

Integrating Approaches to Word Representation

Yuval Pinter *

September 2021

The problem of representing the atomic elements of language in modern neural learning systems is one of the central challenges of the field of natural language processing. I present a survey of the distributional, compositional, and relational approaches to addressing this task, and discuss various means of integrating them into systems, with special emphasis on the word level and the out-of-vocabulary phenomenon.

1 Introduction

The mission of natural language processing (NLP) as a computational research field is to enable machines to function in human-oriented environments where language is the medium of communication. We want them to understand our utterances, to connect these utterances with the objects and concepts of the surrounding world, to produce language which is meaningful to us and helps us navigate a task or satisfy an emotional need. Over the years of its existence, the mainstream of NLP has known shifts motivated by developments in computation, in linguistics, in foundational artificial intelligence, and in learning theory. Since the mid-2010's, the clear dominant framework for tackling NLP tasks, and an undeniably powerful one, has been that of deep neural networks (DNNs). This connectionist approach was originally motivated by the workings of the human brain, but has since developed its own characteristics, and formed a well-defined landscape for exploration which includes constraints stemming from the fundamental properties of its design.

This survey focuses on one of these built-in constraints, which I believe to be central to DNNs in the context of natural language, and specifically of text processing, namely that of **representations**. DNNs “live” in metric space: their operation manipulates real numbers organized into vectors and matrices, propagating function applications and calculated values within instantiations of pre-defined architectures.

*uvp@cs.bgu.ac.il

This mode of existence is very well-suited to problem domains that inhabit their own metric space, like the physical realms of vision and sound. In stark contrast to these, the textual form of linguistic communication is built atop a discrete alphabet and hinges on notions such as symbolic semantics, inconsistent compositionality, and the arbitrariness of the sign (de Saussure, 1916). The example in (1) exhibits all of these: the symbol *dog* refers to two distinct objects bearing no semantic resemblance; *large* and *white* each describe the (canine) dog’s physical properties, while *dining* categorizes the table based on its function, and *hot* does not modify (the second) *dog* at all, but rather joins it to denote a distinctive atomic concept.

(1) *The large white dog ate the hot dog left on the dining table.*

Given these properties of language, it is far from straightforward to decide the means by which to transform raw text into an input for a neural NLP system tasked with a goal which requires a grasp on the overall communicative intent of the text, such that this initial representation does not lose basic semantics essential to the eventual outcome. This transformation process is known as *embedding*, after which its artifacts are themselves known as **embeddings**, often used synonymously in context with “vectors” or “distributed representations”. Indeed, the choice for default representations has known several shifts within the short DNN era, motivated in part by advances in computational power but also by a collective coming to terms with the limitations of the preceding methods.

The great challenge of representation is compounded by the unboundedness of it all — human concept space is ever-expanding, and each new concept may be assigned an arbitrary sign (e.g., *zoomer*); within an existing concept space, associations capable of inspiring new utterances occupy a combinatorial magnitude which is essentially infinite; and even the form-meaning relationship itself exhibits malleability by humans’ interaction with text input devices and various cognitive biases.¹ Each of these sources of expansion weighs any proposed representational method with the additional burden of generalizing to novel inputs while maintaining consistency in the manner by which they are represented in the system. In the NLP literature, the surface manifestation of the expanding spaces of concept and form, and of the more locally-constrained disparity between text available at different points in time of a model’s training and deployment, is known as the **out-of-vocabulary** problem, and the unseen surface forms themselves are termed OOVs.

In this survey, I consider three central approaches to representing the fundamental units of natural language text in its input stage and the consequences of each approach’s selection on the goals of the systems they are applied in. The first, most popular, and most successful one when used in isolation, is the **distributional** approach where the representation function is trained to embed textual units which

¹As a case in point, over the course of writing this survey I have manually added dozens of new terms to the Overleaf editor’s spell-check dictionary, two in the referring sentence alone.

appear in similar contexts close to each other in vector space. The second is the **compositional** approach which seeks to assemble embeddings for workable textual units by breaking them down into more fundamental elements and applying functions over their own representations, less committed to semantic guarantees. The last is the **relational** approach which makes use of large semantic structures curated manually or in a semi-supervised fashion, leveraging known connections between text and concepts and among concepts in order to create embeddings manifesting humans' notions of "meaning". The **OOV problem** features heavily in the motivation and analysis of the work presented, as it presents challenges to each of the approaches described, yet the exact definition of vocabularies and OOV-ness themselves are challenged by the advent of NLP systems that have become mainstream following the processed described in this work, namely **contextualized subword embeddings**.

2 The Atoms of Language

Natural language is ultimately a system for conveying meaning, information, and social cues from the realm of human experience into a discrete linear form by encoding them as auditory, visual, and/or textual symbols, which are then iteratively composed into more complex units. In order to process such a system's outputs by computational means, it seems fitting to identify those symbols which carry the basic units of meaning, and then find the proper ways to map those meanings into representations for a program which can compose them. The first step, that of identifying linguistic atoms, proves to be a formidable challenge. From the surface output perspective, the common wisdom is that the basic semantic unit of language is what is known as a **morpheme**. The English word *unbelievable*, for example, is composed of a stem morpheme *believe*, a semantic-syntactic suffix *-able* recasting the verb into an adjective pertaining to potential, and a semantic prefix *un-* denoting negation. But this morpheme = atom stipulation is not unassailable. Processes below the morpheme level have been documented across languages, for example the sound symbolism phenomenon known as phonaesthesia, where arbitrary sound patterns correlate with a concept or conceptual properties, such as /gl/ in the English light/shine-related words *glow*, *glitter*, and *glare* (Blake, 2017). Less arbitrarily, patterns and even individual sounds in names are known to evoke semantic qualities based on their acoustic properties (Köhler, 1947; Bergh et al., 1984). In English-language informal communication modes, writers sometimes employ the practice of expressive lengthening, where a single character in a word is repeated in order to amplify its referent's extension on some scale. For example, *loooooong* would be used to describe a particularly long object or period of time. In addition to these sub-morpheme phenomena, the morpheme symbolism and the atoms of our conceptual space relate at neither a univalent nor a one-to-one relation. Certain stem morphemes, like *star*, denote multiple types of concepts or objects (**polysemy** and **homonymy**), while some concepts may be referred to using different morphemes like the relevant meanings of *room* and *space* (**synonymy**). The suffix *-s*

can denote both a third-person present verb or a plural noun (**polyexponence**), and both are replaced by *-es* under certain local conditions (**flexivity**).

Theoretical quibbles notwithstanding, NLP is a practical field, and from its nascence it was clear that finding the most appropriate way to break text down to its purest elements should not set back our efforts to perform sequence-level tasks and develop useful applications. Thus, concessions must be made in the form of selecting a unit easily extractable from text and working with it. This necessity coincides with the reality of having English as the overwhelmingly central target of NLP applications and easiest source of data. The focus on a language with mostly isolating morphology, where morphemes often occupy distinct word forms that are related through sentence-level syntax, conspired with the technical ease of detecting whitespace in text and led to an inevitable starting point for the community in using the **space-delimited word** as the basic unit of text analysis.² The very name of the fundamental bag-of-words approach (BoW) illustrates the implicit synonymy of “word” and “basic unit of representation” in NLP jargon. Although subword- and multiword-level systems were designed and developed outside this paradigm, mostly citing a non-English motivation, when the neural revolution came the predominant methods again anchored the field to the space-delimited word as the atom.

The most obvious advantage of this approach is its simplicity, considering how difficult it is in practice to extract correct sub-word morphemes directly from text. Historically-entrenched orthographic conventions and local-context phonological processes lead to phenomena such as variance in morpheme form at different instantiations, such as the disappearance of the stem’s final *e* in the *unbelievable* or the *s-t* alteration in derivations like *Mars-Martian*, making a deterministic mapping from surface form to morpheme sequence impossible. The lack of overt textual marking of morpheme boundaries (except for the uncommon case of hyphenation) also leads to ambiguous segmentation in words like *unionize*, and the general property of our sound and writing systems’ inventory being relatively small leads to the incidence of affix-identical sequences in single-morpheme words like *reply* (cf. *shortly*) and *bring* (cf. *lying*). Automatic detection of morphemes can be achieved today by unsupervised data-driven systems like Morfessor (Creutz and Lagus, 2002, 2007), which rely on large amounts of training data and provide no guarantee to finding the true morphemes in all cases or downstream applications.

²I will continue throughout to use “space-delimited” to describe a family of simple string tokenization techniques which typically also include minimal heuristics for punctuation separation and a handful of language-specific rules like separating English contractions based on a short closed list, in partial accommodation of the difference between grammatical words and orthographic words (Dixon et al., 2002).

3 Neural Representations

The idea of breaking down concepts in language into numerically-valued axes has played a role in the formation of the modern research landscape in linguistics. Osgood (1952) proposed a low-dimensional space in which nominal objects and concepts are represented by values associated with characteristics which may describe them, such that “eager” and “burning” share a value along the *weak* \leftrightarrow *strong* dimension, while differing along the *cold* \leftrightarrow *hot* dimension. The values were elicited from human subjects.

Scaling this very linguistically-motivated approach manually over an entire language is at the very least impractical, and over the years some relaxations of this scheme to define representations for words which are **distributed** along dimensions gave rise to more automation-friendly processing techniques. Most crucial was the realization that the individual dimensions in the representation space do not have to be meaningful in and of themselves. Liberating the dimensions from their labels allowed the number of dimensions to be governed by concerns of data availability and computational memory and power, rather than by the precision of our semantic theory and ontological thoroughness; it allows for the discovery of unnamed but possibly useful similarities and distinctions between concepts; and it “leaves room” for new properties to be learned if, for example, a domain shift occurs during the process of applying an embedding-based system to a downstream task.

Embedding concepts into a “blank” vector space using learning methods turns the implied causal direction that motivated Osgood’s framework on its head: instead of creating the embeddings based on what we know about language and the relations between concepts, the latter become the proxy target by which we can measure whether or not the embeddings learned by our model are useful to us. Starting with an arbitrary metric space with well-known properties such as \mathbb{R}^d becomes a great advantage, as the space comes with metrics and operations which are easy to conceptualize and imagine as the necessary proxies.³ As the formative instance of this realization served the ability to score the relative directionality of two vectors using the cosine similarity function, which can be compared to annotations in word similarity resources such as WordSim-65 (Rubenstein and Goodenough, 1965), where human subjects were asked to score word pairs without the hassle of decomposing them into their semantic properties first. Metric space also affords the intuitive parallelogram metaphor of word analogy, haunting every introductory text and presentation on embeddings with the equation $\text{king} - \text{man} + \text{woman} \approx \text{queen}$.

³One heroic departure from the shackles of euclidean space is the line of work on embeddings in hyperbolic space (Nickel and Kiela, 2017), touted as a more suitable representation framework for hierarchical structures, including the semantic structure of a language.

4 Distributional Semantics

The development of the distributed view of representation for linguistic objects accompanied the rise of methodologies making use of the distributional hypothesis, traditionally attributed to Harris (1954) and framed as “you shall know a word by the company it keeps”. The maximalist interpretation of this adage as “a word is defined by applying a combination function to the set of its contexts”, used pre-modern-neurally in influential methods such as Brown Clustering (Brown et al., 1992), is an appealing principle to the embedding movement for good reason: breaking words down into contexts provides us with just the distributed fixed dimensions we seek. Once we decide exactly what “context” means to us, we can programatically extract all contexts for all target words given only a corpus, and base our latent dimensions (whose number is limited to hundreds or thousands for practical reasons) on them. The two methods which ended up dominating the distributional embeddings landscape share a definition of context, essentially “words that appear near the target word”, but translate this decision into embedding differently. In SkipGram (Mikolov et al., 2013a), dimension significance is built “bottom-up” from a random initialization and a traversal of the corpus; in GloVe (Pennington et al., 2014), dimensions are the result of an implicit reduction of the full $V \times V$ co-occurrence matrix, where V is the number of words in our vocabulary. The former approach was inspired by early embedding systems (Bengio et al., 2003) developed around the task of language modeling, which is defined with an expectation based in distributional signals, while the latter has origins in latent semantic analysis (LSA; Deerwester et al., 1990). Evaluation on intrinsic tasks such as similarity datasets and analogy benchmarks (e.g., Finkelstein et al., 2001; Mikolov et al., 2013b; Hill et al., 2015) cemented distributional word embeddings as the representation go-to and an accessible replacement to one-hot encodings for a host of applications, while performance on **downstream** tasks within deep learning systems advanced the understanding of the utility that **pre-training** can afford end-to-end systems which include an embedding layer (Collobert and Weston, 2008; Collobert et al., 2011).

5 Out-of-Vocabulary Words

The choice of space-delimited words as the basic unit for representation, and the large resource investment necessary to pre-train a distributional model over a large corpus, in both money and time, create a situation where vectors can mostly be trusted **as long as the words they represent are present in the pre-training corpus**. The models so far discussed have no intrinsic ability to represent words not present in their lookup table, or out-of-vocabulary, or **OOVs** (Brill, 1995; Brants, 2000; Plank, 2016; Heigold et al., 2017; Young et al., 2018). Empirical analyses such as the one in Pinter et al. (2017) show that indeed, the overwhelming majority of downstream datasets contain words not present in the pre-training corpora. Pinter et al. (2020a) present a diachronical

dataset showcasing the volume of novel terms entering a large, steady daily publication in English over time; but even a snapshot of a language at a given moment contains unlimited domain-specific terms, morphological derivations, named entities, potential loanwords, typographical errors, and other sources of OOVs which would appear very reasonably in text analysis tasks and which the downstream model should be given the faculty to handle. In fact, according to Kornai (2002), statistical reasoning leads us to conclude that languages have an infinite vocabulary. But even if a language’s word set were finite, and all present in some corpus, practical memory and lookup constraints would still limit embedding tables to non-exhaustive vocabularies.

To overcome the intrinsic limits of corpus-learned embedding tables, the distributional system has begotten some heuristics that try and initialize embeddings for OOVs beyond the trivial random initialization fallback. If one were to stay true to Firth’s maxim, one possible strategy would be to keep SkipGram’s context embedding table as well as the main table (for “target” words), and initialize OOV embeddings based on the context in which they are first encountered (Horn, 2017). This approach has not caught on, and instead most practitioners took to the use of a special <UNK> embedding, named as an abbreviation of *unknown* (Bengio et al., 2003). In a pre-training stage, such an embedding is learned by replacing a small percentage of the corpus with a dedicated <UNK> token, thus gaining at least some prior for an initialization, in some sense an average over possible contexts for encountering *any* word. This approach is brutally simplistic; it assumes not only that all novel words are representable using the same approximation technique, but that they are all *exactly the same*. The first assumption alone is easy to dispute: a careful observation of any taxonomy of word formation processes (Lieber, 2005; Plag, 2018) suggests that embedding new words into an existing space must involve considering multiple approaches in parallel.

- Words created by processes at the multi-word level, such as compounding or blending, require means of extracting the underlying constructed words and composing the semantic contribution from each word. For example, *brunch* is a blend of *breakfast* and *lunch*; a reasonable initial embedding can be the mean vector for these two words, hopefully keeping it at a high similarity with other meals and the appropriate time of day.
- Words that are inflections of known words, for example *ameliorating*, can benefit from a morphological analysis which finds its stem and syntactic suffix, placing the new vector at the sum of the verb *ameliorate* and the generalized notion of *-ing* verbs, if one is realized in the embedding space (arguably, in a good space it should at least be reliably extractable).
- Novel named entities such as *Lyft* or *SARS-COV-2*, more often than not, reflect arbitrary naming practices and cultural primitives, and even recognition of their type (person / organization / location, etc.) might well be impossible without access to knowledge bases covering the appropriate domain, noting explicitly where in concept space the novel word should be embedded.

- Some OOVs are the result of unpredictable subword processes such as typographical errors (typos) and stylistic variation, like the aforementioned expressive lengthening. In such cases, it is sometimes best to opt out of creation of a new embedding at all and simply map the new form to the existing embedding of its intended canonical word form. This choice will depend on the intended application; in certain cases like sentiment analysis, the stylistic information itself is essential.
- Loanwords like *vespa* originate in a different language than the one the embedding was produced for, but in some cases we have access to an embedding space for the origin language and a function which translates between the two languages' space. A system which can detect the word and its origin, perhaps overcoming processes like writing-system transliteration and phonological adaptation, can start by embedding the target language word in a position projected from the source language's embedding for the equivalent word form.

This is not a comprehensive list. More types of novel words are identified in [Pinter et al. \(2020a\)](#), and not all suggestions in the taxonomy above correspond to actual existing work. Limiting this discussion to a strict interpretation of written-form uniqueness also prevents us from considering as OOVs concepts which are spelled in the same way as other words, either by chance (homography, for example *row* as a line or a fight), by naming (e.g., *Space Force*), or by processes such as zero-derivation (the verb *smoke*, derived from the noun). In languages other than English, some OOV-creating forces may be more dominant in word formation than in English. Morphologically-rich languages, as one edge case, feature large percentages of OOVs in novel texts for a given task's text size compared to English, and this property is often compounded by the fact that many of these are low-resource languages, possessing a relatively small corpus-extracted vocabulary to begin with.

The richness and unpredictability of the OOV problem calls for complementing the word representation systems obtained distributionally with additional approaches, which is the focus of this survey.

6 Subword Compositionality

The first approach considered is an attempt to break the space-delimited word paradigm and get at the finer atomic units of meaning, which can then either be used as the fundamental representation layer, or induce better representations at the word level. This perspective, known as the **compositional** approach, is inspired mostly by the cases where insufficient generalizations are made for cases of morphological word formation processes. Under the compositional framework, an ideal representation for *unbelievable* can be obtained by (1) detecting its three morphological components *un-*, *believe*, and *-able*, (2) querying reliable representations learned for each of them, distributionally or

otherwise, and (3) properly assembling them via some appropriate function.⁴

Each of these three steps is a challenge in itself and open to various implementational approaches. Learning representations for subword units is usually done by considering the subword elements in unison with the full word while applying a distributional method (e.g., [Bojanowski et al., 2017](#)), but some have opted for pre-processing the pre-training corpus such that only lemma forms exist as raw text and the other tokens are explicit representations of the morphological attributes attached to each lemma ([Avraham and Goldberg, 2017](#); [Tan et al., 2020](#)), inducing the production of more consistent vocabularies. Others yet leave the learning to the downstream task itself, feeding off the backpropagated signal from the training instances ([Sutskever et al., 2011](#); [Ling et al., 2015](#); [Lample et al., 2016](#); [Garneau et al., 2019](#)); while others train a compositional network based on the word embedding table in an intermediate phase between pre-training and downstream application ([Pinter et al., 2017](#); [Zhao et al., 2018](#)). The composition function from subwords to the word level is also open to many different approaches: prior work has opted for construction techniques as diverse as using the subword strings as one-hot entries to represent the words themselves ([Huang et al., 2013](#)); summing morpheme embeddings to produce word embeddings ([Botha and Blunsom, 2014](#)); traversing a possibly deep morphological parse tree using a recursive neural network ([Luong et al., 2013](#)); positing probabilistic word embeddings for which the morpheme embeddings act as a prior distribution ([Bhatia et al., 2016](#)); side-by-side training of both word-level and character-level modules followed by concatenating the resulting representations, to allow the downstream model to learn from both levels independently and control the interaction terms directly ([Plank et al., 2016](#)); assembling a hierarchical recurrent net that progressively encodes longer portions of text in each layer ([Chung et al., 2019](#)); or dispensing with the word level altogether and just representing text with a single atomic layer of characters or subwords ([Sennrich et al., 2016](#)).

Most challenging of all is the detection of the subwords themselves. As noted above, morphemes are hard to detect from the surface form of a word. For the default setting where no curated resources exist to allow correct morpheme extraction from a word’s form, as is the case in nearly all languages in the world, the mainstream of compositional representation research has centered on the raw character sequence, the unarguable atom of text,⁵ which is used either via direct operation or as a basis for heuristics that define subword units based on statistical objectives. The great advantage of using characters or primitive character n-grams as the atomic

⁴I will use the term **subword** to denote textual units which are largely between the character level and the word level, when no guarantee of their morphological soundness is attempted. In appropriate contexts, this can also denote word-long or character-long elements which are nevertheless obtained by a subword tokenizer.

⁵At least in languages using the Latin script, like English. Chinese text analysis has benefitted from decomposing characters into strokes or radicals; Hebrew and Arabic include diacritical marks that are not character-intrinsic; and elsewhere, treatment of individual bytes from the Unicode representation of characters has also shown merit.

unit for the model (Santos and Zadrozny, 2014; Kim et al., 2016; Wieting et al., 2016; Bojanowski et al., 2017; Peters et al., 2018) is that it rids us of the need to explicitly designate morphemes altogether; the challenge is to still capture the information they convey, somehow. In contrast, heuristically learning a subword vocabulary from information-theoretic notions (Sennrich et al., 2016; Kudo and Richardson, 2018) or character-sequence unigram distribution (Kudo, 2018) may find us many true morphemes, but there is no guarantee of either precision or recall: corpus collection effects are significant in determining the ultimate vocabulary, orthographic norms may still obfuscate many useful generalized morphemes, and many frequent character sequences may enter the subword vocabulary as the result of coincidental quirks. For example, the character sequence *eva* might contribute to the representation of *unbelievable*, passing along signals learned from unrelated words such as *Eva* or *evaluate*. The ever-growing popularity of systems which use such vocabularies in conjunction with the null composition function that ignores sub-word hierarchy and passes the downstream model embeddings corresponding to the raw subword sequence (see §8) prevents any possibility of correcting incorrect subword tokens at the word level: in this scenario, the next processing layer of the model will use the embedding for *eva* as if it were part of the input equally important to a frequent word like *house*.

7 Relational Semantics

Another way to complement distributionally-trained embeddings is to incorporate signals from curated type-level **relational** resources. The prominent category of such resources is semantic graphs, such as WordNet (Fellbaum, 1998) and BabelNet (Navigli and Ponzetto, 2010), which encode the structural qualities of language as a representation of human knowledge. The core goal of semantic graphs is to describe connections between referents in the perceived and conceived world, and to this end they make an explicit distinction between words as character sequences and an internal semantic primitive which we can call **concepts**. Concepts form the chief node type in the semantic graph, connected by individual edges typed into relations such as hypernymy (*elm* “is a” *tree*) or meronymy (*branch* “is part of a” *tree*), as well as linguistic facts about concept names (*shop.verb* “is derivationally related to” *shop.noun*) which make use of the word-form partition of the graph’s node set. In similar vein, relations which straddle the divide between form and function, like synonymy, are extractable from the bipartite subgraph relating word forms and their available meanings.

In the context of language representation, these structures offer a notion of atom-icity stemming from our conceptual primitives, an attractive premise. They may not answer all needs arising from inflectional morphology (since syntactic properties do not explicitly denote concepts) or some of the other word formation mechanisms, but the rich ontological scaffolding offered by the graph and the prospects of assigning separate embeddings for homonyms in a model-supported manner, assuming sense can be disambiguated in usage, seems much “cleaner” than relying on large

corpora and heuristics to statistically extract linguistic elements and their meaning. In addition to this conceptual shift, as it were, the graph structure itself provides a learning signal not present in linear corpus text, relating the basic units to each other through various types of connections and placing all concepts within some quantifiable relation of each other (within each connected component, although lack of any relation path is also a useful signal). The structure can also occupy the place of the fragile judgment-based word similarity and analogy benchmarks, allowing more exact, refined, well-defined relations to be used for both learning the representations and evaluating them. Methods which embed nodes and relations from general graph structures before even considering any semantics attached to individual nodes and edges, like Node2vec (Grover and Leskovec, 2016) and graph convolutional nets (Schlichtkrull et al., 2017), indeed serve as a basis and inspiration for many of the works in this space.

The fundamentally different manner in which the relational paradigm is complementary to the distributional one in contrast with the compositional one has bearing on the OOV problem, which can be viewed from several perspectives. First is the potential of semantic graphs to improve representation of words that are rare or not present in a large corpus used to initialize distributional embeddings. This has proven to be a powerful direction by methods such as retrofitting (Faruqui et al., 2015), where embeddings of related concepts are pushed together in a post-processing learning phase, showcasing WordNet’s impressive coverage of English domain-specific taxonomies such as classical natural sciences. Elsewhere, properly modelling hypernymy, for example, has been found to help understand text with rare words whose hypernyms are well-represented in the pre-training corpus (Shwartz et al., 2017).⁶ Still, semantic graphs provide only a partial solution to the overall goal of OOV impact mitigation, given their limited scope and heavy reliance on expert annotation.

From the other direction, systems relying on semantic graphs for applications such as question answering and dialogue generation are likely to encounter “OOVs” of their own, i.e. words and concepts not present in the underlying graph. Unlike the corpus-OOV problem, which cannot be quantified convincingly without selecting a specific downstream task first, coping with graph-OOVs can be examined through tasks intrinsic to the graph structure itself. One such task is **relation prediction**, where we assume a concept has a known connection with *some* other concept, and need to figure out which one. Depending on our perspective, either the source or target of the relation may be the OOV concept; for example, on first encounter of the concept *indian lettuce*, we wish to know its hypernym from our set of known concepts. This task is also useful for a similar class of graphs known as **knowledge graphs** (KGs),

⁶A tangential but noteworthy approach considers relations that are not curated in large graphs, but rather corpora annotated for inter-word relations such as syntactic dependencies (Madhyastha et al., 2016). Their system creates a mapping between a distributionally-obtained embedding table and one trained on the annotated parses, and generalizes this mapping to words which are now out-of-vocabulary for a further downstream task (e.g., sentiment analysis). In this case, the reference vocabulary (for defining OOV-ness) is not the unsupervised corpus, but rather an intermediate downstream task.

such as Freebase (Bollacker et al., 2008)⁷ and WikiData (Vrandečić and Krötzsch, 2014), which differ from semantic graphs in several aspects. While WordNet curates connections between semantic concepts and dictionary entries, including certain aspects of the physical world (e.g. “an elm is a tree”), KGs focus on real-world entities and often time-sensitive encyclopedic knowledge (e.g. “Satya Nadella is the CEO of Microsoft”). WordNet is a manually-crafted resource created by language and domain experts, whereas many KGs are either crowdsourced or automatically extracted from databases and large text corpora. As a result, KGs are typically disconnected, shallow, and sparse, boasting areas of hubness and areas of isolation; this contrasts with semantic graphs, where systematic connectedness and hierarchy have been observed (Sigman and Cecchi, 2002). KGs are also distinguished by the richness of their relation type variety, in the hundreds or thousands, compared to WordNet’s 18 relation types (including seven pairs of relations reciprocal to each other). Nevertheless, much of the work on the relation prediction problem has been developed and evaluated on both semantic and knowledge graphs, as well as on derived tasks like **graph completion**, where the entirety of a node’s connections are to be inferred at once, imitating real-world scenarios of knowledge discovery.

Over the years, distributional methods have been used to feed increasingly complex neural nets predicting relations by embedding both concept nodes and relation edges based on corpus-trained tables, to a large degree of success (e.g. Nickel et al., 2011; Socher et al., 2013; Bordes et al., 2013; Yang et al., 2014; Toutanova and Chen, 2015; Neelakantan et al., 2015; Ji et al., 2015; Shi and Wenginger, 2017; Dettmers et al., 2018; Nathani et al., 2019). The basic idea calls for embedding concepts into a metric space and modeling relations by some operator that induces a score for an embedding pair input, either by translating the concept vectors, combining them via bilinear operators, projecting them onto a “scoring scale”, or designing an intricate deep system that finds complex relationships. While these systems achieve impressive results, they all build on an implicit assumption that relation prediction is a **strictly local** task: the fit of an edge can be estimated from the nodes it connects and the intended label alone. In KGs, where structure is of secondary concern, this assumption may go a long way before its limitations stress out performance; in the much more structure-crucial semantic graphs, it is increasingly likely that connections are predicted which should not be permissible from enforceable structural constraints alone, e.g. that the hypernym graph cannot contain cycles. Some systems indeed go beyond the individual edge to embed and predict relations, for example the idea of a path prediction task (Guu et al., 2015) which demands more structure reliance, or embedding methods leveraging local neighborhoods of relation interactions and automatic detection of relations from syntactically parsed text in an iterative manner (Riedel et al., 2013; Toutanova et al., 2015; Schlichtkrull et al., 2017). Others have constructed prediction models where an adversary produces examples which violate structural constraints such as symmetry and transitivity (Minervini and Riedel, 2018). Pinter and Eisenstein (2018) present a sys-

⁷Now defunct.

tem which improves WordNet prediction by augmenting the distributionally-obtained signal with features (motifs) representing the global structure of the semantic edifice. In addition to the task benefit, the emerging feature weights lead to discovery of some general properties of English semantics.

8 Contextualized Representations

Recent developments in NLP have brought about a shift in the balance depicted so far with respect to the atomic level chosen to represent language in applications and the approaches taken to create these representations. Advances in multi-task learning and transfer learning, both in non-neural NLP and in non-NLP deep methods, matured well enough to allow deep NLP to use them effectively as well. The increase of available computation power and the extreme utility found to lie in recurrent nets, most notably the Long Short-Term Memory cell (LSTM; Hochreiter and Schmidhuber, 1997), led to a series of works suggesting the incorporation of instance-specific context into the feature extraction part of a model, before applying any task-specific elements, beginning with simple prediction tasks (Melamud et al., 2016), followed by near-full coverage of core NLP (Peters et al., 2018). The next step was to continue training the shared-architecture context learner, which we can now safely call a language model, during the downstream step, in a process known as fine-tuning (Howard and Ruder, 2018). Design and processing power considerations, but also downstream performance, fueled the shift (Radford et al., 2018) from recurrent net infrastructure to transformer models (Vaswani et al., 2017), which in turn facilitated another major conceptual innovation where autoregressive token prediction was replaced by masked language modeling, where sequence-medial tokens are hidden from the representation layer and must be predicted based on the remaining context (Devlin et al., 2019; Liu et al., 2019). Throughout this evolution, one main principle remained stable: the language prediction task acts as the pre-training step, providing a scaffolding model which is capable of representing tokens within a sequence at a level of effectiveness that allows downstream tasks to begin training with meaningful **contextualized** representations. The heart of contextualization lies in the distributional approach.

The design of these pre-training tasks meant they can no longer tolerate OOV tokens at the rate encountered by static embedding algorithms, as that might render the models unusable for any words that appear in context with OOVs downstream, rather than just the OOVs themselves. On the other hand, the prediction layer creates a computational bottleneck which scales with the size of the vocabulary, since every token must be available for prediction at all model steps. Therefore, these models resorted to compositional techniques for the bottom layer where the input sequence is processed into tokens. The character convolution net selected for ELMo (Peters et al., 2018) did not gain traction, possibly because it didn't provide an adequate method for predicting text from the output layer, and so subsequent models, particularly those relying on transformers, operate over a sequence of equal-status tokens, each

representing a word or a subword, from a mid-size vocabulary (tens of thousands) built in a pre-training phase using statistical heuristic techniques mentioned in §6. These models inherit the problems endemic to these methods like inadequacy for certain OOV classes, morphological unsoundness, and length-imbalance; as well as issues like the added burden they impose on already limited-length token sequences. Common wisdom seems to hold that they make up for these shortcomings within the depths of their fully-connected transformer layers, and end up with satisfactory top-layer representations. Recent work challenging these models with truly novel word forms suggest otherwise (Pinter et al., 2020a,b), while work on either incorporating the compositional signal into subword-vocabulary transformers (Ma et al., 2020; Aguilar et al., 2020; El Boukkouri et al., 2020; Pinter et al., 2021), or replacing the subwords with characters or bytes altogether (Clark et al., 2021; Xue et al., 2021), is rapidly gaining traction as well.

Acknowledgments

This survey is an adapted version of the introduction my PhD thesis. I thank my committee for helping to shape it: my advisor, Jacob Eisenstein; Mark Riedl, Dan Roth, Wei Xu, and Diyi Yang.

References

- Aguilar, G., McCann, B., Niu, T., Rajani, N., Keskar, N., and Solorio, T. (2020). Char2subword: Extending the subword embedding space from pre-trained models using robust character compositionality. *arXiv preprint arXiv:2010.12730*.
- Avraham, O. and Goldberg, Y. (2017). The interplay of semantics and morphology in word embeddings. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 422–426, Valencia, Spain. Association for Computational Linguistics.
- Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. (2003). A neural probabilistic language model. *The journal of machine learning research*, 3:1137–1155.
- Bergh, B. G. V., Collins, J., Schultz, M., and Adler, K. (1984). Sound advice on brand names. *Journalism Quarterly*, 61(4):835–840.
- Bhatia, P., Guthrie, R., and Eisenstein, J. (2016). Morphological priors for probabilistic neural word embeddings. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 490–500, Austin, Texas. Association for Computational Linguistics.
- Blake, B. J. (2017). Sound symbolism in english: Weighing the evidence. *Australian Journal of Linguistics*, 37(3):286–313.

- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2017). Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Bollacker, K., Evans, C., Paritosh, P., Sturge, T., and Taylor, J. (2008). Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250.
- Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., and Yakhnenko, O. (2013). Translating embeddings for modeling multi-relational data. In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 26*, pages 2787–2795. Curran Associates, Inc.
- Botha, J. A. and Blunsom, P. (2014). Compositional morphology for word representations and language modelling. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Brants, T. (2000). Tnt: a statistical part-of-speech tagger. In *Proceedings of the sixth conference on Applied natural language processing*, pages 224–231. Association for Computational Linguistics.
- Brill, E. (1995). Transformation-based error-driven learning and natural language processing: A case study in part-of-speech tagging. *Computational Linguistics*, 21(4):543–565.
- Brown, P. F., Della Pietra, V. J., deSouza, P. V., Lai, J. C., and Mercer, R. L. (1992). Class-based n -gram models of natural language. *Computational Linguistics*, 18(4):467–480.
- Chung, J., Ahn, S., and Bengio, Y. (2019). Hierarchical multiscale recurrent neural networks. In *5th International Conference on Learning Representations, ICLR 2017*.
- Clark, J. H., Garrette, D., Turc, I., and Wieting, J. (2021). Canine: Pre-training an efficient tokenization-free encoder for language representation. *arXiv preprint arXiv:2103.06874*.
- Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(76):2493–2537.
- Creutz, M. and Lagus, K. (2002). Unsupervised discovery of morphemes. In *Proceedings of the ACL-02 Workshop on Morphological and Phonological Learning*, pages 21–30. Association for Computational Linguistics.

- Creutz, M. and Lagus, K. (2007). Unsupervised models for morpheme segmentation and morphology learning. *ACM Transactions on Speech and Language Processing (TSLP)*, 4(1):1–34.
- de Saussure, F. (1916). *Cours de linguistique générale* (roy harris, trans.). London: Duckworth.
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American society for information science*, 41(6):391–407.
- Dettmers, T., Pasquale, M., Pontus, S., and Riedel, S. (2018). Convolutional 2d knowledge graph embeddings. In *Proceedings of the 32th AAAI Conference on Artificial Intelligence*.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Dixon, R. M., Aikhenvald, A. Y., et al. (2002). Word: a typological framework. *Word: A cross-linguistic typology*, 1:41.
- El Boukkouri, H., Ferret, O., Lavergne, T., Noji, H., Zweigenbaum, P., and Tsujii, J. (2020). CharacterBERT: Reconciling ELMo and BERT for word-level open-vocabulary representations from characters. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6903–6915, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Faruqui, M., Dodge, J., Jauhar, S. K., Dyer, C., Hovy, E., and Smith, N. A. (2015). Retrofitting word vectors to semantic lexicons. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1606–1615, Denver, Colorado. Association for Computational Linguistics.
- Fellbaum, C. (1998). *WordNet: An Electronic Lexical Database*. Bradford Books.
- Finkelstein, L., Gabrilovich, E., Matias, Y., Rivlin, E., Solan, Z., Wolfman, G., and Ruppin, E. (2001). Placing search in context: The concept revisited. In *Proceedings of the 10th international conference on World Wide Web*, pages 406–414.
- Garneau, N., Leboeuf, J.-S., Pinter, Y., and Lamontagne, L. (2019). Attending form and context to generate specialized out-of-vocabulary words representations. *arXiv preprint arXiv:1912.06876*.
- Grover, A. and Leskovec, J. (2016). node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864.

- Guu, K., Miller, J., and Liang, P. (2015). Traversing knowledge graphs in vector space. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 318–327, Lisbon, Portugal. Association for Computational Linguistics.
- Harris, Z. S. (1954). Distributional structure. *Word*, 10(2-3):146–162.
- Heigold, G., Neumann, G., and van Genabith, J. (2017). How robust are character-based word embeddings in tagging and mt against wrod scrambling or randdm nouse? *arXiv preprint arXiv:1704.04441*.
- Hill, F., Reichart, R., and Korhonen, A. (2015). SimLex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4):665–695.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Horn, F. (2017). Context encoders as a simple but powerful extension of word2vec. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 10–14, Vancouver, Canada. Association for Computational Linguistics.
- Howard, J. and Ruder, S. (2018). Universal language model fine-tuning for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339, Melbourne, Australia. Association for Computational Linguistics.
- Huang, P.-S., He, X., Gao, J., Deng, L., Acero, A., and Heck, L. (2013). Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, pages 2333–2338.
- Ji, G., He, S., Xu, L., Liu, K., and Zhao, J. (2015). Knowledge graph embedding via dynamic mapping matrix. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 687–696, Beijing, China. Association for Computational Linguistics.
- Kim, Y., Jernite, Y., Sontag, D., and Rush, A. M. (2016). Character-aware neural language models. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*.
- Köhler, W. (1947). *Gestalt psychology*, 2nd edn new york. NY: Liveright Publishing Corporation.[Google Scholar].
- Kornai, A. (2002). How many words are there? *Glottometrics*, 4:61–86.
- Kudo, T. (2018). Subword regularization: Improving neural network translation models with multiple subword candidates. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 66–75, Melbourne, Australia. Association for Computational Linguistics.

- Kudo, T. and Richardson, J. (2018). SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., and Dyer, C. (2016). Neural architectures for named entity recognition. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 260–270, San Diego, California. Association for Computational Linguistics.
- Lieber, R. (2005). English word-formation processes. In *Handbook of word-formation*, pages 375–427. Springer.
- Ling, W., Dyer, C., Black, A. W., Trancoso, I., Fernandez, R., Amir, S., Marujo, L., and Luís, T. (2015). Finding function in form: Compositional character models for open vocabulary word representation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1520–1530, Lisbon, Portugal. Association for Computational Linguistics.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Luong, T., Socher, R., and Manning, C. (2013). Better word representations with recursive neural networks for morphology. In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 104–113, Sofia, Bulgaria. Association for Computational Linguistics.
- Ma, W., Cui, Y., Si, C., Liu, T., Wang, S., and Hu, G. (2020). CharBERT: Character-aware pre-trained language model. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 39–50, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Madhyastha, P. S., Bansal, M., Gimpel, K., and Livescu, K. (2016). Mapping unseen words to task-trained embedding spaces. In *Proceedings of the 1st Workshop on Representation Learning for NLP*, pages 100–110, Berlin, Germany. Association for Computational Linguistics.
- Melamud, O., Goldberger, J., and Dagan, I. (2016). context2vec: Learning generic context embedding with bidirectional LSTM. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 51–61, Berlin, Germany. Association for Computational Linguistics.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. In *Proceedings of International Conference on Learning Representations*.

- Mikolov, T., Yih, W.-t., and Zweig, G. (2013b). Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 746–751, Atlanta, Georgia. Association for Computational Linguistics.
- Minervini, P. and Riedel, S. (2018). Adversarially regularising neural NLI models to integrate logical background knowledge. In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 65–74, Brussels, Belgium. Association for Computational Linguistics.
- Nathani, D., Chauhan, J., Sharma, C., and Kaul, M. (2019). Learning attention-based embeddings for relation prediction in knowledge graphs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4710–4723, Florence, Italy. Association for Computational Linguistics.
- Navigli, R. and Ponzetto, S. P. (2010). BabelNet: Building a very large multilingual semantic network. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 216–225, Uppsala, Sweden. Association for Computational Linguistics.
- Neelakantan, A., Roth, B., and McCallum, A. (2015). Compositional vector space models for knowledge base inference. In *2015 aai spring symposium series*.
- Nickel, M. and Kiela, D. (2017). Poincaré embeddings for learning hierarchical representations. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6341–6350.
- Nickel, M., Tresp, V., and Kriegel, H.-P. (2011). A three-way model for collective learning on multi-relational data. In *ICML*, volume 11, pages 809–816.
- Osgood, C. E. (1952). The nature and measurement of meaning. *Psychological bulletin*, 49(3):197.
- Pennington, J., Socher, R., and Manning, C. (2014). GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Pinter, Y. and Eisenstein, J. (2018). Predicting semantic relations using global graph properties. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1741–1751, Brussels, Belgium. Association for Computational Linguistics.

- Pinter, Y., Guthrie, R., and Eisenstein, J. (2017). Mimicking word embeddings using subword RNNs. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 102–112, Copenhagen, Denmark. Association for Computational Linguistics.
- Pinter, Y., Jacobs, C. L., and Bittker, M. (2020a). NYTWIT: A dataset of novel words in the New York Times. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6509–6515, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Pinter, Y., Jacobs, C. L., and Eisenstein, J. (2020b). Will it unblend? In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1525–1535, Online. Association for Computational Linguistics.
- Pinter, Y., Stent, A., Dredze, M., and Eisenstein, J. (2021). Learning to look inside: Augmenting token-based encoders with character-level information. *arXiv preprint arXiv:2108.00391*.
- Plag, I. (2018). *Word-formation in English*. Cambridge University Press.
- Plank, B. (2016). What to do about non-standard (or non-canonical) language in nlp. *arXiv preprint arXiv:1608.07836*.
- Plank, B., Søgaard, A., and Goldberg, Y. (2016). Multilingual part-of-speech tagging with bidirectional long short-term memory models and auxiliary loss. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 412–418, Berlin, Germany. Association for Computational Linguistics.
- Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. (2018). Improving language understanding by generative pre-training. *Technical report, OpenAI*.
- Riedel, S., Yao, L., McCallum, A., and Marlin, B. M. (2013). Relation extraction with matrix factorization and universal schemas. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 74–84, Atlanta, Georgia. Association for Computational Linguistics.
- Rubenstein, H. and Goodenough, J. B. (1965). Contextual correlates of synonymy. *Communications of the ACM*, 8(10):627–633.
- Santos, C. D. and Zadorozny, B. (2014). Learning character-level representations for part-of-speech tagging. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1818–1826.
- Schlichtkrull, M., Kipf, T. N., Bloem, P., van den Berg, R., Titov, I., and Welling, M. (2017). Modeling relational data with graph convolutional networks. *stat*, 1050:17.

- Sennrich, R., Haddow, B., and Birch, A. (2016). Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Shi, B. and Weninger, T. (2017). Proje: Embedding projection for knowledge graph completion. In *AAAI*, volume 17, pages 1236–1242.
- Shwartz, V., Santus, E., and Schlechtweg, D. (2017). Hypernyms under siege: Linguistically-motivated artillery for hypernymy detection. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 65–75, Valencia, Spain. Association for Computational Linguistics.
- Sigman, M. and Cecchi, G. A. (2002). Global organization of the wordnet lexicon. *Proceedings of the National Academy of Sciences*, 99(3):1742–1747.
- Socher, R., Chen, D., Manning, C. D., and Ng, A. (2013). Reasoning with neural tensor networks for knowledge base completion. In *Advances in Neural Information Processing Systems*, pages 926–934.
- Sutskever, I., Martens, J., and Hinton, G. E. (2011). Generating text with recurrent neural networks. In *ICML*.
- Tan, S., Joty, S., Varshney, L., and Kan, M.-Y. (2020). Mind your inflections! Improving NLP for non-standard Englishes with Base-Inflection Encoding. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5647–5663, Online. Association for Computational Linguistics.
- Toutanova, K. and Chen, D. (2015). Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd Workshop on Continuous Vector Space Models and their Compositionality*, pages 57–66, Beijing, China. Association for Computational Linguistics.
- Toutanova, K., Chen, D., Pantel, P., Poon, H., Choudhury, P., and Gamon, M. (2015). Representing text for joint embedding of text and knowledge bases. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1499–1509, Lisbon, Portugal. Association for Computational Linguistics.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Vrandečić, D. and Krötzsch, M. (2014). Wikidata: A free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85.
- Wieting, J., Bansal, M., Gimpel, K., and Livescu, K. (2016). Charagram: Embedding words and sentences via character n-grams. In *Proceedings of the 2016 Conference on*

Empirical Methods in Natural Language Processing, pages 1504–1515, Austin, Texas. Association for Computational Linguistics.

Xue, L., Barua, A., Constant, N., Al-Rfou, R., Narang, S., Kale, M., Roberts, A., and Raffel, C. (2021). Byt5: Towards a token-free future with pre-trained byte-to-byte models. *arXiv preprint arXiv:2105.13626*.

Yang, B., Yih, W.-t., He, X., Gao, J., and Deng, L. (2014). Embedding entities and relations for learning and inference in knowledge bases. *arXiv preprint arXiv:1412.6575*.

Young, T., Hazarika, D., Poria, S., and Cambria, E. (2018). Recent trends in deep learning based natural language processing. *IEEE Computational Intelligence Magazine*, 13(3):55–75.

Zhao, J., Mudgal, S., and Liang, Y. (2018). Generalizing word embeddings using bag of subwords. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 601–606, Brussels, Belgium. Association for Computational Linguistics.