

ABSTRACT

Traditionally, verbs like *base* have combined with the preposition *on* to express a meaning of derivation (*based on*). However, many writing in a US context have noticed the rapid rise of *based off (of)* alongside *based on* (Curzan 2013, Behrens 2014, Janda 2020). In this paper, we document the relative increase of *off* in two English-language corpora in the verb *base* and six other verbs. The results show a clear real-time trend of increasing use of *off*, with some differences in the course of the change across different verbs. We also see an increase in use of *off* in apparent time, which we infer from the topical organization of comments in one of our corpora, the social media site Reddit.

Keywords: Variation, Prepositions, Lexical diffusion, Apparent-time study, Corpus methodology

1 INTRODUCTION

2 This paper studies variation between *on* and *off* as the prepositional complement
3 of a select set of English verbs. One verb in which the variation has been well
4 documented is *base*; (1) gives examples of the variants.¹

5 (1) *Base on/off*

6 (a) I replied to your comment because you ***based it on*** a bunk article.

7 (b) So you didn't ***base it off of*** what the OP [original poster] said, you ***based***
8 ***it off of*** something in your head [...]

9 The *Oxford English Dictionary* (s.v. *base*) gives only examples with *on* (or *upon*)
10 complements, dating back to 1776. But the variation demonstrated in (1) has re-
11 ceived some attention in the linguistic literature, much of it observing rapid change
12 in progress. Janda (2020) finds examples of forms like *based off (of)* from as early
13 as 1980, but dates his first encounter with the *off* variant to ca. 2000, and suggests
14 rapid change thereafter:

15 [W]ithin a few years, the strength and breadth of this construction (in
16 the sense of characterizing almost everyone below a certain age) had
17 become evident. (596)

18 This is confirmed by Curzan (2013), who finds that *based off of* is rare in the Corpus
19 of Contemporary American English (Davies, 2008–) but growing: in Google Books
20 Ngram Corpus data from 2000 (Michel et al. 2011), *based on* outnumbers *based off*

¹All of the numbered examples provided in this paper are from the r/Parenting subreddit of the
Reddit corpus described in Section 3.1 unless a different subreddit source is noted.

21 *of* by 100,000:1, but by 2008, this has fallen to 10,000:1. Janda (2020: p. 597) finds
 22 that Google hits containing *based on* outnumber those containing *based off (of)* at a
 23 ratio of only 163:1, and the raw numbers of *based off (of)* hits are high, exceeding
 24 50 million. Finally, Behrens (2014) provides a more qualitative assessment of the
 25 growing popularity of *based off of* (as opposed to *based on*):

26 As of this writing, I hear it and see it written all the time from my
 27 students and from my younger colleagues; my older colleagues dismiss
 28 the structure as just plain wrong. (67)

29 Anecdotally, we note that the *off* variant is used in pop-up text in Google Sheets as
 30 of July 2024 (Figure 1).

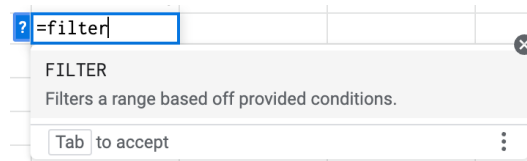


Figure 1: *Based off* in pop-up text in Google Sheets.

31 Janda (2020: p. 597) observes that this variation between *on* and *off* can be
 32 found with other verbs, namely *derive*, *ground*, *justify*, *predicate*, *draw*, *go*, and
 33 *live*. Examples of this variation in still other verbs are provided in (2)–(7).

34 (2) *Build on/off*

35 (a) Peaceful parent, happy siblings should give you a good foundation to
 36 ***build on***.

37 (b) We don't use all of it, but it gave us a good foundation to ***build off of***.

38 (3) *Capitalize on/off*

39 (a) Those people are only *capitalizing on* parents that are unprepared and
40 worried about disappointing their kids.

41 (b) They just don't seem to have a problem *capitalizing off of* our goodwill
42 and never reciprocating, so that's the problem.

43 (4) *Feed on/off*

44 (a) Ignore the tantrums, she's *feeding on* them.

45 (b) It sounds like he's *feeding off* your stress.

46 (5) *Profit on/off*

47 (a) You have legal rights since they are *profiting on* your son[']s image, no
48 matter how little he may have been involved.

49 (b) If you are *profiting off* my image without my knowledge in a space that
50 isn't public, then you owe me compensation.

51 (6) *Survive on/off*

52 (a) For three months we *survived on* our credit cards, then when the credit
53 ran out, we burned through the savings we had set aside for remodeling.

54 (b) We *survived very well off* my dad[']s salary so it wasn't for the money,
55 just for something to do.

56 (7) *Thrive on/off*

57 (a) She also was likely traumatized by the repeated moving and insecure liv-
58 ing arrangements - know that adage kids *thrive on* consistency?

59 (b) Remember children *thrive off of* consistency, that is how they feel safe
60 and calm.

61 Examples (1)–(7) demonstrate that the variation occurs in a variety of tenses and
62 aspects, with and without intervening object pronouns and adverbs.

63 Our contributions in this paper are as follows. First, we present novel quantita-
64 tive evidence for variation and change in the prepositional complements of the verbs
65 given in (1)–(7). We provide real-time evidence for increased use of *off* (*of*) in in-
66 formal written language, drawing on data from the online discussion forum Reddit,
67 and in the Corpus of Contemporary American English (Davies, 2008–). This builds
68 on Janda’s and Curzan’s corpus studies by presenting data on verbs beyond *base*,
69 and by including *off* with and without the following *of*. Second, we present a proof
70 of concept for an innovative methodology that uses the structure of Reddit to infer
71 the demographics of authors – in particular, their age – thus also providing apparent-
72 time evidence for increased use of *off* (*of*). Although the Reddit corpus has been
73 used for sociolinguistic research before (e.g. Brook & Blamire 2023), its potential
74 for inferring demographics is previously unexploited. Thus, our study suggests that
75 Reddit, whose enormous size makes it a valuable potential source of sociolinguistic
76 data, can be used (with caution) to study demographic factors like age (and, likely,
77 geography) despite lacking overt demographic metadata.

78 2 VARIATION AND CHANGE IN PREPOSITIONS ELSEWHERE IN ENGLISH

79 Variation and change in prepositions has been attested elsewhere in English. The
80 *based on/off* variable is reminiscent of variation in the complement of *different*, for
81 which *from*, *than*, and *to* are all attested, with geographical and social condition-
82 ing of their use (Iyeiri, Yaguchi & Okabe 2004, Mair 2007). Behrens & Mercer
83 (2007), Behrens (2014), and Schlüter (2022) give additional examples of preposi-

84 tion variation in fixed expressions in English, some of which can be observed in the
85 writing of contemporary native American English speakers – such as *have concerns*
86 *on* (standardly *about*) and *look forward for* (standardly *to*) – and others of which
87 show regional variation – such as *chat with*, which skews North American, versus
88 *chat to*, which skews British. We know of no indication, however, that these cases
89 of prepositional variation are showing anything like the rapid change observed for
90 *base*.

91 The *off* variant of our variable is implicated in another case of prepositional
92 variation: the variable presence of *of* after *off*. This variation is fairly widespread
93 in English, appearing not just with *off* (e.g. *get off (of) the bus*, *the islands off (of)*
94 *the coast*), but also with other prepositions and words that take *of*-headed com-
95 plements: *out (of) the window*, *all (of) the children*, *not that big (of) a deal* (Est-
96 tling 1999, 2000, Nylund & Seals 2010, Vartiainen & Höglund 2020). This vari-
97 ation between *of* and \emptyset shows social and geographical conditioning, though the
98 specifics depend on the particular construction: for instance, after *off*, the use of
99 *of* is deemed non-standard and prescribed against in formal writing (Vartiainen &
100 Höglund 2020), while after *out*, *of* is favored in formal written language (Estling
101 1999). We do not speculate on the social correlates or diachronic trajectory of the
102 *of* variant in the *based off* construction, instead grouping together *off* and *off of*
103 variants. Still, an interesting direction for future work could be to examine their
104 patterning separately, given the rarity of variationist studies that examine ternary
105 variables like *based on* – *based off* – *based off of* (MacKenzie 2020).

106 In fact, there are even more combinations of prepositions possible in the con-
107 struction under study. Both *on* and *off* appear sporadically in our corpus preceding

108 *from*:

109 (8) (a) And how much of ANS2 *builds on from* prior knowledge from ANS1?

110 (r/UCDavis)

111 (b) His test[s] are very straight forward, heavily *based off from* his lectures.

112 (r/UniversityOfHouston)

113 As far as we are aware, this combination of prepositions has not been previously
114 remarked upon, and it is quite rare – we group the 34 such sentences in our data
115 under the broader umbrellas of *on* and *off*.

116 3 METHODS

117 3.1 *The variable and the data sources*

118 We examine variation and change in the prepositional complements of seven verbs
119 (*base, build, capitalize, feed, profit, survive, and thrive*) in two corpora. These verbs
120 were selected through manual inspection of tokens of verbs appearing with both
121 *on* and *off* in a small sample corpus of posts from the social media website Reddit
122 (Chang et al. 2020, Baumgartner 2019), described in more detail below. Verbs were
123 selected primarily for practical purposes – for example, verbs that yielded too many
124 irrelevant tokens (such as *live*, whose combination with *off* and *on* often expresses
125 a location, like *Many students live on Fifth Street*) were not included. This list of
126 verbs is intended to be representative, not exhaustive: our goal is to show that the
127 shift is occurring in at least this handful of verbs.

128 Of these seven verbs, none is attested in the Oxford English Dictionary (OED)
129 with an *off* complement, but five are attested with *on* complements, confirming that

130 they traditionally take *on* in the standard language: *base*, *build*, *capitalize*, *feed*, and
131 *thrive*. Of the remaining two, *survive* is not shown combining with any prepositions
132 in the OED, but is attested with *on* in Google Ngrams (Michel et al. 2011) and our
133 data. *Profit* is perhaps the outlier among our verbs: the OED lists it as combining
134 with other prepositions, *by*, *of*, and *from*, whose usage rates exceed those of *on* in
135 Google Ngrams at most time points. However, *profit on* is attested fairly robustly
136 in Google Ngrams from 1800, and appears in our data, as in (5). Accordingly, we
137 include it in our data (and return to its special status in Section 5).

138 Because our studied variable is infrequent (other than with *base*), we prioritized
139 large data sets. Our data come from two sources, chosen for their size, their ease of
140 use, and their ability to provide real- and apparent-time data. The first is a corpus of
141 posts from Reddit (Chang et al. 2020, Baumgartner 2019), a news and discussion
142 website divided into topic-specific ‘subreddits’, such as ‘r/linguistics’, a forum for
143 discussion of topics and questions related to linguistics, and ‘r/Legomarket’, a fo-
144 rum where users coordinate buying, selling, and swapping LEGO products. Within
145 a subreddit, discussions are grouped into threads: for instance, r/linguistics contains
146 discussion threads devoted to specific academic articles and weekly Q&A threads
147 where users are encouraged to ask and answer linguistics-related questions. Our
148 Reddit data ranges from 2009 to 2018, though data before 2012 is sparse. The
149 Reddit corpus contains over 7 billion utterances – that is, post submissions and
150 comments (Baumgartner et al. 2020).

151 For this study, we selected posts from subreddits comprising three rough ‘age
152 cohorts’: college, pregnancy, and young parent. Our college cohort data set includes
153 posts from the r/college subreddit and subreddits from individual colleges (Ding

154 2018). The other cohorts comprise posts from r/BabyBumps (a pregnancy-related
155 forum) and r/Parenting, respectively. These age cohorts are intended to show the
156 presence of change in apparent time: we presume that participants in college-related
157 forums tend to be younger than those in pregnancy forums, who in turn tend to
158 be somewhat younger than participants in parenting discussions. These subreddits
159 included 19.4 million utterances (of which 13 million were on college subreddits)
160 with a total of 993 million words (547 million from college subreddits).²

161 Our second data source is the Corpus of Contemporary American English (COCA,
162 Davies, 2008–), which includes approximately one billion words from 1990–2019
163 in eight genres across formal written language (academic texts, newspapers, mag-
164 azines, fiction), online written language (websites, blogs), television and movie
165 subtitles, and spoken language (unscripted conversations from television and radio
166 programs).

167 Our COCA data uses the magazines, newspapers, and spoken language genres;
168 each includes approximately 125 million words. Two web-based genres (web and
169 blog posts) were excluded because they are all indexed to 2012. The other genres
170 were excluded after preliminary searches showed very little use of the *off* variant.

171 Of our two data sources, Reddit plays the primary role. It has important ad-
172 vantages: it is very large and is written in more informal language, meaning that it
173 contains many tokens of our verb–preposition constructions (three times as many as
174 COCA; see Section 3.2 for precise counts). The division into subreddits also allows
175 us to sample from (presumed) different demographics.

²These numbers count punctuation marks as separate words and thus overstate the size of the corpus slightly.

176 Reddit also has downsides as a source of sociolinguistic data. Its text is not
177 lemmatized, so we cannot restrict our searches to verbal forms only, increasing the
178 false positive rate (although we took measures to mitigate this; see Section 3.2). The
179 geographical distribution of the Reddit data is also difficult to determine: Reddit
180 draws users from around the world, and the college subreddits include colleges
181 from outside the US (Ding 2018). In addition, the Reddit data falls within a narrow
182 window of time, primarily 2012–2018.

183 These limitations of Reddit lead us to caution in using it to study variation in
184 American English. Accordingly, we conduct a parallel study in COCA. Although
185 COCA also does not contain sociodemographic information, its texts are all Amer-
186 ican English and its data is generally high-quality. COCA also has part-of-speech
187 tagging, which in theory allows us to target verbal forms (however, the tagger some-
188 times misclassifies nouns like *building* as verbs; see Section 3.2). In addition, the
189 greater time scale of COCA (stretching back to 1990) allows us to observe variation
190 in the use of *off* for longer, and before use of *based off (of)* began to become salient
191 – around 2000, according to Janda (2020) and Curzan (2013).

192 At the same time, COCA has disadvantages compared to Reddit. Much of its
193 text, even in the genres chosen, is more formal and edited. While we do look for
194 genre differences in COCA, there is no expected demographic difference (and thus,
195 no apparent-time interpretation) between the genres.

196 Thus, our main, more interesting findings are in the Reddit data. The COCA
197 data is interpreted primarily as a sanity check on the Reddit results: since the two
198 data sets produce qualitatively similar results, we conclude that the Reddit data is
199 sensible.

200 3.2 Data extraction

201 We searched for various constructions including one of our seven verbs followed
202 by the preposition *on* or *off*. These two components could be adjacent or separated
203 by a nominal phrase³ and/or one or more adverbs. In order to distinguish verbal
204 constructions (e.g. *was based on it, will profit off it*) from non-verbal constructions
205 (e.g. *a class based on it, make a profit off it*), we also tracked instances in which
206 the verb was preceded by an auxiliary like forms of *be*, again with possible adverbs
207 intervening. This allowed us to control for part of speech when modeling (see
208 Section 3.3).

209 Although some of the verbs being studied (*build, feed, and survive*) optionally
210 take an overt object, they most reliably show the desired variation when intransitive.
211 The verb *build* shows *on/off* variation only in its metaphorical meaning, which the
212 OED defines as ‘to establish, develop, or construct (something abstract, such as a
213 system of thought or belief, a reputation, a relationship, etc.)’ (e.g. *The new law is*
214 ***built on solid legal principles***). However, transitive or passive uses of *build on/off*
215 more often involve physical construction, meaning that false-positive sentences like
216 *The first buildings were **built on** campus in 1812* are very common, especially on the
217 college subreddits. On the other hand, intransitive *build on/off*, as exemplified by
218 (2) above, is exclusively metaphorical. Similarly, intransitive *feed* shows variation
219 in preposition whether metaphorical (as in (4) above) or literal (e.g. *Some birds **feed***
220 ***off insects***), while transitive *feed* includes many more irrelevant examples showing

³Nominal phrases could be composed of a stand-alone *pronoun* or a sequence of words centered around a noun, where optional components are in parentheses: (*article/determiner/possessive pronoun*) (*numeral*) (*adjective(s)*) *noun*.

221 no variation: *feeding my baby on the couch, off my plate*, and so on. Likewise,
222 *survive* can take an object in the meaning desired (e.g. *I survived my pregnancy*
223 *on plain pasta*), but such cases have higher rates of false positives because the
224 prepositional phrase can be part of the object (like *The king survived the attempt on*
225 *his life*). To limit ourselves to a consistent construction that yields the most reliable
226 data, we exclude tokens with an intervening object (indicative of a transitive verb)
227 for all verbs except *base* (which is only ever used transitively).

228 Data from Reddit was retrieved from the Pushshift.io Reddit Corpus (Baum-
229 gartner 2019) through ConvoKit (Chang et al. 2020). We used a Python script to
230 search for the sequences described in the previous paragraph. The Reddit Corpus
231 is not lemmatized, so we conducted a string-based search using lists of forms ac-
232 cording to their parts of speech in CELEX (Baayen, Piepenbrock & Gulikers 1995).
233 Thus, for example, hits for the verb *survive* included the words *survive, survived,*
234 *survives, and surviving*.

235 This type of search naturally yields false positives. To investigate how many,
236 we looked at a sample of 100 sentences for a number of configurations based on
237 verb form, presence of a direct object (for *base*), and presence of an auxiliary (if
238 a given category had fewer than 100 sentences, we looked at all of them). This
239 sample revealed several frequent undesired prepositional phrases, which we filtered
240 out of our data: *on/off campus* (with up to three words intervening to account for
241 phrases results like *on the main campus*, extremely common in college subreddits),
242 *on X's own* (with one word between the preposition and *own*, most common with
243 *survive* and *thrive*), and *on demand/a schedule/a routine* (commonly used to discuss
244 feeding practices in the pregnancy and parenting subreddits).

245 Configurations that still yielded rates of false positives above 13% in the sample
246 were removed as well. These included:

- 247 • *bases* not followed by an object (often nominal: *the **bases on** the baseball*
248 *field*)
- 249 • *building* (often nominal: *a **building on** campus*)
- 250 • *built* (often passive: *a community **built on** respect, housing **built on** the quad*)
- 251 • most forms of *feed* without auxiliaries (often nominal: *a **feed on** YouTube*) or
252 with passive auxiliaries (reliably passive: *the students were **fed on** junk food*)
- 253 • *profit or profits* (often nominal: *make a **profit off** his image*) unless preceded
254 by an auxiliary (e.g. *the school doesn't **profit off** of certain classes*)

255 While the remaining data does still have a small proportion of false positives,
256 we do not think they substantially skew the results. In fact, the results are quite
257 robust to the presence of false positives: earlier versions of the data set, with fewer
258 tokens removed, yielded very similar results.

259 In COCA, each word is tagged with its lemma and part of speech. Our COCA
260 search included forms of our seven verbs tagged as verbs. While we expected that
261 COCA would have fewer false positives, this turned out to not always be true, and
262 we removed tokens with *on X's own* and several configurations that had false posi-
263 tive rates above 15%. As with the Reddit corpus, some of these had nouns misclas-
264 sified as verbs (*building, profits*). The tagger classifies *fed* as either a past participle
265 or a past-tense form; the former were removed, as they were more often passive.

266 The tagger was not so accurate with *built*, so all sentences with this form were re-
 267 moved, whether it was tagged as a participle or past tense. Sentences with *survived*,
 268 *surviving*, and *thriving* were also removed due to false positives stemming from lo-
 269 cational prepositional phrases (e.g. *public education is **thriving on** the West Coast*).
 270 Ironically, this seems to be an issue specific to COCA because its texts are *more*
 271 *diverse* than our Reddit data, where many of the locational prepositional phrases
 272 for *survive* and *thrive* involved mentions of campus and were thus easy to filter out.

273 Token counts by corpus and lemma are provided in Table 1.

	<i>base</i>	<i>build</i>	<i>capitalize</i>	<i>feed</i>	<i>profit</i>	<i>survive</i>	<i>thrive</i>	total
Reddit	133 675	3497	803	533	242	1498	1252	141 500
COCA	41 744	1648	1770	1644	150	355	976	48 287

Table 1: Token counts by corpus and lemma.

274 3.3 Statistical analysis

275 The data was analyzed using logistic regression in R (R Core Team 2023). For each
 276 corpus, we fitted three regressions involving different subsets of verbs and factors to
 277 account for the differences between *base* and other verbs: first of all, *base* dwarfs
 278 the other verbs in frequency, comprising 94% of tokens for Reddit and 86% for
 279 COCA. Second, as described in Section 3.2, *base* is transitive (appearing either
 280 with an object or as a passive), while the others are intransitive in our data set. This
 281 difference in syntactic construction makes it difficult to compare *base* with the other
 282 verbs.

283 Our dependent variable is preposition, coded as a binary between *on* (marked
 284 as 0) and *off* (marked as 1). The sequence *off of* is classified as *off* and is quite

285 common: *off* without *of* is only slightly more frequent than *off of* in Reddit and
286 3.2 times more frequent than *off of* in COCA).

287 The regressions were fitted using using R's *buildmer* package (Voeten 2023)
288 through forward stepwise comparison; factors were only included in the model if
289 they significantly improved its fit and improved its Akaike Information Criterion
290 (AIC), which penalizes model complexity (additional model factors). Factors that
291 improved the model but were problematic due to sparse data were removed.

292 The first regression for each corpus looks only at *base*; the second looks at the
293 remaining verbs, which are obligatorily intransitive. The third regression compares
294 the verbal passive construction *be based* (that is, *based* preceded by a passive aux-
295 iliary) to the other intransitive verbs. This filtering increases parallels between *base*
296 and the other verbs by removing two configurations in which *base* regularly appears
297 but the other verbs do not: transitive uses in which *base* is separated from its prepo-
298 sition by an overt object (like *the professor based the textbook on his lectures*) and
299 adjectival passives that are not fully verbal in structure (like *a textbook based on*
300 *the professor's lectures*).

301 The models included the following key factors:

- 302 • Year (centered around the median year with substantial data, 2015 for Reddit
303 and 2005 for COCA)
- 304 • Source/genre: college (baseline) vs. pregnancy vs. parenting for Reddit, mag-
305 azine (baseline) vs. news vs. spoken for COCA
- 306 • Verb: not used in first regression (*base* only), sum-coded in second regression
307 (all verbs other than *base*), and dummy-coded in third regression (all verbs,

308 with *base* as baseline)

309 Two-way interaction terms between these three factors were considered as candi-
310 dates for the models.

311 A number of morphosyntactic factors were also considered. These factors dif-
312 fered slightly according to the properties of the model, as follows. The regressions
313 for *base* had a candidate factor comparing passives with intervening adverbs to pas-
314 sives with no interveners on the one hand and to actives with intervening objects
315 and, optionally, adverbs on the other. As shown in Table 2, the factors of voice
316 and presence of an intervener are confounded, in that active sentences must have
317 interveners;⁴ accordingly, these two factors were combined into a single factor to
318 avoid an unbalanced combination of factors in the models.

	active	passive
no intervener	—	<i>based on</i>
intervener	<i>based it (mostly) on</i>	<i>based entirely on</i>

Table 2: Interaction of voice and presence of intervener for *base*.

319 The other two regressions for each corpus, which included only verbs without
320 direct objects, had a candidate factor marking the presence of an intervening adverb,
321 again to test for an effect of verb–preposition adjacency.

322 The regression for verbs other than *base* also included a factor comparing unin-
323 flected verb forms (*capitalize*) to third-person singular (*capitalizes*), past/participle

⁴The corpora do contain a small number of tokens coded as active without interveners. Some of these involve extraction of the object (e.g. *He presented papers that he **based on** his research*), while many are typos where the last letter of *based* is omitted. All such tokens were removed from our data.

324 (*capitalized*), and progressive (*capitalizing*) forms; this factor was excluded from
325 the regression with all verbs (*base* included) because the form of *base* tokens in
326 this regression was uniformly *based*. Finally, all of the regressions had a candidate
327 factor marking whether the object of the preposition was definite (that is, beginning
328 with *the*).

329 Output for all models can be found in the Appendix.

330 4 RESULTS

331 We test the following hypotheses on the Reddit data:

- 332 1. A real-time shift toward *off*: the proportion of *off* is increasing year-by-year.
- 333 2. An apparent-time shift toward *off*: the proportion of *off* hits in the college-
334 age cohort will be higher than that of the pregnancy-age cohort, which will in
335 turn be higher than that of the parent-age cohort.

336 We find both of these hypotheses to be confirmed in the Reddit data set, which
337 we present first. Afterward, we replicate the real-time trend in the COCA data
338 set, as well as many of the comparisons between *base* and the other verbs. As
339 described in Section 3.1, the COCA data set is smaller but has better tagging and
340 metadata. Replication of the Reddit results in COCA strengthens our confidence
341 that the Reddit data is giving us a real signal despite its shortcomings (e.g. dialectal
342 heterogeneity).

343 4.1 Reddit

344 Figure 2 shows the rate of use of *off* (as opposed to *on*), aggregated across all seven
 345 verbs studied (*base, build, capitalize, feed, profit, survive, thrive*), over nine years
 346 of real time in the Reddit corpus. Though *off* is the minority variant, it shows a
 347 steady, linear rise from 7% to 10% over the decade. The effect of year is significant
 348 ($p \leq .004$) in all three regressions (*base* alone, all other verbs, all verbs combined).

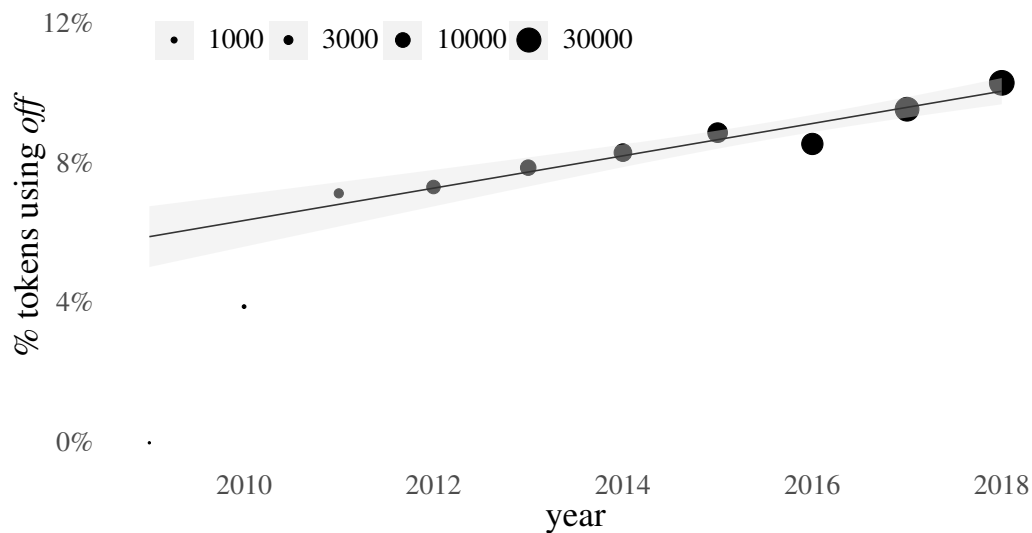


Figure 2: Rate of use of *off* in real time, aggregated over all seven verbs studied, Reddit data.

349 Figure 3 adds an apparent-time perspective to Figure 2 by plotting the college,
 350 pregnancy, and parenting cohorts separately. We see a neat cohort effect: college
 351 posters have the highest rate of *off*, parenting posters have the lowest, and preg-
 352 nancy posters are in the middle. In all three Reddit regression models, the differ-

ences between the cohorts are generally significant ($p \leq .04$).⁵ Figure 3 suggests that the difference between the college and pregnancy cohorts is equivalent to about five years of real time (about two percentage points, approximately half the rise shown by the aggregated data over the decade studied), and the difference between the pregnancy and parenting cohorts is similar, if somewhat larger.

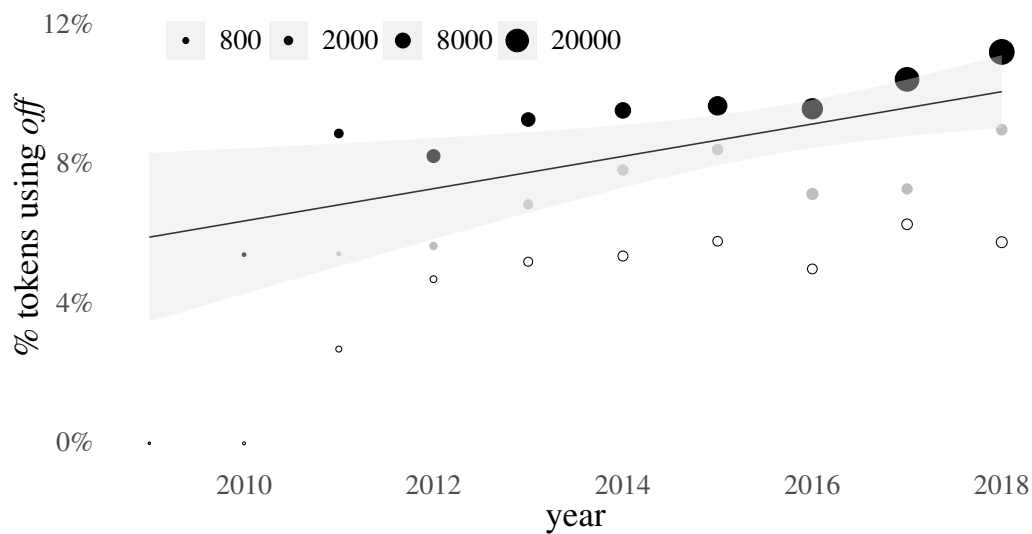


Figure 3: Rate of use of *off* in real time, aggregated over all seven verbs studied, by subreddit type (college vs. pregnancy vs. parenting), Reddit data.

The apparent-time effect can also be derived from the regression models. Table 3 shows the coefficients for year and subreddit type for the three models based on the Reddit data. The coefficient for year represents the estimated yearly change

⁵The effect of pregnancy in the regression containing intransitive verbs has $p = .04$; all others have $p < .001$. The estimated marginal means (cf. Lenth 2023) between pregnancy and the other two cohorts in the two regressions containing non-*base* verbs are also not significant, likely due to the presence of an interaction term between verb and subreddit type.

361 in use of *off*, while the coefficients for pregnancy and parenting compare those
 362 respective cohorts with the baseline, college (in the models, these coefficients are
 363 negative, since *off* is used less frequently in these subreddits than in college subred-
 364 dits). Dividing the subreddit type coefficient by the year coefficient thus gives an
 365 estimate of the apparent-time effect of subreddit types. Indeed, the first two models
 366 yield plausible results, indicating that posters in pregnancy and parenting subreddits
 367 are 5–7 and 10–15 years older than posters in college subreddits, respectively. The
 368 estimated apparent-time effects of the model comparing passive *based* to the other
 369 verbs has a much larger estimated effect (15 and 33 years, respectively); however,
 370 this is likely due to the substantially lower coefficient size for year, which in turn
 371 may reflect the fact that much of the weight of year in this model is caught up in its
 372 interaction with verb.

model	year	pregnancy	parenting	pregnancy/year	parenting/year
<i>base</i>	0.054	0.362	0.818	6.70	15.15
intransitive verbs	0.072	0.394	0.755	5.47	10.49
all verbs	0.029	0.423	0.962	14.59	33.17

Table 3: Absolute value of coefficients for year and subreddit type for Reddit models, with their quotients (interpretable as apparent-time differences).

373 Finally, Figure 4 shows verb-by-verb data for verbs other than *base*. Since 94%
 374 of tokens are *base*, the real- and apparent-time patterns for this verb are largely
 375 captured by the aggregated patterns in Figure 3, and including them in Figure 4
 376 would lead to an issue in depicting the differences in scale. At the baseline of
 377 2015, the verb *capitalize* has a significantly lower rate of *off* than *base* ($\beta = -1.55$,
 378 $p < .001$), while all other verbs have significantly higher rates of *off* than *base*
 379 ($p < .001$). For *profit*, in particular, the rate of *off* is very high – nearly at ceiling.

380 In addition, most of the verbs show a similar pattern of real- and apparent-
381 time effects as *base*, though often with more noise: real-time increase, with col-
382 lege posters leading parenting and pregnancy posters. The one main exception is
383 *feed*, where the aggregate real-time line in Figure 4 trends *downward*. Indeed, in
384 the model comparing *base* with other verbs, *feed* is the only verb with a negative
385 interaction with year (though it is not significant). The interaction term ($-.061$) is
386 greater in absolute value than the main effect of year ($.029$), meaning that the model
387 suggests that use of *off* is decreasing for *feed* year-by-year (not just increasing more
388 slowly than *base*). This downward trend seems to be concentrated in a larger num-
389 ber of tokens of *feed on* than expected in the last couple of years in parenting and
390 pregnancy forums. While we have no explanation for this distribution, we note that
391 these forums include frequent discussion of babies' feeding habits. Many of the
392 false positives (including the common *feed on demand*, see Section 3.2) have been
393 successfully filtered out, but some remain. The relatively small number of tokens
394 and issue with specialized vocabulary mean that this verb's results should be taken
395 with a grain of salt. There is one verb with a significant interaction term with year:
396 the rate of *off* increased significantly more quickly for *survive* than *base*.

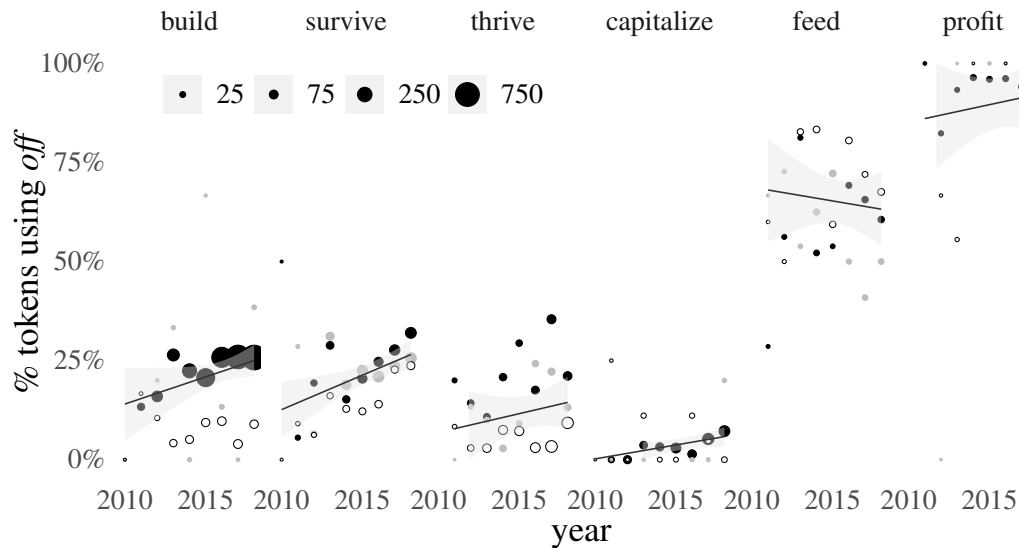


Figure 4: Rate of use of *off* in real time for verbs other than *base*, split by verb (ordered by frequency in the studied corpus) and by subreddit type (college vs. pregnancy vs. parenting), Reddit data.

397 4.2 COCA

398 Figure 5 shows the rate of *off* (as opposed to *on*), with or without a following
 399 preposition, across all verbs in real time from 1990 to 2019 in COCA. Compared
 400 to Reddit, *off* is much less common in COCA: even in 2019, the rate of *off* only
 401 reaches about 3%, compared to 10% in Reddit. However, there is a clear trend
 402 upwards, as *off* appeared well below 1% of the time in 1990. The effect of year is
 403 significant ($p < .001$) in all three regressions.

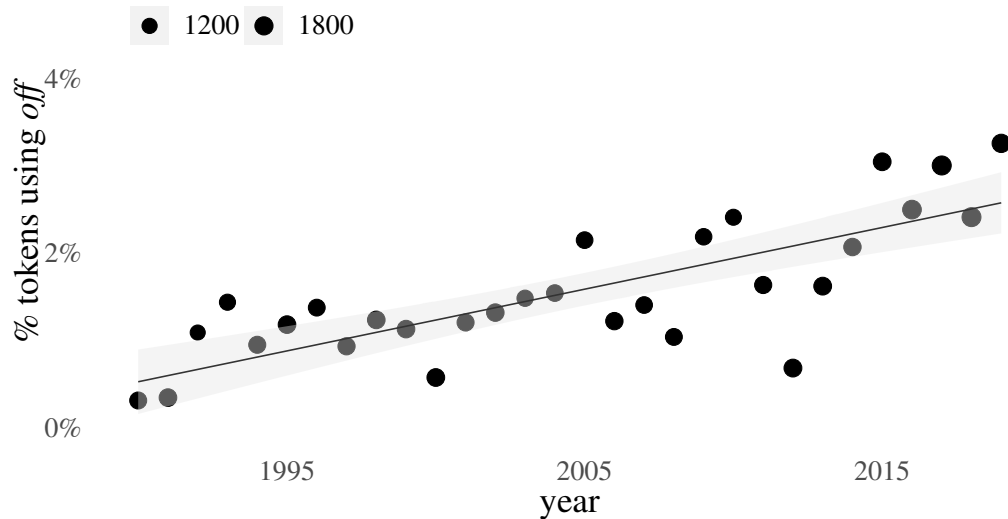


Figure 5: Rate of use of *off* in real time, aggregated over all seven verbs studied, COCA data.

404 Figure 6 splits the data according to text type: magazines, newspapers, and spo-
 405 ken language. Here we do not see the same stark pattern as in Reddit: the three text
 406 types are intermingled and seem quite similar on visual inspection. The statistical
 407 models do detect significant differences. In the model limited to *base*, *off* appears
 408 in spoken language more often than in magazines ($p < .001$), but there is no sig-
 409 nificant difference between magazines and newspapers; comparison of estimated
 410 marginal means finds a significant difference between spoken language and news-
 411 papers ($p < .001$). The difference is equivalent to 5.6 years of real time, given that
 412 the coefficient for spoken language is 5.6 times greater than the year coefficient;
 413 however, there is no reason to suspect that this corresponds to any apparent-time
 414 difference, especially because the difference is not consistent year-over-year as it
 415 is with subreddit types in the Reddit data. In the model including verbs other than

416 *base*, both newspapers and spoken language have higher use of *off* than magazines
 417 ($p < .001$ for both); in fact, *off* appears *more often* in newspapers than spoken lan-
 418 guage, though the difference is not significant. However, as we will see below, this
 419 effect seems to be located in specific verbs; this model does not include an interac-
 420 tion term between verb and text type because the categorical patterning of *survive*
 421 in spoken language (76 tokens, all with *on*) throws off the confidence intervals for
 422 all of the spoken-language interaction terms.

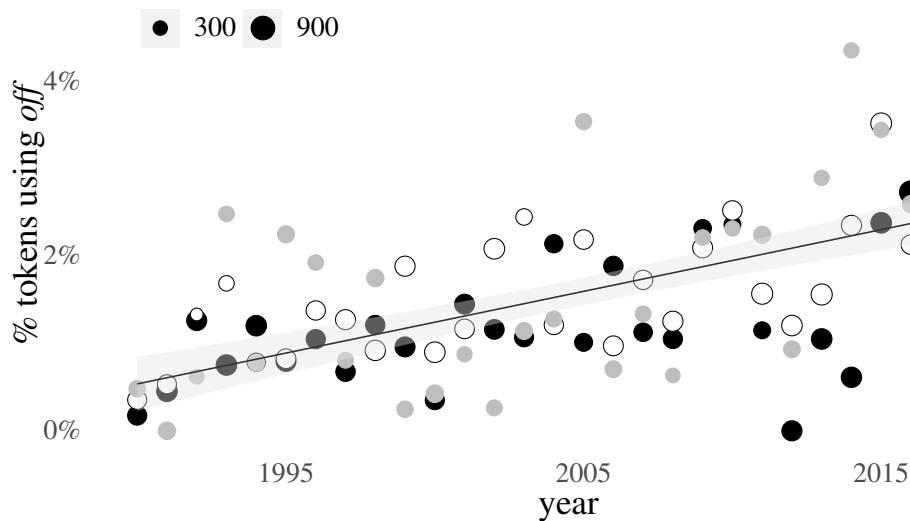


Figure 6: Rate of use of *off* in real time, aggregated over all seven verbs studied, by text type (magazines vs. newspapers vs. **spoken language**), COCA data.

423 Finally, Figure 7 shows verb-by-verb data for verbs other than *base*. We see
 424 that *capitalize* and *survive* very rarely, if ever, take *off*, while *profit* almost always
 425 takes *off*. Meanwhile, *feed* shows a stark genre difference in its use of *off*: maga-
 426 zines have a low rate of *feed off*, while newspapers and especially spoken language

427 show much higher rates. From reading the graph, we would expect that the dif-
428 ferences in text type should be concentrated in their interaction with verb, and this
429 is what we see in the regression comparing passive *be based* to the other verbs.
430 According to this model, *off* is used more often in spoken language than maga-
431 zines, though not quite significantly so ($\beta = .75, p = .054$), while there is almost
432 no difference between newspapers and magazines ($\beta = .02, p = .960$). However,
433 looking at the interaction term, *feed off* appears significantly and substantially more
434 often in newspapers than in magazines ($\beta = 1.54, p = .001$). Inspection of the rel-
435 evant cases reveals no obvious pattern explaining this effect. The verb *profit* has
436 significant interactions as well: *profit off* appears significantly *less* often in news-
437 papers ($\beta = -2.11, p = .023$) and spoken language ($\beta = -2.38, p = .008$) than
438 in magazines, in which we find 42 tokens of *profit off* and only two of *profit on*.
439 Finally, *capitalize off* is more common in spoken language than magazines, though
440 the effect does not reach significance ($\beta = 2.06, p = .071$).

441 The verb-by-verb results in COCA are qualitatively similar to those of Reddit:
442 *feed* has a somewhat higher rate of *off* than most of the verbs, suggesting that the
443 high rate of *off* for *feed* is not due to idiosyncrasies of the data source. Likewise,
444 *profit* appears almost entirely with *off*. Since COCA has a much lower rate of *off*
445 in general, the very low rate of *off* for *capitalize* does not stand out from that of the
446 other verbs; its difference from *be based* is not significant. The other verbs appear
447 with *off* significantly more than *be based* ($p \leq .04$).

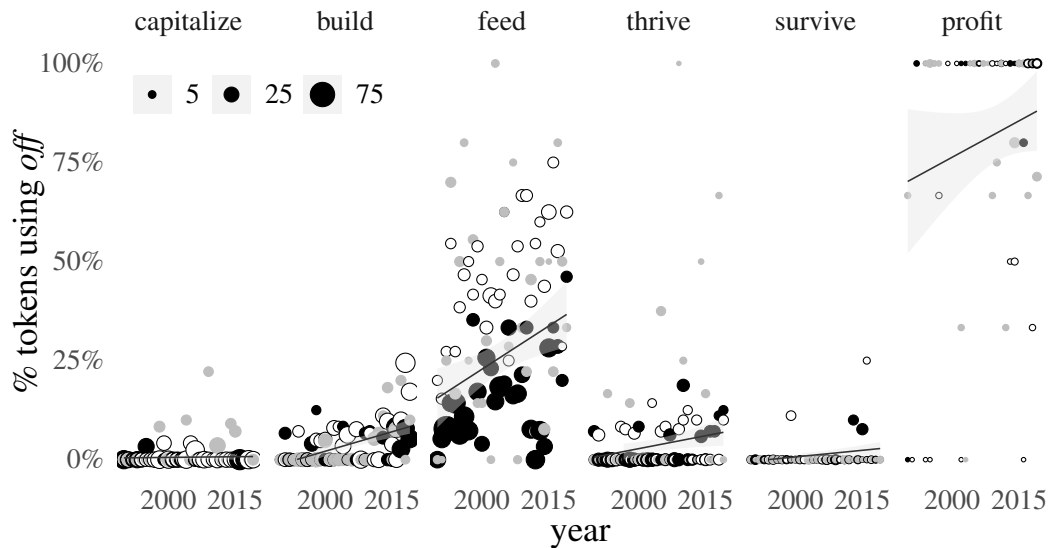


Figure 7: Rate of use of *off* in real time for verbs other than *base*, split by verb (ordered by frequency in the studied corpus) and by text type (magazines vs. newspapers vs. **spoken language**), COCA data.

448 4.3 Morphosyntactic factors

449 In the previous sections, we discussed extrinsic and lexical factors affecting the
 450 choice of preposition: real time, subreddit/genre, and verb. Some of our regressions
 451 also showed significant morphosyntactic effects: verb form, presence of material
 452 intervening between verb and preposition, and whether the prepositional object is
 453 definite.

454 The regression with *base* on Reddit data showed that active uses of this verb
 455 are significantly more likely to take *off* than passive uses, and the difference is very
 456 large ($\beta = 1.96$, $p < .001$; for comparison, the effect size of real time is $\beta = .05$
 457 per year). Other morphosyntactic effects did not substantially improve the model

458 and were not added. The *base* model for COCA does not include this factor, be-
459 cause it is categorical in the COCA data: *off never* occurs in passive uses of *base*
460 with an adverbial intervener, whereas in active uses and in passive uses without an
461 intervener, *off* is merely very rare.

462 The Reddit regression with the intransitive verbs (that is, all except *base*) shows
463 an effect of verb form: *off* occurs significantly more often with progressive forms
464 (e.g. *surviving*) than uninflected forms (e.g. *survive*; $\beta = .56, p < .001$). Verb form
465 does not improve the COCA model and is not added to it.

466 Another detectable syntactic effect in the regressions including verbs other than
467 *base* is that *off* is used less often when an adverbial intervenes between verb and
468 preposition (e.g. *survives mostly on*) than when verb and preposition are adjacent.
469 This factor is significant in the Reddit model comparing *base* to other verbs ($\beta =$
470 $-.30, p < .001$), though it is not added to the Reddit model comparing the non-*base*
471 verbs to one another. It is significant in both of these models for the COCA data.

472 Finally, the two COCA models including verbs other than *base* have one more
473 significant syntactic factor: *off* is used more often when its complement is definite
474 (starts with *the*). This factor is not added to the Reddit models.

475 4.4 Summary

476 Our findings can be summarized as follows:

- 477 • The *off* variant is steadily increasing in real time, in both corpora. The effect
478 is strongest for *base*, which is the most frequent verb, but is present for others
479 too.
- 480 • The Reddit corpus shows an apparent-time increase of the *off* variant as well.

481 While there is some difference in the rate of *off* between different text types
482 in COCA, we do not see an analogous steady, consistent gap.

483 • Different verbs are at different points in the change toward *off*. In both cor-
484 pora, *profit* and *feed* take higher rates of *off* than others, with *profit* almost
485 categorically taking *off*.

486 • The use of *off* is also influenced by internal morphosyntactic factors. Most
487 consistently, *off* is less common when the verb and preposition are separated
488 by an adverb.

489 5 DISCUSSION

490 The previous section confirmed that all seven verbs studied here take *off* comple-
491 ments. The general trend, with few exceptions, is that all verbs are changing toward
492 use of *off* in both real and, where data is available, apparent time.

493 5.1 Main findings

494 There are two questions we want to address here. The first is: why is this change
495 happening? The second is: why is it happening in these verbal constructions specif-
496 ically? After all, *off* is not replacing *on* across the English language in general:
497 not when *on* is used with its core physical meaning, nor when it is used in other
498 metaphorical ways, such as *airing on television* or *kept on file*.

499 To answer the first question, we turn to a suggestion from Janda 2020. In ex-
500 plaining the rise of *based off (of)*, Janda proposes:

501 [I]f one derives something from a source, then a crucial pathway be-

502 tween them leads from the source to the derivative; something takes off
503 from – or is taken off (of) – the source and travels – or is brought –
504 to/as the derivative. [...] Yet *basing* or *being based ON* portrays the
505 implied motion as oriented in [the opposite] direction, and thus sounds
506 more like planting a flagstaff downward into the ground. (597)

507 In other words, *off* suggests extraction, while *on* suggests foundation. Perhaps,
508 then, the shift from *on* to *off* is a change from a metaphor of foundation to one
509 of extraction. That, then, suggests an answer to the second question: the verbs
510 undergoing this change are those that are compatible with this ‘extraction’ meaning.

511 Somewhat speculatively, we observe that there may be a correlation between the
512 strength of a verb’s association with these meanings of extraction and/or foundation
513 and its likelihood of change. We found in both corpora that *profit* had the highest
514 rates of *off* by far: even going back to the 1990s in the COCA data, *profit off* well
515 exceeds 50% *off* usage (19/25 tokens in that decade). Data from the Google Books
516 Ngram Corpus from 2019, plotted in Figure 8, likewise show that *profit off* started
517 gaining ground in the 1990s. This suggests that *profit* was an earlier shifter than the
518 other verbs, not categorically different from them. Indeed, the lexical semantics of
519 *profit* are particularly well-suited for a metaphor of extraction.

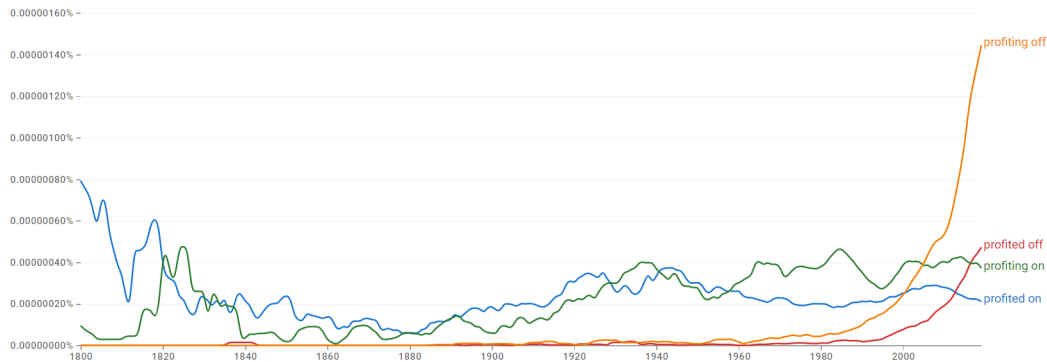


Figure 8: Rates of *on* and *off* with forms of *profit* in Google Book Ngram data from 2019.

520 By contrast, *base* shows lower rates of *off* than any other verb except *capital-*
 521 *ize*, a pattern that holds in both corpora. Again, the lexical semantics of *base* are
 522 particularly compatible with a metaphor of foundation, perhaps leading it to have
 523 resisted shifting longer than the others. (On the other hand, the even lower rate of
 524 *off* with *capitalize* may be a reflection of its rather more educated nature.)

525 We note also that *base* is by far the most frequent verb in both data sets. The
 526 resistance of highly frequent forms to change is well known from work on mor-
 527 phological changes such as analogical leveling (e.g. Hooper 1976). Still, despite
 528 its overall low rate of *off*, *base* is making up for lost time: in COCA, at least, it is
 529 changing toward *off* faster than all the other verbs in our study (though this differ-
 530 ence is not significant for two verbs).

531 All told, the picture here suggests a change that has started in different verbs
 532 at different times, and is progressing for different verbs at different rates. This is
 533 reminiscent of lexical diffusion in phonology: a change starts in one environment
 534 (in this case, potentially *profit*) and then gradually expands to others (Wang 1969).

535 The lexically specific nature of this change raises questions about whether any

536 social evaluation associated with the change is similarly lexically specific. As noted
537 in Section 1, *based off* has reached the level of salience associated with prescriptive
538 judgments of incorrectness. It is likely that *base* has attracted overt comment due
539 to its high frequency and its rapid rate of change. But given that the other six verbs
540 studied here are also changing in the same direction, do they share the same social
541 evaluation? The question of whether the social evaluation of a variant extends to
542 all environments in which that variant surfaces, or whether social evaluation may
543 interact with internal (linguistic) constraints on variation, is longstanding in soci-
544 olinguistics, going back at least to Weiner & Labov 1983. This could be examined
545 in follow-up work on this change.

546 5.2 Text sources and apparent time

547 Both the Reddit and the COCA corpus are divided up into three different text
548 sources. In the former, examples were drawn from subreddits themed around col-
549 lege, pregnancy, and parenting. We use this as a proxy for apparent time, under the
550 assumption that posters on college subreddits are younger (and thus more advanced
551 in the change) than posters on pregnancy subreddits, who are in turn younger than
552 people posting about parenting. This assumption was borne out: qualitatively, Fig-
553 ure 3 shows a fairly consistent difference between the three subreddit types that
554 looks equivalent to the difference of a few years; this impression is largely con-
555 firmed quantitatively as well by the models, as shown in Table 3.

556 The COCA corpus included texts from magazines, newspapers, and spoken lan-
557 guage. This difference, in contrast, is not expected to correspond to apparent time,
558 since it reflects differences in the mode of production rather than the demograph-

559 ics of people producing the texts. Indeed, Figure 6 presents a stark contrast to the
560 Reddit data in Figure 3: the rates of *off* among the different text types are quite in-
561 termingled and close together, and certainly do not show the lockstep pattern of the
562 Reddit data. The COCA models generally show that *off* is more common in spoken
563 language than magazines – which is understandable, given that the latter likely go
564 through more editing. However, the lack of a sharp cohort effect in COCA further
565 reinforces the apparent-time interpretation of the Reddit data.

566 5.3 Morphosyntactic factors

567 The regressions showed evidence of a number of internal grammatical factors influ-
568 encing the use of *off* – in particular, *off* is used less frequently when an adverb or
569 adverbial phrase intervenes between the verb and the preposition (like *based almost*
570 *entirely on*). There is no obvious explanation for this effect. If the shift in prepo-
571 sitions is lexically driven – that is, mediated by each individual lexical item – then
572 intervening material between verb and preposition may make the ties between the
573 two weaker and cause reversion to a default preposition – which, at least for now,
574 is more likely to be *on*.

575 6 CONCLUSION

576 The purpose of this study was twofold. First, we aimed to capture an in-progress
577 shift in the prepositional complement of verbs like *base* from *on* to *off*. While this
578 change has been previously documented for both *base* and other verbs, this is the
579 first study that systematically investigates the change in progress across multiple
580 verbs, text types, and morphosyntactic contexts. Our results are clear: *off* is used

581 across many different verbs, and its use is increasing in both real and apparent
582 time. Moreover, this change shows no sign of stopping and looks to be picking
583 up other verbs in its path as well: examples of *off* can be found even with verbs
584 like *depend* and *rely*, which the authors of this paper, who generally accept the
585 tokens in our corpora and likely produce *off* fairly regularly as well, find crashingly
586 bad.⁶ Thus, our study lays important groundwork for future study of this linguistic
587 variable, including both sociolinguistic factors that we did not study systematically
588 (region, gender, etc.) and a closer look at internal linguistic factors, including those
589 we studied and those we did not. For example, one direction for future research
590 is to investigate geographical patterning of this variation in a large data set with
591 geographical metadata, such as a Twitter-based corpus.

592 The second main purpose of our study was methodological. Reddit represents
593 an enormous body of informal text that could serve as a valuable resource for socio-
594 and other linguistic research. However, its users are anonymous and we have no de-
595 mographic information about them. Thus, we use this study as a proof of concept
596 for the efficacy of the Reddit corpus. Its results are qualitatively similar to those of

⁶Examples, taken from outside the subreddits studied in this paper, are shown below.

- (i) Edit: forgot to mention on the TooGoodToGo app you can get a bunch of bagels for \$3.99
depending off the bagel shop. There are also many other cool findings, so you should check it
out. (r/AskNYC)

- (ii) They **relied off of** my written statement more than anything, because I have issues talking about
any of it, so the statement will be important if you have issues talking about it as well.
(r/Veterans)

597 a more cultivated corpus, COCA: even though some of the tokens in the Reddit cor-
598 pus were undoubtedly made by non-native speakers, it is reliable enough to display
599 broad trends. Moreover, the organization of Reddit into subreddits, which are often
600 very specific, allows us to innovatively approximate its users' demographic proper-
601 ties – in this case, age. This analytical move was successful: the effects of subreddit
602 were interpretable in terms of time and yielded sensible results well within the range
603 of our expectations from the inferred demographic correlates of subreddit. Thus, we
604 hope that this study will serve as inspiration for future use and exploration of Reddit
605 as a source of sociolinguistic data, both general and demographically specific.

606 A MODEL OUTPUT

607 In Tables 5 and 8, the factor representing the six verbs is sum-coded: the five factors
608 compare each of the first five verbs to the grand mean (the mean of the means of the
609 dependent variable – in this case, likelihood of *off* – for each verb). The estimate
610 for the sixth verb, *profit*, is the negative sum of all five factors. Thus, in the tables
611 below, we provide an estimate for *profit* (the negative sum of the five factors) but
612 not a standard error or *p* value, since it is not represented by a separate factor in
613 the model. Each of the five factors, in turn, includes some influence of *profit* in
614 addition to the listed verb. Thus, the estimates for sum-coded factors are rather
615 easier to interpret than their corresponding standard error and *p* values.

616 Terms are listed in the order in which they are added to the model, roughly
617 corresponding with importance. Low *p* values are marked as follows: *** for $p <$
618 $.001$, ** for $.001 \leq p < .01$, * for $.01 \leq p < .05$, . for $.05 \leq p < .1$.

619 A.1 Reddit models

	β	SE	p
Intercept	-2.44	.07	<.001 ***
Voice and intervener (default: passive, intervener)			
active, intervener	1.96	.07	<.001 ***
passive, no intervener	-0.03	.07	.694
Subreddit type (default: college)			
pregnancy	-0.36	.03	<.001 ***
parenting	-0.82	.04	<.001 ***
Year	0.05	.01	<.001 ***

Table 4: Coefficients for model with all tokens of *base*, Reddit data.

	β	SE	p
Intercept	-0.75	.08	<.001 ***
Verb (compared to grand mean)			
build	-0.46	.08	<.001 ***
survive	-0.47	.10	<.001 ***
thrive	-0.52	.12	<.001 ***
capitalize	-2.66	.18	<.001 ***
feed	1.17	.15	<.001 ***
profit	2.93	—	—
Subreddit type (default: college)			
pregnancy	-0.39	.19	.039 *
parenting	-0.75	.14	<.001 ***
Verb form (default: base)			
third-person singular	-0.01	.07	.844
progressive	0.56	.12	<.001 ***
past	-0.07	.14	.639
Year	0.07	.01	<.001 ***
Verb * Subreddit type (compared to grand mean)			
build * pregnancy	0.29	.32	.365
survive * pregnancy	0.18	.22	.430
thrive * pregnancy	-0.24	.27	.378
capitalize * pregnancy	0.51	.64	.432
feed * pregnancy	0.11	.27	.674
profit * pregnancy	-0.85	—	—
build * parenting	-0.61	.21	.004 **
survive * parenting	0.23	.21	.269
thrive * parenting	-0.88	.22	<.001 ***
capitalize * parenting	1.09	.41	.007 **
feed * parenting	1.14	.23	<.001 ***
profit * parenting	-0.97	—	—

Table 5: Coefficients for model with intransitive tokens of verbs other than *base*, Reddit data.

	β	SE	p
Intercept	-1.86	.02	<.001 ***
Verb (default: base)			
build	0.66	.06	<.001 ***
survive	0.72	.11	<.001 ***
thrive	0.57	.13	<.001 ***
capitalize	-1.55	.24	<.001 ***
feed	2.34	.16	<.001 ***
profit	4.40	.27	<.001 ***
Subreddit type (default: college)			
pregnancy	-0.42	.07	<.001 ***
parenting	-0.96	.08	<.001 ***
Year	0.03	.01	.004 **
Intervener (default: no)			
yes	-0.30	.08	<.001 ***
Verb * Subreddit type (default: base * college)			
build * pregnancy	0.31	.32	.341
survive * pregnancy	0.28	.15	.064 .
thrive * pregnancy	-0.15	.24	.545
capitalize * pregnancy	0.56	.76	.464
feed * pregnancy	0.29	.24	.220
profit * pregnancy	-0.90	.72	.207
build * parenting	-0.42	.21	.050 *
survive * parenting	0.44	.20	.031 *
thrive * parenting	-0.69	.23	.003 **
capitalize * parenting	1.29	.48	.007 **
feed * parenting	1.49	.24	<.001 ***
profit * parenting	-0.71	.58	.221
Verb * Year (default: base)			
build * year	0.03	.03	.180
survive * year	0.07	.03	.044 *
thrive * year	0.09	.05	.055 .
capitalize * year	0.18	.10	.068 .
feed * year	-0.06	.05	.193
profit * year	0.01	.11	.920

Table 6: Coefficients for model with passive tokens of *base* and intransitive tokens of other verbs, Reddit data.

620 A.2 COCA models

	β	SE	p
Intercept	-7.13	.24	<.001 ***
Year	0.17	.02	<.001 ***
Text type (default: magazine)			
newspaper	0.06	.25	.826
spoken	0.94	.22	<.001 ***

Table 7: Coefficients for model with all tokens of *base*, COCA data.

	β	SE	p
Intercept	-3.60	.16	<.001 ***
Verb (compared to grand mean)			
capitalize	-2.50	.28	<.001 ***
build	-1.01	.20	<.001 ***
feed	1.95	.15	<.001 ***
thrive	-0.59	.20	.004 **
survive	-1.98	.56	<.001 ***
profit	4.13	—	—
Text type (default: magazine)			
newspaper	1.23	.11	<.001 ***
spoken	0.95	.14	<.001 ***
Year	0.07	.02	<.001 ***
Prepositional object (default: indefinite)			
definite	0.51	.11	<.001 ***
Intervener (default: no)			
yes	-1.38	.41	.001 ***
Verb * Year (compared to grand mean)			
capitalize * year	-0.05	.03	.121
build * year	0.04	.02	.054 .
feed * year	-0.03	.02	.051 .
thrive * year	0.00	.02	.889
survive * year	0.05	.06	.393
profit * year	-0.01	—	—

Table 8: Coefficients for model with intransitive tokens of verbs other than *base*, COCA data.

	β	SE	p
Intercept	-6.67	.38	<.001 ***
Verb (default: base)			
capitalize	-0.03	1.07	.974
build	2.51	.48	<.001 ***
feed	4.89	.39	<.001 ***
thrive	2.79	.49	<.001 ***
survive	1.91	.94	.042 *
profit	9.49	.84	<.001 ***
Text type (default: magazine)			
newspaper	0.02	.43	.960
spoken	0.75	.39	.054 .
Year	0.16	.03	<.001 ***
Intervener (default: no)			
yes	-1.53	.42	<.001 ***
Prepositional object (default: indefinite)			
definite	0.48	.11	<.001 ***
Verb * Text type (default: base * magazine)			
capitalize * newspaper	0.96	1.23	.434
build * newspaper	0.77	.52	.140
feed * newspaper	1.54	.45	.001 ***
thrive * newspaper	0.35	.61	.565
survive * newspaper	0.36	1.10	.741
profit * newspaper	-2.11	.93	.023 *
capitalize * spoken	2.06	1.14	.071 .
build * spoken	-0.53	.56	.344
feed * spoken	0.30	.42	.487
thrive * spoken	0.60	.59	.311
survive * spoken	-12.76	259.81	.960
profit * spoken	-2.38	.90	.008 **
Verb * Year (default: base)			
capitalize * year	-0.14	.05	.002 **
build * year	-0.05	.03	.109
feed * year	-0.12	.03	<.001 ***
thrive * year	-0.09	.04	.012 *
survive * year	-0.04	.09	.658
profit * year	-0.12	.04	.002 **

Table 9: Coefficients for model with passive tokens of *base* and intransitive tokens of other verbs, COCA data.

REFERENCES

621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643

Baayen, R. H., R. Piepenbrock & L. Gulikers. 1995. *The CELEX Lexical Database (CD-ROM)*. Philadelphia.

Baumgartner, Jason. 2019. Reddit corpus. <https://files.pushshift.io/reddit/> (16 July, 2024).

Baumgartner, Jason, Savvas Zannettou, Brian Keegan, Megan Squire & Jeremy Blackburn. 2020. The Pushshift Reddit dataset. *Proceedings of the International AAAI Conference on Web and Social Media* 14. 830–839.

Behrens, Susan J. 2014. *Understanding language use in the classroom: a linguistic guide for college educators*. Toronto: Multilingual Matters.

Behrens, Susan J. & Cindy Mercer. 2007. The style of which this is written: neutralization of prepositions in English. *NADE Digest* 3(2). 47–58.

Brook, Marisa & Emily Blamire. 2023. Language play is language variation: quantitative evidence and what it implies about language change. *Language* 99(3). 491–530.

Chang, Jonathan P., Caleb Chiam, Liye Fu, Andrew Wang, Justine Zhang & Cristian Danescu-Niculescu-Mizil. 2020. ConvoKit: a toolkit for the analysis of conversations. In Olivier Pietquin, Smaranda Muresan, Vivian Chen, Casey Kennington, David Vandyke, Nina Dethlefs, Koji Inoue, Erik Ekstedt & Stefan Ultes (eds.), *Proceedings of the 21th annual meeting of the Special Interest Group on Discourse and Dialogue*, 57–60. Association for Computational Linguistics.

Curzan, Anne. 2013. Based off of what? <https://www.chronicle.com/blogs/linguafranca/based-off-of-what> (6 August, 2024).

- 644 Davies, Mark. 2008–. The Corpus of Contemporary English (COCA). <https://www.english-corpora.org/coca/> (16 July, 2024).
645
- 646 Ding, Karl. 2018. GitHub - karlding/college-subreddits: an eventually-complete list
647 of college (and university) subreddits. <https://github.com/karlding/college-subreddits> (10 June, 2021).
648
- 649 Estling, Maria. 1999. Going out (of) the window? *English Today* 15(3). 22–27.
- 650 Estling, Maria. 2000. Competition in the wastebasket: a study of constructions
651 with *all*, *both* and *half*. In Christian Mair & Marianne Hundt (eds.), *Corpus*
652 *linguistics and linguistic theory: papers from the Twentieth International Con-*
653 *ference on English Language Research on Computerized Corpora (ICAME 20),*
654 *Freiburg im Breisgau 1999*, 103–116. Amsterdam, Atlanta: Rodopi.
- 655 Hooper, Joan B. 1976. Word frequency in lexical diffusion and the source of mor-
656 phophonological change. In William M. Christie Jr. (ed.), *Current progress in*
657 *historical linguistics: Proceedings of the Second International Conference on*
658 *Historical Linguistics, Tucson, Arizona, 12–16 January 1976*, 96–105. Amster-
659 dam: North Holland.
- 660 Iyeiri, Yoko, Michiko Yaguchi & Hiroko Okabe. 2004. *To be different from* or *to be*
661 *different than* in present-day American English? *English Today* 20(3). 29–33.
- 662 Janda, Richard D. 2020. Perturbations, practices, predictions, and postludes in a
663 bioheuristic historical linguistics. In Richard D. Janda, Brian D. Joseph & Bar-
664 bara S. Vance (eds.), *The handbook of historical linguistics*, vol. II, chap. 24,
665 523–650. Hoboken, NJ: Wiley Blackwell.

- 666 Lenth, Russell V. 2023. *emmeans: Estimated Marginal Means, aka Least-Squares*
667 *Means*. R package version 1.9.0. [https://CRAN.R-project.org/package=](https://CRAN.R-project.org/package=emmeans)
668 [emmeans](https://CRAN.R-project.org/package=emmeans).
- 669 MacKenzie, Laurel. 2020. Comparing constraints on contraction using Bayesian
670 regression modeling. *Frontiers in Artificial Intelligence: Language and Compu-*
671 *tation* 3. 58.
- 672 Mair, Christian. 2007. British English/American English grammar: convergence in
673 writing—divergence in speech? *Anglia* 125(1). 84–100.
- 674 Michel, Jean-Baptiste, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew
675 K Gray, The Google Books Team, Joseph P. Pickett, Dale Hoiberg, Dan Clancy,
676 Peter Norvig, Jon Orwant, Steven Pinker, Martin A. Nowak & Erez Lieberman
677 Aiden. 2011. Quantitative analysis of culture using millions of digitized books.
678 *Science* 331(6014). 176–182.
- 679 Nylund, Anastasia & Corinne Seals. 2010. “It’s not that big (of) a deal”: the soci-
680 olinguistic conditioning of inverted degree phrases in Washington, DC. *Univer-*
681 *sity of Pennsylvania Working Papers in Linguistics* 16(2). 133–140.
- 682 R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. R
683 Foundation for Statistical Computing. Vienna. <https://www.R-project.org/>.
- 684 Schlüter, Julia. 2022. Language corpora and the teaching and learning of English
685 as an international language. In Marcus Callies & Stefanie Hehner (eds.), *Pluri-*
686 *centric languages and language education: pedagogical implications and inno-*
687 *vative approaches to language teaching*, 166–189. London: Routledge.

- 688 Vartiainen, Turo & Mikko Höglund. 2020. How to make new use of existing re-
689 sources: tracing the history and geographical variation of *off of*. *American Speech*
690 95(4). 408–440.
- 691 Voeten, Cesko C. 2023. *buildmer: Stepwise Elimination and Term Reordering for*
692 *Mixed-Effects Regression*. R package version 2.9. [https://CRAN.R-project.org/](https://CRAN.R-project.org/package=buildmer)
693 [package=buildmer](https://CRAN.R-project.org/package=buildmer).
- 694 Wang, William S-Y. 1969. Competing changes as a cause of residue. *Language*
695 45(1). 9–25.
- 696 Weiner, E. Judith & William Labov. 1983. Constraints on the agentless passive.
697 *Journal of Linguistics* 19(1). 29–58.