

# LETTER AND WORD FREQUENCIES

REX GOOCH

Letchworth Garden City, Herts, England

rexgooch@ntlworld.com

Counts of letters and words have been made from various bodies of text over the years. Perhaps the earliest well-known corpus was that assembled by W.N. Francis and H. Kucera of Brown University (hence known as the Brown corpus). It contains one million words of American English from 1961 in 15 categories of 500 texts. If it has a use today, it is because matching studies were done of British English (LOB corpus) and Indian English (Kolhapur Corpus) among others. In *Word Ways*, I have preferred to use the Collins COBUILD corpus, known as the Bank of English®: it runs to hundreds of millions of words of English text from British, US, Australian, and Canadian sources (including textbooks, novels, newspapers, guides, magazines, and websites). The corpus has been automatically word-class tagged, and a 200-million-word corpus has been parsed. It is at <http://www.cobuild.collins.co.uk/>. Recently, I have used the British National Corpus, which is a 100 million word collection of samples of written and spoken language from a wide range of sources, designed to represent a wide cross-section of current British English, both spoken and written (available from <http://www.natcorp.ox.ac.uk/>, the Oxford University website). There is a Baby version, and the original Brown corpus, tagged, is also available.

In the February issue of *Word Ways* (page 41), Anil remarked that the order of popularity of letters varied greatly depending on whether the letters are counted from running text, or from a list of the most popular words. I believe he was talking about column 6 below, to which I have added actual frequencies. He was mainly comparing the order with that in column 4, the Linotype machine. The Linotype order is very close to that of the BNC running text in column 3. There is a significant difference in *order* between letters and running text. However, in all cases E is dominant, followed by the same cluster of letters — differences in frequency are small in this cluster, so letters easily change places, and the differences in *frequency* are not truly significant. I thought that perhaps any difference between running text and words might be due to a difference between the most popular words (often repeated in running text), and others. I therefore decided to count the letters in the most popular 999 words in the BNC, then in the next 85,801: the results in columns 1 and 2 do in my view show significant differences, eg in the more popular words (which are shorter), E is more frequent, and H less so. Similarly, T seems to be more common in running text. I leave the reader to browse the results further. Column 5 is from the same source as Anil's column 6, but has no plural words or words with common suffixes: the lack of terminal -S and -ING are just apparent.

Finally, I include the top 100 words from the BNC, and from the American Heritage Word Frequency Book (Carroll, Davies and Richman, 1971). It draws upon 5 million running words used in US schools. The numbers are the number of occurrences per 10,000 words for the BNC.

I trust the tables will also serve as a reference for *Word Ways* contributors.

Parts per thousand for individual letters in various selections of text

1	2	3	4	5	6
E 134	E 108	E 125	E	E 114	E 117
A 73	A 85	T 92	T	A 85	I 86
T 73	I 83	A 81	A	I 79	S 81
R 73	S 76	O 76	O	R 75	A 79
O 71	R 73	I 73	I	T 74	R 74
N 70	N 72	N 70	N	O 71	N 74
S 67	T 65	S 63	S	N 64	T 67
I 65	O 64	R 61	H	S 55	O 59
L 50	L 60	H 54	R	L 55	L 53
D 41	C 41	L 41	D	C 47	C 41
C 39	D 38	D 39	L	U 36	D 38
U 33	U 32	C 31	U	P 32	U 32
H 32	M 30	U 28	C	M 32	G 28
M 31	G 28	M 24	M	D 31	P 27
P 28	P 27	F 22	F	H 27	M 27
G 24	H 26	P 20	G	G 23	H 22
Y 21	B 20	G 20	Y	B 21	B 20
W 18	Y 18	W 20	P	Y 20	Y 17
F 17	F 14	Y 19	W	F 14	F 14
B 14	K 11	B 16	B	V 10	V 11
V 12	V 11	V 10	V	W 9	K 9
K 9	W 10	K 7	K	K 8	W 9
J 1	Z 4	X 2	X	X 3	Z 5
X 1	X 3	J 2	J	Z 2	X 3
Q 1	J 2	Q 1	Q	Q 2	J 2
Z 0	Q 2	Z 1	Z	J 1	Q 2

1. top 999 words in BNC (5.5 letters per word)

2. next 85801 words in BNC (7.7 letters per word)

3. running text, BNC

4. Linotype

5. 18584 Common Base words

6. 45406 Common Words

BNC	AmHer	BNC	AmHer		
THE	651	the	NO	22	out
OF	309	of	SAID	22	them
AND	282	and	WHO	21	then
TO	269	a	MORE	21	she
A	226	to	ABOUT	20	many
IN	198	in	UP	19	some
THAT	117	is	THEM	18	so
IT	114	you	SOME	18	these
IS	105	that	COULD	17	would
WAS	97	it	HIM	17	other
I	95	he	INTO	17	into
FOR	89	for	ITS	17	has
ON	76	was	THEN	16	more
YOU	73	on	TWO	16	her
HE	71	are	OUT	16	two
BE	69	as	TIME	16	like
WITH	68	with	LIKE	16	him
AS	54	his	ONLY	16	see
BY	53	they	MY	16	time
AT	50	at	DID	15	could
HAVE	49	be	OTHER	14	no
ARE	49	this	ME	14	make
THIS	48	from	YOUR	14	than
NOT	48	I	NOW	14	first
BUT	47	have	OVER	13	been
HAD	46	or	JUST	13	its
HIS	45	by	MAY	13	who
THEY	45	one	THESE	13	now
FROM	43	had	NEW	13	people
SHE	40	not	ALSO	13	my
WHICH	39	but	PEOPLE	13	made
OR	39	what	ANY	13	over
WE	37	all	KNOW	12	did
AN	36	were	VERY	12	down
THERE	34	when	SEE	12	only
HER	34	we	FIRST	12	way
WERE	33	there	WELL	12	find
ONE	30	can	AFTER	12	use
DO	29	an	SHOULD	11	may
BEEN	28	your	THAN	11	water
ALL	27	which	WHERE	11	long
THEIR	27	their	BACK	10	little
HAS	27	said	HOW	10	very
WOULD	26	if	GET	10	after
WILL	26	do	MOST	10	words
WHAT	26	will	WAY	10	called
IF	24	each	DOWN	10	just
CAN	24	about	OUR	9	where
WHEN	22	how	MADE	9	most
SO	22	up	GOT	9	know