

Zipf’s law holds for phrases, not words

Jake Ryland Williams,^{1,*} Paul R. Lessard,^{2,†} Suma Desu,^{3,‡} Eric M. Clark,^{1,§}
James P. Bagrow,^{4,5,¶} Christopher M. Danforth,^{1,**} and Peter Sheridan Dodds^{1,††}

¹*Department of Mathematics & Statistics, Vermont Complex Systems Center,
Computational Story Lab, & the Vermont Advanced Computing Core,
The University of Vermont, Burlington, VT 05401.*

²*Department of Mathematics, University of Colorado, Boulder CO, 80309*

³*Center for Computational Engineering, Massachusetts Institute of Technology, Cambridge, MA, 02139*

⁴*Computational Story Lab, Vermont Advanced Computing Core,
& the Department of Mathematics and Statistics, University of Vermont, Burlington, VT, 05401*

⁵*Vermont Complex Systems Center, University of Vermont, Burlington, VT, 05401*

(Dated: March 5, 2015)

With Zipf’s law being originally and most famously observed for word frequency, it is surprisingly limited in its applicability to human language, holding over no more than three to four orders of magnitude before hitting a clear break in scaling. Here, building on the simple observation that phrases of one or more words comprise the most coherent units of meaning in language, we show empirically that Zipf’s law for phrases extends over as many as nine orders of rank magnitude. In doing so, we develop a principled and scalable statistical mechanical method of random text partitioning, which opens up a rich frontier of rigorous text analysis via a rank ordering of mixed length phrases.

PACS numbers: 89.65.-s,89.75.Da,89.75.Fb,89.75.-k

INTRODUCTION

Over the last century, the elements of many disparate systems have been found to approximately follow Zipf’s law—that element size is inversely proportional to element size rank [1, 2]—from city populations [2–4], to firm sizes [5], and family names [6]. Starting with Mandelbrot’s optimality argument [7], and the dynamically growing, rich-get-richer model of Simon [3], strident debates over theoretical mechanisms leading to Zipf’s law have continued until the present [8–11]. Persistent claims of uninteresting randomness underlying Zipf’s law [8] have been successfully challenged [9], and in non-linguistic systems, good evidence supports Simon’s model [3, 12, 13] which has been found to be the basis of scale-free networks [14, 15].

For language, the vast majority of arguments have focused on the frequency of an individual word which we suggest here is the wrong fundamental unit of analysis. Words are an evident building block of language, and we are naturally drawn to simple counting as a primary means of analysis (the earliest examples are Biblical concordances, dating to the 13th Century). And while we have defined morphemes as the most basic meaningful ‘atoms’ of language, the meaningful ‘molecules’ of language are clearly a mixture of individual words and phrases. The identification of meaningful phrases, or multi-word expressions, in natural language poses one of the largest obstacles to accurate machine translation [16]. In reading the phrases “New York City” or “Star Wars”, we effortlessly take them as irreducible constructions, different from the transparent sum of their parts. Indeed, it

is only with some difficulty that we actively parse highly common phrases and consider their individual words.

While partitioning a text into words is straightforward computationally, partitioning into meaningful phrases would appear to require a next level of sophistication requiring online human analysis. But in order to contend with the increasingly very large sizes and rapid delivery rates of important text corpora—such as news and social media—we are obliged to find a simple, necessarily linguistically naive, yet effective method.

A natural possibility is to in some way capitalize on N -grams, which are a now common and fast approach for parsing a text. Large scale N -gram data sets have been made widely available for analysis, most notably through the Google Books project [17]. Unfortunately, all N -grams fail on a crucial front: in their counting they overlap, which obscures underlying word frequencies. Consequently, and crucially, we are unable to properly assign rankable frequency of usage weights to N -grams combined across all values of N .

Here, we introduce ‘random partitioning’, a method that is fast, intelligible, scalable, and sensibly preserves word frequencies: i.e., the sum of sensibly-weighted partitioned phrases is equal to the total number of words present. As we show, our method immediately yields the profound basic science result that phrases of mixed lengths, as opposed to just individual words, obey Zipf’s law, indicating the method can serve as a profitable approach to general text analysis. To explore a lower level of language, we also partition for sub-word units, or graphemes, by breaking words into letter sequences. In the remainder of the paper, we first describe random partitioning and then present results for a range of texts.

TEXT PARTITIONING

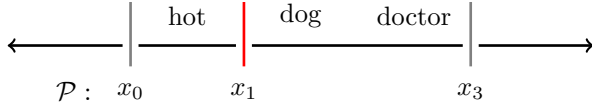
To begin our random partitioning process, we break a given text T into clauses, as demarcated by standard punctuation (other defensible schemes for obtaining clauses may also be used), and define the length norm, ℓ , of a given clause t (or phrase, $s \in S$) as its word count, written $\ell(t)$. We then define a partition, \mathcal{P} , of a clause t to be a sequence of the boundaries surrounding its words:

$$\mathcal{P} : x_0 < \dots < x_{\ell(t)}, \quad (1)$$

and note that $x_0, x_{\ell(t)} \in \mathcal{P}$ for any \mathcal{P} , as we have (a priori) the demarcation knowledge of the clause. For example, consider the highly ambiguous text:

“Hot dog doctor!”

Forgoing punctuation and casing, we might attempt to break the clause down, and interpret through the partition:



i.e., $\mathcal{P} = \{x_0, x_1, x_3\}$, which breaks the text into phrases, “hot” and “dog doctor”, and assume it as reference to an attractive veterinarian (as was meant in [18]). However, depending on our choice, we might have found an alternative meaning:

- hot dog; doctor: A daring show-off doctor.
 - : One offers a frankfurter to a doctor.
- hot; dog doctor: An attractive veterinarian (vet).
 - : An overheated vet.
- hot dog doctor: A frank-improving condiment.
 - : A frank-improving chef.
- hot; dog; doctor: An attractive vet of canines.
 - : An overheated vet of canines.

Note in the above that we (as well as the speaker in [18]) have allowed the phrase “dog doctor” to carry idiomatic meaning in its non-restriction to canines, despite the usage of the word “dog”.

Now, in an ideal scenario we might have some knowledge of the likelihood for each boundary to be “cut” (which would produce an ‘informed’ partition method), but for now our goal is generality, and so we proceed, assuming a uniform boundary-cutting probability, q , across all $\ell(t) - 1$ word-word (clause-internal) boundaries of a clause, t . In general, there are $2^{\ell(t)-1}$ possible partitions of t involving $\frac{1}{2}\ell(t)(\ell(t) + 1)$ potential phrases. For each integral pair i, j with $1 \leq i < j \leq \ell(t)$, we note that the probability for a randomly chosen partition

of the clause t to include the (contiguous) phrase, $t_{i\dots j}$, is determined by successful cutting at the ends of $t_{i\dots j}$ and failures within (e.g., x_2 must *not* be cut to produce “dog doctor”), accommodating for $t_{i\dots j}$ reaching one or both ends of t , i.e.,

$$P_q(t_{i\dots j} | t) = q^{2-b_{i\dots j}}(1-q)^{\ell(s)-1} \quad (2)$$

where $b_{i\dots j}$ is the number of the clause’s boundaries shared by $t_{i\dots j}$ and t . Allowing for a phrase $s \in S$ to have labeling equivalence to multiple contiguous regions (i.e., $s = t_{i\dots j} = t_{i'\dots j'}$, with $i, j \neq i', j'$) within a clause e.g., “ha ha” within “ha ha ha”, we interpret the ‘expected frequency’ of s given the text by the double sum:

$$f_q(s | T) = \sum_{t \in T} f_q(s | t) = \sum_{t \in T} \sum_{s=t_{i\dots j}} P_q(t_{i\dots j} | t). \quad (3)$$

Departing from normal word counts, we may now have $f_q \ll 1$, except when one partitions for word ($q = 1$) or clause ($q = 0$) frequencies. When weighted by phrase length, the partition frequencies of phrases from a clause sum to the total number of words originally present in the clause:

$$\ell(t) = \sum_{1 \leq i < j \leq \ell(t)} \ell(t_{i\dots j}) P_q(t_{i\dots j} | t), \quad (4)$$

which ensures that when the expected frequencies of phrases, s , are summed (with the length norm) over the whole text:

$$\sum_s \ell(s) f_q(s | T) = \sum_{t \in T} \ell(t) f(t), \quad (5)$$

the underlying mass of words in the text is conserved (see SI-2 for proofs of Eqs. 4 and 5). Said differently, phrase partition frequencies (random or otherwise) conserve word frequencies through the length norm ℓ , and so have a physically meaningful relationship to the words on “the page.”

STATISTICAL MECHANICAL INTERPRETATION

Here, we focus on three natural kinds of partitions: $q=0$: clauses are partitioned only as clauses themselves; $q=\frac{1}{2}$: what we call ‘pure random partitioning’—all partitions of a clause are equally likely; $q=1$: clauses are partitioned into words.

In carrying out pure random partitioning ($q=\frac{1}{2}$), which we will show has the many desirable properties we seek, we are assuming all partitions are equally likely, reminiscent of equipartitioning used in statistical mechanics [19]. Extending the analogy, we can view $q=0$ as a zero temperature limit, and $q=1$ as an infinite temperature one. As an anchor for $f_{\frac{1}{2}}$, we note that words that appear once within a text—hapax legomena—will have $f_q \in \{\frac{1}{4}, \frac{1}{2}, 1\}$ (depending on clause boundaries), on the order of 1 as per standard word partitioning.

EXPERIMENTS AND RESULTS

Before we apply the random partition theory to produce our generalization of word count, f_q , we will first examine the results of applying the random partition process in a ‘one-off’ manner. We process through the clauses of a text once, cutting word-word boundaries (and in a parallel experiment for graphemes, cutting letter-letter boundaries within words) uniformly at random with probability $q = \frac{1}{2}$.

In Fig. 1A, we present an example ‘one-off’ partition of the first few lines of Charles Dickens’ “Tale of Two Cities” We give example partitions at the scales of clauses (red), pure random partition phrases (orange), words (yellow), pure random partition graphemes (green), and letters (blue). In Fig. 1B, we show Zipf distributions for all five partitioning scales. We see that clauses ($q=0$) and pure random partitioning phrases ($q=\frac{1}{2}$) both adhere well to the pure form of $f \propto r^{-\theta}$ where r is rank. For clauses we find $\theta \simeq 0.78$ and for random partitioning, $\theta \simeq 0.98$ (see supplementary material for measurement details and for examples of other works of literature). The quality of scaling degrades as we move down to words and graphemes with the appearance of scaling breaks [21–23]. Scaling vanishes entirely at the level of letters.

Moving beyond a single work, we next summarize findings for a large collection of texts [25] in Fig. 2A, and compare the Zipf exponent θ for words and pure random $q=\frac{1}{2}$ ‘one-off’ partitioning for around 4000 works of literature. We plot the corresponding marginal distributions in Fig. 2B, and see that clearly $\theta \lesssim 1$ for $q=\frac{1}{2}$ phrases, while for words, there is a strong positive skew with the majority of values of $\theta > 1$. These steep scalings for words (and graphemes), $\theta > 1$, are not dynamically accessible for Simon’s model [10].

Leaving aside this non-physicality of Zipf distributions for words and concerns about breaks in scaling, we recall that Simon’s model connects the rate, α , at which new terms are introduced, to θ in a simple way: $1 - \alpha = \theta$ [3]. Given frequency data from a pure Simon model, the word/phrase introduction rate is determined easily to be $\alpha = N/M$, where N is the number of unique words/phrases, and M is the sum total of all word/phrase frequencies. We ask how well works of literature conform to this connection in Fig. 2C, and find that words (green dots) do not demonstrate any semblance of a linear relationship, whereas phrases (blue dots) exhibit a clear, if approximate, linear connection between $1 - \alpha$ and θ .

Despite this linearity, we see that a pure Simon model fails to accurately predict the phrase distribution exponent θ . This is not surprising, as when $\alpha \rightarrow 0$, an immediate adherence to the rich-get-richer mechanism produces a transient behavior in which the first few (largest-count) word varieties exist out of proportion to the eventual scaling. Because a pure Zipf/Simon distribution preserves

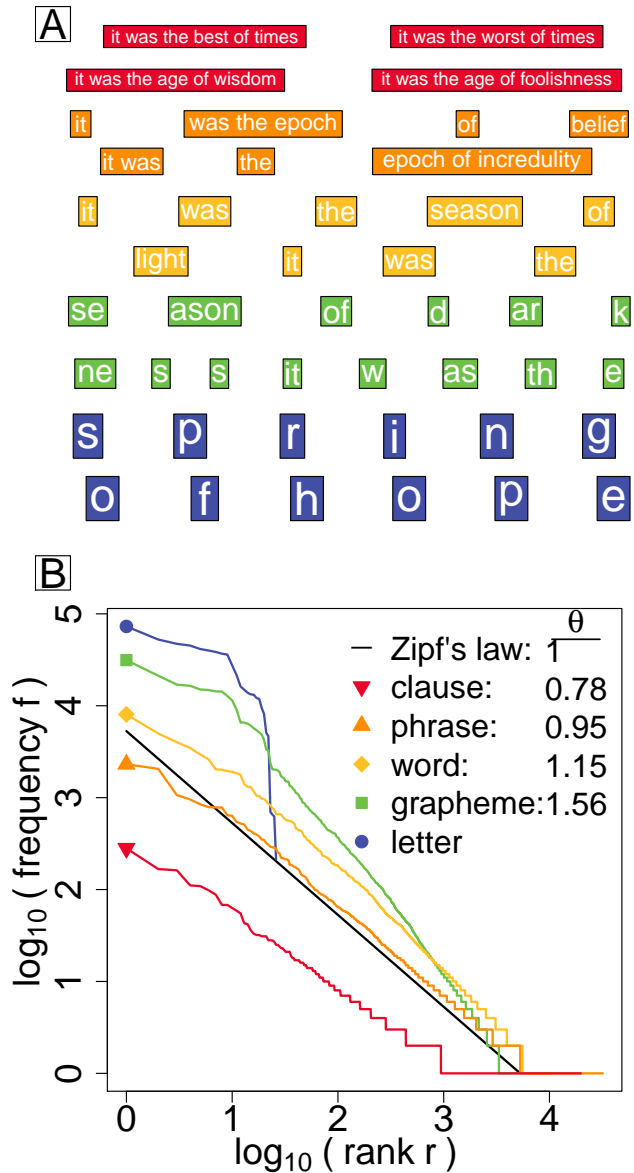


FIG. 1. **A.** Partition examples for the start of Charles Dickens’ “Tale of Two Cities” at five distinct levels: clauses (red), pure random partitioning phrases ($q = \frac{1}{2}$, orange), words (yellow), pure random partitioning graphemes ($q = \frac{1}{2}$, green), and letters (blue). The specific phrases and graphemes shown are for one realization of pure random partitioning. **B.** Zipf distributions for the five kinds of partitions along with estimates of the Zipf exponent θ when scaling is observed. No robust scaling is observed at the letter scale. The colors match those used in panel A, and the symbols at the start of each distribution are intended to strengthen the connection to the legend. See Ref. [20] and supplementary material for measurement details.

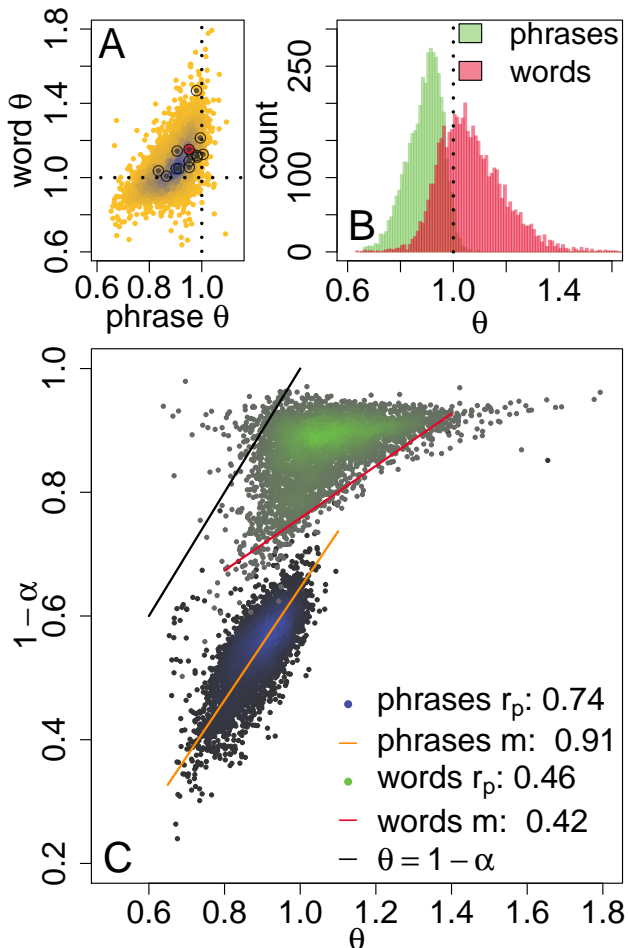


FIG. 2. **A.** Density plot showing the Zipf exponent θ for ‘one-off’ randomly partitioned phrases and word Zipf distributions ($q=1$ and $q=\frac{1}{2}$) for around 4000 works of literature. We indicate “Tale of Two Cities” by the red circle, and with black circles, we represent measurements for 14 other works of literature analyzed further in the supplementary material. **B.** Histograms of the Zipf exponent θ for the same set of books (marginal distributions for **A**). Phrases typically exhibit $\theta \leq 1$ whereas words produce unphysical $\theta > 1$, according to Simon’s model **C.** Test of Simon’s model’s analytical connection $\theta = 1 - \alpha$, where θ is the Zipf exponent and α is the rate at which new terms (e.g., graphemes, words, phrases) are introduced throughout a text. We estimate α as the number of different words normalized by the total word volume. For both words and phrases, we compute linear fits using Reduced Major Axis (RMA) regression [24] to obtain slope m , along with the Pearson correlation coefficient r_p . Words (green) do not exhibit a simple linear relationship whereas phrases do (blue), albeit clearly below the $\alpha = 1 - \theta$ line in black.

$\theta = 1 - \alpha$, we expect that a true, non-transient power-law consistently makes the underestimate $1 - N/M < \theta$.

Inspired by our results for one-off partitions of texts, we now consider ensembles of pure random partitioning for larger texts. In Fig. 3, we show Zipf distributions of

expected partition frequency, f_q , for $q=\frac{1}{2}$ phrases for four large-scale corpora: English Wikipedia, the New York Times (NYT), Twitter, and music lyrics (ML), coloring the main curves according to the length of a phrase for each rank. For comparison, we also include word-level Zipf distributions ($q=1$) for each text in gray, along with the canonical Zipf distribution (exponent $\theta=1$) for reference.

We observe scalings for the expected frequencies of phrases that hover around $\theta = 1$ for over a remarkable 7–9 orders of magnitude. We note that while others have observed similar results by simply combining frequency distributions of N -grams [26], these approaches were unprincipled as they over-counted words. For the randomly partitioned phrase distributions $f_{\frac{1}{2}}$, the scaling ranges we observe persist down to 10^{-2} , beyond the hapax legomena, which occur at frequencies greater than 10^{-1} . Such robust scaling is in stark contrast to the very limited scaling of word frequencies (gray curves). For pure word partitioning, $q=1$, we see two highly-distinct scaling regimes exhibited by each corpus, with shallow upper (Zipf) scalings at best extending over four orders of magnitude, and typically only three. (In a separate work, we investigate this double scaling finding evidence that text-mixing is the cause [23].)

For all four corpora, random partitioning gives rise to a gradual interweaving of different length phrases when moving up through rank r . Single words remain the most frequent (purple), typically beginning to blend with two word phrases (blue) by rank $r = 100$. After the appearance of phrases of length around 10–20, depending on the corpus, we see the phrase rank distributions fall off sharply, due to long clauses that are highly unique in their construction (upper right insets).

In the supplementary material, we provide structured tables of example phrases extracted by pure random partitioning for all four corpora (Tabs. S1–S4), along with complete phrase data sets. As with standard N -grams, the texture of each corpus is quickly revealed by examining phrases of length 3, 4, and 5. For example, the second most common phrases of length 5 for the four corpora are routinized phrases: “the average household size was” (EW), “because of an editing error” (NYT), “i uploaded a youtube video” (TW), and “na na na na na” (ML). By design, random partitioning allows us to quantitatively compare and sort phrases of different lengths. For music lyrics, “la la la la la” has an expected frequency similar to “i don’t know why”, “just want to”, “we’ll have”, and “whatchu” (see Tab. S4), while for the New York Times, “the new york stock exchange” is comparable to “believed to have” (see Tab. S2).

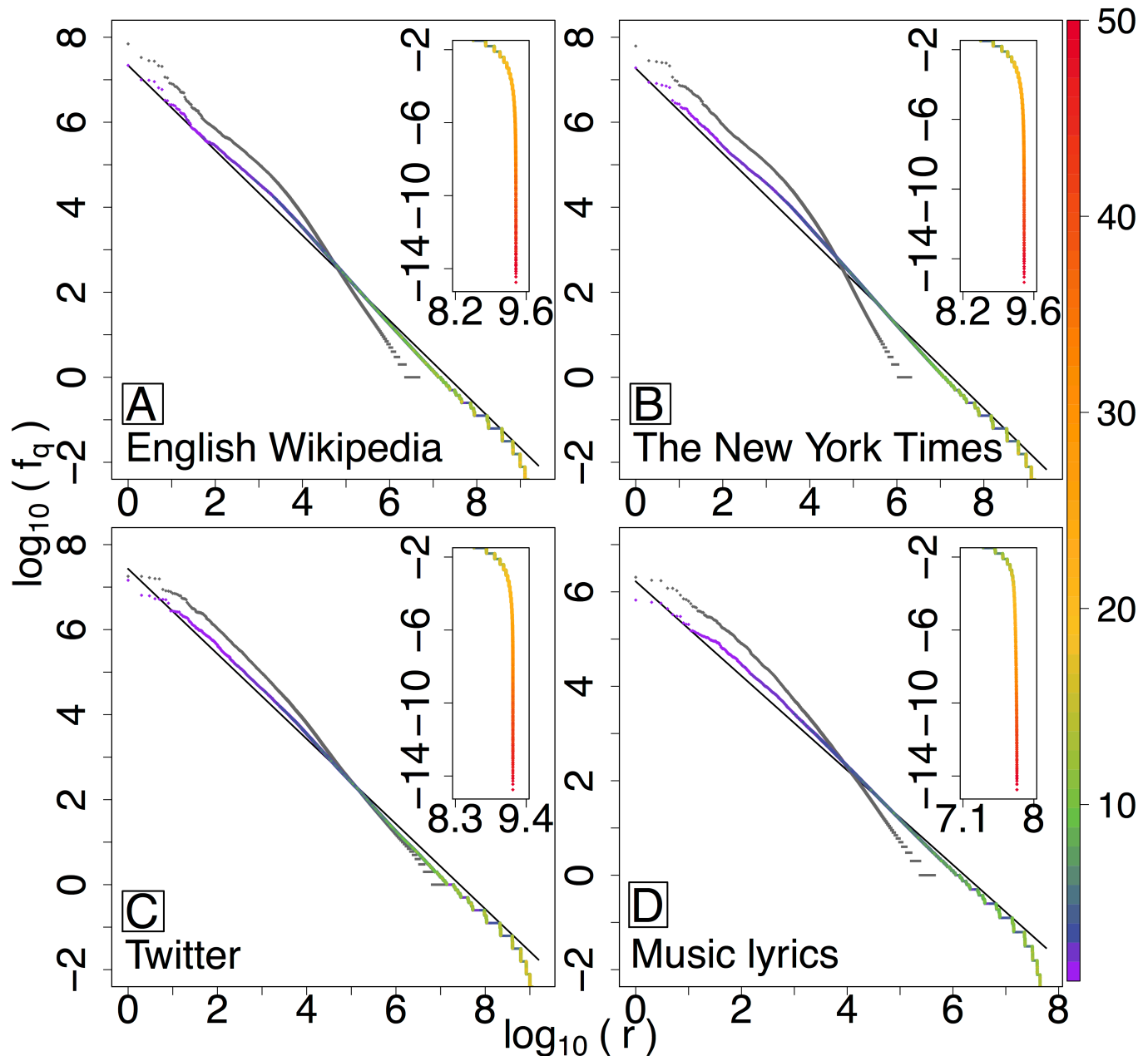


FIG. 3. Random partitioning distributions ($q=\frac{1}{2}$) for the four large corpora: (A) Wikipedia (2010); (B) The New York Times (1987–2007); (C) Twitter (2009); and (D) Music Lyrics (1960–2007). Top right insets show the long tails of random partitioning distributions, and the colors represent phrase length as indicated by the color bar. The gray curves are standard Zipf distributions for words ($q=1$), and exhibit limited scaling and with clear scaling breaks. See main text and Tabs. S1–S4, for example phrases.

DISCUSSION

The phrases and their effective frequencies produced by our pure random partitioning method may serve as input to a range of higher order analyses. For example, information theoretic work may be readily carried out,

context models may be built around phrase adjacency using insertion and deletion, and specific, sentence-level partitions may be realized from probabilistic partitions.

While we expect that other principled, more sophisticated approaches to partitioning texts into rankable mixed phrases should produce Zipf’s law spanning simi-

lar or more orders of magnitude in rank, we believe random partitioning—through its transparency, simplicity, and scalability—will prove to be a powerful method for exploring and understanding large-scale texts.

To conclude, our results reaffirm Zipf’s law for language, uncovering its applicability to a vast lexicon of phrases. Furthermore, we demonstrate that the general semantic units of statistical linguistic analysis can and must be phrases—not words—calling for a reevaluation and reinterpretation of past and present word-based studies in this new light.

The authors are grateful for the computational resources provided by the Vermont Advanced Computing Core which was supported by NASA (NNX 08A096G). CMD was supported by NSF grant DMS-0940271; PSD was supported by NSF CAREER Award #0846668.

* jake.williams@uvm.edu
 † paul.lessard@boulder.edu
 ‡ sdesu@mit.edu
 § eric.clark@uvm.edu
 ¶ james.bagrow@uvm.edu
 ** chris.danforth@uvm.edu
 †† peter.dodds@uvm.edu

- [1] G. K. Zipf, *The Psycho-Biology of Language*, patterns (Houghton-Mifflin, New York, NY, 1935).
- [2] G. K. Zipf, *Human Behaviour and the Principle of Least-Effort*, patterns (Addison-Wesley, Cambridge, MA, 1949).
- [3] H. A. Simon, *Biometrika* **42**, 425 (1955).
- [4] M. Batty, *Science Magazine* **319**, 769 (2008).
- [5] R. Axtell, *Science* **293**, 1818 (2001).
- [6] D. H. Zanette and S. C. Manrubia, *Physica A* **295**, 1 (2001).
- [7] B. B. Mandelbrot, in *Communication Theory*, edited by W. Jackson (Butterworth, Woburn, MA, 1953) pp. 486–502.
- [8] G. A. Miller, *American Journal of Psychology* **70**, 311 (1957).
- [9] R. Ferrer-i Cancho and B. Elvevåg, *PLoS ONE* **5**, e9411 (2010).
- [10] R. M. D’Souza, C. Borgs, J. T. Chayes, N. Berger, and R. D. Kleinberg, *Proc. Natl. Acad. Sci.* **104**, 6112 (2007).
- [11] B. Coromina-Murtra and R. Solé, *Physical Review E* **82**, 011102 (2010).
- [12] S. Bornholdt and H. Ebel, *Phys. Rev. E* **64**, 035104(R) (2001).
- [13] T. Maillart, D. Sornette, S. Spaeth, and G. von Krogh, *Phys. Rev. Lett.* **101**, 218701 (2008).
- [14] D. J. de Solla Price, *J. Amer. Soc. Inform. Sci.* **27**, 292 (1976).
- [15] A.-L. Barabási and R. Albert, *Science* **286**, 509 (1999).
- [16] I. A. Sag, T. Baldwin, F. Bond, A. A. Copestake, and D. Flickinger, in *Proceedings of the Third International Conference on Computational Linguistics and Intelligent Text Processing*, CICLing ’02 (Springer-Verlag, London, UK, 2002) pp. 1–15.
- [17] Google Labs ngram viewer. Available at <http://ngrams.googlelabs.com/>. Accessed May 15, 2014.
- [18] Cougar Town, season 4 episode 4, “I should have known it,” (Broadcast January 29, 2013), <http://www.imdb.com/title/tt2483134/>.
- [19] N. Goldenfeld, *Lectures on Phase Transitions and the Renormalization Group*, general, *Frontiers in Physics*, Vol. 85 (Addison-Wesley, Reading, Massachusetts, 1992).
- [20] A. Clauset, C. R. Shalizi, and M. E. J. Newman, *SIAM Review* **51**, 661 (2009).
- [21] R. Ferrer-i Cancho and R. V. Solé, *Journal of Quantitative Linguistics* **8**, 165 (2001).
- [22] M. Gerlach and E. G. Altmann, *Phys. Rev. X* **3**, 021006 (2013).
- [23] J. R. Williams, J. P. Bagrow, C. M. Danforth, and P. S. Dodds, “Text mixing shapes the anatomy of rank-frequency distributions: A modern zipfian mechanics for natural language,” (2015), <http://arxiv.org/abs/1409.3870>.
- [24] J. M. V. Rayner, *J. Zool. Lond. (A)* **206**, 415 (1985).
- [25] Project Gutenberg, (2010), <http://www.gutenberg.org>.
- [26] L. Q. Ha, E. I. Sicilia-Garcia, J. Ming, and F. J. Smith, in *Proceedings of the 19th International Conference on Computational Linguistics (COLING)* (2002) pp. 315–320.
- [27] E. Sandhaus, “The New York Times Annotated Corpus,” Linguistic Data Consortium, Philadelphia (2008).
- [28] Twitter, (2009), twitter API. <http://dev.twitter.com/>.
- [29] P. S. Dodds and C. M. Danforth, *Journal of Happiness Studies* (2009), doi:10.1007/s10902-009-9150-9.
- [30] Wikipedia, (2010), <http://dumps.wikimedia.org/enwiki/>.

SI-1: MATERIALS AND METHODS

To obtain the results in Fig. 2, we utilize the maximum likelihood estimation (MLE) procedure developed in [20]. In applying this procedure to clause and phrases distributions, several quantities are generally considered:

- $\hat{\theta}$: Zipf exponent estimate.
- r_{\max} : upper cutoff in rank r determined by MLE procedure.
- D : Kolmogorov-Smirnov (KS) statistic.
- p -value determined by the MLE procedure (note that higher is better in that the null hypothesis is more favored).
- $1-\alpha$: Estimate of Zipf exponent θ based on Simon's model [3] where α is the introduction rate of new terms. We estimate α as the number of unique terms (N) divided by the total number of terms (M).

which we report for 14 famous works of literature in SI-3.

In Fig. 2C we measure covariation between regressed values of $\hat{\theta}$ and the Simon model prediction $1-\alpha$. Since both are subject to measurement error ($\hat{\theta}$ is a regressed quantity and α is only coarsely approximated by N/M), we adhere to Reduced Major Axis regression [24], which produces equivalent results upon interchanging x and y variables, and hence guarantees that no information is assumed or lost when we place $\hat{\theta}$ as the x -variable).

To produce the rank-frequency distributions in Fig. 3 and words in tables S1–S4, we apply the random partition process to several large corpora from a wide scope of content. These corpora are: twenty years of New York Times articles (NYT, 1987–2007) [27], approximately 4% of a year's tweets (Twitter, 2009) [28], music lyrics from thousands of songs and authors (Lyrics, 1960–2007) [29], and a collection of complete Wikipedia articles (Wikipedia, 2010) [30]. In Fig. 2 we also use a subset of more than 4,000 books from the Project Gutenberg eBooks collection (eBooks, 2012) [25] of public-domain texts.

SI-2: PROOF OF f_q WORD CONSERVATION

In the body of this document we claim that the random partition frequencies of the phrases within a text T conserve the text's underlying mass of words, M_T . This claim relies on the fact that the partition frequencies of phrase-segments, $t_{i\dots j}$, emerging from a single clause, t , preserve its word mass, $\ell(t)$. We represented this by the summation presented (Eq. 4) in the body of this document, which is equivalent to, $f_q(S | t)E_S[\ell(s) | t]$, i.e.,

the total number of words represented by the frequency of appearance of all phrases generated by the q -partition:

$$\begin{aligned} f(S | t)E_S[\ell(s) | t] &= \sum_{s \in S} \ell(s) f_q(s | t) \\ &= \sum_{s \in S} \sum_{s=t_{i\dots j}} \ell(t_{i\dots j}) P_q(t_{i\dots j} | t) \\ &= \sum_{1 \leq i < j \leq \ell(t)} \ell(t_{i\dots j}) P_q(t_{i\dots j} | t), \end{aligned} \quad (6)$$

which we now denote by $M(S | t)$ for brevity. For convenience, we now let $n = \ell(t)$ denote the clause's length and observe that for each phrase-length $k < n$ there are two single-boundary phrases having partition probability $q(1-q)^{k-1}$, and $n-k-1$ no-boundary phrases having partition probability $q^2(1-q)^{k-1}$. The contribution to the above sum by all k -length phrases is then given by

$$2kq(1-q)^{k-1} + (n-k-1)kq^2(1-q)^{k-1}. \quad (7)$$

Upon noting the frequency of the single phrase (equal to the clause t) whose length is n , $(1-q)^{n-1}$, we consider the sum over all $k \leq n$,

$$\begin{aligned} M(S | t) &= (1-q)^{n-1} \\ &+ [2q + nq^2] \sum_{k=1}^{n-1} k(1-q)^{k-1} \\ &- q^2 \sum_{k=1}^{n-1} k(k+1)(1-q)^{k-1}, \end{aligned} \quad (8)$$

which we will show equals n . We now define the quantity $x = 1-q$ (the probability that a space remains intact), and in these terms find the sum to be:

$$\begin{aligned} M(S | t) &= nx^{n-1} \\ &+ [2(1-x) + n(1-x)^2] \sum_{k=1}^{n-1} kx^{k-1} \\ &- (1-x)^2 \sum_{k=1}^{n-1} k(k+1)x^{k-1}. \end{aligned} \quad (9)$$

This framing through x affords a nice representation in terms of the generating function

$$f(x) = \frac{1-x^{n+1}}{1-x}, \quad (10)$$

which allows us to express the summations through derivatives of $f(x)$:

$$\begin{aligned} \sum_{k=1}^{n-1} kx^{k-1} &= f'(x) - nx^{n-1}, \text{ and} \\ \sum_{k=1}^{n-1} k(k+1)x^{k-1} &= f''(x), \end{aligned} \quad (11)$$

to find

$$\begin{aligned} M(S | t) &= nx^{n-1} \\ &+ [2(1-x) + n(1-x)^2] (f'(x) - nx^{n-1}) \\ &- (1-x)^2 f''(x). \end{aligned} \quad (12)$$

Substitution of the second derivative term

$$f''(x)(1-x) = 2f'(x) - n(n+1)x^{n-1} \quad (13)$$

then produces the reduced form:

$$\begin{aligned} M(S | t) &= n[f'(x)(1-x)^2 \\ &- (nx^{n+1} - (n+1)x^n)], \end{aligned} \quad (14)$$

into which we substitute the first derivative term

$$f'(x)(1-x)^2 = 1 + nx^{n+1} - (n+1)x^n, \quad (15)$$

to render

$$\begin{aligned} M(S | t) &= n[1 + nx^{n+1} - (n+1)x^n \\ &- (nx^{n+1} - (n+1)x^n)] = n, \end{aligned} \quad (16)$$

which proves Eq. 4. Putting this together into a sum over all clauses, we see proof of Eq. 5 naturally follows:

$$\begin{aligned} \sum_{s \in S} \ell(s) f_q(s | T) &= \sum_{t \in T} \sum_{s \in S} \ell(s) f_q(s | t) \\ &= \sum_{t \in T} M(S | t) = \sum_{t \in T} \ell(t). \end{aligned} \quad (17)$$

SI-3: PARAMETERS FOR WELL-KNOWN TEXTS

Below are tables showing fits of Zipf's exponent, $\hat{\theta}$, for 14 famous works of literature, along with details of the maximum likelihood estimation (MLE) procedure in [20]. The quantities used in these table are described in SI-1, Materials and Methods.

A Tale of Two Cities

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.783	3	0.0124	0.961	0.176
phrase	0.951	3	0.00742	0.772	0.603
word	1.15	4	0.0077	0.811	0.925
grapheme	1.56	4	0.0146	0.359	0.986

Moby Dick

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.296	1	0.0192	0	0.154
phrase	0.902	3	0.0132	0.0626	0.576
word	1.05	7	0.00986	0.61	0.912
grapheme	1.42	13	0.0109	0.953	0.986

Great Expectations

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.301	1	0.0199	0	0.186
phrase	0.995	5	0.0164	0.225	0.622
word	1.21	4	0.00943	0.526	0.938
grapheme	1.66	3	0.0147	0.181	0.988

Pride and Prejudice

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	1	3	0.0204	0.911	0.172
phrase	0.983	3	0.0148	0.149	0.617
word	1.11	18	0.0201	0.662	0.947
grapheme	1.43	24	0.0226	0.698	0.989

Adventures of Huckleberry Finn

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.881	4	0.0192	0.977	0.197
phrase	0.98	3	0.0119	0.385	0.625
word	1.47	1	0.0183	0.83	0.94
grapheme	1.66	6	0.0239	0.203	0.987

Alice's Adventures in Wonderland

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.707	2	0.0198	0.711	0.191
phrase	0.906	2	0.0108	0.687	0.555
word	1.14	6	0.0353	0.105	0.899
grapheme	1.19	49	0.0338	0.972	0.975

The Adventures of Tom Sawyer

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.321	1	0.0208	0	0.188
phrase	1.01	6	0.0173	0.826	0.555
word	1.12	3	0.0162	0.108	0.893
grapheme	1.51	4	0.0134	0.683	0.978

The Adventures of Sherlock Holmes

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.308	1	0.0231	0	0.191
phrase	0.952	4	0.0093	0.892	0.586
word	1.09	9	0.0144	0.733	0.921
grapheme	1.44	12	0.0191	0.663	0.983

Sense and Sensibility

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.274	1	0.0176	0	0.142
phrase	0.982	3	0.00945	0.611	0.614
word	1.12	20	0.017	0.907	0.946
grapheme	1.41	28	0.0264	0.584	0.989

Leaves of Grass

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.486	2	0.00768	0.783	0.0717
phrase	0.865	3	0.00971	0.463	0.543
word	1.01	6	0.0095	0.78	0.886
grapheme	1.39	7	0.0131	0.692	0.981

Oliver Twist

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.93	3	0.0152	0.808	0.242
phrase	0.962	3	0.00945	0.439	0.622
word	1.13	8	0.0118	0.695	0.931
grapheme	1.52	7	0.0153	0.521	0.987

Ulysses

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.34	1	0.0192	0	0.193
phrase	0.912	4	0.0062	0.854	0.551
word	1.05	5	0.00773	0.515	0.887
grapheme	1.48	4	0.00874	0.61	0.983

Frankenstein; Or, The Modern Prometheus

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.257	1	0.0121	0	0.0741
phrase	0.834	2	0.0085	0.55	0.532
word	1.04	5	0.0215	0.057	0.906
grapheme	1.31	12	0.019	0.682	0.982

Wuthering Heights

level	$\hat{\theta}$	r_{\max}	D	p -value	$1 - \alpha$
clause	0.927	3	0.0217	0.751	0.178
phrase	0.952	7	0.0104	0.978	0.581
word	1.06	10	0.0163	0.533	0.917
grapheme	1.54	5	0.0165	0.345	0.984

SI-4: PHRASE FREQUENCY TABLES

The following tables contain selected phrases extracted by random partitioning for the four corpora examined in the main text. We provide complete phrase lists in csv format along with other material online at: <http://www.uvm.edu/storylab/share/papers/williams2014a/>.

rank	order=1	order=2	order=3	order=4	order=5
1	i (668838.75)	in the (28174.25)	i love you (2566.75)	la la la la (514.06)	la la la la (184.89)
2	you (600813.50)	and i (25040.88)	i don't know (2094.00)	i don't want to (315.31)	na na na na (93.98)
3	the (576318.50)	i know (17993.00)	i want to (1750.06)	na na na na (281.78)	on and on and on (48.28)
4	and (440698.25)	you know (16977.75)	la la la (1449.50)	in love with you (237.28)	i want you to know (47.70)
5	to (330196.75)	i don't (16237.12)	i want you (1229.00)	i want you to (227.75)	you know what i mean (45.64)
6	me (305085.75)	on the (14977.12)	you and me (1159.00)	i don't know what (201.38)	don't know what to do (45.22)
7	a (301126.50)	if you (13856.62)	i don't want (1105.88)	i don't know why (187.59)	oh oh oh oh (40.80)
8	it (219505.25)	to me (13048.50)	i know you (1086.00)	oh oh oh oh (181.59)	da da da da (40.41)
9	my (205611.00)	to the (12940.75)	i need you (1065.12)	i want to be (172.69)	do do do do (40.02)
10	in (203916.25)	to be (12614.00)	and i know (1051.62)	know what to do (144.06)	one more chance at love (35.66)
11	that (150464.50)	i can (12372.12)	i don't wanna (914.00)	what can i do (141.41)	i don't want to be (35.38)
12	of (149402.75)	and the (11679.88)	i got a (904.25)	yeah yeah yeah (138.19)	in the middle of the (34.66)
13	on (143576.50)	but i (11512.50)	i know that (903.00)	you don't have to (137.38)	i don't give a fuck (33.81)
14	your (135024.00)	of the (11239.88)	you know i (902.69)	i close my eyes (130.31)	yeah yeah yeah yeah (33.05)
15	but (132235.00)	i can't (10372.88)	i can see (872.62)	you want me to (129.19)	i don't know what to (32.39)
16	all (124985.00)	for you (10147.75)	and i don't (844.81)	you make me feel (128.31)	all i want is you (31.78)
17	so (121375.75)	when i (10046.38)	in your eyes (844.06)	i just want to (128.00)	you know i love you (26.88)
18	no (116877.00)	come on (9924.25)	i don't care (832.06)	da da da da (123.78)	the middle of the night (26.73)
19	we (113865.25)	you can (9686.00)	and if you (825.94)	if you want to (123.06)	the rest of my life (26.34)
20	is (11375.25)	i got (9577.88)	the way you (824.94)	come back to me (121.56)	no no no no (26.11)
21	for (108828.50)	in my (9473.12)	all the time (817.62)	in the middle of (119.16)	at the end of the (25.30)
22	oh (107477.25)	all the (9467.25)	na na na (790.38)	and i don't know (118.72)	i wanna be with you (22.77)
23	be (107432.75)	i want (9396.50)	don't you know (766.62)	let me tell you (117.66)	all i wanna do is (22.44)
24	love (104438.50)	that i (9190.88)	this is the (766.25)	give it to me (111.97)	no matter what i do (22.41)
25	it's (99026.75)	i am (9141.88)	can't you see (761.19)	you are the one (111.94)	the way you love me (21.42)
26	now (95016.75)	and you (9048.75)	you love me (753.44)	do do do do (111.28)	no matter what you do (21.36)
27	don't (94956.00)	i was (9028.12)	oh oh oh (749.56)	i love you so (111.16)	what you do to me (20.83)
28	yeah (92807.00)	tell me (8783.50)	i wanna be (744.50)	all i want is (109.81)	when i close my eyes (20.31)
29	when (91600.75)	like a (8614.12)	you know that (714.38)	how does it feel (109.69)	and i don't know why (20.09)
30	with (90323.75)	the way (8512.38)	you want to (709.62)	know what i mean (109.12)	let me be the one (19.86)
31	what (90190.50)	to you (8289.50)	you don't know (707.62)	no no no no (104.03)	the end of the day (18.64)
32	this (90120.00)	when you (8157.62)	in my heart (693.69)	to be with you (100.81)	in the name of love (18.50)
33	know (89600.00)	if i (7941.50)	you and i (691.50)	i don't wanna be (97.50)	lemme see you drip sweat (18.00)
34	like (84259.00)	in a (7893.38)	you make me (675.19)	and on and on (96.47)	i like the way you (17.91)
35	just (83346.75)	my heart (7882.88)	if you want (663.81)	the end of the (94.66)	it's been a long time (17.89)
36	baby (83182.75)	for me (7880.50)	yeah yeah yeah (662.38)	i wish i could (93.09)	till the end of time (17.67)
37	do (81926.00)	this is (7754.62)	don't want to (654.62)	don't give a fuck (92.94)	i wish that i could (17.61)
38	up (81529.00)	for the (7570.88)	want to be (624.56)	can you feel it (91.88)	if you want me to (17.47)
39	if (74941.25)	let me (7539.25)	in my life (622.44)	the way i feel (91.00)	see it in your eyes (17.20)
40	chorus (72833.00)	with you (7482.62)	if i could (619.25)	i don't know how (90.47)	no matter what they say (16.78)
41	can (67057.50)	i need (7424.62)	you know what (615.06)	gon play with it (90.00)	and i don't know what (16.73)
42	down (66636.75)	with me (7386.00)	what you want (605.19)	you know that i (89.84)	let me hear you say (16.70)
43	get (63408.50)	you are (7208.25)	i used to (604.88)	at the end of (89.38)	i look into your eyes (16.70)
44	time (62579.50)	i wanna (7083.00)	on and on (595.94)	can you hear me (89.06)	i love the way you (16.64)
45	out (62562.50)	what you (6949.00)	i see you (592.88)	want you to know (88.38)	and i don't want to (16.45)
46	go (62101.75)	love you (6900.38)	in the sky (587.75)	out of my mind (86.62)	when i think of you (16.38)
47	quot (61793.50)	the world (6774.62)	in the air (584.06)	i need to know (86.56)	i look in your eyes (16.31)
48	got (60347.00)	do you (6733.50)	what to do (577.12)	all i wanna do (84.03)	the end of the world (16.16)
49	one (59306.50)	from the (6679.88)	all night long (558.19)	on the other side (83.88)	when the sun goes down (16.11)
50	see (58662.50)	want to (6649.88)	i know i (557.00)	do you love me (83.72)	still in love with you (16.02)
100	that's (28709.75)	in love (4324.38)	i just want (441.88)	that you love me (61.00)	you want me to do (11.83)
150	always (17981.50)	i won't (3225.00)	make me feel (369.31)	take a look at (51.09)	the end of the line (9.78)
200	en (13668.25)	without you (2692.50)	for you to (308.31)	you make me wanna (43.81)	that's the way it goes (8.77)
250	side (10606.00)	when i'm (2280.75)	who i am (277.81)	the rest of my (39.50)	the way that you do (7.88)
300	words (8896.00)	so long (2050.12)	on the wall (254.81)	open up your eyes (36.66)	i want to see you (7.23)
350	coming (7424.50)	have a (1815.25)	no one else (236.19)	get out of my (34.09)	makes the world go round (6.59)
400	ground (6669.25)	that's the (1645.12)	that's what i (218.94)	i don't want no (31.16)	tell me what you need (6.22)
450	death (5688.75)	then you (1506.12)	come back to (206.38)	don't mean a thing (29.25)	hey hey hey hey hey (5.81)
500	slow (5006.25)	i try (1382.25)	just want to (194.12)	goes on and on (27.84)	my my my my my (5.44)
600	cut (3808.00)	here i (1196.62)	i see your (172.44)	me like you do (25.44)	hey ladies drop it down (5.00)
700	grew (3091.25)	love with (1066.25)	in the game (158.62)	in front of you (23.47)	don't know if i can (4.61)
800	shut (2569.75)	my hands (969.25)	not the same (145.62)	you broke my heart (21.91)	it's been so long since (4.30)
900	doo (2167.75)	i tell (879.75)	yes i am (134.25)	me what you want (20.78)	you were the only one (4.06)
1000	seven (1898.75)	s a (802.88)	it was the (126.00)	all that i want (19.69)	just the way it is (3.88)
1500	food (1140.25)	am the (562.75)	a whole lot (95.38)	i wanna thank you (15.59)	mean a thing to me (3.20)
2000	fields (776.75)	caught up (434.12)	give me love (79.25)	got nothing to say (13.09)	a shoulder to cry on (2.78)
2500	vie (575.50)	saturday night (352.00)	yes you are (68.12)	know that you can (11.62)	was it good for you (2.48)
3000	compromise (451.00)	of things (295.38)	all about the (60.12)	is how we do (10.34)	right round like a record (2.27)
3500	couch (363.00)	the white (254.50)	think you can (53.75)	joy to the world (9.38)	your love would be untrue (2.08)
4000	pu (301.25)	they see (223.75)	i can fly (49.00)	if i don't get (8.69)	he was the only one (1.94)
4500	collect (254.75)	we'll have (197.62)	you said you'd (44.81)	give it all to (8.09)	you that we won't stop (1.81)
5000	product (219.25)	you drive (179.12)	want to hold (41.25)	wanna get with you (7.59)	cut me down to size (1.72)
6000	whatchu (169.50)	where you're (149.50)	take a breath (35.94)	your eyes on me (6.78)	round the ole oak tree (1.56)
7000	battered (135.25)	a plane (128.62)	right here in (32.00)	i wish i may (6.19)	move on down the line (1.44)
8000	verloren (111.25)	step out (111.88)	of all that (29.00)	we can make love (5.69)	bow wow wow yippie yo (1.33)
9000	nt (93.25)	fuck what (99.12)	be waiting for (26.44)	who the fuck are (5.25)	to warm a lonely night (1.25)
10000	honda (79.75)	you should've (88.38)	that what you (24.38)	like a loaded gun (4.88)	ain't that what you said (1.19)
15000	fuma (43.75)	little angel (57.50)	i wouldn't mind (17.44)	it's better this way (3.75)	on christmas day in the (0.94)
20000	cooper (28.50)	the undertow (42.00)	the wrong place (13.69)	since she left me (3.09)	and let it all go (0.80)
25000	fishy (20.25)	a major (32.88)	for one last (11.38)	and maybe you can (2.66)	no matter how far away (0.70)
30000	illtown (15.25)	alright baby (27.00)	you should try (9.75)	it take to make (2.34)	t want french fried potatoes (0.62)
35000	ndelo (12.00)	loud enough (22.62)	never give it (8.56)	i came to bring (2.12)	what you gave to me (0.58)
40000	rees (9.75)	view mirror (19.50)	to me a (7.62)	things i'm gonna do (1.94)	love is out the door (0.53)
45000	metaphoric (8.00)	your concern (17.12)	roll roll roll (6.88)	gotta say too much (1.75)	non ci sono solo io (0.50)
50000	memorizing (6.75)	the cancer (15.25)	on the eyes (6.25)	lay on the floor (1.62)	set the floor on fire (0.48)
60000	ajai (5.00)	an' then (12.38)	keep my eye (5.31)	give up on yourself (1.44)	right here next to you (0.44)
70000	aleiki (4.00)	cats be (10.38)	no more rummin' (4.62)	there's only one god (1.31)	gates of the seven seals (0.38)
80000	saatanan (3.25)	blif ik (8.88)	we'll show them (4.12)	skies from now on (1.19)	we don't even have to (0.38)
90000	sauber (2.75)	yo tell (7.75)	time you say (3.69)	it comes to that (1.09)	ooh when you walk by (0.34)
100000	mosques (2.25)	believe anymore (6.88)	seemed so right (3.38)	but if i leave (1.00)	van de kille stemmen die (0.31)

TABLE S4. Example phrases for Music Lyrics extracted by random partitioning.