

A Probabilistic Approach to Linguistic Analysis in Machine Translation Output Evaluation

Olivier Gouirand

Laboratoire BABEL, Faculté de Lettres et Sciences Humaines
Université du Sud Toulon-Var
BP 20132, F-83957 La Garde (France)
olivier.gouirand@univ-tln.fr

Abstract

In order to overcome some shortcomings of previous machine translation (MT) evaluation methods, mainly pertaining to the difficulty of measuring actual linguistic quality (if possible in an automated way), a metric based on the distributional features of grammatical categories within the scope of the *Law on Anomalous Numbers* was devised : the text is tagged for parts-of-speech and these are computed into a positional distribution. The deviation of the latter from said law indicates translation errors, which alone is useful in diagnostic evaluations but is also of importance to theoretical linguists.

1 Introduction

The MT evaluation community has contributed a great number of approaches over the years but there are two problems which hamper progress in the field: first, generic evaluation techniques have always been considered difficult to achieve due to the specificity of user needs (King & Falkedal, 1990:1 ; Spark & Galliers, 1996:70) ; second, measuring the essential feature of MT i.e. text quality is a very complex task (Hovy et al., 2002:5) due to the many possible planes of analysis that it implies. The response to the former has been to design adaptative frameworks like EAGLES (EWG, 1996) or ISLE (described in Hovy et. al., 2002:4), the risk being to minimize the performance of core technologies factor as set forth by DARPA (White, O'Connell & O'Mara, 1994:intro-

Natural language generated by means of computer software indeed disrupts an equilibrium that lies in the syntax-semantics relations, because a text may seem well-formed while being nonsense all the same. The reason is grammatical categories are primitive constructs controlled by semantic categories and structures into which they fit. Studying MT output is about identifying software issues but also raising interesting questions for linguists to address : why is it that parts-of-speech are organized (of course not necessarily in a uniform way across languages) along the lines of other human and natural constructs, as demonstrated herein? What is their ontological status? etc.

duction). As to the latter, the quest for objective metrics - preferably automatable ones - that would grasp text quality has lead to mainly statistical criteria, the risk being to abstract away from linguistic quality. For example, in the BLEU algorithm (Papinemi, 2001:2-3), the distance between any translation and a set of reference translations stems from a calculus involving n -gram similarity (modified unigram precision) deemed to mirror translation quality criteria adequacy and fluency. It only rests on word counts though, so the linguistic level cannot be directly invoked.

The present paper seeks to cope with both problems by asserting the priority of linguistic quality in MT output assessment over any other criterion because that is fundamentally what the end user needs. To pursue this aim, a method which would incorporate linguistic content into a generalizable

framework is needed, other than statistical equalizations of human judgements *à la* ALPAC (1966).

What follows is basically the result of an investigation into the very nature of traditional grammatical categories in the wake of Aristotle, viz. nouns as "substance", determiners as "quantity", adjectives as "quality", verbs as "action / affection", just to name a few. We aim at demonstrating that the distribution of categories may be generalizable - although categories are not strictly mappable -, at least to a pair of natural languages, namely English and French ; that could be useful to better understand the ties between natural languages out of a linguistic standpoint (which is of paramount importance in translation) and shedding light upon phenomena best interpreted at the higher levels of cognition, or rather, semantics following Morris's classification (1938:6).

We shall also see that they clearly fit into the more complex scheme of the syntax-semantics interface, in the sense of words' distributional and semantic properties tallying. Indeed, one of the hurdles of MT, thus its evaluation, has been the gap between syntax and semantics, the former having historically prevailed over the latter, making full-blown linguistic evaluations troublesome, not to mention the other major stumbling block, i.e. the lack of expected results for any one translation, in other words a "correct" or "yardstick" translation.

2 Setup and experimenting with categorial distributions

Our experiment started with an in-depth experiment conducted on four English-French parallel manually translated corpora, each about 180,000 words, drawn from various subject areas and contained in the well-documented MCI/ECI corpus (ACL, 1994) [subcorpora "ILO" and "CCITT"], as well the *Hansard* (Parliament of Canada, 1999) and the *Annual Reports of the UN Secretary General* (UN, 1993 ; UN, 1995-2000). Technically, the first stage of it was devoted to cleaning up the texts from unnecessary information (empty lines, doublets and the like) as well marking sentence boundaries, a requirement of the piece of software carrying out the actual linguistic work. A robust tagger previously trained on the *Penn Treebank* and the *Trésor de la Langue Française* corpora (WinBrill© by ATILF, formerly InaLF, 1998) was used to assign grammatical tags to parts-of-speech.

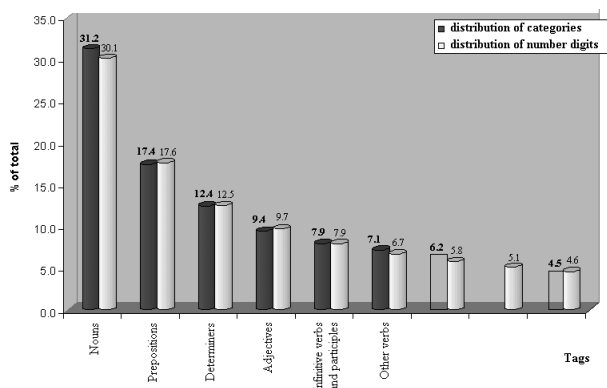
After stripping the texts of its lexical information, the tags alone were fed into a standard spreadsheet application, each sentence forming a row. Then the columns were scanned for specific tags and these were totalled (see Figure I). The resulting positional "grammatical n-grams" formed the data, the interpretation of which follows. It should be emphasized that these are linguistic in the sense that every part-of-speech mirrors some syntactic relation : for example, a determiner is only such with respect to the "determined", an adjective (*adjectivum nomen*) relates to a noun, an adverb to a verb and so forth.

The starting point in the interpretation of the results could be the well-known fact that word distributions comply with Zipf's law (Zipf, 1965:38), that is : the rank of a word is inversely proportional to its frequency and their product is constant. But what is less known is that, above a threshold of about 2,000 tokens (although we first noticed this on a manually counted 500-word text, spurring our initial hypothesis), the distribution of parts-of-speech in English and French matches another statistical discovery, the *Law on Anomalous Numbers* (Benford, 1938). This law, which states that the distribution D of number digits equals $\log(1+1/D)$, applies to all non-random (thus not to lotteries and the like) and non-constrained numbers (thus not to supermarket prices ending with 9 for marketing reasons etc.), ranging from land areas to stock-market prices, regardless of unit used. It is readily and successfully used in accounting to trace frauds but otherwise limited to a narrow audience despite its mathematical demonstration by Terence Hill in 1996, probably because it encompasses other laws in a disturbing fashion ; the history of science has been an array of laws where some paradigms seem to be both the product of man's analysis of his world and the world's laws, which some scientists as Thom (1990:344) acknowledge in the biological nature of language e.g. noun as *salience* (Thom, 1988:23).

In our experiment, the average distribution of grammatical categories in the English set of texts from the corpora is : nouns 31.23 %, prepositions 17.39 %, determiners 12.44 %, adjectives 9.43 %, infinitive verbs and participles 7.92 %, other verbs (very weakly correlated to the others positionally) 7.15 % and presumably two groups made up of some wh- words and conjunctions as well as one comprising pronouns and adverbs, account for 6.18 and 4.48 % respectively, see Table I below.

Table I

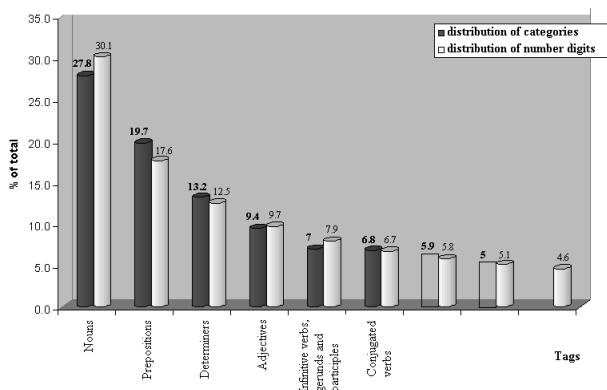
Compared distribution of grammatical categories and number digits (English set)



In the French texts, the distribution is : nouns 27.83 %, prepositions 19.73 %, determiners 13.19 %, adjectives 9.41 %, infinitive verbs and participles 6.96 %, conjugated verbs 6.80 % and presumably two groups made up of pronouns, as well as one comprising conjunctions, account for 5.91 and 4.96 % respectively, see Table II below.

Table II

Compared distribution of grammatical categories and number digits (French set)



The lower part of the spectrum is nonetheless kept truncated because intervals in such a distribution are too close to discriminate between aggregate "minor" categories with certainty, unless finer-grained groupings emerge from future research.

These results were checked for reliability : deviations were within acceptable limits and no data relative to the categories presented above that were below the confidence level of 1,000 due to the large number of tokens for each type. Also, only the most strongly correlated groupings were made. When it comes to informativity, the set underwent a subsequent series of processes consisting in a machine translation of all source texts, then the post-processing as described above and the com-

parison with the manually translated target texts. Interestingly, with a first-generation MT system (Systran), the distribution was distorted significantly : 30 % less proper nouns in English and 22 % in French, more than a fifth less present tense verbs and participles in English, a quarter and a sixth more determiners and prepositions respectively in English, roughly the same loss for French pronouns and other verbs (see Table III and IV). To take just one example for the sake of brevity, Systran has a persistent bug : proper nouns are treated as common nouns when there's a match, as in "Mr Pat O'Brien" translated as *M. tapotent O'Brien* ; that explains the aforesaid deflated proper nouns.

Table III

Compared distribution of grammatical categories (English)

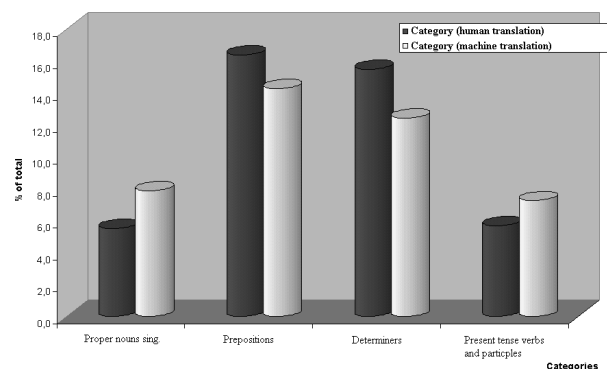
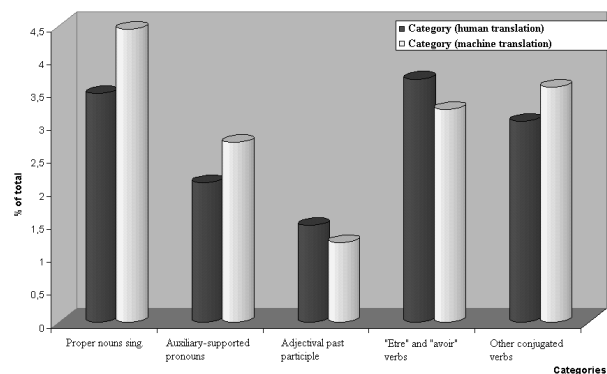


Table IV

Compared distribution of grammatical categories (French)



Since the various subsets of the corpus were coherent in our initial experiment, the distortion in the distribution of categories can only mean that it is sensitive to Natural Language Generation (NLG). What's more, the syntactic correctness was preserved despite altered relative values for each category and obviously with an altered sense, as usual with MT, owing to missing or incomplete semantic interpretation (Gouirand, 2002:75).

3 Discussion

Going back to Benford's law, the correlation with Zipf's law is that they are both based upon the principle of invariance (Zipf's law is insensitive to tagging while discriminating meaning, see Wilks 1996:6) and that there must be an underlying law which can account for the distribution. It has been claimed that lexical distributions relate to Zipf's law by virtue of word length ; however, it is easy to show that inflectional grammars such as German or Swedish have a different bearing on word lengths from those of, say, English or French. So the underlying law remains implicit and our duty as linguists is to make it explicit.

Yet the analogy is not complete, for one would not expect Benford's law to be relevant to linguistics as it conveys plain numerical rather than lexical information. The numbering of tags can be posited as arbitrary : the fact that some plot of land measures n acres doesn't tell much about what n actually means, it is simply the outcome of a composite geometrical law expressing surface areas. The numbers are merely arbitrary because whatever their name is, they always retain the underlying law's features, which is precisely why linguistic tags also apply : as Saussure (1916:98-100) showed at length, linguistic signs are arbitrarily chosen, they are signs pointing to "something else", at another level ; thus it is an underlying law which steers the said distribution.

Then, it is our contention that such an underlying law must be of semantic nature since the ancillary experiment set forth above clearly demonstrated that data-processing the corpora with machine translation software having a very limited semantic competence distorted the distribution. In other words, human processing of grammatical categories bears on ontological categories in the realm of semantics (Jackendoff, 1990:23) and therefore keeps the distribution within the boundaries of the law, whereas MT systems don't, insofar as there is no real semantic component attached to them (when they have one, they may to some extent).

Another proof of the validity of our model is given by feeding the text tags into a concordancer. Instead of using words as queries, we used our tags and then we reordered left-hand and right-hand side tags in decreasing order for each corpora.

As seen from Figure II and III (which only display the three most probable tags because as one moves

down, volatility rises sharply), nouns (NNx in English/SBx in French), prepositions (IN/PREP) and determiners (DT/DTN) are once again found to be the most frequent top-level n -grams (ca. 70 percent of total) and roughly distributed in a way similar to categories, which makes sense : let us assume that tags were evenly scattered over all categories ; their count alone would create a vertical hierarchy within each category in building n -grams but this simply does not happen because the distribution takes place horizontally instead, meaning that the distribution is linked to relations among categories, which was the point to be made here.

Stating that syntactic constructs are driven by semantics means that ontological and/or biological nature of grammatical categories and their interplay can be investigated within the scope of an integrated theory of meaning where pre-verbal features of language are the driving force. In the domain of conceptual semantics (Jackendoff, 1983:57), syntactic structures parallel conceptual structures in that ontological categories (such as THING, EVENT, STATE, ACTION, PLACE, PATH, PROPERTY, AMOUNT) are the primitives associated with syntactic *constituents* ; by the way, some of the former are quite similar to Aristotle's findings, suggesting a continuum in the approach. However, in the internal workings of a MT system, invoking a concept from the lexical stratum is not straightforward, because lexemes may refer to several concepts, so clear-cut one-to-one correspondences are out of reach, at least for a computer program, so future research in lexicology should focus on a better mapping of lexical units into semantic categories.

For the French school of linguistics (Culioli, 1999:63), there is also a pre-verbal (and definitively pre-lexical) level of primitive relations where a distinction between syntax semantics is pointless. Hence the *notion* is a set of physical and cultural features that are inserted in a given context (Culioli, 1990:50). For Thom (1988:206), a concept is workable only through categorial operations which constrain the extensional perimeter centered around a prototypical core.

Another assumption is that the number of workable categories cannot be large or the system's output would result in unlimited semiosis and cognitive burdening. This is in line with Aristotle's initial formalization of categories but also with most speakers phenomenological experience.

The current state of this research allows an evaluation based on the distribution of grammatical *n*-grams that is useful to identify major translation failures but for a finer grained pinpointing, the definition of meaningful relations between syntax and semantics, parts-of-speech and meaning needs to be addressed. These relationships are to be grasped not only at the level of lexical entries analysable as semantic categories (with several impediments that are not yet solved) but also as ontological constructs or notional domains which are modelled as networks or semantic spaces. The overall sense then depends on such features but also on information conveyed by the discursive context as argumentative underpinning, demonstrated by Anscombre and Ducrot (1983:86) as well as Anscombre (1995:32), which provides directions for future research in the field.

Also, the question of whether the distributional features of grammatical categories in English and French are valid in other languages is still open. If they were to be ported, the metric would then be generic and so would the linguistic implications. There are obstacles, however : quality parallel corpora and robust parallel taggers for "exotic" languages (and even less exotic ones for that matter) are scarce, so the task of investigating the matter would be immensely time-consuming and it would require extensive linguistic expertise, let alone the "toolbox" problem. This is because some languages are often presented as resistant to some theories, for example Basque, which lacks a verbal category : the verbal function is still to be found though, in some other category from where it can be extracted ; In a language such as Swedish, determiners are embedded in nouns and can thus be backtracked. So applying our model to other languages wouldn't necessarily weaken it, only such an endeavour would require substantial efforts.

References

- ALPAC. 1966. *Language and Machines. Computers in Translation and Linguistics*. A Report by the Automatic Language Processing Advisory Committee, Division of Behavioural Sciences, National Academy of Sciences. Publication 1416. National Research Council, Washington, DC.
- Association for Computational Linguistics. 1994. *European Corpus Initiative (ECI)*. Multilingual Corpus 1(MCI) [CD-ROM]. ACL, Geneva/Edinburgh.
- Jean-Claude Anscombre and Oswald Ducrot. 1983. *L'argumentation dans la langue*. Mardaga, Brussels, Belgium.
- Jean-Claude Anscombre. 1995. *Théorie des topoï*. Kimé, Paris, France.
- Frank Benford. 1938. The Law of Anomalous Numbers. *Proc. Amer. Phil. Soc.*, 78:551-554.
- Antoine Culioli. 1990. *Pour une linguistique de l'énonciation. Opérations et représentations*, volume 1. Ophrys, Paris, France.
- Antoine Culioli. 1999. *Pour une linguistique de l'énonciation. Formalisation et opérations de repérage*, volume 2. Ophrys, Paris, France.
- EWG. 1996. *EAGLES Evaluation Group*. Final Report. Octobre 1996 version. Center foer Sprogteknologi, Copenhagen, Denmark.
- Olivier Gouirand. 2002. Quelques réflexions sur la traduction automatique comme médiation. *Terminologie et Traduction (T&T)*. La revue des services linguistiques des institutions européennes, 2/2002:72-83. Office des Publications officielles des Communautés européennes, Luxembourg.
- E. Hovy, M. King and A. Popescu-Belis. 2002. In Introduction to MT Evaluation. *Workbook of the LREC 02 Workshop on Machine Translation : Human Evaluators Meet Automated Metrics*, Las Palmas, Spain, 2002, pp. 1-7.
- Institut National de la Langue Française (InaLF), Massachusetts Institute of Technology (MIT) and University of Pennsylvania. 1993-98. *WinBrill tagger and lemmatizer*. ATILF, Nancy, France.
- Ray Jackendoff. 1983. *Semantics and Cognition*. MIT Press, Cambridge, Mass.
- Ray Jackendoff. 1990. *Semantic Structures*. MIT Press, Cambridge, Mass.
- Margaret King and Kirsten Falkedal. 1990. Using Test Suites in Evaluation of Machine Translation Systems. *Proceedings of COLING '90*, Helsinki, Finland, 1990.
- Charles W. Morris. 1938. *Foundations of The Theory of Signs*. University of Chicago Press, Chicago, IL.
- K. Papineni, S. Roukos, T. Ward et al. 2001. *BLEU : a Method for Automatic Evaluation of Machine Translation*. Computer Science IBM Research Report RC22176 (W0109-022). IBM Research Division, Yorktown, NY.
- Parliament of Canada. 1999. *Debates of the House of Commons of Canada (Hansard), Dec. 1,2,3 & 10*,

1999. Parliament of Canada, Ottawa. <http://www.parl.gc.ca/cgi-bin/36/pb_chb_hou_deb.pl?e> [page viewed April 28, 2000].
- Ferdinand de Saussure. 1916. *Cours de linguistique générale*. Payot, Paris, France.
- Karen Sparck-Jones and Julia Galliers. 1996. *Evaluating Natural Language Processing*. Springer, Berlin, Germany.
- René Thom. 1988. *Esquisse d'une sémiophysique*. Interéditions, Paris, France.
- René Thom. 1990. *Apologie du logos*. Hachette, Paris, France.
- United Nations (UN). 1993. *Annual Report of the Secretary-General on the Work of the Organization*. UN, New York, NY. <<http://www-rali.iro.umontreal.ca/arc-a2/BAF/>> [page viewed on Aug 15, 2000].
- UN. 1995. *Annual Report of the Secretary-General on the Work of the Organization*. UN, New York, NY. <<http://www.un.org/Docs/SG/SG-Rpt/>> [page viewed on Aug 15, 2000].
- UN. 1996. *Annual Report of the Secretary-General on the Work of the Organization*. UN, New York, NY. <<http://www.un.org/Docs/SG/sgrepot.htm>> [page viewed on Aug 15, 2000].
- UN. 1997. *Annual Report of the Secretary-General on the Work of the Organization*. UN, New York, NY. <<http://www.un.org/Docs/SG/Report97/97con.htm>> [page viewed on Jun 27, 2000].
- UN. 1998. *Annual Report of the Secretary-General on the Work of the Organization*. UN, New York, NY. <<http://www.un.org/Docs/SG/Report98/con98.htm>> [page viewed on Jun 27, 2000].
- UN. 1999. *Annual Report of the Secretary-General on the Work of the Organization*. UN, New York, NY. <<http://www.un.org/Docs/SG/Report99/toc.htm>> [page viewed on Aug 9, 2000].
- UN. 2000. *Millennium Report of the Secretary-General*. United Nations, New York, NY. <http://www.un.org/millennium/sg/report/> [page viewed on Aug 9, 2000].
- J. White, T. O'Connell, F. O'Mara. 1994. *ARPA Machine Translation Program : 3Q94 Evaluation*. Georgetown University and Advanced Research Projects Agency, Washington, DC.
- Yorick Wilks. 1996. *The grammar of sense: Is word-sense tagging much more than part-of-speech tagging?* Research Memoranda CS-96-05. University of Sheffield, Sheffield, UK..
- George K. Zipf. 1965. *Psycho-Biology of Languages*. MIT Press, Cambridge, Mass.

APPENDIX

Figure I : illustration of spreadsheet processing

A420 = ABR										A448 = ADJ									
A	B	C	D	E	F	G	H	I		A	B	C	D	E	F	G	H	I	
419	SUB	DTN	SBC	ECJ	ADV	PREP	SBC	PREP	DTN	419	SUB	DTN	SBC	ECJ	ADV	PREP	SBC	PREP	DTN
420	ABR	SBP	SBP	REL	ACJ	VPAR	VNCF	DTN	SBC	420	ABR	SBP	SBP	REL	ACJ	VPAR	VNCF	DTN	SBC
421	ADV	DTN	DTN	SBC	PREP	DTN	SBP	ACJ	VPAR	421	ADV	DTN	DTN	SBC	PREP	DTN	SBP	ACJ	VPAR
422	DTN	SBC	ECJ	ADJ1P	ADV	DTC	SBC	PREP	SBC	422	DTN	SBC	ECJ	ADJ1P	ADV	DTC	SBC	PREP	SBC
423	PREP	PRO	REL	VCJ	DTN	SBC	PREP	DTN	SBC	423	PREP	PRO	REL	VCJ	DTN	SBC	PREP	DTN	SBC
424	PRV	ACJ	VPAR	ADV	DTC	SBC	PREP	DTN	SBC	424	PRV	ACJ	VPAR	ADV	DTC	SBC	PREP	DTN	SBC
425	PRV	ACJ	VPAR	DTN	SBC	DTC	SBC	DTC	SBC	425	PRV	ACJ	VPAR	DTN	SBC	DTC	SBC	DTC	SBC
426	PREP	PRO	PRV	ADV	ECJ	ADV	ADJ	PREP	VNCF	426	PREP	PRO	PRV	ADV	ECJ	ADV	ADJ	PREP	VNCF
427	DTN	SBC	PREP	SBC	SUB	PRV	VCJ	VNCF	VCJ	427	DTN	SBC	PREP	SBC	SUB	PRV	VCJ	VNCF	VCJ
428	SBP	ADJ	PREP	DTN	SBC	PREP	SBC	ADJ	COO	428	SBP	ADJ	PREP	DTN	SBC	PREP	SBC	ADJ	COO
429	DTN	SBC	PREP	DTN	SBC	VCJ	SUB	SBC	DTN	429	DTN	SBC	PREP	DTN	SBC	VCJ	SUB	SBC	DTN
430	DTN	SBC	PRV	ACJ	VPAR	DTC	SBC	PREP	DTN	430	DTN	SBC	PRV	ACJ	VPAR	DTC	SBC	PREP	DTN
431	PREP	PRO	REL	VCJ	DTN	SBC	ADJ2PA	DTC	SBC	431	PREP	PRO	REL	VCJ	DTN	SBC	ADJ2PA	DTC	SBC
432	PREP	DTN	SBC	PREP	PRO	DTN	SBC	REL	VCJ	432	PREP	DTN	SBC	PREP	PRO	DTN	SBC	REL	VCJ
433	PREP	SBC	ACJ	VPAR	DTN	SBC	PRV	ADV	PRV..	433	PREP	SBC	ACJ	VPAR	DTN	SBC	PRV	ADV	PRV..
434	DTN	SBC	VCJ	DTN	SBC	DTC	SBC	PREP	REL	434	DTN	SBC	VCJ	DTN	SBC	DTC	SBC	PREP	REL
435	DTN	SBC	VCJ	DTN	SBC	PREP	VNCF	DTN	SBC	435	DTN	SBC	VCJ	DTN	SBC	PREP	VNCF	DTN	SBC
436	SBC	ADJ	PREP	DTN	SBC	PREP	SBC	ADJ	COO	436	SBC	ADJ	PREP	DTN	SBC	PREP	SBC	ADJ	COO
437	PRV	PRV	VCJ	PREP	SBC	SBP	SBP	SBP	SBP	437	PRV	PRV	VCJ	PREP	SBC	SBP	SBP	SBP	SBP
438	DTN	SBC	VCJ	ADV	SUB	SBC	SBP	SBP	SBP	438	DTN	SBC	VCJ	ADV	SUB	SBC	SBP	SBP	SBP
439	DTN	SBC	VCJ	SUB	DTN	SBC	ADV	ACJ	ADV	439	DTN	SBC	VCJ	SUB	DTN	SBC	ADV	ACJ	ADV
440	DTN	SBC	VCJ	PREP	ADJ	DTC	SBC	PREP	VNCF	440	DTN	SBC	VCJ	PREP	ADJ	DTC	SBC	PREP	VNCF
441	DTN	SBC	VCJ	ADV	SUB	PRV	PRV..	ACJ	SBC	441	DTN	SBC	VCJ	ADV	SUB	PRV	PRV..	ACJ	SBC
442	DTN	SBC	VCJ	SUB	PREP	DTN	SBC	SUB	DTN	442	DTN	SBC	VCJ	SUB	PREP	DTN	SBC	SUB	DTN
443	SBP	SBP	SBP	ACJ	EPAR	ADJ1PA	PUL	SUB		443	SBP	SBP	SBP	ACJ	EPAR	ADJ1PA	PUL	SUB	
444	DTN	SBC	VCJ	PREP	DTN	SBC	SUB	DTN	SBC	444	DTN	SBC	VCJ	PREP	DTN	SBC	SUB	DTN	SBC
445	PRV	VCJ	ADV	SUB	DTN	SBC	SBP	SBP	SBP	445	PRV	VCJ	ADV	SUB	DTN	SBC	SBP	SBP	SBP
446	ADV	DTN	SBC	VCJ	SUB	DTN	SBC	ADV	ACJ	446	ADV	DTN	SBC	VCJ	SUB	DTN	SBC	ADV	ACJ
447	PREP	SBC	ADJ	DTN	SBC	PREP	ADJ	SBC	SUB	447	PREP	SBC	ADJ	DTN	SBC	PREP	ADJ	SBC	SUB
448	ADJ	SBC	DTN	SBC	VCJ	ADV	SUB	DTN	SBC	448	ADJ	SBC	DTN	SBC	VCJ	ADV	SUB	DTN	SBC
449	DTN	SBC	VCJ	DTN	SBC	PREP	PRV	VNCF	DTN	449	DTN	SBC	VCJ	DTN	SBC	PREP	PRV	VNCF	DTN

B2 = =NB.SI(Données!A\$1:A\$50000;"ABR")																
V	Pos./Cat.	ABR	ADJ	ADJ1PA	ADJ2PA	ADV	CAR	COO	DTN	DTC	FGW	INJ	PFX	PREP	PRV	
2	1	14	10	0	9	117	59	17	875	35	0	0	0	513	178	
3	2	0	32	0	28	58	69	8	357	117	0	0	0	144	23	
4	3	0	98	6	25	88	93	10	196	90	0	0	0	159	45	
5	4	0	63	27	34	147	49	10	388	107	0	0	0	247	29	
6	5	0	85	30	30	83	82	17	518	93	0	0	0	284	42	
7	6	0	88	19	28	58	70	21	381	87	0	0	0	303	38	
8	7	0	141	23	39	65	75	31	342	117	0	1	0	275	71	
9	8	0	98	26	38	67	104	43	356	139	0	0	0	309	33	
10	9	0	104	30	55	63	78	38	378	124	0	0	0	293	36	
11	10	0	142	28	43	56	81	53	310	150	0	0	0	304	33	
12	11	0	109	33	52	59	75	64	321	120	0	1	0	317	41	
13	12	0	114	29	34	50	87	55	355	117	0	0	0	341	46	
14	13	0	138	32	35	79	76	65	315	121	0	0	0	270	32	
15	14	0	109	30	60	47	70	54	302	144	0	0	0	349	33	
16	15	0	121	37	32	64	68	67	318	119	1	0	0	300	25	
17	16	0	145	24	39	64	75	59	272	115	0	0	0	313	38	
18	17	0	125	33	52	63	94	53	297	99	0	0	0	301	33	
19	18	0	118	31	36	80	67	80	269	93	0	0	0	298	35	
20	19	0	132	33	32	68	67	83	301	104	2	0	0	281	35	

Figure II : Category n-grams for English

N	N2	N3	N4	Left	NODE	Right	N	N2	N3	N4	N	N2	N3	N4	Left	NODE	Right	N	N2	N3	N4	N	N2	N3	N4	Left	NODE	Right	N	N2	N3	N4
1174	2070	1657	1696	DT	jj	NN	1475	2858	2119	2596	438	126	685	196	NNP	vbd	VBN	604	108	337	219	3276	3725	3263	3383	NN	in	DT	4501	4424	3850	3467
690	1127	669	1237	IN	jj	NNS	869	1621	1023	1891	435	86	377	186	NN	vbd	IN	307	65	321	199	1492	1604	1274	1742	NNS	in	NN	1357	1877	1244	1390
174	377	376	573	CC	jj	IN	217	388	359	573	313	66	157	163	NN	vbd	DT	303	53	210	158	1151	1594	148	847	VBN	in	NNP	875	1022	1101	1237
152				JJ	jj	TO	171				287				PRP	vbd	RB	266				904				NNP	in	JJ	651			
152				VB	jj	JJ	152				170				WDT	vbd	JJ	131				369				VBD	in	NNS	502			
131				VBD	jj	CC	146				85				RB	vbd	TO	113				302				RB	in	CD	451			
97				RB	jj	DT	79				75				WP	vbd	NN	75				294				CD	in	PRP\$	347			
93				NN	jj	NNP	48				68				CC	vbd	NNP	40				292				VZ	in	IN	286			
91				VBN	jj	CD	14				61				CD	vbd	NNS	32				296				IN	in	VBS	253			
82				PRP\$	jj	RB	13				38				EX	vbd	VBG	32				281				VB	in	PRP	231			

Figure III : Category n-grams for French

N 2 3 4 Left NODE Right N N2 N3 N4													N N2 N3 N4 Left NODE Right N N2 N3 N4													N N2 N3 N4 Left NODE Right N N2 N3 N4													N N2 N3 N4 Left NODE Right N N2 N3 N4												
3368	6027	2378	3773	SBC	adj	PREP	1181	2636	1302	1473	686	262	642	388	SBC	acj	VPAR	907	143	919	811	4733	6469	4599	4638	SBC	prep	DTN	5008	4326	4324	4463																			
515	907	979	752	DTN	adj	SBC	908	1434	1038	1342	231	108	361	349	PRV	acj	EPAR	504	143	350	288	901	1026	878	1034	ADJ	prep	SBC	2402	3705	2351	2597																			
415	591	455	616	PREP	adj	DTN	634	738	459	893	231	75	213	207	ADV	acj	ADV	314	139	222	250	675	632	685	563	ADJ2PA	prep	VNCFG	833	855	1195	1146																			
340				SBP	adj	COO	513				220				SBP	acj	DTN	118				489				ADJ1PA	prep	ADV	313																						
185				ADJ	adj	DTN	305				167				REL	acj	APAR	58				460				COO	prep	PRO	265																						
174				ADV	adj	ADJ	185				148				ADJ	acj	PREP	53				425				VCJ	prep	SBP	244																						
140				DTN	adj	VCJ	153				70				PRV++	acj	SBC	18				303				SBP	prep	CAR	194																						
116				COO	adj	ACJ	148				49				COO	acj	DJ2PAR	17				291				ADV	prep	REL	138																						
89				ECJ	adj	DJ2PAR	143				35				CAR	acj	PRO	6				283				VNCFG	prep	ADV	137																						
81				CAR	adj	ADV	140				29				ADJ2PA	acj	COO	4				227				VPAR	prep	ADV	91																						