

Word frequency effect and response quality: A comparison of four questionnaire versions on a web panel

Ana Slavec and Vasja Vehovar
University of Ljubljana

3MC, conference
Chicago, July 2016

Outline

1. Background
2. Description of linguistic resources used
3. Previous studies
4. Present study: split-ballot experiment
5. Conclusions, limitations and plans for future research

1. Research question

- Each question can be worded/translated in several way:
How to choose the most optimal wording?
- Example (translation from Slovenian to English)
 - “Zbirati” (verb)

zbirati to collect, to gather; to assemble; to rally; to accumulate; ~ (denarna) sredstva to collect (financial) resources; on zbira znamke he collects stamps; ~ se to gather, to assemble, to rally, to converse
 - [MADE UP QUESTION ITEM] “Skrbi me, da vlada zbira preveč informacij o ljudeh, kot sem jaz.”
 - [ONE OF THE SEVERAL POSSIBLE ENGLISH TRANSLATIONS] “I am concerned that the governemnt is collecting too much information about people like me.”
 - **Why not assembling, gathering or other synonymous word?**

Word frequency effect

- Words commonly used in daily speech are recognised and processed more quickly than less commonly used words (Howes and Solomon 1951; Broadbent 1967)
- **Unfamiliar words** as one of the psycholinguistic determinants of question difficulty
- Word frequencies in text corpora as possible estimate of word familiarity:
 - Lower Frequency Wording (LFW)
 - Higher Frequency Wording (HFW)

2. Text corpora

www.natcorp.ox.ac.uk/corpus/index.xml?ID=products

BRITISH NATIONAL CORPUS

▷ BNC ▷ corpus

BOOKMARK

search site

BNC Products

What is the BNC?
How the BNC was created
The BNC in numbers

When users obtain a BNC product, they agree to the licence which gives them the right to hold and use a copy of the corpus. A corpus is a dataset which can be used in many different ways, and we regret that the University of Oxford is not able to offer support to users of the corpus. Funding for the development and support of the corpus ended many years ago, but the corpus has been created in such a way that it

corpus.byu.edu/coca/

CORPUS OF CONTEMPORARY AMERICAN ENGLISH

450 MILLION WORDS, 1990-2012 [DOWNLOAD ALL 190,000 TEXTS]

EMAIL
PASSWORD
(HELP) LOG IN (REGISTER)

DISPLAY

LIST CHART KWIC COMPARE

SEARCH STRING

WORD(S)

COLLOCATES

POS LIST

SECTIONS SHOW

1 IGNORE

SPOKEN
FICTION
MAGAZINE
NEWSPAPER
ACADEMIC

2 IGNORE

SPOKEN
FICTION
MAGAZINE
NEWSPAPER
ACADEMIC

SORTING AND LIMITS

SORTING

MINIMUM 10

CLICK TO SEE OPTIONS

There are a wide range of additional resources that are based on the BYU corpora:

Full-text	Download 440 million words of full-text data for COCA (190,000 texts), or 1.8 billion words for GloWbE (1,800,000 texts). With this data, you will have the texts from the corpora on your own computer , rather than having to use the web interface.
Wikipedia corpus (NEW)	Quickly and easily create "virtual" corpora from the 4.4 million articles of Wikipedia (1.9 billion words) on almost any topic -- biology, investments, cars, Buddhism, etc. Search these virtual corpora, compare them to each other, and create keyword/frequency lists from your corpora.
Word and Phrase (analyze texts)	Enter entire texts and see detailed frequency information on the words in the text, and create word lists based on your text. Click through the words to see detailed information on any word. Highlight phrases in your text and have it search for related phrases in COCA.
Word and Phrase (frequency lists)	Search and browse the most complete frequency dictionary of English. See detailed information (all on one page) -- definition, frequency by genre, collocates (nearby words), concordance lines, synonyms, and Wordnet-related words, all with useful links from one resource to another.
Word Frequency	You can also download lists showing the frequency of the top 60,000 lemmas by genre (and sub-genre). Free list of the top 5,000 lemmas in COCA. Download the 100,000 integrated word list from COCA, COHA, BNC, and SOAP -- the largest, corrected frequency list of English.
Collocates	Download lists with the top 200-300 collocates (nearby words) for 60,000 different lemmas -- 4,300,000 node/collocate pairs in all.
N-grams	Download free lists containing the top 1,000,000 2-grams (two word sequences), 3-grams, 4-grams, and 5-grams in COCA. There are also other lists that contain the

INTRODUCTION

[Help / information / contact](#)

[WHERE SHOULD I START?]

[COMPARE TO OTHER CORPORA / ARCHITECTURES]

The Corpus of Contemporary American English (COCA) is the largest freely-available corpus of English, and the only **large and balanced** corpus of American English. The corpus was created by **Mark Davies** of **Brigham Young University**, and it is used by **tens of thousands** of users every month (linguists, teachers, translators, and other researchers). COCA is also related to **other large corpora** that we have created.

The corpus contains more than **450 million words** of text and is equally divided among spoken, fiction, popular magazines, newspapers, and academic texts. It includes **20 million words each year from 1990-2012** and the corpus is also updated regularly (the most recent texts are from Summer 2012). Because of its design, it is perhaps the only corpus of English that is suitable for looking at **current, ongoing changes** in the language (see the **2011 article in *Literary and Linguistic Computing***).

The interface allows you to search for **exact words or phrases, wildcards, lemmas, part of speech, or any combinations of these**. You can **search for surrounding words (collocates)** within a ten-word window (e.g. all nouns somewhere near *faint*, all adjectives near *woman*, or all verbs near *feelings*), which often gives you good insight into the meaning and use of a word.

The corpus also allows you to easily **limit searches by frequency and compare the frequency** of words, phrases, and grammatical constructions, in at least two main ways:

- By genre: comparisons between spoken, fiction, popular magazines, newspapers, and academic, or even between sub-genres (or domains), such as movie scripts, sports magazines, newspaper editorial, or scientific journals
- Over time: compare different years from 1990 to the present time

Collocations in Sketch Engine



user: Ms. Ana Slavec corpus: [enTenTen \[2012\]](#)

- Concordance
- Word List
- Word Sketch
- Thesaurus
- Find X
- Sketch-Diff
- Corpus Info
- [?](#)

Change options

participation/involvement

enTenTen [2012] freqs = [386,641](#) | [278,334](#)

participation 6.0 4.0 2.0 0 -2.0 -4.0 -6.0 involvement

and/or	80,545	56,488	0.2	0.2
non-participation	57	0	4.5	--
inclusiveness	57	0	4.4	--
equality	207	14	4.3	0.6
attendance	1,091	104	6.5	3.2
attainment	94	9	4.4	1.3
inclusion	505	69	5.5	2.7
representation	448	67	4.4	1.7
involvement	1,367	258	6.0	3.6
transparency	607	125	6.0	3.8
citizenship	222	50	4.8	2.8
openness	111	27	4.4	2.6
collaboration	641	196	4.8	3.1
accountability	366	110	5.2	3.6
cooperation	531	179	4.8	3.3
sponsorship	168	55	4.4	2.9
empowerment	355	121	6.1	4.7

subject_of	28,736	21,404	0.1	0.1
decline	60	0	2.3	--
exemplify	14	0	2.1	--
entail	21	0	2.0	--
energize	11	0	1.9	--
lag	10	0	1.6	--
outweigh	12	0	1.6	--
bolster	9	0	1.5	--
amount	12	0	1.4	--
constitute	30	13	1.3	0.1
enrich	47	21	3.0	1.9
evidence	54	33	3.8	3.1
characterise	19	14	2.7	2.4
wane	13	11	2.6	2.5
vary	79	75	1.3	1.3
characterize	50	53	2.2	2.3
cease	0	16	--	1.3

adj_subject_of	6,919	3,679	0.2	0.1
contingent	12	0	2.7	--
compulsory	18	0	2.6	--
anonymous	26	0	2.6	--
conditional	11	0	2.6	--
open	357	0	2.0	--
confidential	13	0	1.4	--
free	468	0	1.4	--
voluntary	740	19	7.0	1.7
mandatory	178	10	4.7	0.5
optional	117	9	4.2	0.5
welcome	94	20	2.7	0.5
vital	139	74	2.3	1.4
invaluable	15	11	1.7	1.3
essential	237	190	2.1	1.7
key	129	112	1.2	1.0
critical	137	147	1.7	1.8

WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

Noun

- [S: \(n\)](#) **impact** (the striking of one body against another)
- [S: \(n\)](#) **impact**, [wallop](#) (a forceful consequence; a strong effect) *"the book had an important impact on my thinking"; "the book packs a wallop"*
- [S: \(n\)](#) [impingement](#), [encroachment](#), **impact** (influencing strongly) *"they resented the impingement of American values on European culture"*
- [S: \(n\)](#) [shock](#), **impact** (the violent interaction of individuals or groups entering into combat) *"the armies met in the shock of battle"*

Verb

- [S: \(v\)](#) **impact** (press or wedge together; pack together)
- [S: \(v\)](#) [affect](#), **impact**, [bear upon](#), [bear on](#), [touch on](#), [touch](#) (have an effect upon) *"Will the new rules affect me?"*

<http://wordnetweb.princeton.edu/perl/webwn>

<http://globalwordnet.org/wordnets-in-the-world/>

Example

	COCA	enTenTen
assembling	1481	32.420
collecting	7537	208.996
gathering	11133	340.008
assembling information	7	83
collecting information	154	4417
gathering information	312	7542

Previous studies on effects of low frequency wordings

Indicator	Studies
Gaze times	<u>Longer</u> (Inhoff and Reyner 1986; Jurafsky 2003; Lenzner et al. 2011)
Response times	<u>Longer</u> (Lenzner et al. 2010) No sig. effect found (Slavec and Vehovar 2015)
Drop-out rate	No sig. effect found (Lenzner et al. 2010) <u>Higher but small effect</u> (Slavec and Vehovar 2015)
Item non-response	No sig. difference (Lenzner et al. 2010; Slavec and Vehovar 2015)
Satisficing	No sig. difference (Lenzner et al. 2010; Slavec and Vehovar 2015)
Subjective evaluation of difficulty	<u>Moderate effecte for the difficulty of understanding and small effect for the difficulty of providing answers</u> (Slavec and Vehovar 2015)

Present study: the questionnaire

- P1. In general, how well do you think the United States government is doing in reducing the threat of terrorism?
- P2. How worried are you that there will soon be another terrorist attack in the United States?
- P3. Do you think the use of torture against suspected terrorists in order to gain important information can ever be justified?
- P4. you completely agree, mostly agree, mostly disagree, or completely disagree with this statement?
 - I often worry about the chances of a nuclear attack by terrorists.
 - Freedom of speech should not extend to groups that are sympathetic to terrorists.
 - The police should be allowed to search the houses of people who might be sympathetic towards terrorists without a court order.
 - The government's anti-terrorism policies have gone too far in restricting the average person's civil liberties.
 - I am concerned that the government is collecting too much information about people like me.
- P5. As you may know, the United States government has a policy that it NEVER pays ransom money for hostages held by terrorist groups. Overall, do you approve or disapprove of this policy?
- P6. statement comes closer to your own views even if neither is exactly right? Please select:
 - Some religions are more prone to violence than others.
 - All religions are about the same when it comes to violence.
- P7. Which statement comes closer to your own views even if neither is exactly right? Please select:
 - The Islamic religion is more likely to encourage violence among its believers.
 - The Islamic religion does not encourage violence more than others.
- P8. How concerned, if at all, are you about Islamic extremism around the world these days?

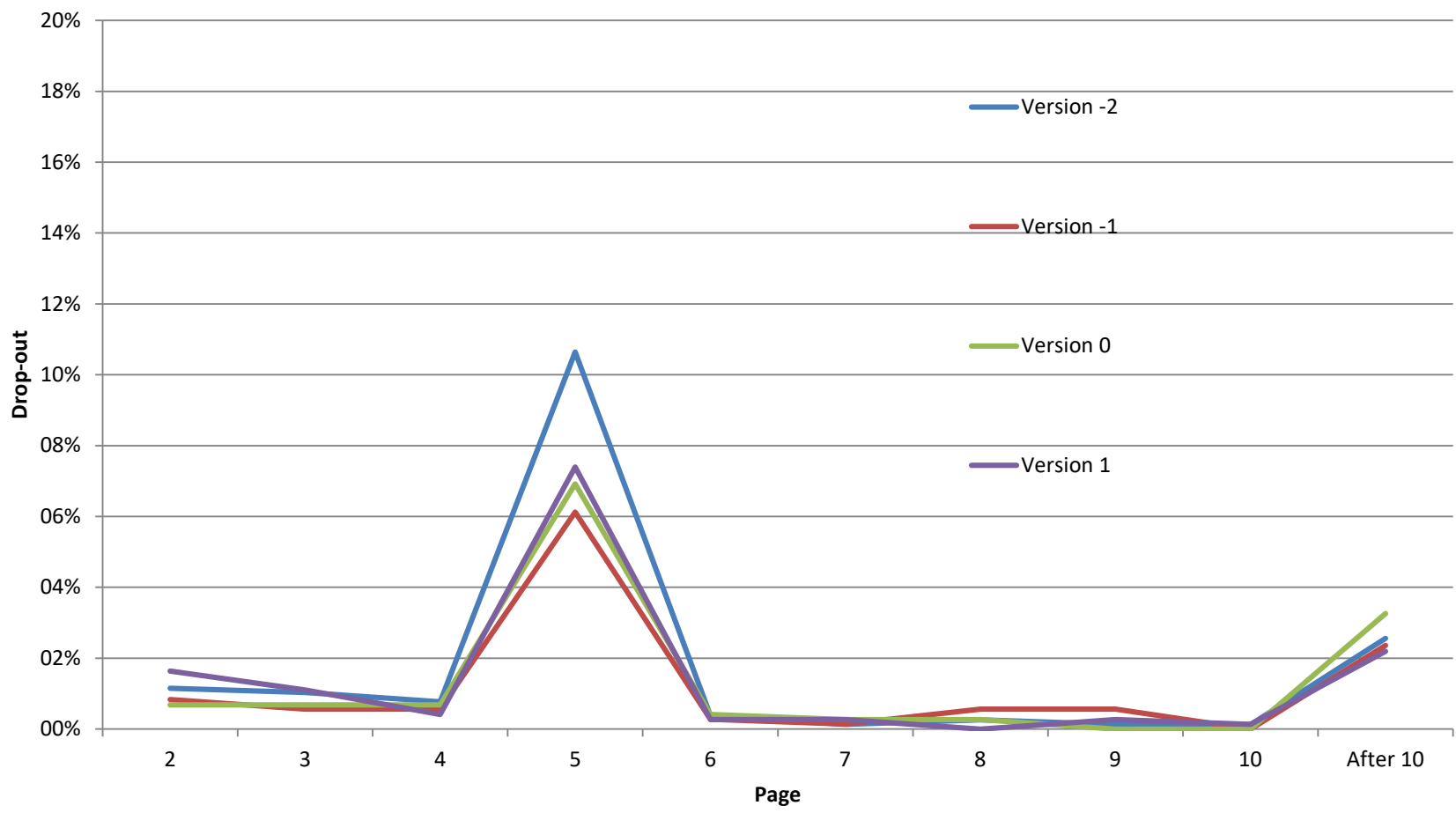
Comparison of four questionnaire versions

Differences		Version -2	Version -1	Version 0	Version 1
Number of changes		34	16	-	11
Median WF ratio	String	13.3	15.4	-	3.0
	Single word	8.2	8.2	-	2.6
Max WF ratio (+ examples)	String	7.240 Court → Tribunal	258 Encourage → Boost	-	25.4 Sympathetic to → Support
	Single word	169 Too far → Excessively	44 Reckon → Consider	-	497 Sympathetic to → Support

Data collection

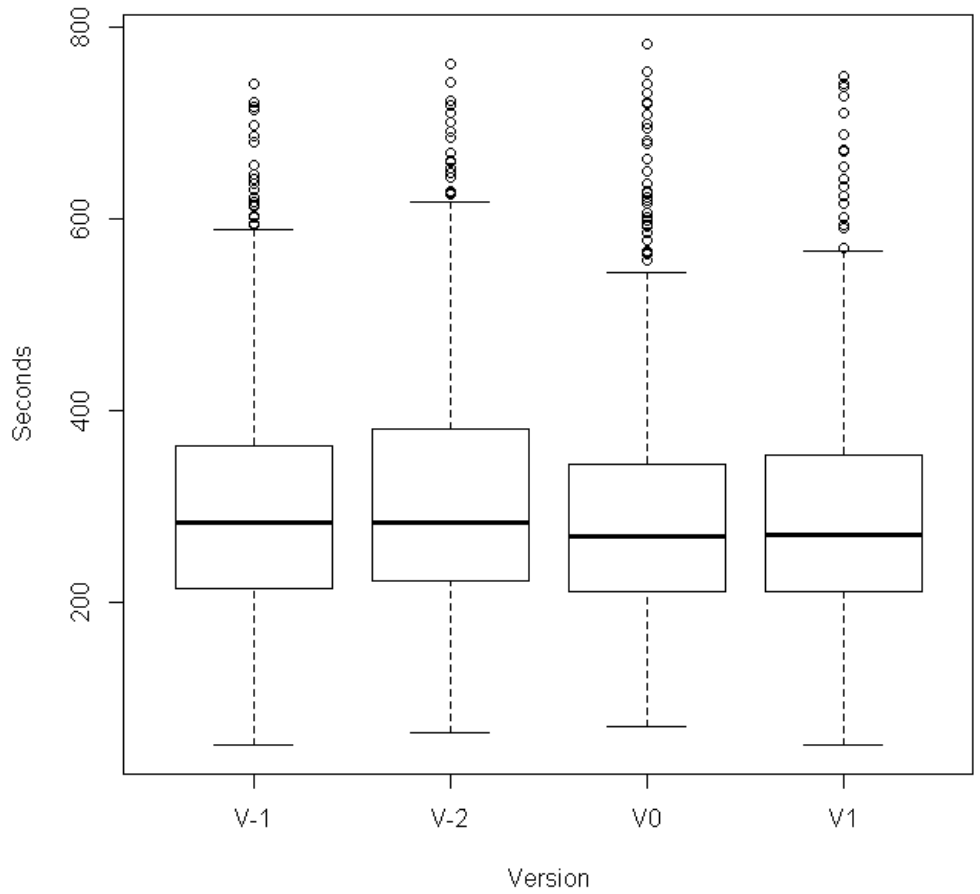
- Survey Monkey Audience Panel
 - Non-probability online panel recruited from a diverse population of Survey Monkey website visitors
 - Non-cash incentives (charitable donations)
- Sample of 2,4000 units (600 for each group)
- US residents, 18+
- Consistent socio-demographic structure across groups
- October 1-2 2015

Drop-out



RQ indicators	Version -2 N=780	Version -1 N=719	Version 0 N=739	Version 1 N=730
Drop-out	17%	12%	13%	13%

Response times



	Version -2	Version -1	Version 0	Version 1
Median time	4 m 53s	4m 50s	4m 35s	4m 37s
Median time (< 13 min)	4m 44s	4m 43s	4m 29s	4m 31s

% of DK answers

- P3. Do you think the use of torture against suspected terrorists in order to gain important information can ever be vindicated/justified?

	Version -2	Version -1	Version 0	Version 1
% DK	11.3 %	11.0 %	7.6 %	8.2 %

How much did you enjoy completing the questionnaire?

S1	Version -2	Version -1	Version 0	Version 1
N	656	642	649	640
1 - Not at all well	13.6%	13.2%	12.6%	12.5%
2 - A little	22.6%	18.7%	16.2%	20.3%
3 - A moderate amount	35.4%	41.1%	39.8%	38.3%
4 - A lot	16.5%	15.1%	18.2%	16.9%
5 - A great deal	12.0%	11.8%	13.3%	12.0%
Average	4.6	4.6	4.7	4.7

How difficult was to interpret the meaning of questions in this questionnaire?

S2	Version -2	Version -1	Version 0	Version 1
N	655	642	649	640
1 - Extremely difficult	1.1%	1.2%	0.8%	0.8%
2 - Very difficult	2.4%	1.1%	1.1%	1.1%
3- Moderately difficult	8.1%	6.5%	6.6%	5.9%
4 - Slightly difficult	15.3%	15.3%	9.7%	10.8%
5 - Not difficult at all	73.1%	75.9%	81.8%	81.4%
Average	4.8	3.9	3.0	4.6

How difficult was to interpret the meaning of questions in this questionnaire?

S3	Version -2	Version -1	Version 0	Version 1
N	649	635	641	633
10+ words	2.5%	2.4%	0.9%	1.6%
5-9 words	4.5%	2.5%	1.7%	2.7%
4 words	2.8%	2.0%	1.9%	1.1%
3 words	3.8%	3.9%	2.9%	2.5%
2 words	9.4%	7.5%	3.4%	4.6%
1 word	11.1%	10.4%	7.1%	7.1%
0 words	66.1%	71.2%	82.0%	80.4%
Average	1.3	1.1	0.6	1.0

Controlling for gender, education and language

- Men, the less educated and non-native speakers found the questionnaire more difficult than women, the more educated and native speakers
- Controlling the association between questionnaire difficulty and version (-2, -1, 0 and 1)
 - Language: association only for native speakers
 - Education: association only for those educated
 - Gender: association for both genders but weaker

Conclusions

- The worst version has a higher drop-out rate than the other three
- After removing outliers, the response time longer only for the worst version
- Except for one item, no effects on DK rates
- Respondents in the two worst versions found the questionnaire more difficult and reported a higher number of words that were at least a little difficult to understand
- Interaction with gender, education and language

Study limitations and potentials for future research

- Study Limited to case studies and selected examples
- Not all cases were pure synonyms
- Not all response quality indicators could be measured
- Future:
 1. Integration of language resources in questionnaire development tools
 2. Additional case studies and a meta-analytic approach to discover key factors that affect response quality



University of Ljubljana
Faculty *of Social Sciences*



CENTRE
FOR SOCIAL
INFORMATICS

Feedback and questions welcome.

Ana.Slavec@fdv.uni-lj.si

Expert evaluations

Text corpora frequency estimates (string)	Expert evaluations (median)	Context
Threat > Menace > Danger	Threat > Danger > Menace	... of terrorism
Upset > Concerned > ...	Concerned ≈ Worried	How ...
Justified > Legitimate > ...	Justified > Excused > ...	Ever ...
Risk > Chances > ...	Risk > Chances > of attack
Support > Sympathetic to	Support > Sympathetic to	... terrorists
Restricting > Curtailing	Limiting ≈ Restricting > liberties
Gathering > Collecting > ...	Gathering ≈ Collecting > information
Ransom > Demanded Hostages > Sureties	Ransom > Demanded Hostages > Sureties	... money for ...
Prone > Inclined	Prone ≈ Inclined	... to violence
Promote > Encourage > ...	Promote ≈ Encourage	... violence
Concerned > ...	Concerned ≈ Worried	... about extremism



Cognitive interviews

- Participants asked to paraphrase question or to define a certain item
- Half were assigned a LFW and half a HFW
- Level of match:
 - High: careful/cautious, threat/meance, and ransom/demanded m.
 - Medium: sympathetic/support, collecting/gathering, prone/inclined to, and chances/risk.
 - Low: worried/apprehensive, justified/vindicated, restricting/limiting, ecourage/promote, and concerned/preoccupied.
- When presented with a low-frequency wording, respondents used its high-frequency alternative.