

Advances in how programs treat natural language words have a big impact in AI.

BY NOAH A. SMITH

Contextual Word Representations: Putting Words into Computers

THIS ARTICLE AIMS to tell the story of how we put words into computers. It is part of the story of the field of natural language processing (NLP), a branch of artificial intelligence.^a It targets a wide audience with a basic understanding of computer programming, but avoids a detailed mathematical treatment, and it does not present any algorithms. It also does not focus on any particular NLP application, such as translation, question answering, or information extraction. The ideas presented here were developed by many researchers over many decades, so the citations are not exhaustive but rather direct the reader to a handful of papers that are, in the author's view, seminal. After reading this article, you should have a general understanding of word vectors (also known as word embeddings): why they exist, what problems they solve, where they come from, how they have changed over time, and what open questions exist about them.

a Recommended NLP textbooks: Jurafsky and Martin¹⁹ and Eisenstein.¹⁰

There are two ways to talk about words:

- A *word token* is a word observed in a piece of text. In some languages, identifying the boundaries of the word tokens is a complicated procedure (and speakers of the language may not agree on the “correct” rules for splitting text into words), but in English we tend to use whitespace and punctuation to delimit words, and in this article we assume this problem, known as tokenization, is “solved.” For example, the first sentence of this paragraph is typically tokenized as follows (with the end-of-sentence punctuation treated as a separate token):

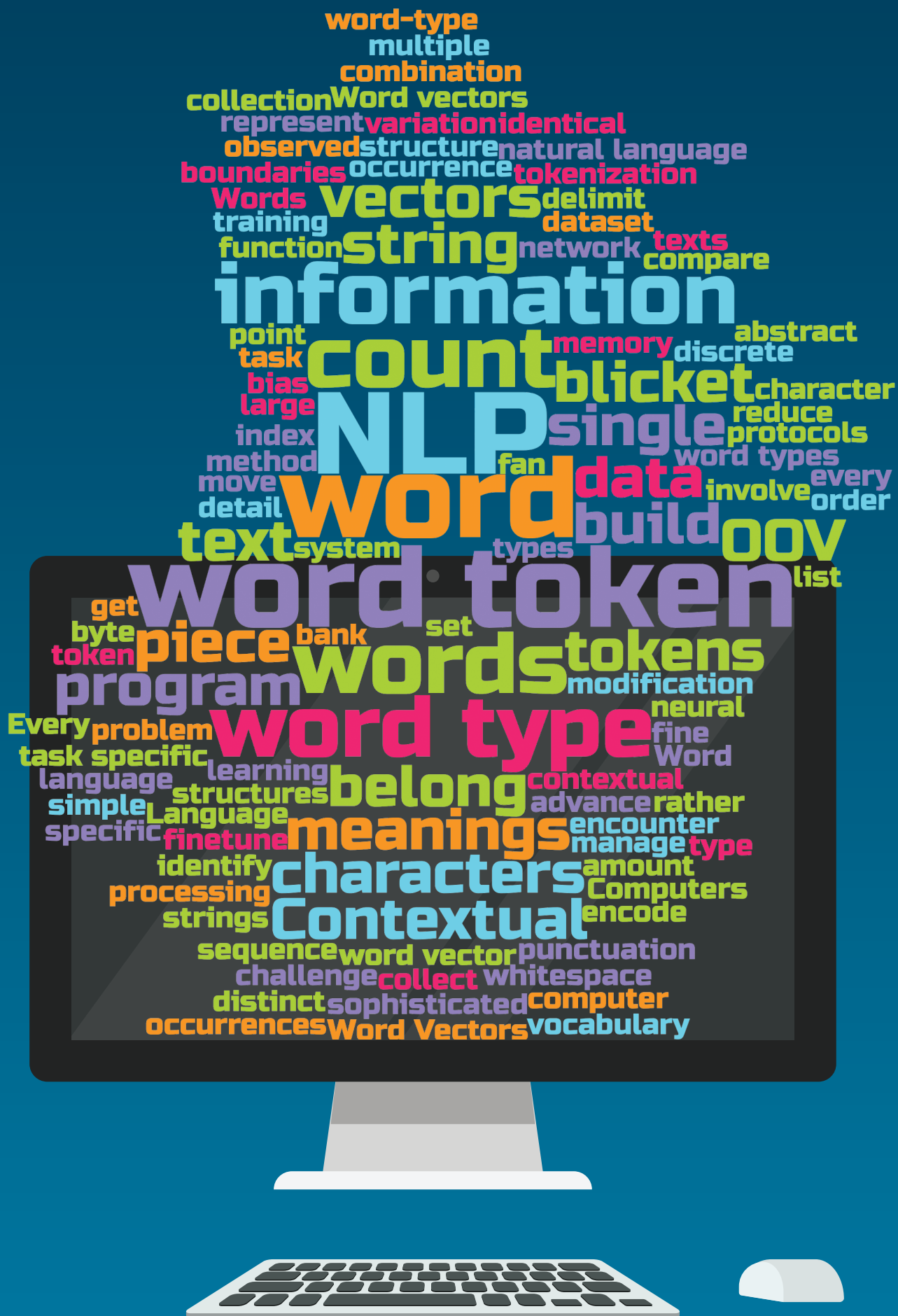
A word token is a word observed in a piece of text.

There are 13 tokens in the sentence.

- A *word type* is a distinct word, in the abstract, rather than a specific instance. Every word token is said to “belong” to its type. In the example, there are only 11 word types, since the two instances of word share the same type, as do the two instances of *a*. (If we ignored the distinction between upper- and lower-case letters, then there would only be 10 types, since the first word *A* would have the same type as the fifth and ninth words.) When we count the occurrences of a vocabulary word in a collection of texts (known as a corpus, plural corpora), we are counting the tokens that belong to the same word type.

» key insights

- Even at the most fundamental level of words, representing the meaning of natural language text computationally is a difficult challenge.
- Different words' meanings can be more or less similar. Continuous vectors have been used to effectively capture this property, and large collections of text have made it possible to automate discovery of many aspects of word meaning similarity. Conventionally, every word in the vocabulary got a single, fixed vector.
- A word's meaning can also vary greatly based on the contexts it appears in. The latest advances recognize and learn about this variation using classical tools from NLP and ML. These methods have shown large improvements across many benchmarks.



Discrete Words

In a computer, the simplest representation of a piece of text is a sequence of characters (depending on the encoding, a character might be a single byte or several). A word type can be represented as a string (ordered list of characters) but comparing whether two strings are identical is costly.

Not long ago, words were usually *integerized*, so that each word type was given a unique (and more or less arbitrary) nonnegative integer value. This had the advantages that every word type was stored in the same amount of memory, and array-based data structures could be used to index other information by word types (like the string for the word, or a count of its tokens, or a richer data structure containing detailed information about the word's potential meanings). The vocabulary could be continuously expanded as new word types were encountered (up to the range of the integer data type, over four billion for 4-byte unsigned integers). And, of course, testing whether two integers are identical is very fast.

The integers themselves did not mean anything; the assignment might be arbitrary, alphabetical, or in the order word tokens were observed in a reference text corpus from which the vocabulary was derived (that is, the type of the first word token observed would get 0, the type of the second word token would get 1 if it was different from the first, and so on). Two word types with related meanings might be assigned distant integers, and two “adjacent” word types in the assignment might have nothing to do with each other. The use of integers is only a convenience following from the data types available in the fashionable programming languages of the day; for example, in Lisp “gensym” would have served the same purpose, although perhaps less efficiently. For this reason, we refer to integer-based representations of word types as *discrete* representations.

Words as Vectors

To see why NLP practitioners no longer treat word types as discrete, it's useful to consider how words get used in NLP programs. Here are some examples:

► Observing a word token in a given document, use it as evidence to help predict a category for the document.

The idea that words can be more or less similar is critical when we consider that NLP programs, by and large, are built using supervised machine learning.

For example, the word *delightful* appearing in a review of a movie is a cue the reviewer might have enjoyed the film and given it a positive rating.^b

► Observing a word token in a given sentence, use it as evidence to predict a word token in the translation of the sentence. For example, the appearance of the word *cucumber* in an English sentence is a cue that the word *concombre* might appear in the French translation.

► Conversely, given the full weight of evidence, choose a word type to write as an output token, in a given context.

In each of these cases, there is a severe shortcoming to discrete word types: information about how to use a particular word as evidence, or whether to generate a word as an output token, cannot be easily shared across words with similar properties. As a simple example, consider filling in the blank in the following sentence:

S. will eat anything, but V. hates ____

Given your knowledge of the world, you are likely inclined to fill in the blank with high confidence as a token of a type like *peas*, *sprouts*, *chicken*, or some other mass or plural noun that denotes food. These word types share something (together with the other words for food), and we would like for a model that uses words to be able to use that information.^c To put it another way, our earlier interest in testing whether two words are identical was perhaps too strict. Two non-identical words may be more or less similar.

The idea that words can be more or less similar is critical when we consider that NLP programs, by and large, are built using supervised machine learning, that is, a combination of examples demonstrating the inputs and

^b Context matters. “The most delightful part of seeing this movie was the popcorn” would signal just the opposite. See Pang and Lee²⁷ for a detailed treatment of the problems of sentiment and opinion analysis.

^c One situation where this lack of sharing is sorely noticed is the case of new words, sometimes called “out of vocabulary” (OOV) words. When an NLP program encounters an OOV word token, say *blicket*, what should it do? By moving away from discrete words, we have managed to reduce the occurrence of truly OOV word types by collecting information about an increasingly large set of words in advance of building the NLP program.

outputs to a task (at least one of which consists of words) and a mechanism for generalizing from those input-output pairings. Such a mechanism should ideally exploit similarity: anything it discovers about one word should transfer to similar words.

Where might this information about similarity come from? There are two strands of thought about how to bring such information into programs. We might trace them back to the rationalist and empiricist traditions in philosophy, though I would argue it's unwise to think of them in opposition to each other.

One strand suggests that humans, especially those trained in the science of human language, know this information, and we might design data structures that encode it explicitly, allowing our programs to access it as needed. An example of such an effort is WordNet,¹³ a lexical database that stores words and relationships among them such as synonymy (when two words can mean the same thing) and hyponymy (when one word's meaning is a more specific case of another's). WordNet also explicitly captures the different senses of words that take multiple meanings, such as *fan* (a machine for blowing air, or someone who is supportive of a sports team or celebrity). Linguistic theories of sentence structure (syntax) offer another way to think about word similarity in the form of categories like “noun” and “verb.”

The other strand suggests the information resides in artifacts such as text corpora, and we can use a separate set of programs to collect and suitably organize the information for use in NLP. With the rise of ever-larger text collections on the Web, this strand came to dominate, and the programs used to draw information from corpora have progressed through several stages, from count-based statistics, to modeling using more advanced statistical methods, to increasingly powerful tools from machine learning.

From either of these strands (or, more commonly in practice, by intertwining them), we can derive a notion of a word type as a vector instead of an integer.^d In doing so, we can choose the

dimensionality of the vector and allocate different dimensions for different purposes. For example:

- Each word type may be given its own dimension and assigned 1 in that dimension (while all other words get 0 in that dimension). Using dimensions only in this way, and no other, is essentially equivalent to integerizing the words; it is known as a “one hot” representation, because each word type's vector has a single 1 (“hot”) and is otherwise 0.

- For a collection of word types that belong to a known class (for example, days of the week), we can use a dimension that is given binary values. Word types that are members of the class get assigned 1 in this dimension, and other words get 0.

- For word types that are variants of the same underlying root, we can similarly use a dimension to place them in a class. For example, in this dimension, *know*, *known*, *knew*, and *knows* would all get assigned 1, and words that are not forms of *know* get 0.

- More loosely, we can use surface attributes to “tie together” word types that look similar; examples include capitalization patterns, lengths, and the presence of a digit.

- If word types' meanings can be mapped to magnitudes, we might allocate dimensions to try to capture these. For example, in a dimension we choose to associate with “typical weight” *elephant* might get 12,000 while *cat* might get 9. Of course, it's not entirely clear what value to give *purple* or *throw* in this dimension.

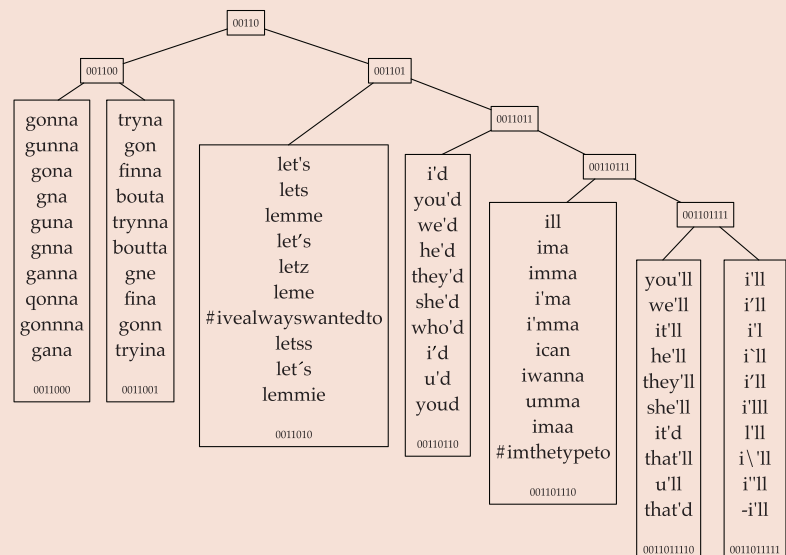
Examples abound in NLP of the allocation of dimensions to vectors representing word types (either syntactic, like “verb,” or semantic, like “animate”), or to multiword sequences (for example, *White House* and *hot dog*). The technical term used for these dimensions is *features*. Features can be designed by experts, or they can be derived using automated algorithms. Note that some features can be calculated even on out-of-vocabulary word types. For example, noting the capitalization pattern of characters in an out-of-vocabulary word might help a system guess whether it should be treated like a person's name.

Words as Distributional Vectors: Context as Meaning

An important idea in linguistics is that words (or expressions) that can be

Figure 1. Example Brown clusters.

These were derived from 56M tweets, see Owoputi et al.²⁶ for details. Shown are the 10 most frequent words in clusters in the section of the hierarchy with prefix bit string 001110. Intermediate nodes in the tree correspond to clusters that contain all words in their descendants. Note that differently spelled variants of words tend to cluster together, as do words that express similar meanings, including hashtags. The full set of clusters can be explored at http://www.cs.cmu.edu/~ark/TweetNLP/cluster_viewer.html. Note that there are several Unicode characters that are visually similar to the apostrophe, resulting in different strings with similar usage.



^d A vector is a list, usually a list of numbers, with a known length, which we call its dimensionality. It is often interpreted and visualized as a direction in a Euclidean space.

used in similar ways are likely to have related meanings¹⁴ (consider our day of the week example). In a large corpus, we can collect information about the ways a word type w is used, for

example, by counting the number of times it appears near every other word type. When we begin looking at the full distribution of contexts (nearby words or sequences of words) in a corpus

where w is found, we are taking a distributional view of word meaning.

One highly successful approach to automatically deriving features based on this idea is *clustering*; for example, the Brown et al.⁴ clustering algorithm automatically organized words into clusters based on the contexts they appear in, in a corpus. Words that tended to occur in the same neighboring contexts (other words) were grouped together into a cluster. These clusters could then be merged into larger clusters. The resulting hierarchy, while by no means identical to the expert-crafted data structure in WordNet, was surprisingly interpretable and useful (an example is shown in Figure 1). It also had the advantage that it could be rebuilt using any given corpus, and every word observed would be included. Hence, suitable word clusters could be built separately for news text, or biomedical articles, or tweets.

Another line of approaches started by creating word vectors in which each dimension corresponded to the frequency the word type occurred in some context.⁸ For instance, one dimension might correspond to *the* and contain the number of times the word occurred within a small window of a token *the*. Contextual patterns on the left or the right, and of varying distances and lengths, might be included. The result is a vector perhaps many times longer than the size of the vocabulary, in which each dimension contains a tiny bit of information that may or may not be useful. An example is shown in Figure 2. Using methods from linear algebra, aptly named *dimensionality reduction*, these vectors could be compressed into shorter vectors in which redundancies across dimensions were collapsed.

These reduced-dimensionality vectors had several advantages. First, the dimensionality could be chosen by the NLP programmer to suit the needs of the program. More compact vectors might be more efficient to compute with and might also benefit from the lossiness of the compression, since corpus-specific “noise” might fall away. However, there is a trade-off; longer, less heavily compressed vectors retain more of the original information in the distributional vectors. While the individual dimensions of

Figure 2. Example calculation of word vectors.

We consider three-word types occurring a science news story ([https:// bit.ly/2B9uaKr](https://bit.ly/2B9uaKr)): astronomers, bodies, and objects. The table above shows the frequency of each word occurring within two positions on either side of the word whose vector we are constructing, giving three (vertical) word vectors with 34 visible dimensions (zeroes and other dimensions not shown). Do you expect bodies to be more similar to astronomers or to objects? The calculation beneath is the cosine similarity score, applied to each word paired with the others (and itself, always giving similarity of one). In this tiny corpus, bodies are closer to objects than either is to astronomers.

Contextual Word Representations:

context words	v(astronomers)	v(bodies)	v(objects)
't			1
,		2	1
.	1		1
1			1
And			1
Belt			1
But	1		
Given			1
Kuiper			1
So	1		
and		1	
are		2	1
between			1
beyond		1	
can			1
contains		1	
from	1		
hypothetical			1
ice		1	
including		1	
is	1		
larger		1	
now	1		
of	1		
only			1
out		1	
potential		1	
the	1		1
these		2	1
they	1		
think	2		
those			1
thought		2	
what	1		

$$\text{cosine_similarity}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \cdot \|\mathbf{v}\|}$$

	astronomers	bodies	objects
astronomers	$\frac{14}{\sqrt{14} \cdot \sqrt{14}} = 1$	$\frac{0}{\sqrt{24} \cdot \sqrt{14}} = 0$	$\frac{1+1}{\sqrt{14} \cdot \sqrt{16}} \approx 0.134$
bodies		$\frac{24}{\sqrt{24} \cdot \sqrt{24}} = 1$	$\frac{2+2+2}{\sqrt{24} \cdot \sqrt{16}} \approx 0.306$
objects			$\frac{16}{\sqrt{16} \cdot \sqrt{16}} = 1$

the compressed vectors are not easily interpreted, we can use well-known algorithms to find a word's nearest neighbors in the vector space, and these were often found to be semantically related words, as one might hope.

Indeed, these observations gave rise to the idea of *vector space semantics*,³³ in which arithmetic operations were applied to word vectors to probe what kind of “meanings” had been learned. Famously, analogies like “*man is to woman as king is to queen*” led to testing whether $v(\text{man}) - v(\text{woman}) = v(\text{king}) - v(\text{queen})$. Efforts to design word vector algorithms to adhere to such properties followed.

The notable disadvantage of reduced-dimensionality vectors is that the individual dimensions are no longer interpretable features that can be mapped back to intuitive building blocks contributing to the word's meaning. The word's meaning is distributed across the whole vector; for this reason, these vectors are sometimes called *distributed representations*.^e

As corpora grew, scalability became a challenge, because the number of observable contexts grew as well. Underlying all word vector algorithms is the notion that the value placed in each dimension of each word type's vector is a parameter that will be optimized, alongside all the other parameters, to best fit the observed patterns of the words in the data. Since we view these parameters as continuous values, and the notion of “fitting the data” can be operationalized as a smooth, continuous objective function, selecting the parameter values is done using iterative algorithms based on gradient descent. Using tools that had become popular in machine learning, faster methods based on stochastic optimization were developed. One widely known collection of algorithms is available as the *word2vec* package.²⁴ A common pattern arose in which industry researchers with large corpora and powerful computing infrastructure would construct word vectors using an established (often expensive) iterative method, and then

publish the vectors for anyone to use.

There followed a great deal of exploration of methods for obtaining distributional word vectors. Some interesting ideas worth noting include:

- When we wish to apply neural networks to problems in NLP (see Figure 3), it's useful to first map each input word token to its vector, and then “feed” the word vectors into the neural network model, which performs a task like translation. The vectors can be fixed in advance (or pretrained from a corpus, using methods like those above, often executed by someone else), or they can be treated as parameters of the neural network model, and adapted to the task specifically.⁶ Finetuning refers to initializing the word vectors by pretraining, then adapting them through task-specific learning algorithms. The word vectors can also be initialized to random values, then estimated solely through task learning, which we might call “learning from scratch”^f

- We can use expert-built data structures like WordNet as additional in-

put to creating word vectors. One approach, *retrofitting*, starts with word vectors extracted from a corpus, then seeks to automatically adjust them so that word types that are related in WordNet are closer to each other in vector space.¹¹

- We can use bilingual dictionaries to “align” the vectors for words in two languages into a single vector space, so that, for example, the vectors for the English word type *cucumber* and the French word type *concombre* have a small Euclidean distance.¹² By constructing a function that reorients all the English vectors into the French space (or vice versa), researchers hoped to align *all* the English and French words, not just the ones in the bilingual dictionary.

- A word's vectors are calculated in part (or in whole) from its character sequence.²¹ These methods tend to make use of neural networks to map arbitrary-length sequences into a fixed-length vector. This has two interesting effects: in languages with intricate word formation systems (morphology),^g variants

f The result of vectors learned from scratch for an NLP task is a collection of distributed representations that were derived from something other than distributional contexts (the task data).

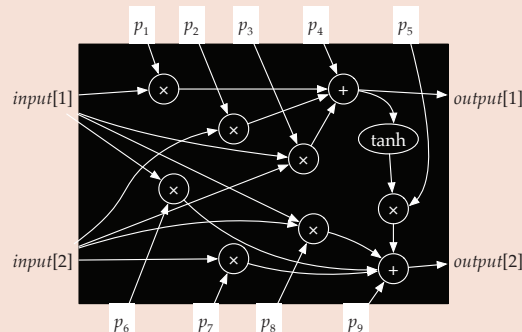
g For example, the present tense form of the French verb *manger* is *mange*, *manges*, *mangeons*, *mangez*, or *mangent*, depending on whether the subject is singular or plural, and first, second, or third person.

Figure 3. A simple neural network.

A neural network is a function from vectors to vectors. A very simple example is a function from two-dimensional inputs to two-dimensional outputs, such as:

$$\begin{aligned} \text{output}[1] &= p_1 \times \text{input}[1] + p_2 \times \text{input}[2] + p_3 \times \text{input}[1] \times \text{input}[2] + p_4 \\ \text{output}[2] &= p_5 \times \tanh(\text{output}[1]) + p_6 \times \text{input}[1] + p_7 \times \text{input}[2] \\ &\quad + p_8 \times \text{input}[1] \times \text{input}[2] + p_9 \end{aligned}$$

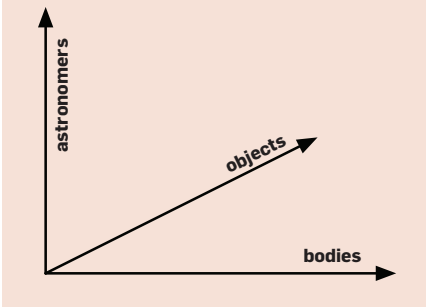
Neural networks are almost always defined in terms of parameters, here denoted by p_1, \dots, p_9 , which are automatically chosen using standard machine learning algorithms. Typically, they include at least one transformation that is not linear (for example, the hyperbolic tangent).



The illustration displays the toy neural network as a computation graph (inside the black box), with inputs on the left, outputs on the right, and each parameter corresponds to a gray box placed along the top and bottom. Round nodes inside the box correspond to intermediate operations (addition, multiplication, and tanh).

e Though distributional information is typically used to build distributed vector representations for word types, the two terms are not to be confused and have orthogonal meanings!

Figure 4. Approximate visualization of the relationships between the three-word vectors calculated in Figure 2.



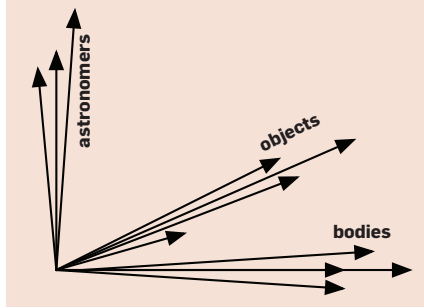
of the same underlying root may have similar vectors; and, differently spelled variants of the same word will have similar vectors. This kind of approach was quite successful for social media texts, where there is rich spelling variation. For example, these variants of the word *would*, all attested in social media messages, would have similar character-based word vectors because they are spelled similarly: *would*, *wud*, *wld*, *wuld*, *wouldd*, *woud*, *wudd*, *whould*, *wouldl*, and *wOuld*.

Contextual Word Vectors

We started this discussion by differentiating between word tokens and word types. All along, we have assumed each word type was going to be represented using a fixed data object (first an integer, then a vector) in our NLP program. This is convenient, but it makes some assumptions about language that do not fit with reality. Most importantly, words have different meanings in different contexts. At a coarse-grained level, this was captured by experts in crafting WordNet, in which, for example, *get* is mapped to over 30 different meanings (or senses). It is difficult to obtain widespread agreement on how many senses should be allocated to different words, or on the boundaries between one sense and another; word senses may be *fluid*.^h Indeed, in many NLP programs based on neural networks, the very first thing that happens is that each word token's type vector is passed into a function that transforms it based on the words in

^h For example, the word *bank* can refer to the side of a river or to a financial institution. When used to refer to a blood bank, we can debate whether the second sense is evoked or a third, distinct one.

Figure 5. Hypothetical visualization of contextual vectors for tokens of astronomers, bodies, and objects from figures 2 and 4.



its nearby context, giving a new version of the word vector, now specific to the token in its particular context. In our example sentence earlier, the two instances of *a* will therefore have different vectors, because one occurs between *is* and *word* and the other occurs between *in* and *piece*; for example, compare figures 4 and 5.

With hindsight, we can now see that by representing word types independent of context, we were solving a problem that was more difficult than it needed to be. Because words mean different things in different contexts, we were requiring that type representations capture *all* of the possibilities (for example, the 30 meanings of *get*). Moving to word token vectors simplifies things, asking the word token representation to capture only what a word means *in this context*. For the same reasons the collection of contexts a word type is found in provide clues about its meaning(s), a particular token's context provides clues about its specific meaning. For instance, you may not know what the word *blicket* means, but if I tell you that I ate a strawberry blicket for dessert, you likely have a good guess.ⁱ

Returning to the fundamental notion of similarity, we would expect words that are similar to each other to be good substitutes for each other. For example, what are some good substitutes for the word *gin*? This question is difficult to answer about the word type (WordNet tells us that *gin* can refer to a liquor for drinking, a trap for hunting, a machine for separating seeds from cotton fibers, or a card game), but easy in a given context (for example, “I use two parts gin to one part vermouth.”). Indeed, *vodka* might

ⁱ Though such examples abound in linguistics, this one is due to Chris Dyer.

even be expected to have a similar contextual word vector if substituted for *gin*.^j

ELMo, which stands for “embeddings from language models,”²⁸ brought a powerful advance in the form of word *token* vectors—that is, vectors for words in context, or contextual word vectors—that are pretrained on large corpora. There are two important insights behind ELMo:

- If every word token is going to have its own vector, then the vector should depend on an arbitrarily long context of nearby words. To obtain a “context vector,” we start with word type vectors, and pass them through a neural network that can transform arbitrary-length sequences of left- and/or right-context word vectors into a single fixed-length vector. Unlike word type vectors, which are essentially lookup tables, contextual word vectors are built from both type-level vectors and neural network parameters that “contextualize” each word. ELMo trains one neural network for left contexts (going back to the beginning of the sentence a token appears in) and another neural network for right contexts (up to the end of the sentence). Longer contexts, beyond sentence boundaries, are in principle possible as well.

- Recall that estimating word vectors required “fitting the data” (here, a corpus) by solving an optimization problem. A longstanding data-fitting problem in NLP is language modeling, which refers to predicting the next word given a sequence of “history” words (briefly alluded to in our filling-in-the-blank example). Many of the word (type) vector algorithms already in use were based on a notion fixed-size contexts, collected across all instances of the word type in a corpus. ELMo went farther, using arbitrary-length histories and directly incorporating the language models known at the time to be most effective (based on recurrent neural networks³²). Although recurrent networks were already widely used in NLP (see Goldberg¹⁵ for a thorough introduction), training them as language models, then using the context vectors they provide for each word token as pretrained word (token) vectors was novel.

It's interesting to see how the ideas around getting words into computers


^j The author does not endorse this substitution in actual cocktails.

have come full circle. The powerful idea that text data can shed light on a word's meaning, by observing the contexts in which a word appears, has led us to try to capture a word token's meaning primarily through the specific context it appears in. This means that every instance of *plant* will have a different word vector; those with a context that look like a context for references to vegetation are expected to be close to each other, while those that are likely contexts for references to manufacturing centers will cluster elsewhere in vector space. Returning to the example in Figure 2, while instances of *bodies* in this article will likely remain closer to *objects* than to *astronomers*, in a medical news story, we might expect instances of *bodies* to be closer to humans (and by extension *astronomers*).


Why is this advance so exciting? Whether the development of contextual word vectors completely solves the challenge of ambiguous words remains to be seen. New ideas in NLP are often tested on benchmark tasks with objectively measurable performance scores. ELMo was shown to be extremely beneficial in NLP programs that:

- ▶ answer questions about content in a given paragraph (9% relative error reduction on the SQuAD benchmark),
- ▶ label the semantic arguments of verbs (16% relative error reduction on the Ontonotes semantic role labeling benchmark),
- ▶ label expressions in text that refer to people, organizations, and other named entities (4% relative error reduction on the CoNLL 2003 benchmark), and
- ▶ resolve which referring expressions refer to the same entities (10% relative error reduction on the Ontonotes coreference resolution benchmark).

Gains on additional tasks were reported by Peters et al.²⁸ and later by other researchers. Howard and Ruder¹⁸ introduced a similar approach, ULMFiT, showing a benefit for text classification methods. A successor approach, bidirectional encoder representations from transformers (BERT⁹) that introduced several innovations to the learning method and learned from more data, achieved a further 45% error reduction (relative to ELMo) on the first task and 7% on the second. On the SWAG benchmark,



An important idea in linguistics is that words (or expressions) that can be used in similar ways are likely to have related meanings.



recently introduced to test grounded commonsense reasoning,³⁵ Devlin et al.⁹ found that ELMo gave 5% relative error reduction compared to non-contextual word vectors, and BERT gave another 66% relative to ELMo. A stream of papers since this article was conceived have continued to find benefits from creative variations on these ideas, resulting in widely adopted models like GPT-2,³⁰ RoBERTa,²³ T5,³¹ XLM,²⁰ and XLNet.³⁴ It is rare to see a single conceptual advance that consistently offers large benefits across so many different NLP tasks.

At this writing, there are many open questions about the relative performance of the different methods. A full explanation of the differences in the learning algorithms, particularly the neural network architectures, is out of scope for this introduction, but it's fair to say that the space of possible learners for contextual word vectors has not yet been fully explored; see Peters et al.²⁹ for some exploration. Some of the findings on BERT suggest that the role of finetuning may be critical; indeed, earlier work that used pre-trained language models to improve text classification *assumed* finetuning was necessary.⁷ While ELMo is derived from language modeling, the modeling problem solved by BERT (that is, the objective function minimized during estimation) is rather different.^k The effects of the dataset used to learn the language model have not been fully assessed, except for the unsurprising pattern that larger datasets tend to offer more benefit.

Cautionary Notes

Word vectors are biased. Like any engineered artifact, a computer program is likely to reflect the perspective of its builders. Computer programs that are built from data will reflect what is in the data—in this case, a text corpus. If the text corpus signals associations between concepts that reflect cultural biases, these associations should be expected to persist in the word vectors and any system that uses them. Hence, it is not surprising that NLP programs that use corpus-derived word vectors associate,

^k BERT pretraining focuses on two tasks: prediction of words given contexts on both sides (rather than one or the other) and predicting the words in a sentence given its preceding sentence.

for example, *doctor* with male pronouns and *nurse* with female ones. Methods for detecting, avoiding, and correcting unwanted associations are an active area of research.^{3,5} The advent of contextual word vectors offers some possibility of new ways to avoid unwanted generalization from distributional patterns.

Language is a lot more than words. Effective understanding and production of language is about more than knowing word meanings; it requires knowing how words are put together to form more complicated concepts, propositions, and more. This is not nearly the whole story of NLP; there is much more to be said about approaches to dealing with natural language syntax, semantics, and pragmatics, and how we operationalize tasks of understanding and production that humans perform into tasks for which we can attempt to design algorithms. One of the surprising observations about contextual word vectors is that, when trained on very large corpora, they make it easier to disambiguate sentences through various kinds of syntactic and semantic parsing; it is an open and exciting question how much of the work of understanding can be done at the level of words in context.

NLP is not a single problem. While the gains above are quite impressive, it's important to remember that they reflect only a handful of benchmarks that have emerged in the research community. These benchmarks are, to varying degrees, controversial, and are always subject to debate. No one who has spent any serious amount of time studying NLP believes they are "complete" in any interesting sense. NLP can only make progress if we have ways of objectively measuring progress, but we also need continued progress on the design of the benchmarks and the scores we use for comparisons. This aspect of NLP research is broadly known as evaluation and includes both human-judgment-based and automatic methods. Anyone who is an enthusiast of NLP (or AI, more generally) should take the time to learn how progress is measured and understand the shortcomings of evaluations currently in use.

What's Next

Over the next few years, I expect to see new findings that apply variations on contextual word vectors to new problems and that explore modifications to the learning methods. For example, build-

ing a system might involve sophisticated protocols in which finetuning and task-specific training are carried out on a series of dataset/task combinations. Personally, I'm particularly excited about the potential for these approaches to improve NLP performance in settings where relatively little supervision is available. Perhaps, for example, ELMo-like methods can improve NLP for low-resource genres and languages.²⁵ Likewise, methods that are computationally less expensive have the potential for broad use (for example, Gururangan et al.¹⁷). I also expect there will be many attempts to characterize the generalizations that these methods are learning (and those that they are not learning) in linguistic terms; see for example Goldberg¹⁶ and Liu et al.²²

Further Reading

An introductory guide to linguistics for those interested in NLP is provided by Bender¹ and Bender and Lascarides.² A more thorough mathematical treatment of the topics noted here is given in chapter 14 of Eisenstein.¹⁰ For contextual word vectors, the original papers are recommended at this writing.^{9,28}

Acknowledgments. Thanks to Oren Etzioni, Chris Dyer, and students from the University of Washington's Winter 2019 CSE 447 class. NSF grant IIS-1562364 supported research related to this article. C

References

1. Bender, E.M. *Linguistic Fundamentals for Natural Language Processing: 100 Essentials from Morphology and Syntax*. Morgan & Claypool, 2013.
2. Bender, E.M. and Lascarides, A. *Linguistic Fundamentals for Natural Language Processing II: 100 Essentials from Semantics and Pragmatics*. Morgan & Claypool, 2019.
3. Bolukbasi, T., Chang, K., Zou, J.Y., Saligrama, V., and Kalai, A.T. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. In *Proceedings of Advances in Neural Information Processing Systems*, 2016, 4349–4357.
4. Brown, P.F., Desouza, P.V., Mercer, R.L., Della Pietra, V.J., and Lai, J.C. Class-based n-gram models of natural language. *Computational Linguistics* 18, 4 (1992), 467–479.
5. Caliskan, A., Bryson, J.J., and Narayanan, A. Semantics derived automatically from language corpora contain human biases. *Science* 356, 6334 (2017), 183–186.
6. Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. 2011. Natural language processing (almost) from scratch. *J. Machine Learning Research* 12 (2011), 2493–2537.
7. Dai, A.D. and Le, Q.V. Semisupervised sequence learning. In *Proceedings of Advances in Neural Information Processing Systems*, 2015.
8. Deerwester, S.C., Dumais, S.T., Landauer, T.K., Furnas, G.W., and Harshman, R.A. Indexing by latent semantic analysis. *J. Amer. Soc. Information Science* 41, 6 (1990), 391–407.
9. Devlin, J., Chang, M.W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of 2019 NAACL*.
10. Eisenstein, J. *Introduction to Natural Language Processing*. MIT Press, 2019.

11. Faruqui, M., Dodge, J., Jauhar, S.K., Dyer, C., Hovy, E., and Smith, N.A. Retrofitting word vectors to semantic lexicons. In *Proceedings of 2016 NAACL*.
12. Faruqui, M. and Dyer, C. Improving vector space word representations using multilingual correlation. In *Proceedings of 2014 EACL*.
13. Fellbaum, C. (Ed.). *WordNet: An Electronic Lexical Database*. MIT Press, 1998.
14. Firth, J.R. A synopsis of linguistic theory 1930–1955. *Studies in Linguistic Analysis*. Blackwell, 1957, 1–32.
15. Goldberg, Y. *Neural Network Methods for Natural Language Processing*. Morgan & Claypool, 2017.
16. Goldberg, Y. Assessing BERT's syntactic abilities. 2019; arXiv:1901.05287.
17. Gururangan, S., Dang, T., Card, D., and Smith, N.A. Variational pretraining for semi-supervised text classification. In *Proceedings of 2019 ACL*.
18. Howard, J. and Ruder, S. Universal language model finetuning for text classification. 2018; arXiv:1801.06146.
19. Jurafsky, D. and Martin, J.H. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition* (3rd Ed.). Prentice Hall, forthcoming.
20. Lample, G. and Conneau, A. Cross-lingual language model pretraining. In *Proceedings of 2019 NeurIPS*.
21. Ling, W. et al. Finding function in form: Compositional character models for open vocabulary word representation. In *Proceedings of 2015 EMNLP*.
22. Liu, N.F., Gardner, M., Belinkov, Y., Peters, M.E. and Smith, N.A. Linguistic knowledge and transferability of contextual representations. In *Proceedings of 2019 NAACL*.
23. Liu, Y. et al. RoBERTa: A robustly optimized BERT pretraining approach. 2019; arXiv:1907.11692.
24. Mikolov, T., Chen, K., Corrado, G., and Dean, J. Efficient estimation of word representations in vector space. In *Proceedings of 2013 ICLR*.
25. Mulcaire, P., Kasai, J., and Smith, N.A. Polyglot contextual representations improve crosslingual transfer. In *Proceedings of 2019 NAACL*.
26. Owoputi, O., O'Connor, B., Dyer, C., Gimpel, K., Schneider, N., and Smith, N.A. Improved part-of-speech tagging for online conversational text with word clusters. In *Proceedings of 2013 NAACL*.
27. Pang, B. and Lee, L. *Opinion Mining and Sentiment Analysis. Foundations and Trends in Information Retrieval*, Vol. 2. Now Publishers, Inc. 2008.
28. Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. Deep contextualized word representations. In *Proceedings of 2018 NAACL*.
29. Peters, M.E., Neumann, M., Zettlemoyer, L., and Yih, W. Dissecting contextual word embeddings: Architecture and representation. In *Proceedings of 2018 EMNLP*.
30. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D. and Sutskever, I. Language models are unsupervised multitask learners. 2019.
31. Raffel, C. et al. Exploring the limits of transfer learning with a unified text-to-text transformer. 2019; arXiv:1910.10983.
32. Sundermeyer, M., Schütter, R., and Ney, H. LSTM neural networks for language modeling. In *Proceedings of 2012 Interspeech*.
33. Turney, P.D. and Pantel, P. From frequency to meaning: Vector space models of semantics. *J. Artificial Intelligence Research* 37, 1 (2010), 141–188.
34. Yang, Z., Dai, Y., Yang, Y., Carbonell, J., Salakhutdinov, R. and Le, Q.V. XLNet: Generalized autoregressive pretraining for language understanding. In *Proceedings of Advances in Neural Information Processing Systems*, 2019.
35. Zellers, R., Bisk, Y., Schwartz, R., and Choi, Y. SWAG: A large-scale adversarial dataset for grounded commonsense inference. In *Proceedings of 2018 EMNLP*.

Noah A. Smith (nasmith@cs.washington.edu) is a professor in the School of Computer Science and Engineering at the University of Washington, and senior research manager for the AllenNLP team at the Allen Institute for Artificial Intelligence, Seattle, WA, USA.

© 2020 ACM 0001-0782/20/6 \$15.00.



Watch the author discuss this work in the exclusive *Communications* video. <https://cacm.acm.org/videos/contextual-word-representations>