهواء . هذا هو قنديل أم هاشم الموجوه منه وقة القوى ، ذابلة الأعين ، يلبس كل منهم ما قدر عليه ، أو إن شئت : فها وقعت عليه يده من شيء فهو لابسه . نداءات الباعة كلها نغم حزين .

[...]

ما هذا الظلم الخفى الذى يشكون منه؟ وما هذا العبء الذى يجثم على الصدور جميعها؟ ومع ذلك فعلى الوجوه كلها نوع من الرضا والقناعة. ما أسهل ما ينسون! تتناول أيد كثيرة قروشاً وملاليم قليلة ليس هنا قانون ومعيار وسعر، بل عرف وخاطر وفصال وزيادة فى الكيل أو طبة فى الميزان. وقد يكون الكيل مدلساً والميزان مغشوشاً، كل بالبركة ، صفوف تستند إلى جدار الجامع جالسة على الأرض ، وبعضهم يتوسد الرصيف. خليط من رجال ونساء وأطفال ، لا تدرى من أين جاعوا ولا كيف سيختفون ، ثمار سقطت من شجرة الحياة فتعفنت فى كنفها . هنا مدرسة الشحاذين . حامل كيس اللقم يثقل الحمل ظهرة ينادى :

[...]

والشابة التي تنبت فجأة وسط الحارة عارية أو شبه عارية :

#### [...]

صوبها الصارخ يجذب الوجوه للنوافذ ، وعيناها الساحرتان تستهويان المطلات ، فتمطر عليها أكوام من الخرق ورث الثياب في لحظة واحدة تذوب وتختفى ، فلا تدرى أطارت ، أم ابتلعتها الأرض فغارت . وهذا بائع الدقة الأعمى الذى لا يبيعك إلا إذا بدأته السلام و أقرأك وراءه الصيغة الشرعية للبيع والشراء .

ينقضى النهار فيودع كرش الطرشجى بقية براميله ، وتترك أقدام الخراط عملها اليومى وأدواتها ، لتعود بصاحبها إلى الدار . لا يزال الترام هنا وحشاً مفترساً له فى كل يوم ضحية غريزة . يتقدم المساء ينعشه نسيم ذو دلال . تسمع من القهاوى ضحكات غضة وأخرى غليظة «حشاشى» . وإذا دلفت من الميدان إلى مدخل شارع مراسينه سمعت ضجيج السكارى فى خمارة أنسطاسى التى يلقبها أهل الحى بفكاهتهم خمارة أنست » . يخرج منها سكير هائج يتطوح ويتعرض للهارة :

#### [...]

أشباح الميدان الحزينة المتعبة يحركها الآن نوع من البهجة والمرح ليس فى الدنيا هم . والمستقبل بيد الله تتقارب الوجوه بود ، وينسى الوجيع شكايته . ويبذر الرجل آخر نقوده فى الجوزة أو الكتشينة وليكن ما يكون : تقل أصوات اصطدام كفف الموازين ، وتختفيى عربات اليد ، وتطفأ الشموع داخل المشنات ، عندئذ تنتهى جولة إسهاعيل فى الميدان . هو خبير بكل ركن وشبر وحجر ، لا يفاجئه نداء بائع ، ولا ينبهم عليه مكانه . تلفه الجموع فيلتف معها كقطرة المطريلقمها المحيط . صور متكررة متشابهة اعتادها فلا تجد فى روحه أقل مجاوبة لا يتطلع ولا يمل . لا يعرف الرضا ولا الغضب إنه ليس منفصلا عن الجمع حتى تتبينه عينه . من يقول له إن كل ما يسمعه ولا يفطن له من الأصوات . وكل ما نقع عليه عينه ولا يراه من الأشباح ، لها كلها مقدرة عجيبة على التسلل إلى القلب ، والنفوذ إليه خفية ، والاستقرار فيه ، والرسوب فى أعماقه فتصبح فى كل يوم عوامه . أما الآن فلا تمتاز نظرته بأية حياة . . نظرة سليمة ، كل عملها أن تبصر .

21A

#### PHRASE EXTRAITE DU TEXTE.

فهذه الأرغفة تعد وتملأ بالفول النابت وتخرج بها أم محمد تحملها في مقطف على رأسها : ما تهل في الميدان حتى تختطف الأرغفة ، ويختفى المقطف و تطير ملاءتها ، وترجع خجلة تتعثر في أذيالها غاضبة ضاحكة من جشع شحاذي السيدة وتصير حادثتها فكاهة الأسرة بضعة أيام يتندرون

### PHRASE ANALYSÉE PAR Gn3(c,c,c)

Nº	Mot arabe	Analyse	Racine	Schème
1	فهذه	coor/ه:R1 /غ:R2 /ه:R3	هذه	فعل
1	فهذه	e:R1 ه', R2:ه'/ 1:82 غ'/ 1:84 غار: R1:ه	فهذ	فعل
2	الأرغفة	sgfm:ة/ R3:ف/ R2:غ/ R1:ر/ipref! / h:art:	رغف	أفعلة
3	تعل	R3:د/ R2:ع/ R1:ت/	تعد	فعل
5	بالفول	R3:ل/erprep/ \:pref/ ن-/ R1:ل/prep/ \:pref! ب-/	لفل	افعول
6	النابت	R3:ت/ R2:ن/ R1:ن/ R1:ن/ art:ن/ :art:	نبت	فاعل
7	وتخرج	R3:ج/ R2:خ/pref:ت/coor:و/	خرج	تفعل
10	محمد	R3:د/ R2:م/ R1:ح/pref:م/	حمد	مفعل
11	تحملها	post/ ۱:post: ر/ R1:م/ R1:ح/pref:ت/	حمل	تفعل
13	مقطف	R3:ف/ R2:ق/pref:ق/pref:م/	قطف	مفعل
17	تهل	R1 (ن/ R2:ه/ R1:ت/	تهل	فعل
19	الميدان	duma:ن/l:duma:ن/R3 / 1:ch2:e/ اي/ R1:م/ 1:pref/	لد	افعيل
21	تختطف	R3:ف/ R2:ف/ :R1:ت/ R1:خ/pref:ت/	خطف	تفتعل
22	الأرغفة	sgfm:ة/ R3:ف/ R2:ف/ R1:ر/erf:أ / art:ل/ art:	رغف	أفعلة
24	المقطف	R3:ف/ R2:ظ/ R1:ق/pref:م/ art:ل/ art:	قطف	مفعل
26	تطير	R3:ر/2:inf2:ي/ R2:ك/ R1:ت/	تطر	فعيل
28	وترجع	R3:ج/ R1:ر/pref:ت/coor:و/	رجع	تفعل
29	خجلة	:R3 :sgfm:ق/ R2:خ/ R1:خ/	خجل	فعلة
33	غاضبة	R3 :sgfm:أ / R1:ف/ 1:inf1:غ/	غضب	فاعلة
34	ضاحكة	R1 / \:inf1/ :R2:خض/ R3:sgfm	ضحك	فاعلة
36	جشع	R3:ع/ R2:ج/	جشع	فعل
37	شحاذي	suff: ک/R2: خ/ R1:ش/	شحذ	فعالى
37	شحاذي	duma:ی/ R2:خ/ 1:inf2: ( R2:خ/ R1:ش/	شحذ	فعال
37	شحاذي	post: الـ R2: ج/ 1:inf3: الـ R1: ش/	شحذ	فعال
37	شحاذي	nisb:/ R2 / \:inf2/3:R3 (:inf2) اش/	شحذ	فعال
37	شحاذي	plma:// :R1:ح/ 1:inf2:/غ:R1:ش/	شحذ	فعال
38	السيدة	sinf2/۵:R3:« R3:«gfm:» / ۱:pref/ن/:R1:س/ R1:	لسد	افعيلة
39	وتصير	R3:ر/inf2:ي/ R2:ص/ R1:ت/coor:و/	تصر	فعيل
40	حادثتها	sgfm/: post/ \:post : ر/11: R1: \ 1: 1/1: R1: ح/	حدث	فاعلة
41	فكاهة	sgfm:«/ 1:e/2 / \:inf2/،:R3 (1:e/2:R2 / 1:e/2:e/2:e/2:e/2:e/2:e/2:e/2:e/2:e/2:e/2	فكه	فعالة
43	بضعة	sgfm:ة/ R3:ع/ R2:ض/ R1:ب/	بضع	فعلة

## PHRASE ANALYSÉE PAR Gn3(c,s,c)

N°	Mot arabe	Analyse	Racine	Schème
5	بالفول	R3:ل/ R2:و/ R1:ف/ art:ل/ R2:erep:/	فول	فعل
19	الميدان	duma:ن/ R1:ن/ R2:د/ R2:ي/ R1:م/ :art ال :art ال :art ال :art ال	مید	فعل
26	تطير	R3:ر/ R2:ي/ R1:ط/pref:ت/	طير	تفعل
32	أذيالها	r2 / \!int /ن:R3 (:int / :int	ذيل	فعال
32	أذيالها	post/ :post/ !:post: ر/R2:(/ R2:د/Post: د/pref:	ذيل	أفعال
38	السيدة	R1 :art /ن:art /	سيد	فعلة
39	وتصير	R3:ر/ R2:ي/ R1:ص/pref:ت/coor:و/	صير	تفعل
41	فكاهة	sgfm:ة/ R2: / \:R1 / \:R2:ف/	کاه	فعلة