



Quelques aspects de la perception de la richesse lexicale

Thoiron Philippe, Arnaud Pierre J. L.

Pour citer cet article

Thoiron Philippe, Arnaud Pierre J. L., « Quelques aspects de la perception de la richesse lexicale », *Cycnos*, vol. 8. (Apparences textuelles et réalité linguistique), 2008, mis en ligne en juillet 2008.
<http://epi-revel.univ-cotedazur.fr/publication/item/392>

Lien vers la notice <http://epi-revel.univ-cotedazur.fr/publication/item/392>
Lien du document <http://epi-revel.univ-cotedazur.fr/cycnos/392.pdf>

Cycnos, études anglophones

revue électronique éditée sur épi-Revel à Nice

ISSN 1765-3118 ISSN papier 0992-1893

AVERTISSEMENT

Les publications déposées sur la plate-forme épi-revel sont protégées par les dispositions générales du Code de la propriété intellectuelle. Conditions d'utilisation : respect du droit d'auteur et de la propriété intellectuelle.

L'accès aux références bibliographiques, au texte intégral, aux outils de recherche, au feuilletage de l'ensemble des revues est libre, cependant article, recension et autre contribution sont couvertes par le droit d'auteur et sont la propriété de leurs auteurs. Les utilisateurs doivent toujours associer à toute unité documentaire les éléments bibliographiques permettant de l'identifier correctement, notamment toujours faire mention du nom de l'auteur, du titre de l'article, de la revue et du site épi-revel. Ces mentions apparaissent sur la page de garde des documents sauvegardés ou imprimés par les utilisateurs. L'université Côte d'Azur est l'éditeur du portail épi-revel et à ce titre détient la propriété intellectuelle et les droits d'exploitation du site. L'exploitation du site à des fins commerciales ou publicitaires est interdite ainsi que toute diffusion massive du contenu ou modification des données sans l'accord des auteurs et de l'équipe d'épi-revel.

EPI-REVEL

Revue électronique de l'Université Côte d'Azur

Quelques aspects de la perception de la richesse lexicale

Philippe Thoiron

C.R.T.T. - Université Lumière-Lyon 2

Pierre J.L. Arnaud

C.R.T.T. - Université Lumière-Lyon 2

I. Introduction

On parle souvent de "richesse lexicale" : c'est un concept qui plaît*. Il est immédiatement saisissable par la métaphore qui le nomme. Cette facilité d'abord le dessert probablement car elle dispense d'une définition exhaustive qu'il n'est d'ailleurs pas facile de donner de façon satisfaisante (cf. Dugast 1978). A défaut d'une définition stricte rangée au rang des accessoires (Thoiron 1988), chacun accorde à ce concept les attributs qui lui semblent les plus convenables et les plus appropriés à son champ d'investigation. Certains y ont recours dans le cas d'études stylistiques, d'autres, en linguistique appliquée, l'utiliseront comme critère d'évaluation de la compétence lexicale acquise en L2 (compte rendu dans Arnaud, à paraître). Nous reviendrons sur ces divergences qui sont évidemment dues à la complexité d'un problème à première vue simple, mais le travail que nous présentons ici concerne un aspect à notre connaissance relativement inexploré de la richesse lexicale, à savoir sa perception par les lecteurs¹.

Définir et mesurer la richesse lexicale est une chose, en effet, mais déterminer si les destinataires du message y sont sensibles en est une autre. Une démarche expérimentale, parce qu'elle rend possible l'isolation et le contrôle des variables, permet d'apporter un début de réponse à ce problème de la sensibilité des destinataires. La question de recherche qui organise les deux expériences présentées ci-dessous est la suivante : confrontés à deux versions du même texte différant par leur richesse lexicale, des sujets y repéreront-ils des faits liés à la richesse ?

II. La richesse "lexicale"

II.1. Le concept : qu'entend-on par richesse lexicale ?

Le concept, fort utilisé par les stylostatisticiens, a fait l'objet de plusieurs études dont les plus récentes sont celles de Ménard (1983) et Thoiron, Labbé et Serant (1988). Si des divergences subsistent entre les chercheurs du domaine, il existe néanmoins un consensus qui repose sur quelques distinctions fondamentales, mais laisse subsister un flou terminologique qui peut parfois se révéler dangereux.

Tout d'abord, la richesse lexicale est un concept quantitatif : les jugements qualitatifs ne sont pas considérés comme pertinents. Ceci est, en particulier, manifeste lorsqu'on oppose, comme il convient, richesse lexicale et rareté lexicale. Ménard (1983:10) déplore leur confusion qui obscurcit souvent les débats sur le vocabulaire. Bernet (1988:2-3) rappelle que, si "le besoin d'évaluer l'étendue et la richesse du vocabulaire ne date pas d'hier", il convient de bannir les impressions et les idées préconçues. Celles-ci sont souvent liées à "l'illusion de la variété et de la richesse" qui "peut reposer sur la présence de mots rares ou de termes techniques".

Parler de critères quantitatifs pour fonder le concept de richesse lexicale, c'est indiquer que ce qui compte, et ce que l'on compte, c'est la proportion des vocables différents dans un texte. Selon la tradition des spécialistes français (en particulier à la suite de Muller 1964 : 235), on

¹ Les tests de perception de fréquence (cf. Schneider 1978), bien que concernant une problématique différente, ne sont évidemment pas étrangers au domaine de la présente étude.

note, pour un texte T, N le nombre de ses mots-occurrences (donc sa longueur) et V le nombre de ses vocables différents (donc l'étendue de son vocabulaire). C'est que, toujours à la suite de Muller, on distingue soigneusement entre le vocabulaire "réalité tangible et chiffrable" et le lexique qui a "un caractère virtuel" (Ménard 1983 :13).

Une autre distinction doit être établie entre richesse lexicale absolue et richesse lexicale relative. Parler de richesse lexicale absolue, c'est espérer comme Guiraud (1954), pouvoir trouver un indice absolu permettant de fixer un "degré zéro en dessous duquel le vocabulaire est pauvre" (Ménard 1983:12). Cette démarche critiquée par beaucoup a conduit à l'abandon du concept lui-même et l'on s'accorde maintenant à admettre le caractère relatif de la richesse lexicale (Muller 1979:282). On ne dit plus d'un texte qu'il est riche en soi mais qu'il est plus riche que tel autre.

On aura relevé, à travers ces quelques mises au point, que ce qui est en vérité objet d'étude et objet de consensus parmi les spécialistes, c'est *stricto sensu* la richesse du *vocabulaire*. C'est par commodité (mais aussi par abus de langage) que l'on parle de richesse lexicale, car, en fait, ce n'est pas le lexique, au sens que Muller donne à ce terme, qui est l'objet de la mesure. On n'attache généralement aucune importance à ce glissement terminologique en stylostatistique littéraire parce qu'il ne génère aucune confusion : on sait bien que c'est de la richesse des oeuvres qu'il s'agit dans la grande majorité des études. On verra plus loin que, s'agissant de l'étude de la *perception* des textes, la terminologie doit être fixée avec plus de rigueur.

II.2. La mesure de la richesse lexicale et les problèmes qu'elle pose

La base des études de richesse lexicale ainsi fixée, restent les problèmes de sa mesure. Le rapport V/N, intuitivement perçu comme un bon estimateur, est tellement sensible aux variations de N qu'il ne peut guère être utilisé qu'avec des textes de longueurs très voisines. Tous les spécialistes se sont donc efforcés de trouver soit des indices (Brunet avec w (Brunet 1978), Dugast avec U (Dugast 1979), par exemple) soit des méthodes (raccourcissement de textes, caractérisation quantitatives de la distribution des classes de fréquences, etc.) pour contourner cet obstacle. Il semble bien qu'avec les travaux les plus récents (cf. par exemple bibliographie dans Thoiron et al. 1988) on soit maintenant en mesure de comparer finement la richesse lexicale de textes de longueurs très différentes. On a recours, pour cela, à la distribution des fréquences, ou, en d'autres termes, à l'ensemble des entiers formé par les effectifs des vocables ayant la même fréquence dans le texte. Ces effectifs sont notés classiquement V_1 (nombre de vocables de fréquence 1), V_2 , ..., V_i , ... Enfin, la relation entre mesure de la richesse lexicale et taux de répétitivité lexicale a fait l'objet de plusieurs études récentes : Lafon (1984) s'intéresse à la localisation des formes répétées; Serant et Thoiron (1988) cherchent à inclure dans les mesures de la richesse lexicale les éléments relatifs à la topographie des vocables répétés.

II.3. La richesse lexicale dans le cadre du présent travail

En optant pour la relativité du concept de richesse lexicale, on peut mieux poser le problème de la relation richesse lexicale/rareté. On juge en effet de deux textes et de leurs V respectifs sans se préoccuper du lexique L dont les vocabulaires sont issus. On n'ignore pas le lien entre L et V, on choisit de ne pas en tenir compte. C'est d'ailleurs une des raisons pour lesquelles l'adjectif lexical fonctionne sans inconvénient comme équivalent à qui concerne le vocabulaire.

La rareté d'un lexème ne peut être mesurée qu'en termes de fréquence dans L. Or, on a exprimé (Muller 1968 :136,141; 1979 :440) les plus expresses réserves sur la notion de fréquence en langue. Mais, si, là encore, on renonce à une mesure absolue pour se satisfaire d'une mesure relative, mais opératoire, on peut admettre qu'on parle de rareté relative de

lexèmes. On rejoint ainsi le concept de fréquence intuitive qui semble bien fonctionner chez les locuteurs natifs (cf. Schneider 1978 pour la fréquence expérimentale et Arnaud 1989 pour la fréquence d'usage)². Il est alors possible de dire que tel lexème est plus, ou moins, fréquent que tel autre.

Le présent travail s'intéresse à la perception de faits relatifs à la variété, à l'étendue, des vocables d'un texte (ou de plusieurs) mais aussi à leur rareté. Il convient donc de construire un système où cohabitent et interagissent richesse lexicale (relative) et rareté lexématique. La richesse lexicale est mesurée par le biais d'un dénombrement de vocables différents, au sein du texte lui-même, la rareté lexématique est évaluée par référence aux résultats d'un autre dénombrement issu d'un ensemble extérieur au texte et considérablement plus grand que lui (un corpus donc). Pour la langue française, on admettra que le dépouillement du TLF constitue le meilleur dénombrement disponible et que les fréquences relevées permettent d'établir une hiérarchisation entre les lexèmes, à défaut de fournir une valeur absolue et indiscutable.

Il importe, pour une bonne hygiène méthodologique, de bien saisir la différence de nature entre les diverses comparaisons qui peuvent être conduites. Comparer des valeurs de V entre elles, c'est ne pas sortir de l'univers du discours. Introduire des critères de rareté, c'est pénétrer dans l'univers de la langue, fût-ce par la porte étroite/dérobée d'une estimation sur corpus.

Il nous a semblé que, pour satisfaire aux exigences de cette hygiène, il serait préférable de préciser la terminologie que nous utiliserons, ainsi que ses implications. *Dans le cadre de cette étude* nous distinguerons entre richesse lexicale, richesse du vocabulaire et rareté lexématique. La richesse lexicale est un ensemble, qui englobe richesse du vocabulaire et rareté lexématique. Ce qui est étudié ici c'est l'ensemble lui-même. Il nous apparaît, en effet, que, dans un premier temps, il n'est pas utile de distinguer plus avant. Bien entendu, si les résultats obtenus y invitent, il faudra poursuivre en sériant les questions.

III. Expériences

III.1. Texte et modifications

Pour modifier la richesse du vocabulaire d'un texte, on peut donc intervenir sur plusieurs paramètres : sur N sans toucher à V, ou sur V sans toucher à N, ou bien encore sur les deux simultanément. La distribution des fréquences peut être affectée aussi : on peut transférer un vocable de la classe de fréquence i à une classe des vocables de fréquence $(i \pm d)$, où d est un entier quelconque. Enfin, on peut modifier la topographie des formes répétées sans changer les effectifs des classes.

Dans le cas présent, nous avons délibérément évité de modifier sensiblement N, notre contrainte initiale étant même de le laisser invariant³. Ainsi les seules variations considérées concernent-elles le vocabulaire du texte. Ces modifications sont, elles aussi, possibles de plusieurs façons. On peut fort bien procéder à des changements vocable pour vocable de telle manière que le total V reste lui aussi invariant. Les modifications sont alors de nature strictement qualitative. On se rendra vite compte que la richesse du vocabulaire est inchangée et que c'est la rareté lexématique qui a varié. On peut aussi procéder à des changements qui influenceront sur la richesse du vocabulaire et sur la rareté lexématique conjointement. C'est d'ailleurs, le résultat, linguistiquement quasi-inévitable, de toute modification du vocabulaire qui vise à produire un texte appauvri. On comprendra facilement qu'en changeant un vocable

2 Par ailleurs, le lien entre la fréquence des mots d'un texte et sa lisibilité est connu depuis longtemps (cf. Klare 1968)

3 Dans un nombre très limité de cas (remplacement d'un verbe intransitif par un verbe transitif : *recourir* à par *donner* (une explication) nous avons dû déroger.

pour un autre dont la rareté lexématique est moindre on a plus de chances de réemployer un lexème déjà présent dans le texte que dans le cas inverse où le changement aurait pour visée une complexification lexicale, par recours à des lexèmes plus rares forcément.

Dans le cadre expérimental, il nous a semblé préférable d'aller dans le sens de l'appauvrissement du texte plutôt que de son enrichissement (i.e. par le biais de la diminution et non de l'augmentation de la rareté lexématique). Nous avons choisi de remplacer des vocables présents dans le texte par des vocables moins rares chaque fois que c'était possible tout en gardant un texte instinctivement lisible (en fait, 10% du vocabulaire a été ainsi modifié). Pour apprécier le degré de rareté, nous avons utilisé, en dépit des inconvénients liés à son caractère exagérément littéraire, le TLF et ses fréquences pour la "deuxième moitié du XXème siècle" rapportées à 100 millions d'occurrences.

Contrairement à ce qu'on pourrait croire *a priori* ce genre de modifications pose des problèmes redoutables, en eux-mêmes et par leur corrélation. Les difficultés sont d'abord de nature linguistique bien sûr, mais leurs implications statistiques ne sont pas négligeables non plus.

Au plan linguistique, la contrainte de la substitution vocable pour vocable est difficile à respecter. On se rend très vite compte que l'absence de vraie synonymie pose des problèmes qui concernent :

- la substitution vocable monosème / vocable polysème
- les collocations
- les nominalisations.

La question de la monosémie se manifeste très rapidement puisqu'on sait bien qu'il existe une corrélation nette entre rareté et monosémie. En substituant à un lexème du texte un lexème moins rare, on risque d'obtenir un lexème plus polysémique.

Les problèmes sont particulièrement sensibles au niveau des faits de collocation. Tel vocable, inséré dans une collocation (même assez faiblement figée, *animaux stylisés* par exemple), ne peut être facilement échangé au profit d'un autre (*simplifié*) qui, bien que bon équivalent *a priori*, "refuse" de s'intégrer dans le contexte.

L'abandon des nominalisations constitue une modification simple du point de vue lexical. Leur rareté est généralement supérieure à celle des verbes correspondants et la rareté est moindre lorsqu'on substitue la forme verbale. Mais c'est au plan syntaxique que se trouvent les obstacles, prévisibles au demeurant. On peut être amené à modifier plus qu'on ne le souhaite au départ un membre de phrase où une relation de prédication doit être substituée à une relation de détermination (ex. *la théorie est juste* pour *la justesse de la théorie*) et le principe de la substitution lexème pour lexème est mis en échec.

Ces difficultés à trouver des équivalences mono-élémentaires ont des répercussions au plan quantitatif. Certaines sont triviales et ne seront pas développées. Il s'agit bien sûr des jeux sur N et V, donc sur V/N. Plus délicats sont les faits de modification de la répétitivité de certains vocables, donc, à terme, de la topographie du texte. Le remplacement de *accomplir* par *faire* induit une augmentation du taux de répétitivité ainsi que de l'indice de topographie. Or, ces phénomènes statistiques sont à leur tour objets de perception. La répétition d'un vocable de haute fréquence n'est pas perçue de la même manière que celle d'un vocable de fréquence plus faible. La répétition de *accomplir* n'a pas, à cet égard, la même valeur que celle de *faire*. On voit reparaître ainsi la rareté lexématique dans le cadre de la mesure de la topographie des formes répétées.

Nous avons décidé de retenir un document privilégiant le contenu et ne s'écartant pas trop, d'un point de vue typologique, de ce que nos sujets ont l'habitude de lire. Un passage de prose informative, extrait de *Science et Vie* (novembre 1983), traitant du mystère des figures géantes tracées sur le sol près de Nazca au Pérou a été choisi.

La longueur du texte a été adaptée aux conditions expérimentales. Deux expériences ayant été prévues, il fallait que la plus longue des deux puisse se dérouler en 40 minutes maximum, que

les deux versions du même texte puissent être lues et que des réponses (brèves) à trois questions soient formulées. Nous avons choisi le début du texte original l'interrompant lorsque la longueur idoine (700 mots-occurrences environ) a été atteinte. Aucun fragment de texte n'a été omis, à l'exclusion du chapeau de l'article.

Les modifications lexicales ont été faites selon des contraintes simples et peu nombreuses :

- substitution d'un vocable par un vocable moins rare;
- réemploi si possible d'un vocable déjà présent dans le texte.

Beaucoup de nos choix contiennent une part d'arbitraire. On aurait évidemment pu en faire d'autres et il y avait certainement d'autres versions appauvries possibles.

Le lecteur trouvera, en annexe, le texte original et les modifications que nous y avons apportées. Le Tableau 1 permet une comparaison des deux versions sur divers paramètres.

Tableau 1 : Comparaison des deux versions du texte

Version	normale	appauvrie
N	716	714
V	332	305
V/N	0.46	0.43
N _{LEX}	318	321
V _{LEX}	233	206
fréq. moy.	47319	54021

Tableau 3 : Impressions sur les différences

Ordre des textes →	normal	appauvri
réponse ↓	appauvri	normal
Exactes	9	4
Pas nettes	4	2
Fausse	1	1
« différence non repérée »	1	0

Tableau 4 : Repérage de la version la plus riche

Ordre des textes →	normal	appauvri
réponse ↓	appauvri	normal
Exactes	15	7
Fausse	0	2
« non réponse »	0	1

III.2. Expérience 1

La question à laquelle on tente de répondre par la première expérience est la suivante : des sujets à qui on fait lire successivement deux versions du même texte, l'une normale et l'autre appauvrie, sont-ils capables de repérer les différences entre les deux versions ?

Les sujets étaient des étudiants de DEUG L.E.A. (1ère et 2ème année, donc âgés de 18 à 22 ans pour la plupart), convoqués à l'extérieur de leur horaire normal; leur participation était volontaire. Il s'agissait pour les 9/10 de filles.

Les 39 étudiants furent répartis en trois groupes. Le premier groupe reçut, selon la procédure décrite ci-dessous, d'abord le texte normal, puis le texte appauvri. Le deuxième groupe reçut

les textes dans l'ordre inverse afin de contrôler la variable "ordre de lecture". Le reste des étudiants constituait un groupe de contrôle et reçut deux fois le texte normal.

Les trois groupes participèrent ensemble à l'expérience. Les étudiants furent prévenus qu'il s'agissait d'une expérience anonyme sur la lecture, sur les buts de laquelle nous ne pouvions rien leur dire avant qu'elle soit terminée. Il leur fut demandé de lire les textes une seule fois, sans revenir en arrière, "à vitesse normale, comme ils liraient chez eux un magazine". Ils reçurent ensuite une feuille de réponses, identique pour les trois groupes et ne comportant que des emplacements numérotés, à l'exclusion de toute autre indication. Le premier texte fut ensuite distribué, puis, lorsque tous les sujets eurent terminé leur lecture (avec des temps qui, incidemment, allaient du simple au double selon les individus), le second texte. Les deux lectures eurent donc lieu sans que les étudiants sachent ce qui allait leur être demandé. Les textes furent alors ramassés, et les sujets furent invités à répondre sur les feuilles prévues aux questions qui leur étaient lues selon le protocole suivant :

"Certains d'entre vous ont reçu deux fois exactement le même texte (ce qui n'empêche pas les lettres d'identification d'être différentes); les autres ont reçu des textes légèrement différents. Si vous avez l'impression d'avoir lu deux fois le même texte, inscrivez à l'emplacement n°1 les mots 'même texte'. Si par contre vous avez l'impression qu'il y avait des différences entre les deux textes, inscrivez 'différences'."

La répartition des réponses est indiquée au Tableau 2.

Tableau 2 : Réponses sur le statut respectif des deux textes lus

Ordre des textes →	normal	appauvri	normal
réponse ↓	appauvri	normal	normal
« même texte »	3	1	<u>7</u>
« des différences »	15*	10*	3

Les effectifs correspondant aux réponses exactes sont dotés d'un astérisque

Les étudiants ayant répondu "*même texte*" furent priés de ne pas répondre aux questions suivantes. Les autres furent invités à inscrire à l'emplacement n°2 leurs "*impressions sur ces différences entre les deux textes*"; pour le cas où ils s'estimaient incapables d'analyser ces différences, ils devaient répondre "*différence non repérée*". Les réponses fournies ont été triées en "justes" (c'est-à-dire citant expressément le vocabulaire ou donnant des exemples lexicaux), "pas nettes" (portant sur des faits de forme pouvant englober le vocabulaire, mais sans que celui-ci soit mentionné), et "fausses" (exemple: "il n'y a pas de plan"). Elles apparaissent au Tableau 3.⁴

Pour la troisième question, le texte suivant fut lu :

"Nous allons vous mettre sur la voie. L'un des deux textes avait un vocabulaire plus riche que l'autre. Inscrivez la lettre d'identification du texte qui vous a semblé le plus riche. Si vous n'êtes pas capable de répondre à cette question, inscrivez 'non réponse'".

Les effectifs correspondant à ces réponses apparaissent au Tableau 4.

Tableau 4 = page manquante de la version papier (p. 40)

Les données des tableaux ont fait l'objet d'analyses statistiques au moyen de tests de χ^2 , avec correction de Yates pour les petits effectifs (Schwartz 1963 : 99). On a testé:

- Le repérage du statut relatif des deux textes (différence vs. identité). Le χ^2 sur l'ensemble du Tableau 2 (colonnes 1 et 2 regroupées) est égal à 9,12; pour 1 degré de liberté, il est significatif au seuil $\alpha = 0,01$, et permet de dire que, lorsque deux textes différents ont été distribués, les sujets ont perçu l'existence d'une différence.

⁴ Seules les réponses des étudiants ayant répondu correctement à la troisième question sont reprises dans le tableau 3, ce qui entraîne une différence d'effectifs par rapport au tableau 2.

- L'incidence de l'ordre de lecture sur le repérage de la différence entre les textes distribués. Le χ^2 - sur le Tableau 2, colonnes 1 et 2 seulement - très proche de zéro, montre que l'ordre de lecture est sans effet sur ce repérage.
- L'incidence de l'ordre de lecture sur la capacité à repérer le texte le plus riche. Le χ^2 sur le Tableau 4, lignes 1 et 2 seulement, est de 1,10; il est non significatif pour 1 d.d.l. même avec $\alpha = 0,20$, et permet de montrer que les sujets ont été capables de voir quel était le texte le plus riche, qu'ils l'aient lu en premier ou en dernier lieu.

Cependant, même si les sujets ne savaient pas ce qui était recherché dans l'expérience, il est clair que le fait d'avoir à lire deux fois le même texte ou bien deux versions d'un même texte occasionnait une activité ouvertement métalinguistique. Pour cette raison, la seconde expérience fut mise sur pied.

III.3. Expérience 2

Les 53 sujets de cette expérience, semblables à ceux de la première, mais dont aucun n'y avait participé, furent soumis à une procédure très semblable, avec lecture d'instructions préalables, puis distribution d'une feuille de réponses sans indications exploitables. Un seul texte leur fut cependant distribué, à savoir la version normale pour moitié d'entre eux et la version appauvrie pour les autres. Après lecture, les textes furent ramassés; jusqu'à ce point, les sujets ne savaient pas sur quoi portait l'expérience, qui aurait fort bien pu concerner le contenu du texte. Les questions suivantes furent alors posées:

"Du point de vue du style, de la forme, avez-vous remarqué quelque chose de particulier ? Si vous n'avez rien remarqué, inscrivez 'non'; si vous avez remarqué quelque chose, dites quoi ou bien, si vous n'avez pas une idée claire, écrivez "repéré quelque chose, mais je ne sais pas quoi".

Pour les deux versions du texte, les réponses *non* furent majoritaires. Les réponses détaillées produites se révélèrent toutes fausses, portant pour la plupart sur le contenu du texte (malgré les instructions reçues). Les effectifs concernés apparaissent au Tableau 5. Après regroupement des trois dernières lignes, on a testé l'incidence de la richesse sur la capacité à repérer un fait de forme. Avec un χ^2 non significatif (0,55), il est possible de penser que les sujets n'ont pas été sensibles aux faits de forme, c'est-à-dire que, quel que soit le texte lu, ils n'ont rien repéré de particulier.

Tableau 5 : Repérage de faits de forme remarquables dans le texte lu

Texte lu →		normal	appauvri
« rien remarqué »		18	16
Réponse Détaillées → Fournie	réponse juste	0	0
	réponse fausse	6	5
	« ne sait pas quoi »	2	6

Les consignes suivantes furent ensuite données:

"Nous allons vous mettre sur la voie : l'un des textes distribués est une version délibérément appauvrie de l'autre. Maintenant que vous le savez, pensez-vous avoir eu affaire à un texte 'normal' ou à un texte 'pauvre' ?

Les résultats apparaissent au Tableau 6. On a testé alors la capacité des sujets ayant reçu le texte appauvri à le repérer comme tel. Le χ^2 de 0,40, non significatif, permet de croire que, malgré l'information supplémentaire reçue, l'ensemble des sujets prenant part à l'expérience a cru avoir affaire à un texte normal.

Tableau 6 : Repérage du statut du texte lu

Texte lu	normal	appauvri
Réponse :		
« texte normal »	22*	21
« texte pauvre »	4	6*

Les effectifs correspondant aux réponses exactes sont dotés d'un astérisque

Ceci est confirmé par les réponses à la troisième question, qui invitait les sujets ayant répondu *pauvre* à la question 2 à expliciter ce qui leur avait semblé *pauvre* ("*grammaire, vocabulaire, articulations, ...etc.*"). Aucun des étudiants ayant répondu à juste titre "*pauvre*" à la question 2 (les 6 soulignés dans le Tableau 6) ne fournit ici de réponse juste (cinq fournissent des réponses fausses et un se déclare incapable de fournir une réponse précise). La seconde expérience indique donc que les sujets n'ont pas eu de réactions différentes selon la version lue; en d'autres termes, l'appauvrissement lexical du texte est passé inaperçu.

IV. Discussion et conclusion

Avant de tirer des conclusions de ces résultats, il n'est pas inutile de revenir sur les limitations des expériences qui les ont produits. On a vu ci-dessus (3.1) que la méthode d'appauvrissement lexical adoptée, à savoir la substitution vocable pour vocable, avait permis de contrôler, bien sûr la variable *contenu du texte*, mais aussi des variables telles que la complexité syntaxique et la cohésion non-lexicale. On pourrait toutefois reprocher à cette méthode d'avoir induit des effets stylistiques, donc des effets repérables extra ou para-lexicaux, par l'introduction de lexèmes plus fréquents qui pourraient par exemple produire un effet d'abaissement du niveau de langue. On a vu également que le remplacement d'un terme d'une collocation risque de rendre le syntagme résultant plus repérable. Une substitution peut en effet jouer non seulement sur le mot-occurrence visé, mais aussi sur son environnement. Ces effets indésirables apparaissent cependant comme tout à fait marginaux quand on examine la version appauvrie.

Par ailleurs, il ne faudrait pas s'exagérer l'utilité du critère de rareté lexématique retenu. En effet, les listes de fréquences courantes, et celle du TLF en particulier, ne regroupent pas les fréquences des dérivés; or, il y a tout lieu de penser que les dérivés ne sont pas stockés dans le lexique mental indépendamment de leur base (cf. Aitchison 1987 : 107 & sq.). Ceci n'a pu toutefois introduire que des imprécisions mineures, et il est évident que la plus grande fréquence des vocables de la version appauvrie constitue une donnée fiable.

On pourrait en outre objecter que si les sujets de l'Expérience 2 n'ont pas repéré la version appauvrie, c'était parce que celle-ci n'était pas assez pauvre. Or, comme on l'a vu, les spécialistes du domaine ont renoncé à établir des critères absolus de la richesse lexicale d'un texte (au sens où nous l'entendons ici) qui permettraient de fonder un tel jugement, que ce soit en termes de richesse du vocabulaire ou en termes de rareté lexématique, pour lesquelles on ne dispose ni de seuils quantitatifs, ni de textes-étalons auxquels on pourrait comparer les textes expérimentaux. En outre, on remarquera que dans l'Expérience 1, les sujets ont repéré correctement la différence de richesse entre les deux versions.

La lecture d'un texte peut être envisagée comme une interaction entre le lecteur et le texte. Le texte n'apporte des informations au lecteur que si celui-ci apporte à sa lecture sa compétence linguistique et cognitive et sa connaissance du monde, qui permettent un double mouvement de déclenchement et de génération d'hypothèses de bas en haut et de haut en bas. Mais le lecteur apporte également à la lecture ses intérêts et sa motivation. Le texte utilisé a plu aux sujets - certains l'ont indiqué par écrit en répondant à la question sur leurs impressions - et l'intérêt pour son contenu a pu entraîner une grande transparence de la forme, malgré la situation expérimentale. On ne peut exclure que le même dispositif expérimental aboutisse à des résultats différents si on utilise un texte dans lequel la forme serait plus présente, comme un texte ouvertement littéraire.

Nos conclusions seront motivées par deux sources. En premier lieu, les réponses détaillées fournies par les sujets de l'Expérience 1 indiquent que beaucoup d'entre ceux qui ont repéré la différence de richesse lexicale entre les deux textes ont explicitement comparé les vocables correspondants rencontrés lors des deux lectures. Il n'est pas possible de savoir si leur

repérage de la différence repose également sur la richesse du vocabulaire (V/N, rappel), mais la remarque de Bernet (1988) citée plus haut reçoit une justification : certaines occurrences de faible fréquence ou frappantes pour d'autres raisons jouent sans doute un rôle important dans la perception de la richesse lexicale d'un texte. Il n'est pas impossible que celle-ci repose sur l'effort mental entraîné par des vocables d'accès difficile. Nous en concluons que la prise en compte de la rareté lexicématique est nécessaire dès lors que l'on ne s'intéresse plus à l'objet-texte seul, mais aussi à sa réception.

Notre seconde conclusion repose sur la comparaison des résultats des deux expériences, d'où il ressort qu'une pauvreté lexicale repérée en situation de comparaison métalinguistique est passée inaperçue lors d'une lecture simple. On peut en déduire qu'il existe très probablement une grande tolérance des lecteurs vis-à-vis de la pauvreté lexicale, et que celle-ci n'est spontanément repérable que dans des cas extrêmes, devoirs d'écoliers ou productions de locuteurs non-natifs, textes délibérément écrits avec un vocabulaire réduit (par exemple le traité de Nation (1983) sur l'enseignement du vocabulaire anglais⁵) ou bien textes pour apprenants à vocabulaire contrôlé - mais ceux-ci ne s'adressent de toute façon pas à des natifs et, de plus, leur réduction lexicale est clairement annoncée. Le repérage spontané de la richesse (et non plus de la pauvreté) lexicale d'un texte pose d'autres problèmes, et la mise sur pied d'un paradigme expérimental qui permette de distinguer les effets de la richesse du vocabulaire de ceux de la rareté lexicématique de textes naturels constitue un défi à l'astuce des chercheurs.

Aitchison, J. (1987). *Words in the Mind : An Introduction to the Mental Lexicon*. Oxford: Blackwell.

Arnaud, P.J.L. (1989). Estimations subjectives des fréquences des mots. *Cahiers de Lexicologie*. 54, 1. 69-81.

Arnaud, P.J.L. (à paraître). Objective lexical and grammatical characteristics of L2 written compositions and the validity of separate-component tests. in P.J.L. Arnaud et H. Béjoint (eds.). *Vocabulary and Applied Linguistics*. Londres: Macmillan.

Bernet, Ch. (1988). Faits lexicaux, richesse du vocabulaire : Résultats. in Thoiron et al. 1-11.

Brunet, E. (1978). L'analyse statistique du TLF. *Le Français Moderne*. 46, 1. 54-66.

Dugast, D. (1978). Sur quoi se fonde la notion d'étendue théorique du vocabulaire ? *Le Français Moderne*. 46, 1. 25-32.

Dugast, D. (1979). *Vocabulaire et stylistique*. Genève : Slatkine.

----- (1971). *Etudes statistiques sur le vocabulaire français : Dictionnaire des fréquences*. Nancy: CNRS-TLF.

Guiraud, P. (1954). *Les Caractères statistiques du vocabulaire*. Paris : P.U.F.

Klare, G.R. (1968). The role of word frequency in readability. in J.R. Bormuth (ed.). *Readability in 1968*. s.l. (U.S.A.) : National Conference on Research in English.

Lafon, P. (1984). *Dépouillements et statistiques en lexicométrie*. Genève : Slatkine.

Ménard, N. (1983). *Mesure de la richesse lexicale*. Genève : Slatkine.

Muller, Ch. (1964). Calcul des probabilités et calcul d'un vocabulaire. *TraLiLi*. 2. 235-244.

Muller, Ch. (1968). *Initiation à la statistique linguistique*. Paris : Larousse.

Muller, Ch. (1977). *Principes et méthodes de statistique lexicale*. Paris : Hachette.

Muller, Ch. (1979). *Langue française et linguistique quantitative*. Genève : Slatkine.

Nation, I.S.P. (1983). *Teaching and Learning Vocabulary*. Wellington : Victoria University of Wellington, English Language Institute.

Schneider, A. (1978). La fréquence lexicale: Test de perception. *Le Français Moderne*. 46, 1. 6-11.

5 Dont l'un d'entre nous (P. A.) s'était effectivement aperçu en cours de lecture de la pauvreté lexicale (voulue par son auteur en raison du public visé).

Schwartz, D. (1963). *Méthodes statistiques à l'usage des médecins et des biologistes*. Paris: Flammarion.

Serant, D., et Ph.Thoiron (1988). Richesse lexicale et topographie des formes répétées. in Thoiron et al. 125-139.

Thoiron, Ph. (1988). Richesse lexicale et classement de textes. in Thoiron et al. 141-163.

Thoiron, Ph., D.Labbé et D.Serant (1988). *Etudes sur la richesse et la structure lexicales*. Genève: Slatkine.

Texte utilisé pour les expériences. Les occurrences substituées pour la version "appauvrie" sont indiquées entre crochets, sans que toutes les modifications aient été effectuées. Les textes distribués avaient subi toutes les modifications nécessaires.

Le mystère de Nazca

Avec le "secret" de la Grande Pyramide, les terrasses de Baalbek et les mégalithes [pierres] de Carnac ou de Stonehenge, les dessins du désert de Nazca sont un des morceaux de choix des amateurs d'archéologie fantastique. Il faut convenir que même pour les esprits les moins enclins [portés] à vagabonder dans l'imaginaire, ces figures [dessins] étranges [fantastiques] posent une énigme.

Nazca est le nom d'une petite ville qui émerge [apparaît] comme une oasis au milieu des plateaux désolés de la côte méridionale [sud] du Pérou. C'est tout près de là qu'apparaissent, visibles surtout d'avion, des dessins pour la plupart de très grandes dimensions [taille] - quelques-uns mesurent [font] plus de deux cents mètres - éparpillés sur une cinquantaine de kilomètres. Ils ont été tracés [faits] en lignes claires en écartant simplement les gravillons [pierres] du désert, dont la teinte [couleur] est plus foncée que celle du sol sous-jacent [en-dessous] . Presque tous représentent des animaux plus ou moins stylisés - des oiseaux en vol, une araignée, un poisson, un singe - ou bien des bandes [lignes] rectilignes [droites] et des figures [dessins] symboliques, plus difficiles à déchiffrer [comprendre]. Mais, vus du sol, leurs dimensions [tailles] sont telles, que l'on ne comprend pas ces dessins au premier coup d'oeil : pour bien les voir il faut les survoler, et c'est la photographie aérienne qui les a rendus célèbres.

On voit les problèmes qui se posent. Quels artistes inconnus ont donc créé [fait], dans le désert, cette oeuvre [travail] gigantesque [immense] qu'on ne peut admirer [voir] que du ciel ? Dans quel but ? Et quelles techniques ont-ils bien pu mettre en oeuvre pour respecter [suivre] les proportions de leurs modèles, alors que l'homme lui-même ne pouvait embrasser [découvrir] du regard l'ensemble d'un dessin ? Pour faire bref : qui ? pourquoi ? comment ?

Bien des monuments du passé, en Amérique du Sud, restent encore à explorer [découvrir] et à interpréter [comprendre] . Mais les spécialistes de la littérature fantastique vont plus vite que les savants. Il est vrai que pour satisfaire leurs lecteurs, ils n'ont qu'à recourir à [donner] la même explication : les visites de créatures extraterrestres. Sur ce thème, Erich von Däniken publie, en 1970 et 1972, deux ouvrages [livres] dans lesquels il traite à sa façon le problème de Nazca. Sa théorie est simple. Les anciens habitants du Pérou - avant la période de l'empire inca - ont tracé [fait] les dessins de Nazca en obéissant aux instructions [ordres] qui leur étaient données d'en haut par l'équipage d'un vaisseau venu d'un autre monde. L'engin survolait le chantier à basse altitude et guidait les paysans indiens au fur et à mesure qu'ils progressaient [avançaient] au milieu des gravillons [pierres] . Apparemment les extraterrestres n'ignoraient rien de la langue parlée dans le pays. Ou peut-être envoyaient-ils leurs ordres par transmission [message] télépathique ? Quant à la raison pour laquelle ils faisaient exécuter [faire] ces figures [dessins] géantes [immenses], von Däniken en a une toute prête : les tracés [lignes] rectilignes [droites] sont des pistes d'atterrissage, et les dessins, des signaux [messages] destinés aux [pour les] pilotes des soucoupes volantes, comme ceux qu'on voit

aujourd'hui dans les aéroports à l'usage [intention] des navigateurs aériens, ou sur les routes pour les automobilistes.

Maria Reiche, une mathématicienne qui pendant deux ans a travaillé à établir une carte détaillée des dessins de Nazca et s'est employée à [a tout fait pour] les préserver des [éviter leur] destructions, fait remarquer que les "pistes" en question sont débarrassées de [sans aucune] gravillons [pierre], et que le sol dénudé [nu] est plutôt meuble [mou], ce qui ne le rend pas très pratique pour un atterrissage. "J'ai bien peur, dit-elle, que si les soucoupes volantes essaient de s'y poser, elles y restent enfoncées!" Mais après tout, nous ignorons la fiche technique des aéronefs [engins] extraterrestres...

Von Däniken n'a d'ailleurs pas la paternité de son idée. Les pilotes qui, les premiers, avaient observé [vu] les dessins de Nazca du haut des airs, les appelaient par plaisanterie [pour rire] "l'aéroport préhistorique". Ils les comparaient aussi aux fameux canaux de Mars. La seule originalité de von Däniken est de s'être emparé de ces propos [idées] pour les habiller d'une théorie scientifique.

Entre autres explications [théories] du même genre, on doit citer celles qui veulent que les dessins aient été tracés [faits] au laser depuis des vaisseaux [engins] aériens [extra-terrestres] ; qu'ils renferment [contiennent] la clé du trésor des Incas; ou des messages laissés sur la Terre par des extraterrestres, à l'intention d'autres voyageurs de l'espace.

(source : M.Rouzé, Le mystère de Nazca enfin élucidé, *Science et Vie*, 1983, 11, pp.56-61)

* Les auteurs remercient Daniel Serant, professeur à l'Université Claude Bernard-Lyon 1, qui a bien voulu vérifier la validité de leurs analyses statistiques.