

Characterizing Idioms: Conventionality and Contingency

Michaela Socolof^{1,2}, Jackie Chi Kit Cheung^{1,2,3}, Michael Wagner¹, Timothy J. O’Donnell^{1,2,3}

McGill University¹, Quebec AI Institute, Mila², Canada CIFAR AI Chair³

michaela.socolof@mail.mcgill.ca, chael@mcgill.ca,

jcheung@cs.mcgill.ca, timothy.odonnell@mcgill.ca

Abstract

Idioms are unlike most phrases in two important ways. First, words in an idiom have non-canonical meanings. Second, the non-canonical meanings of words in an idiom are contingent on the presence of other words in the idiom. Linguistic theories differ on whether these properties depend on one another, as well as whether special theoretical machinery is needed to accommodate idioms. We define two measures that correspond to the properties above, and we implement them using BERT (Devlin et al., 2019) and XLNet (Yang et al., 2019). We show that English idioms fall at the expected intersection of the two dimensions, but that the dimensions themselves are not correlated. Our results suggest that special machinery to handle idioms may not be warranted.

1 Introduction

Idioms—expressions like *rock the boat*—bring together two phenomena which are of fundamental interest in understanding language. First, they exemplify *non-conventional word meaning* (Weinreich, 1969; Nunberg et al., 1994). The words *rock* and *boat* in this idiom seem to carry particular meanings—something like *destabilize* and *situation*, respectively—which are different from the conventional meanings of these words in other contexts. Second, unlike other kinds of non-conventional word use such as novel metaphor, there is a contingency relationship between words in an idiom (Wood, 1986; Pulman, 1993). It is the specific combination of the words *rock* and *boat* that has come to carry the idiomatic meaning. *Shake the canoe* does not have the same accepted meaning.

In the literature, most discussions of idioms make use of prototypical examples such as *rock the boat*. This obscures an important fact: There is no generally agreed-upon definition of *idiom*;

phrase types such as light verb constructions (e.g., *take a walk*) and semantically transparent collocations (e.g., *now or never*) are sometimes included in the class (e.g., Palmer, 1981) and sometimes not (e.g., Cowie, 1981). This lack of homogeneity among idiomatic phrases has been recognized as a challenge in the domain of NLP, with Sag et al. (2002) suggesting that a variety of techniques are needed to deal with different kinds of multi-word expressions. What does seem clear is that prototypical cases of idiomatic phrases tend to have higher levels of both non-conventional meaning and contingency between words.

This combination of non-conventionality and contingency has led to a number of theories that treat idioms as exceptions to the mechanisms that build phrases compositionally. These theories posit special machinery for handling idioms (e.g., Weinreich, 1969; Bobrow and Bell, 1973; Swinney and Cutler, 1979). An early but representative example of this position is Weinreich (1969), who posits the addition of two structures to linguistic theory: (1) an *idiom list*, where each entry contains a string of morphemes, its associated syntactic structure, and its sense description, and (2) an *idiom comparison rule*, which matches strings against the idiom list. Such theories must of course provide principles for addressing the difficult problem of distinguishing idioms from other instances of non-conventionality or contingency.

We propose an alternative approach, which views idioms not as exceptional, but merely the result of the interaction of two independently motivated cognitive mechanisms. The first allows words to be interpreted in non-canonical ways depending on context. The second allows for the storage and reuse of linguistic structures—not just words, but larger phrases as well (e.g., Di Sciullo and Williams, 1987; Jackendoff, 2002; O’Donnell, 2015). There is disagreement in the literature about the relationship between these two proper-

ties; some theories of representation predict that the only elements that get stored are those with non-canonical meanings (e.g., Bloomfield, 1933; Pinker and Prince, 1988), whereas others predict that storage can happen no matter what (e.g., O’Donnell, 2015; Tremblay and Baayen, 2010). We predict that, consistent with the latter set of theories, neither mechanism should depend on the other.

This paper presents evidence that prototypical idioms occupy a particular region of the space of these two mechanisms, but are not otherwise exceptional. We define two measures, *conventionality*—meant to measure the degree to which words are interpreted in a canonical way, and *contingency*—a statistical association measure meant to capture the degree to which the presence of one word form depends on the presence of another. Our implementations make use of the pre-trained language models BERT (Devlin et al., 2019) and XLNet (Yang et al., 2019). We construct a novel corpus of English phrases typically called idioms, and show that these phrases fall at the intersection of low conventionality and high contingency, but that the two measures are not correlated and there are no clear discontinuities that separate idioms from other types of phrases.

Our experiments also reveal hitherto unnoticed asymmetries in the behavior of head and non-head words of idioms. In idioms, the dependent word (e.g., *boat* in *rock the boat*) shows greater deviation from its conventional meaning than the head.

2 Conventinality and contingency

In this section we describe the motivation behind our two measures and lay out our predictions about their interaction.

Our first measure, *conventionality*, captures the extent to which subparts of a phrase contribute their normal meaning to the phrase. Most of language is highly conventional; we can combine a relatively small set of units in novel ways, precisely because we can trust that those units will have similar meanings across contexts. At the same time, the linguistic system allows structures like metaphors and idioms, which use words in non-conventional ways. Our conventionality measure is intended to distinguish phrases based on how conventional the meanings of their words are.

Our second measure, *contingency*, captures how unexpectedly often a group of words occurs to-

gether in a phrase and, thus, measures the degree to which there is a statistical contingency—the presence of one or more words strongly signals the likely presence of the others. This notion of contingency has also been argued to be a critical piece of evidence used by language learners in deciding which linguistic structures to store (e.g., Hay, 2003; O’Donnell, 2015).

To aid in visualizing the space of phrase types we expect to find in language, we place our two dimensions on the axes of a 2x2 matrix, where each cell contains phrases that are either high or low on the conventionality scale, and high or low on the contingency scale. The matrix is given in Figure 1, with the types of phrases we expect in each cell.

	Low conv.	High conv.
High cont.	Idioms (e.g., <i>raise hell</i>)	Common collocations (e.g., <i>in and out</i>)
Low cont.	Novel metaphors	Regular language use (e.g., <i>eat peas</i>)

Figure 1: Matrix of phrase types, organized by whether they have high/low conventionality and high/low contingency

We expect our measures to place idioms primarily in the top left corner of the space. At the same time, we predict a lack of correlation between the measures and a lack of major discontinuities in the space. We take these predictions to be consistent with theories that factorize the problem into two mechanisms (captured by our dimensions of conventionality and contingency). We contend that this factorization provides a natural way of characterizing not just idioms, but also collocations and novel metaphors, alongside regular language use.

3 Methods

In this section, we describe the creation of our corpus of idioms and define measures of conventionality and contingency. Given that definitions of idioms differ in which phrases in our dataset count as idioms (some would include semantically transparent collocations, others would not), we do not want to commit to any particular definition a priori, while still acknowledging that people share somewhat weak but broad intuitions about idiomaticity. As we discuss below, our idiom dataset consists of phrases that have at some point been called idioms in the linguistics literature.

3.1 Dataset

We built a corpus of sentences containing idioms and non-idioms, all gathered from the British National Corpus (BNC; [Burnard, 2000](#)), which is a 100 million word collection of written and spoken English from the late twentieth century. The corpus we construct is made up of sentences containing *target phrases* and *matched phrases*, which we detail below.

The target phrases in our corpus consist of 207 English phrasal expressions, some of which are prototypical idioms (e.g., *rock the boat*) and some of which are boundary cases that are sometimes considered idioms, such as collocations (e.g., *bits and pieces*). These expressions are divided into four categories based on their syntax: verb object (VO), adjective noun (AN), noun noun (NN), and binomial (B) expressions. Binomial expressions are fixed pairs of words joined by *and* or *or* (e.g., *wear and tear*). The phrases were selected from lists of idioms published in linguistics papers ([Riehemann, 2001](#); [Morgan and Levy, 2016](#); [Stone, 2016](#); [Bruening et al., 2018](#); [Bruening, 2019](#); [Titone et al., 2019](#)). We added the lists to our dataset one-by-one until we had at least 30 phrases of each syntactic type. We chose these four types in advance to investigate a variety of syntactic types to prevent our results from being too heavily skewed by any potential syntactic confounds in particular constructions. The full list of target phrases is given in Appendix A. The numerical distribution of phrases is given in Table 1.

Phrase type	Number of phrases	Example
VO	31	<i>jump the gun</i>
NN	36	<i>word salad</i>
AN	33	<i>red tape</i>
B	58	<i>fast and loose</i>

Table 1: Types, counts, and examples of target phrases in our idiom corpus, with head words bolded

The BNC was constituency parsed using the Stanford Parser ([Manning et al., 2014](#)), then Tregex ([Levy and Andrew, 2006](#)) expressions were used to find instances of each target phrase.

Matched, non-idiomatic sentences were also extracted in order to allow for direct comparison of conventionality scores for the same word in idiomatic and non-idiomatic contexts. To obtain these matches, we used Tregex to find sentences that included a phrase with the same syntactic

structure as the target phrase. Each target phrase was used to obtain two sets of matched phrases: one set where the head word remained constant and one where the non-head word remained constant.¹ For example, to get head word matches of the adjective noun combination *sour grapes*, we found sentences where the lemma *grape* was modified with an adjective other than *sour*. Below is an example of a sentence found by this method:

*Not a **special grape** for winemaking, nor a hidden architectural treasure, but hot steam gushing out of the earth.*

The number of instances of the matched phrases ranged from 29 (the number of verb object phrases with the object *logs* and a verb other than *saw*) to the tens of thousands (e.g., for verb object phrases beginning with *have*), with the majority falling in the range of a few hundred to a few thousand. Issues of sparsity were more pronounced among the target phrases, which ranged from one instance (*word salad*) to 2287 (*up and down*). Because of this sparsity, some of the analyses described below focus on a subset of the phrases.

The syntactic consistency between the target and matched phrases is an important feature of our corpus, as it allows us to compare conventional-ity across semantic contexts while controlling for syntactic structure.

3.2 Conventionality measure

Our measure of conventionality is built on the idea that a word being used in a conventional way should have similar or related meanings across contexts, whereas a non-conventional word meaning can be idiosyncratic to particular contexts. In the case of idioms, we expect that the difference between a word’s meaning in an idiom and the word’s conventional meaning should be large. On the other hand, there should be little difference between the word’s meaning in a non-idiom and the word’s conventional meaning.

Our measure makes use of the language model BERT ([Devlin et al., 2019](#)) to obtain contextualized embeddings for the words in our dataset. BERT was trained on a corpus of English text, both nonfiction and fiction, with the objectives of masked language modeling and next sentence pre-

¹To obtain matched phrases, we follow work such as [Gazdar \(1981\)](#), [Rothstein \(1991\)](#), and [Kayne \(1994\)](#) in treating the first element in a binomial as the head. We discuss this further in Section 6.

diction. For each of our phrases, we compute the conventionality measure separately for the head and non-head words. For each case (head and non-head), we first take the average embedding for the word across sentences *not containing* the phrase. That is, for *rock* in *rock the boat*, we get the embeddings for the word *rock* in sentences where it does not occur with the direct object *boat*. Let O be a set of instances w_1, w_2, \dots, w_n of a particular word used in contexts *other than* the context of the target phrase. Each instance has an embedding $u_{w_1}, u_{w_2}, \dots, u_{w_n}$. The average embedding for the word among these sentences is:

$$\mu_O = \frac{1}{n} \sum_{i=1}^n u_{w_i} \quad (1)$$

We take this quantity to be a proxy for the prototypical, or conventional, meaning of the word. The conventionality score is the negative of the average distance between μ_O and the embeddings for uses of the word across instances of the phrase in question. We compute this as follows:

$$\text{conv}(\text{phrase}) = -\frac{1}{m} \sum_{i=1}^m \left\| \frac{T_i - \mu_O}{\sigma_O} \right\|_2 \quad (2)$$

where T is the embedding corresponding to a particular use of the word in the target phrase, and σ_O is the component-wise standard deviation of the set of embeddings u_{w_i} , and m is the number of sentences in which the target phrase is used.

3.3 Contingency measure

Our second measure, which we have termed *contingency*, refers to whether a particular set of words appears within the same phrase at an unexpectedly high rate. The measure is based on the notion of pointwise mutual information (PMI), which is a measure of the strength of association between two events. We use a generalization of PMI that extends it to sets of more than two events, allowing us to capture the association between phrases that contain more than two words.

The specific generalization of PMI that we use has at various times been called total correlation (Watanabe, 1960), multi-information (Studený and Vejnarová, 1998), and specific correlation (Van de Cruys, 2011).

$$\text{cont}(x_1, x_2, \dots, x_n) = \log \frac{p(x_1, x_2, \dots, x_n)}{\prod_{i=1}^n p(x_i)} \quad (3)$$

For the case of three variables, we get:

$$\text{cont}(x, y, z) = \log \frac{p(x, y, z)}{p(x)p(y)p(z)} \quad (4)$$

To estimate the contingency of a phrase, we use word probabilities given by XLNet (Yang et al., 2019), an auto-regressive language model that gives estimates for the conditional probabilities of words given their context. Like BERT, XLNet was trained on a mix of fiction and nonfiction data. To estimate the joint probability of the words in *rock the boat* in some particular context (the numerator of the expression above), we use XLNet to obtain the product of the conditional probabilities in the chain rule decomposition of the joint. We get the relevant marginal probabilities by using attention masks over particular words, as shown below, where c refers to the context—that is, the rest of the words in the sentence containing *rock the boat*.

$$\begin{aligned} \Pr(\text{boat} \mid \text{rock the}, c) &= \dots \text{rock the } \mathbf{boat} \dots \\ \Pr(\text{the} \mid \text{rock}, c) &= \dots \text{rock } \mathbf{the} \text{ } [_] \dots \\ \Pr(\text{rock} \mid c) &= \dots \mathbf{rock} \text{ } [_] \text{ } [_] \dots \end{aligned}$$

The denominator is the product of the probabilities of each individual word in the phrase, with both of the other words masked out:

$$\begin{aligned} \Pr(\text{boat} \mid c) &= \dots [_] \text{ } [_] \mathbf{boat} \dots \\ \Pr(\text{the} \mid c) &= \dots [_] \mathbf{the} \text{ } [_] \dots \\ \Pr(\text{rock} \mid c) &= \dots \mathbf{rock} \text{ } [_] \text{ } [_] \dots \end{aligned}$$

The conditional probabilities were computed right to left, and included the sentence to the left and the sentence to the right of the target sentence for context. Note that in order to have an interpretable chain rule decomposition for each sequence, we calculate the XLNet-based generalized PMI for the entire string bounded by the two words of the idiom—this means, for example, that the phrase *rock the fragile boat* will return the PMI score for the entire phrase, adjective included.

4 Validation of conventionality measure

Our conventionality measure provides an indirect way of looking at how canonical a word’s meaning is in context. In order to validate that the measure corresponds to an intuitive notion of unusual word meaning, we carried out an online experiment to see whether human judgments of conventionality

correlated with our automatically-computed conventionality scores. The experimental design and results are described below. (Note that our contingency measure directly computes the statistical quantity we want, so validation is not necessary.)

4.1 Human rating experiment

The experiment asked participants to rate the literalness of a word or phrase in context.² We used twenty-two verb object target phrases and their corresponding matched phrases.³ For each target phrase (e.g., *rock the boat*), there were ten items, each of which consisted of the target phrase used in the context of a (different) sentence. Each sentence was presented with the preceding sentence and the following sentence as context, which is the same amount of context that the automatic measure was given. In each item, a word or phrase was highlighted, and the participant was asked to rate the literalness of the highlighted element. We obtained judgments of the literalness of the head word, non-head word, and entire phrase for ten different sentences containing each target phrase.

We also obtained literalness judgments of the head word and entire phrase for phrases matched on the head of the idiom (e.g., verb object phrases with *rock* as the verb and a noun other than *boat* as the object). Similarly, we obtained literalness judgments of the non-head word and the entire phrase for phrases matched on the non-head word of the idiom (e.g., verb object phrases with *boat* as the object and a verb other than *rock*). Participants were asked to rate literalness on a scale from 1 ('Not literal at all') to 6 ('Completely literal'). We chose to use an even number of points on the scale to discourage participants from imposing a three-way partition into 'low', 'neutral', and 'high'. Items were presented using a Latin square design. The experiment was run online using the Prosodylab Experimenter (Wagner, 2021), a JavaScript tool building on jsPsych (De Leeuw, 2015).

Participants were adult native English speakers

²Participants were recruited on Amazon Mechanical Turk and compensated at a rate of \$15/hour. The study was carried out with REB approval.

³We excluded one target phrase from the analyses (*spill the beans*) based on examination of the BERT-based conventionality scores. The verb *spill* used in *spill the beans* scored anomalously high on conventionality; investigation of the target and matched sentences revealed that roughly half of the matched sentences included a different idiom: *spill X's guts*. We checked the rest of our dataset and did not find other instances of this confound.

who gave written informed consent to participate. The experiment took about 10 minutes to complete. The data were recorded using anonymized participant codes, and none of the results included any identifying information. There were 150 participants total. The data from 10 of those participants were excluded due to failure to follow the instructions (assessed with catch trials).

4.2 Results

To explore whether our conventionality measure correlates with human judgments of literalness, we compare the scores to the results from the rating experiment. Ratings were between 1 and 6, with 6 being the highest level of conventionality.

We predicted that the literalness ratings should increase as conventionality scores increased. To assess whether our prediction was borne out, a linear mixed model was fit using the lmerTest (Kuznetsova et al., 2017) package in R (Team, 2017), with conventionality score and highlighted word (head versus non-head) and their interaction as predictors, plus random effects of participant and item.⁴ All random effects were maximal up to convergence. Results are shown in Table 2 in Appendix B. The results confirm our prediction that words that receive higher conventionality scores are rated as highly literal by humans ($\hat{\beta} = 0.185$, $SE(\hat{\beta}) = 0.050$, $p < 0.001$; see Row 2 of Table 2 in Appendix B).

We carried out a nested model comparison to see whether including the BERT conventionality score as a predictor significantly improved the model, and we found that it did. A likelihood ratio test with the above model and one without the BERT conventionality score as a predictor yielded a higher log likelihood for the full model ($\chi^2 = 80.043$, $p < 0.001$).

5 Analyses

In this section we present analyses of our two measures individually, showing that they capture the properties they were intended to capture. We then investigate the interaction between the measures. Section 5.3 evaluates our central predictions.

We predict that the target phrases will score lower on conventionality than the matched phrases, since we expect these phrases to contain words with (often highly) unconventional meanings. We further predict that the target phrases will

⁴Rating~Conv*Head+(1|Item)+(1+Conv||Partp)

have higher contingency scores than the matched phrases, due to all of the target phrases being expressions that are frequently reused. Putting the two measures together, we expect idioms to fall at the intersection of low conventionality and high contingency, but not to show major discontinuities that qualitatively distinguish them from phrases that fall at other areas of intersection.

5.1 Analysis 1: conventionality measure

We find that the target phrases have lower average conventionality scores than the matched phrases, with a difference of -1.654, with $t(145) = -5.829$ and $p < 0.001$. This is consistent with idioms having unconventional word meanings.

5.2 Analysis 2: contingency measure

We find that, averaged across contexts, the target phrases had higher contingency scores, with a difference in value of 2.25 bits, with $t(159) = 8.807$ and $p < 0.001$.

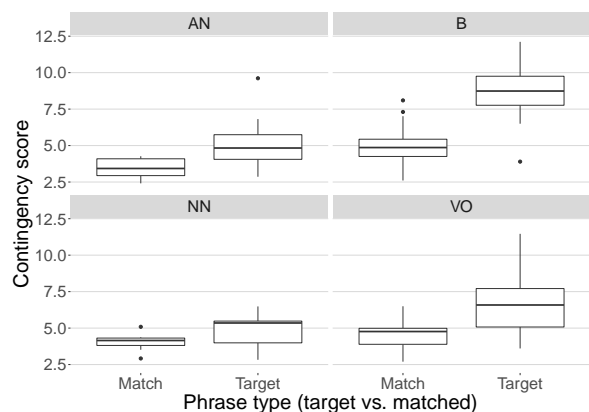


Figure 2: Contingency of target and matched phrases, for phrases with at least 30 instances

Figure 2 shows boxplots of the average contingency score for each phrase type. Since many of the target phrases only occurred in a handful of sentences, we have excluded phrases for which the target or matched sets contain fewer than 30 sentences.⁵ For the most part, there were fewer sentences containing the target phrase than there were sentences containing only the head or only the non-head word in the relevant structural position. This likely explains the greater variance

⁵This threshold was chosen to strike a balance between having enough instances contributing to the average score for each datapoint, and having a large enough sample of phrases. We considered thresholds at every multiple of 10 until we reached one that left at least 100 datapoints remaining.

among the target phrases—the averages are based on fewer data points.

For all syntactic structures, the median contingency score was higher for target phrases than matched phrases. The greatest differences were observed for verb object and binomial phrases.

We fit another mixed effects model to test whether target idioms have higher contingency scores than matched phrases across syntactic classes (AN, B, NN, VO). The model predicts the contingencies for each instance of a phrase used in context, with the target-matched contrast and syntactic class as fixed effects, and random effects for the target-matched pairs.⁶ We find that target phrases have significantly higher contingency scores than matched phrases (see Row 2 of Table 3 of Appendix B).

5.3 Analysis 3: interaction and correlation of measures

Here we show that idioms fall in the expected area of our two-dimensional space, with no evidence of correlation between the measures. Our results provide evidence against the notion of a special mechanism for idioms, whereby conventionality and contingency are expected to covary.

Recall the 2x2 matrix of contingency versus conventionality (Figure 1), where idioms were expected to be in the top left quadrant. Figure 3 shows our results. Since the conventionality scores were for individual words, we averaged the scores of the head word and the primary non-head word (i.e., the verb and the object for verb object phrases, the adjective and the noun for adjective noun phrases, the two nouns in noun noun phrases, and the two words of the same category in binomial phrases). The plot shows the average values of the target and matched phrases.

As discussed above, the target phrases came from lists of idioms in the literature, and thus include a mix of canonical idioms and (seemingly) compositional collocations. We predicted that the target phrases would be distributed between the top two quadrants, with obvious idioms on the top left and collocations on the top right. As a sample, our results placed the following phrases in the top left quadrant: *clear the air*, *bread and butter*, *nuts and bolts*, *red tape*, and *cut corners*. For each of these phrases, the idiomatic meaning cannot be derived by straightforwardly composing the mean-

⁶ $\text{Cont} \sim \text{Target} * \text{Class} + (1 + \text{Target} | \text{Idiom})$

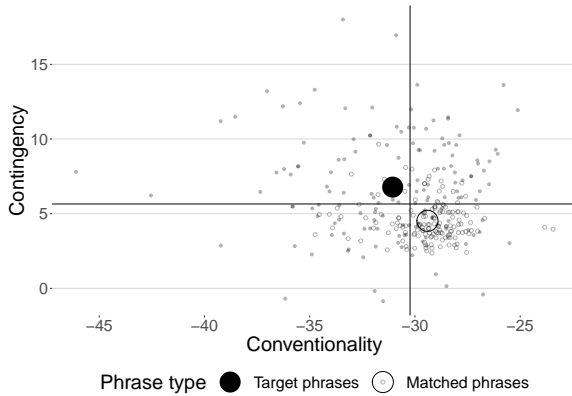


Figure 3: Contingency versus conventionality values of target and matched phrases. Large circles are average values of all target (black) and all matched (white) phrases.

ing of the parts. In the top right quadrant (high conventionality, high contingency), we have *more or less*, *rise and fall*, *back and forth*, and *deliver the goods*. The bottom left quadrant was predicted to contain non-literal phrases whose words are not as strongly associated with one another as those in the most well-known idioms. The phrases in our dataset that fall into this quadrant include *hard sell*, *hit man*, and *cold feet*. A list of which target phrases landed in each quadrant is given in Appendix D.

For the matched phrases, we assumed that the majority were instances of regular language use, so we predicted them to cluster in the bottom right quadrant. Our results are consistent with this prediction. The horizontal and vertical black lines on the plot were placed at the mean values for each measure. Recall that our examples of “regular language use” consist of head-dependent constructions that share one word with an existing idiom. Although obtaining the phrases in this way may have biased our sample of “regular language use” toward similarity with target phrases, the fact that we still see a clear difference between target and matched average values is all the more striking.

Figure 4 shows only the target phrases that received a human annotation of 1 or 2 for head word literality—that is, the phrases judged to be most non-compositional. As expected, the average score for the target phrases moved more solidly into the idiom quadrant.

We also found no evidence of correlation between contingency and conventionality values among the entire set of phrases, target and

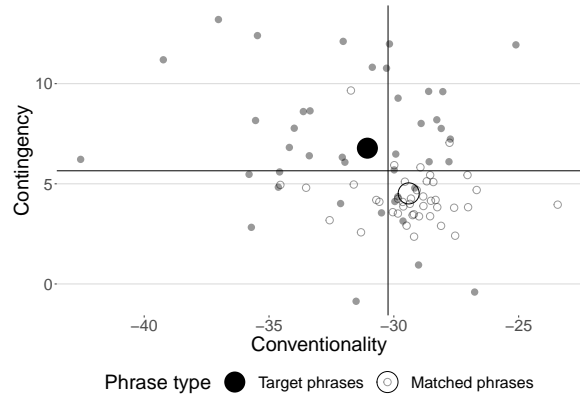


Figure 4: Contingency versus conventionality values of target and matched phrases (for target phrases rated as highly idiomatic). Large circles are average values of all target (black) and all matched (white) phrases.

matched ($r(312) = -0.037$, $p = 0.518$), which is consistent with theories that treat the two properties as independent of each other.

6 Asymmetries between heads and dependents

Our experiments revealed an unexpected but interesting asymmetry between heads and their dependents. Based on conventionality scores, the head word of the target phrases was more conventional on average than the primary non-head word. A two-sample t-test revealed that this difference was significant ($t = 3.029$, $df = 252.45$, $p = 0.0027$). The matched phrases did not show a significant difference between heads and non-heads ($t = 1.506$, $df = 277.42$, $p = 0.1332$).

Figure 5 presents the data in a different way, with target and matched phrases plotted together. The plots show that the variability in overall phrase conventionality, which helps to distinguish idioms and non-idioms, is largely driven by the dependent word (as indicated by the steeper slopes for the non-head effects). This interaction between phrase conventionality and head/non-head is significant (see Row 10 of Table 4 of Appendix B).

In addition, Figure 5 illustrates that this discrepancy between heads and non-heads is largest for verb object phrases. We confirm this by fitting a linear model of word conventionality with predictors for phrase conventionality (average of the component words), head versus non-head word, and syntactic class, plus all interactions, using sum coding to compare factor levels of syntactic class.⁷

⁷ $\text{WordConv} \sim \text{PhraseConv} * \text{Class} * \text{Head}$

The effect of headedness on conventionality scores is significantly greater for verb object phrases than the global effect of headedness (see Panel 4 of Figure 5; Row 14 of Table 4 of Appendix B). We raise the possibility that there is an additive effect of linear order, with conventionality decreasing from left to right through the phrase. For verb object phrases, the two effects go in the same direction, whereas for adjective noun and noun noun phrases, the linear order effect counteracts the headedness effect. We are not aware of any other theory positing the attribution of idiomatic meaning to incremental chunks in this way. Our results suggest that syntactic constituency alone is not enough to explain the observed patterns.

We note that there is disagreement in the literature about whether binomial phrases (which are coordinate structures) contain a head at all. Some proposals treat the first conjunct as the head (e.g., Rothstein, 1991; Kayne, 1994; Gazdar, 1981), while others treat the conjunction as the head or claim that there is no head (e.g., Bloomfield, 1933). We find that in the binomial case, the first conjunct patterns like the heads of the other phrase types, though how much of this effect may be driven by linear order remains unclear. This may provide suggestive converging evidence for the first-conjunct-as-head theory, though further exploration of this idea is needed.

7 Related work

Many idiom detection models build on insights about unconventional meaning in metaphor. A number of approaches use distributional models, such as Kintsch (2000), Utsumi (2011), Sa-Pereira (2016), and Shutova et al. (2012), the latter of which was one of the first to implement a fully unsupervised approach for encoding relationships between words, their contexts, and their dependencies. A related line of work aims to automatically determine whether potentially idiomatic expressions are being used idiomatically or literally, based on contextual information (Katz and Giesbrecht, 2006; Fazly et al., 2009; Sporleder and Li, 2009, 2014). Our measure of conventionality is inspired by the insights of these models; as described in Section 3.2, our measure uses differences in embeddings across contexts.

Meanwhile, approaches to collocation detection have taken a probabilistic or information-theoretic approach that seeks to identify collocation

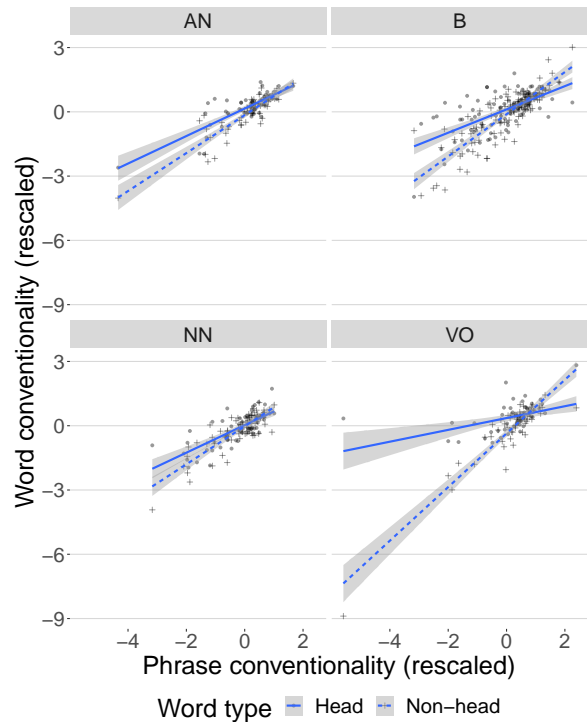


Figure 5: Change in head versus non-head conventionality scores as phrase conventionality increases, for all phrases (target and matched), separated by phrase type (adjective noun, binomial, noun noun, and verb object).

tions using word combination probabilities. PMI is a frequently-used quantity for measuring co-occurrence probabilities (Fano, 1961; Church and Hanks, 1990). Other implementations include selectional association (Resnik, 1996), symmetric conditional probability (Ferreira and Pereira Lopes, 1999), and log-likelihood (Dunning, 1993; Daille, 1996). Like our study, most previous work on idiom and collocation detection focuses specifically on English.

While much of the literature in NLP recognizes that idioms share a cluster of properties, including semantic idiosyncrasy, syntactic inflexibility, and institutionalization (e.g., Sag et al., 2002; Fazly and Stevenson, 2006; Fazly et al., 2009), our approach is novel in attempting to characterize idioms along two orthogonal dimensions that correspond to specific proposals from the cognitive science literature. Our measures may offer a new avenue for tackling automatic idiom detection.

8 Discussion & Conclusion

We investigated whether idioms could be characterized as occupying the intersection between contingency and conventionality, without needing to

appeal to idiom-specific machinery that associates the storage of multi-word expressions with the property of unconventional meaning, as has been proposed in previous work.

When we plotted conventionality and contingency scores against each other, we found that idioms fell, on average, in the area of low conventionality and high contingency, as expected. Regular, non-idiomatic phrases fell in the high conventionality, low contingency area, also as expected. The lack of correlation between the two measures provides support for theories that divorce the notions of conventionality and contingency.

Our results suggest that idioms represent just one of the ways that conventionality and contingency can interact, analogous to collocations or metaphor. We also presented the novel finding that the locus of non-conventionality in idioms resides primarily in the dependent, rather than the head, of the phrase, a result that merits further study.

9 Ethics statement

This paper uses computational tools to argue for a theoretical position about idioms. Our idiom dataset was automatically generated from an existing corpus, and so did not involve data collection from human participants on our part. To validate our conventionality measure, we conducted an additional online experiment with crowdworkers on Amazon Mechanical Turk, for which we obtained REB approval. Details about the participants, recruitment, and consent process are given in Section 4. We note that one limitation of this work is that it only investigates English idioms, potentially contributing to an over-focus on English in this domain.

Acknowledgments

We thank Reuben Cohn-Gordon, Jacob Hoover, Alessandro Sordoni, and the Montreal Computational and Quantitative Linguistics Lab at McGill University for helpful feedback. We also gratefully acknowledge the support of the Natural Sciences and Engineering Research Council of Canada, the Fonds de Recherche du Québec, and the Canada CIFAR AI Chairs Program.

References

Leonard Bloomfield. 1933. *Language*. Henry Holt, New York.

Samuel Bobrow and Susan Bell. 1973. On catching on to idiomatic expressions. *Memory & Cognition*, 1:343–346.

Benjamin Bruening. 2019. Idioms, collocations, and structure: Syntactic constraints on conventionalized expressions. *Natural Language and Linguistic Theory*, 70:491–538.

Benjamin Bruening, Xuyen Dinh, and Lan Kim. 2018. Selection, idioms, and the structure of nominal phrases without classifiers. *Glossa*, 42:1–46.

Lou Burnard. 2000. *The British National Corpus Users Reference Guide*. Oxford University Computing Services.

Kenneth Ward Church and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22–29.

A. P. Cowie. 1981. The treatment of collocations and idioms in learners' dictionaries. *Applied Linguistics*, 1.3:223–235.

Beatrice Daille. 1996. Study and implementation of combined techniques for automatic extraction of terminology. In *The Balancing Act: Combining Symbolic and Statistical Approaches to Language*, pages 49–66. MIT Press, Cambridge, MA.

Josh De Leeuw. 2015. jsPsych: A JavaScript library for creating behavioral experiments in a web browser. *Behavioral Research Methods*, 47(1):1–12.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL-HLT*, 1:4171–4186.

Anna Maria Di Sciullo and Edwin Williams. 1987. *On the Definition of Word*. MIT Press, Cambridge, MA.

Ted Dunning. 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, 19(1):61–74.

Robert Mario Fano. 1961. *Transmission of Information: A Statistical Theory of Communications*. MIT Press, Cambridge, MA.

Afsaneh Fazly, Paul Cook, and Suzanne Stevenson. 2009. Unsupervised type and token identification of idiomatic expressions. *Computational Linguistics*, 35:61–103.

Afsaneh Fazly and Suzanne Stevenson. 2006. Automatically constructing a lexicon of verb phrase idiomatic combinations. In *Proceedings of the 11th Conference of the European Chapter of the ACL (EACL'06)*, pages 337–344. Trento, Italy.

Joaquim Ferreira and Gabriel Pereira Lopes. 1999. A local maxima method and a fair dispersion normalization for extracting multiword units from corpora. In *Sixth Meeting on Mathematics of Language*, pages 369–381.

- Gerald Gazdar. 1981. Unbounded dependencies and coordinate structure. *Linguistic Inquiry*, 12:155–184.
- Jennifer Hay. 2003. *Causes and Consequences of Word Structure*. Routledge, New York, NY.
- Ray Jackendoff. 2002. What’s in the lexicon? In S. Nootboom, F. Weerman, and F. Wijnen, editors, *Storage and Computation in the Language Faculty*. Kluwer Academic Press, Dordrecht.
- Graham Katz and Eugenie Giesbrecht. 2006. Automatic identification of non-compositional multiword expressions using latent semantic analysis. In *Proceedings of the Workshop on Multiword Expressions: Identifying and Exploiting Underlying Properties*, pages 12–19, Sydney, Australia. Association for Computational Linguistics.
- Richard Kayne. 1994. *The Antisymmetry of Syntax*. MIT Press, Cambridge, MA.
- Walter Kintsch. 2000. A computational theory of metaphor comprehension. *Psychonomic Bulletin & Review*, 7:257–266.
- Alexandra Kuznetsova, Per B. Brockhoff, and Rune H. B. Christensen. 2017. lmerTest package: tests in linear mixed effects models. *Journal of Statistical Software*, 82:1–26.
- Roger Levy and Galen Andrew. 2006. Tregex and Tsurgeon: Tools for querying and manipulating tree data structures. In *5th International Conference on Language Resources and Evaluation (LREC 2006)*.
- Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit. In *Association for Computational Linguistics (ACL) System Demonstrations*, pages 55–60.
- Emily Morgan and Roger Levy. 2016. Abstract knowledge versus direct experience in processing of binomial expressions. *Cognition*, 157:384–402.
- Geoffrey Nunberg, Ivan A. Sag, and Thomas Wasow. 1994. Idioms. *Language*, 70:491–538.
- Timothy J. O’Donnell. 2015. *Productivity and reuse in language: A theory of linguistic computation and storage*. MIT Press.
- Frank R. Palmer. 1981. *Semantics*. Cambridge.
- Steven Pinker and Alan Prince. 1988. On language and connectionism: Analysis of a parallel distributed processing model of language acquisition. *Cognition*, 28:73–193.
- Stephen G. Pulman. 1993. The recognition and interpretation of idioms. In C. Cacciari et al., editor, *Idioms—Processing, Structure, and Interpretation*, pages 249–270. Lawrence Erlbaum Associates, Hillsdale, NJ.
- Philip Resnik. 1996. Selectional constraints: An information-theoretic model and its computational realization. *Cognition*, 61:127–159.
- Suzanne Z. Riehemann. 2001. A constructional approach to idioms and word formation. PhD thesis, Stanford University.
- Susan Rothstein. 1991. Heads, projections, and categorial determination. In K. Leffel and D. Bouchard, editors, *Views on phrase structure*, pages 97–112. Kluwer, Dordrecht.
- Fernando Sa-Pereira. 2016. Distributional representations of idioms. Master’s thesis, McGill University.
- Ivan A. Sag, Timothy Baldwin, Francis Bond, Ann Copestake, and Dan Flickinger. 2002. Multiword expressions: A pain in the neck for NLP. In A. Gelbukh, editor, *Computational Linguistics and Intelligent Text Processing. CICLing 2002. Lecture Notes in Computer Science*, volume 2276. Springer, Berlin, Heidelberg.
- Ekaterina Shutova, Tim Van de Cruys, and Anna Korhonen. 2012. Unsupervised metaphor paraphrasing using a vector space model. In *Proceedings of COLING 2012: Posters*, pages 1121–1130.
- Caroline Sporleder and Linin Li. 2009. Unsupervised recognition of literal and non-literal use of idiomatic expressions. In *Proceedings of the 12th Conference of the European Chapter of the ACL*, pages 754–762.
- Caroline Sporleder and Linin Li. 2014. Classifying idiomatic and literal expressions using topic models and intensity of emotions. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2019–2027.
- Megan Schildmier Stone. 2016. The difference between bucket-kicking and kicking the bucket: Understanding idiom flexibility. PhD thesis, University of Arizona.
- Milan Studený and Jirina Vejnarová. 1998. The multi-information function as a tool for measuring stochastic dependence. In *Proceedings of the NATO Advanced Study Institute on Learning in Graphical Models*. Kluwer Academic Publishers, Norwell, MA.
- David Swinney and Anne Cutler. 1979. The access and processing of idiomatic expressions. *Journal of Verbal Learning and Verbal Behavior*, 18:523–534.
- R Core Team. 2017. *R: A Language and Environment for Statistical Computing*. Vienna: R Foundation for Statistical Computing.
- Debra Titone, Kyle Lovseth, Kristina Kasparian, and Mehrgol Tiv. 2019. Are figurative interpretations of idioms directly retrieved, compositionally built, or both? *Canadian Journal of Experimental Psychology/Revue canadienne de psychologie*, 73:216–230.

- Antoine Tremblay and Harald Baayen. 2010. Holistic processing of regular four-word sequences: A behavioral and ERP study of the effects of structure, frequency, and probability on immediate free recall. In D. Wood, editor, *Perspectives on Formulaic Language: Acquisition and Communication*, pages 151–173. The Continuum International Publishing Group, London.
- Akira Utsumi. 2011. Computational exploration of metaphor comprehension processes using a semantic space model. *Cognitive Science*, 35:251–296.
- Tim Van de Cruys. 2011. Two multivariate generalizations of pointwise mutual information. In *Proceedings of the Workshop on Distributional Semantics and Compositionality*. Association for Computational Linguistics, Stroudsburg, PA.
- Michael Wagner. 2021. [Prosodylab experimenter](#).
- Satosi Watanabe. 1960. Information theoretical analysis of multivariate correlation. *IBM Journal of Research and Development*, 4:66–82.
- Uriel Weinreich. 1969. Problems in the analysis of idioms. In J. Puhvel, editor, *Substance and structure of language*, pages 23–81. University of California Press.
- Mary McGee Wood. 1986. *A Definition of Idiom*. University of Birmingham.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. XLNet: Generalized autoregressive pretraining for language understanding. *NeurIPS*, pages 5754–5764.

A

On the following page is a list of the target phrases in our corpus.

Target phrase	Type	Target phrase	Type
deliver the goods	VO	swimming pool	NN
run the show	VO	cash cow	NN
rock the boat	VO	foot soldier	NN
call the shots	VO	attorney general	NN
talk turkey	VO	hit list	NN
cut corners	VO	soup kitchen	NN
jump the gun	VO	bull market	NN
have a ball	VO	boot camp	NN
foot the bill	VO	message board	NN
break the mold	VO	gold mine	NN
pull strings	VO	report card	NN
mean business	VO	comfort food	NN
raise hell	VO	pork barrel	NN
close ranks	VO	flower girl	NN
strike a chord	VO	hit man	NN
cry wolf	VO	blood money	NN
lose ground	VO	cottage industry	NN
make waves	VO	board game	NN
clear the air	VO	death wish	NN
pay the piper	VO	word salad	NN
spill the beans	VO	altar boy	NN
bite the dust	VO	bench warrant	NN
saw logs	VO	time travel	NN
lead the field	VO	love language	NN
take the powder	VO	night owl	NN
buy the farm	VO	life blood	NN
turn tail	VO	road rage	NN
get the sack	VO	light house	NN
hit the sack	VO	bid price	NN
kick the bucket	VO	carrot cake	NN
shoot the bull	VO	command line	NN
		stag night	NN
		husband material	NN

Target phrase	Type	Target phrase	Type
cold feet	AN	by and large	B
green light	AN	more or less	B
red tape	AN	bits and pieces	B
black box	AN	up and down	B
blue sky	AN	rise and fall	B
bright future	AN	sooner or later	B
sour grape	AN	rough and ready	B
green room	AN	far and wide	B
easy money	AN	give and take	B
last minute	AN	time and effort	B
hard heart	AN	pro and con	B
hot dog	AN	sick and tired	B
raw talent	AN	back and forth	B
hard labor	AN	day and night	B
broken home	AN	wear and tear	B
fat chance	AN	nut and bolt	B
dirty joke	AN	tooth and nail	B
happy hour	AN	on and off	B
high time	AN	win or lose	B
rich history	AN	food and shelter	B
clean slate	AN	odds and ends	B
stiff competition	AN	in and out	B
maiden voyage	AN	sticks and stones	B
cold shoulder	AN	make or break	B
clean energy	AN	part and parcel	B
hard sell	AN	loud and clear	B
back pay	AN	cops and robbers	B
deep pockets	AN	short and sweet	B
broken promise	AN	safe and sound	B
dead silence	AN	black and blue	B
blind faith	AN	toss and turn	B
tight schedule	AN	fair and square	B
brutal honesty	AN	heads or tails	B
bright idea	AN	hearts and flowers	B
kind soul	AN	rest and relaxation	B
bruised ego	AN	flesh and bone	B
		life and limb	B
		checks and balances	B
		fast and loose	B
		high and dry	B
		pots and pans	B
		now or never	B
		hugs and kisses	B
		bread and butter	B
		risk and reward	B
		cloak and dagger	B
pins and needles	B	nickel and dime	B
sugar and spice	B	rhyme or reason	B
neat and tidy	B	leaps and bounds	B
step by step	B	live and learn	B
lost and found	B	peace and quiet	B
old and grey	B	song and dance	B

B

Table 2: Model results table with human literalness rating as the dependent variable, using `lmer`

Coefficient	$\hat{\beta}$	$SE(\hat{\beta})$	t	p
Intercept	0.051	0.019	1.655	0.049
Conv	0.185	0.050	3.725	< 0.001
Head(False)	0.015	0.014	1.050	0.147
Conv:Head(False)	0.073	0.053	1.376	0.084

$n = 4945$

Table 3: Model results table for model described in Section 5.2, with contingency score as the dependent variable, using `lmer`

Coefficient	$\hat{\beta}$	$SE(\hat{\beta})$	t	p
Intercept	4.949	0.114	43.379	< 0.001
Target(True)	1.253	0.165	7.587	< 0.001
Class(VO)	-0.195	0.200	-0.975	0.165
Class(AN)	-0.662	0.201	-3.297	< 0.001
Class(B)	1.796	0.179	10.045	< 0.001
Target(True): Class(VO)	0.501	0.303	1.654	0.049
Target(True): Class(AN)	-0.896	0.286	-3.135	< 0.001
Target(True): Class(B)	1.394	0.247	5.641	< 0.001

$n = 99573$

Table 4: Model results table for model described in Section 6, with conventionality score as the dependent variable

Coefficient	$\hat{\beta}$	$SE(\hat{\beta})$	t	p
Intercept	0.163	0.035	4.614	< 0.001
PhraseConv	0.526	0.036	14.453	< 0.001
Class(VO)	0.196	0.065	3.020	0.003
Class(AN)	-0.135	0.063	-2.153	0.032
Class(B)	-0.010	0.064	-0.150	0.881
Head(False)	-0.326	0.050	-6.525	< 0.001
PhraseConv:Class(VO)	-0.250	0.062	-4.043	< 0.001
PhraseConv:Class(AN)	0.117	0.069	1.683	0.093
PhraseConv:Class(B)	0.116	0.068	1.694	0.091
PhraseConv:Head(False)	0.476	0.051	9.247	< 0.001
Class(VO):Head(False)	-0.392	0.092	-4.271	< 0.001
Class(AN):Head(False)	0.271	0.089	3.044	0.002
Class(B):Head(False)	0.019	0.091	0.212	0.832
PhraseConv:Class(VO): Head(False)	0.500	0.087	5.717	< 0.001
PhraseConv:Class(AN): Head(False)	-0.233	0.098	-2.380	0.018
PhraseConv:Class(B): Head(False)	-0.232	0.097	-2.396	0.017

$n = 584$

C

To confirm that our results are not simply an artifact of the dataset we used, we replicated the study on a second dataset, which is the set of phrases used in the idiom detection work of Fazly et al. (2009). We did not have any hand in choosing the phrases in this dataset, and it has very little overlap with our own. We once again fail to find evidence that the two dimensions of conventionality and contingency are correlated with one another in this set of phrases ($r(24) = -0.276$, $p = 0.172$), and we see a similar spread of data across the four quadrants, shown in Figure 6.

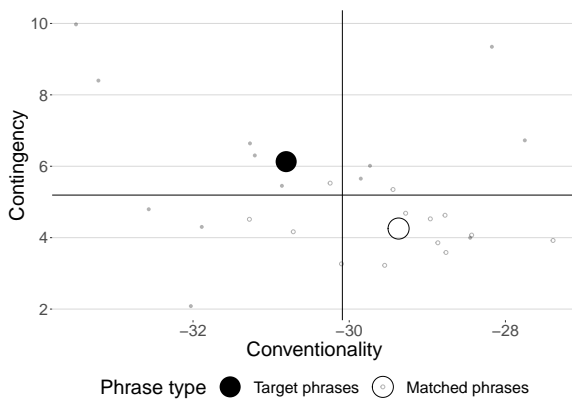


Figure 6: Contingency and conventionality values of target and matched phrases. Large circles are average values of all target (black) and matched (white) phrases.

D

Below is a list of the target phrases that landed in each of the quadrants in Figure 3, for those phrases that occurred at least 30 times in the corpus.

Top left	Top right
black and blue	back and forth
black box	bits and pieces
bread and butter	boot camp
by and large	bright future
call the shots	deep pockets
checks and balances	deliver the goods
clear the air	far and wide
cottage industry	food and shelter
cut corners	heads or tails
day and night	high and dry
foot soldier	more or less
give and take	on and off
gold mine	part and parcel
happy hour	pull strings
have a ball	rise and fall
high time	rock the boat
in and out	run the show
loud and clear	song and dance
make or break	swimming pool
nuts and bolts	up and down
peace and quiet	
red tape	
safe and sound	
sick and tired	
soup kitchen	
sour grapes	
win or lose	

Bottom left	Bottom right
cold feet	blue sky
green light	board game
hard sell	bright idea
hit man	get the sack
hot dog	green room
last minute	hit list
lose ground	report card
mean business	time and effort