

# Signal in Noise: Exploring Meaning Encoded in Random Character Sequences with Character-Aware Language Models

**Mark Bo Chu**

Columbia University

mbc2165@columbia.edu

**Bhargav Srinivasa Desikan**

École Polytechnique Fédérale de Lausanne

bhargav.srinivasadesikan@epfl.ch

**Ethan O. Nadler**

Carnegie Observatories

University of Southern California

enadler@carnegiescience.edu

**Ruggerio L. Sardo**

Sapienza University of Rome

losardor@gmail.com

**Elise Darragh-Ford**

Stanford University

KIPAC & Department of Physics

edarragh@stanford.edu

**Douglas Guilbeault**

University of California, Berkeley

Haas Business School

douglas.guilbeault@berkeley.edu

## Abstract

Natural language processing models learn word representations based on the distributional hypothesis, which asserts that word context (e.g., co-occurrence) correlates with meaning. We propose that  $n$ -grams composed of random character sequences, or *garble*, provide a novel context for studying word meaning both within and beyond extant language. In particular, randomly generated character  $n$ -grams lack meaning but contain primitive information based on the distribution of characters they contain. By studying the embeddings of a large corpus of garble, extant language, and pseudowords using CharacterBERT, we identify an axis in the model’s high-dimensional embedding space that separates these classes of  $n$ -grams. Furthermore, we show that this axis relates to structure within extant language, including word part-of-speech, morphology, and concept concreteness. Thus, in contrast to studies that are mainly limited to extant language, our work reveals that meaning and primitive information are intrinsically linked.

## 1 Introduction

What primitive information do character sequences contain? Modern natural language processing is driven by the *distributional hypothesis* (Firth, 1957), which asserts that the context of a linguistic expression defines its meaning (Emerson, 2020). Because existing words—which represent an extremely small fraction of the space of possible character sequences—appear in context together, the distributional paradigm at this level is limited in

its ability to study the meaning of and information encoded by arbitrary character level  $n$ -grams (word forms). Furthermore, state-of-the-art computational language models operating within the distributional paradigm, such as BERT (Devlin et al., 2019), are mainly trained on extant words. Yet, a plethora of insights into language learning have emerged from inquiries into language beyond extant words, such as the grammatical errors and inference patterns that children exhibit when distinguishing extant words from non-linguistic auditory signals, including emotional expressions, auditory gestures, and other forms of paralinguistic speech (Yang, 2006; Carey, 2000). We therefore propose that character  $n$ -grams (i.e., sequences of alphabetic characters) outside the space of extant language can provide new insights into the meaning of words and how they are represented by these models, beyond that captured by word and subword-based distributional semantics alone. We explore this by studying the embeddings of randomly generated character  $n$ -grams (referred to as *garble*), which contain primitive communicative information but are devoid of meaning, using the CharacterBERT model (El Boukkouri et al., 2020). Such randomly generated character  $n$ -grams are textual analogues of paralinguistic vocalizations—vocal extra-speech sounds and noises.

Our analyses contribute to the growing understanding of BERTology (Rogers et al., 2020) by identifying a dimension, which we refer to as the *information axis*, that separates extant and garble  $n$ -grams. This finding is supported by a Markov

model that produces a probabilistic information measure for character  $n$ -grams based on their statistical properties. Strikingly, this information dimension correlates with properties of extant language; for example, parts of speech separate along the information axis, and word concreteness varies along a roughly orthogonal dimension in our projection of CharacterBERT embedding space. Although the information axis we identify separates extant and randomly generated  $n$ -grams very effectively, we demonstrate that these classes of  $n$ -grams mix into each other in detail, and that *pseudowords*—i.e., phonologically coherent character  $n$ -grams without extant lexical meaning—lie between the two in our CharacterBERT embeddings.

This paper is organized as follows. We first discuss concepts from natural language processing, information theory, and linguistics relevant to our study. We then analyse CharacterBERT representations of extant and randomly generated character sequences and how the relation between the two informs the structure of extant language, including morphology, part-of-speech, and word concreteness. Finally, we ground our information axis in a predictive Markov language model.

## 2 Modeling $n$ -grams Beyond Extant Language

Models in computational linguistics often represent words in a high-dimensional embedding space based on their co-occurrence patterns according to the distributional hypothesis (Landauer and Dumais, 1997; Mikolov et al., 2013). Embeddings that capture the semantic content of extant words are used for many natural language applications, including document or sentence classification (Kowsari et al., 2019), information retrieval and search (Mitra et al., 2018), language modelling and translation (Devlin et al., 2019), language generation (Brown et al., 2020), and more (Jurafsky and Martin, 2021). In these cases, vector operations performed on word embeddings are used for higher-level tasks such as search or classification.

Word embeddings have largely concerned themselves with extant language—that is, commonly used words which carry consistent meaning—and thus cannot represent character  $n$ -grams outside of this space. The few models that encompass *character  $n$ -grams*, which naturally include  $n$ -grams beyond extant words, often use RNNs (Mikolov et al., 2010) or encoder-decoder architectures (Sutskever

et al., 2014) to represent character-level sequences. In parallel, the ubiquitous use of Transformer models has led to studies of their inner representations, weights, and attention mechanism (Rogers et al., 2020; Clark et al., 2019). Most Transformer models are trained using extant words and sub-words, largely focusing on their semantics and syntax; however, some recent models operate at the character level, such as CharacterBERT (El Boukkouri et al., 2020) and CharBERT (Ma et al., 2020). Strikingly, character-level models excel at character-level tasks (e.g., spelling correction; Xie et al. 2016; Chollampatt and Ng 2018) and perform comparably to word-level models at language-modelling tasks (Kim et al., 2016).

Character-level models are therefore an ideal tool for studying the information and meaning encoded in  $n$ -grams beyond the realm of extant language. Given that the current state-of-the-art is driven by Transformer-based models, throughout our study, we use the CharacterBERT model. CharacterBERT is uniquely suited for our study as it uses a CharacterCNN module (Peters et al., 2018) to produce single embeddings for any input token, built as a variant to BERT which relies on sub-word tokenization (El Boukkouri et al., 2020).

## 3 Primitive Information and Meaning Beyond Extant Language

Before presenting our results, we discuss general characteristics of the space beyond extant words; we reiterate that this space is missed by word and sub-word-based models. Due to CharacterBERT’s use of English characters, we restrict our analysis to English character  $n$ -grams, and we study the properties of CharacterBERT embeddings including English-based  $n$ -grams outside of extant language. By studying CharacterBERT’s representations of meaning encoded in  $n$ -grams that do not appear in consistent (or any) context in its training data, our framework goes beyond the traditional distributional hypothesis paradigm. In this way, we seek to understand core properties of information encoded in  $n$ -grams beyond their lexicalized semantics by simultaneously studying  $n$ -grams that contain different types of information.<sup>1</sup>

We use randomly generated character sequences to create  $n$ -grams that contain primitive informa-

<sup>1</sup>In analogy, the theory of ensemble perception in developmental psychology offers a framework to understand the human ability to understand the ‘gist’ of multiple objects at once (Sweeny et al., 2015).

tion but no meaning. We adapt Marr’s notion of primitive visual information for primitive textual information (Marr and Hildreth, 1980), and make the analogue between vision and language because information is substrate independent (Deutsch and Marletto, 2015). In our case, primitive textual information is lower-level communicative information which is present in both text with and without meaning. Being textual, our randomly generated  $n$ -grams are not bound by the constraints of human speech, and may be phonologically impossible; these garble  $n$ -grams may be seen as an example of textual noise.

In the following subsections, we provide three examples of language—distorted speech, paralinguistic, and pseudowords—which motivate our study of character-level embeddings for randomly generated character  $n$ -grams. We then describe the complementary information encoded by word morphology.

### 3.1 Distorted Speech

In popular use, “garble” refers to a message that has been distorted (garbled), such as speech where meaning is corrupted by phonological distortions. For example, the phrase “reading lamp” may become “eeling am” when garbled. Garbled speech contains lesser, or zero, meaning compared to ungarbled speech, but the signal of speech media is nonetheless present as information, which according to Shannon (1951) may contain no meaning at all. Garbled speech satisfies the classical five-part definition of communication provided by Shannon (2001); an *information source* (speaker) can *transmit* (verbalize) an informationally primitive message through the *channel* of speech media through the *receiver* (ears) to the *destination* (listener).

### 3.2 Paralinguistic

Paralinguistic vocalizations are specifically identifiable sounds beyond the general characteristics of speech (Noth, 1990) and present another example of communication beyond lexicalized semantics. Paralinguistic vocalizations include *characterizers*, like moaning; and *segregates*, like “uh-huh” for affirmation. The border between such paralinguistic vocalizations and lexicalized interjections with defined meanings is “fuzzy” (Noth, 1990).

### 3.3 Pseudowords

Pseudowords are phonologically possible character  $n$ -grams without extant lexical meaning. Word-

likeness judgments reveal that human distinctions between pseudowords and phonologically impossible nonwords are gradational (Needle et al., 2020). As a unique informational class, pseudowords have been used in language neuronal activation studies (Price et al., 1996), infant lexical-semantic processing (Friedrich and Friederici, 2005), in poetry through nonsense (Ede, 1975), and in literary analyses (Lecerle, 2012). Pseudowords can also elicit similar interpretations and associations across independent participants (Davis et al., 2019a).

To consider pseudowords generatively, it is helpful to note that an alphabetic writing system covers not only every word but every possible word in its language (Deutsch, 2011); pseudowords can thus be thought of as possible-but-uninstantiated (counterfactual) extant words—e.g., “cyberspace” was a pseudoword before the internet. We embed randomly generated pseudowords into our model to study their information content and relation to both extant words and randomly generated  $n$ -grams.

### 3.4 Morphology

Morphology deals with the systems of natural language that create words and word forms from smaller units (Trost, 1992). Embedding spaces and the distributional hypothesis offer insights into the relationship between character combination, morphology and semantics. Notably, morphological irregularities complicate the statistics of global character-level findings in the embedding space, like through *suppletion*—where word forms change idiosyncratically e.g. *go*’s past tense is *went*, or *epenthesis*—where characters are inserted under certain phonological conditions e.g. fox pluralizes as *foxes* (Trost, 1992); so too do the multiple ‘correct’ spellings of pseudowords under conventional phoneme-to-grapheme mappings (Needle et al., 2020). Distinctions between morphological phenomena can also be hard to define; for example, the boundary between derivation and compounding is “fuzzy” (Trost, 1992).

## 4 Character-Level Language Models for Information Analysis

As described above, state-of-the-art language models serve as a tool to study meaning as it emerges through the distributional hypothesis paradigm. Existing work on the analysis of Transformers and BERT-based models have explored themes we are interested in, such as semantics (Ethayarajh, 2019),

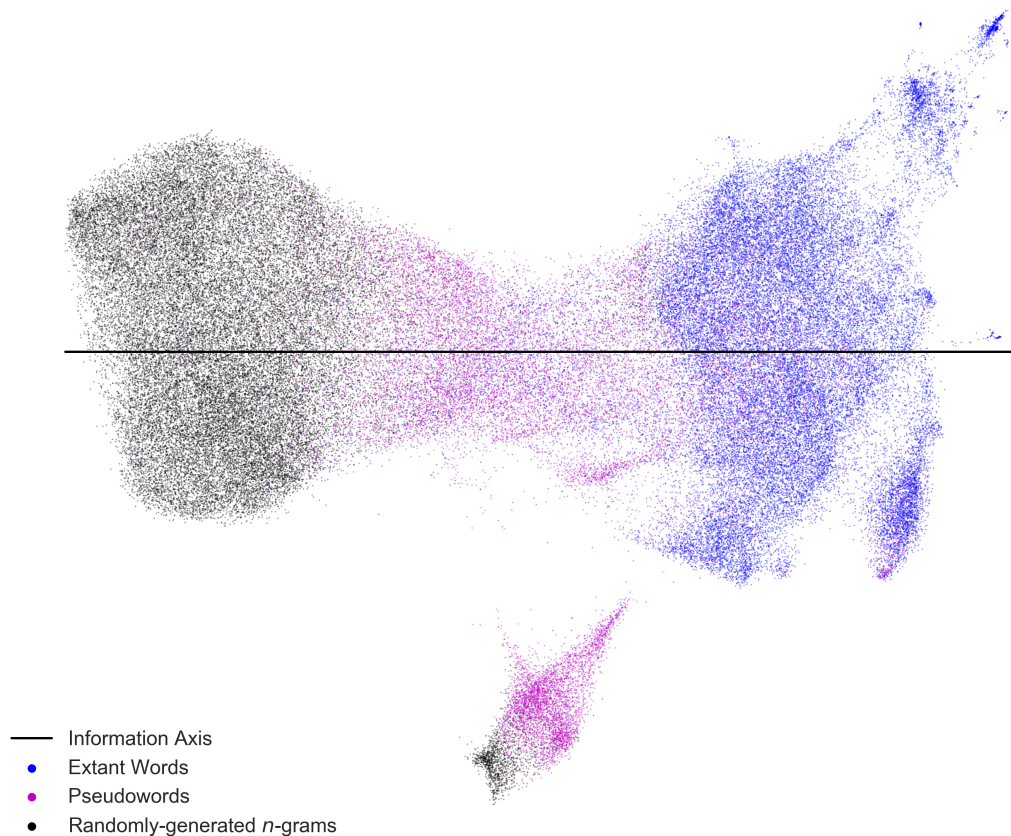


Figure 1: UMAP projection of CharacterBERT embeddings for extant words (blue), pseudowords (magenta), and randomly generated character  $n$ -grams (black). The solid black line shows the information axis that we define in this work. The bottom-most cluster of random and pseudoword character  $n$ -grams is comprised of character  $n$ -grams ending in “s”, and the top-most clusters of extant words are comprised of compound words.

syntax (Goldberg, 2019), morphology (Hofmann et al., 2020, 2021), and the structure of language (Jawahar et al., 2019). However, all of this work limits itself to the focus of extant words due to the word and sub-word-based nature of these models.

We study the structure of the largely unexplored character  $n$ -gram space which includes extant language, pseudowords and garble character  $n$ -grams, seen through the representations created by CharacterBERT, as follows. To explore how the character  $n$ -gram space is structured in the context of character based distributional semantics, we embed 40,000 extant English words, 40,000 randomly generated character  $n$ -grams, and 20,000 pseudowords. We choose the 40,000 most used English words that have been annotated for concreteness/abstractness ratings (Brysaert et al., 2014). Randomly generated character  $n$ -grams are forced to have a string length distribution that matches the corpus of extant words we analyze. To generate pseudowords, we use a popular pseudoword generator.<sup>2</sup>

<sup>2</sup><http://soybomb.com/tricks/words/>

The CharacterBERT (El Boukkouri et al., 2020) general model has been trained on nearly 40 GB of Reddit data using character sequences. We leverage this model to create representations of character  $n$ -grams that may not have been seen in the training data. This allows us to use the resulting 512 dimensional embeddings for exploration via visualisation, topology modelling via distances and projections, and classification error analysis.

#### 4.1 Identifying the Information Axis

To guide our exploration of the high-dimensional topology of the resulting embeddings, we use the UMAP dimensionality reduction technique (McInnes et al., 2018). UMAP creates a low-dimensional embedding by searching for a low-dimensional projection of the data that has the closest possible equivalent fuzzy topological structure as the original representations, thereby preserving both local and global structure. In Appendix A, we demonstrate that our key results are not sensitive to this choice of dimensionality reduction method.

We use the UMAP embeddings to extract an *in-*

Character $n$ -gram type	Information Axis position
Extant	$0.75 \pm 0.12$
Noun	$0.74 \pm 0.12$
Verb	$0.72 \pm 0.09$
Adjective	$0.76 \pm 0.11$
Adverb	$0.87 \pm 0.09$
Pseudoword	$0.50 \pm 0.15$
Random	$0.17 \pm 0.11$

Table 1: Median and standard deviation of minmax-normalized position along the information axis shown in Figure 1, for extant words (including parts of speech), pseudowords, and randomly generated  $n$ -grams.

*formation axis* that captures most variance among extant and randomly generated  $n$ -grams. To assign  $n$ -grams an ‘information axis score,’ we minmax-normalize the UMAP coordinates along this axis. Thus, our information axis establishes a link between extant language and garble, thereby connecting meaning and primitive information. Figure 1 shows how CharacterBERT embeddings of extant, pseudoword, and randomly generated character  $n$ -grams arrange themselves in this space.

#### 4.2 Statistical Properties of $n$ -grams Along the Information Axis

We perform several statistical tests to differentiate between categories of character  $n$ -grams along the information axis. First, Table 1 lists the median and standard deviation of minmax-normalized position along the information axis, demonstrating that extant words, pseudowords, and garble are clearly separated. Note that the scatter within each  $n$ -gram class is much smaller than the distances between classes, indicating that our results are robust to variations in the garble and pseudoword samples.

Next, we use the Kolmogorov-Smirnov (KS; Massey Jr 1951) two-sample test to assess differences between the information axis distributions of our  $n$ -gram classes. All of the KS tests very significantly indicate differences between types of character  $n$ -gram and parts of speech along the information axis ( $p \ll 0.001$ ). Furthermore, the KS statistic score is 0.94 for (extant, random), 0.83 for (pseudoword, random), and 0.70 for (extant, pseudoword), indicating that extant and random  $n$ -grams differ most significantly along the information axis (consistent with Figures 1–2).

#### 4.3 Hyperplane Classifier

The visualisation of the character  $n$ -grams suggests that a hyperplane classifier is suitable for separating

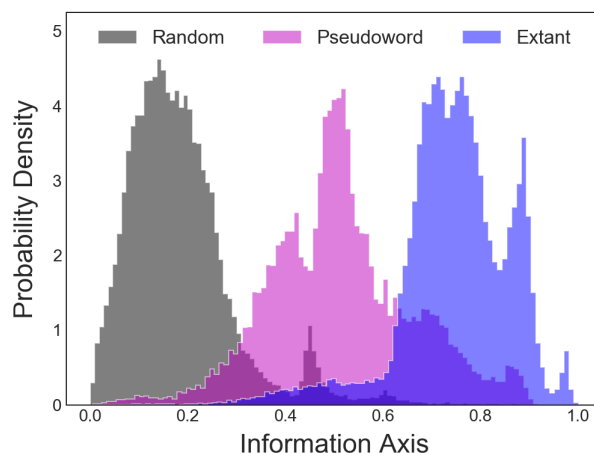


Figure 2: Probability density of CharacterBERT embeddings for extant words (blue), pseudowords (magenta), and randomly generated character  $n$ -grams (black) as a function of minmax-normalized position along the information axis shown in Figure 1.

extant words and garble. We use a support vector machine (Cortes and Vapnik, 1995) trained on half of our 40,000 commonly-used extant words and half of our computer-generated garble to classify unseen extant, garble and pseudoword character  $n$ -grams. We use this method to explore the information axis in the high-dimensional embedding space.

The classifier achieves an accuracy of 98.9% on unseen extant language and garble character  $n$ -grams, suggesting we can learn about the embeddings through error analysis.

In particular, we found similarities among extant words classified as garble. 74.4% (270/363) were compound or derivative words, similar to many extant language terms that lie near the midpoint of the information axis. 19% (69/363) were foreign words like “hibachi” or dialect words like “doohickey.”

The garble classification errors—garble classified as extant language—were in small part due to our randomization method inadvertently creating extant language labelled as garble, accounting for 9.5% (36/377) errors we identify. The garble classified as extant language mostly contained phonologically impossible elements, though some were pseudowords.

When pseudowords were forcibly classified into extant or garble character  $n$ -grams, more pseudowords were classified as extant language than garble (12894 as extant to 7106 as garble). Labelling affirms these intuitions, with pseudowords

like “fought” looking intuitively familiar and being readable. Given CharacterBERT’s massive Reddit training data, typos and localized language may account for the classifier’s tendency to classify pseudowords as extant language. Also, our embedding space only uses the 40,000 most common English words out of 208,000 distinct lexicalized lemma words (Brybaert et al., 2016), which may impact spatial structure if included.

## 5 Structure of Extant Words along the Information Axis

We use this section to discuss the structure of language across the information axis derived from our low-dimensional UMAP space. We structure our analysis across this axis as it organises the relative structure of extant words vs. randomly generated character  $n$ -grams, while also distinguishing internal structure within the extant word space.

### 5.1 Extant vs. Pseudowords vs. Garble

At the scale of global structure, the information axis highlights that extant words are separated from randomly generated character  $n$ -grams (Figure 1). We note that the midpoint of all character  $n$ -gram classes is 0.5 on our information axis. Pseudowords populate the region near the midpoint of the information axis, and also overlap with both extant English and garble character  $n$ -grams (Figure 2). There is no distinct boundary between the three classes of  $n$ -grams, consistent with both morphological descriptions of compound and derivational words and descriptions of paralinguistic as “fuzzy.” This global structure—and the structure internal to extant language (Figure 3)—goes beyond the distributional hypothesis by including  $n$ -grams that do not appear in consistent (or any) contexts, like pseudowords and garble. Pseudowords lie between extant and garble character  $n$ -grams, but there is no distinct boundary between pseudowords and the other classes of  $n$ -grams.

Extant language, pseudoword, and garble regions have different internal structure (Figure 1). The garble region has comparatively less structure than the extant language region, though there is some internal variation, notably a cluster of character  $n$ -grams ending in the character “s” separated from the main garble region. We qualitatively explore the classes of garble and pseudoword embeddings revealed by our analysis in Appendix B, which includes supplementary discussion of the

potential relevance of these findings for linguistic theory.

### 5.2 Parts of Speech and Morphology

In our UMAP projection, detailed structure emerges for extant words split by part-of-speech (Figure 3). In particular KS statistics between all part-of-speech pairs significantly indicate that their distributions differ along the information axis. Furthermore, KS statistic values are 0.12 for (noun, verb), 0.11 for (noun, adjective), 0.64 for (noun, adverb), 0.22 for (verb, adjective), 0.72 for (verb, adverb), and 0.64 for (adjective, adverb). This suggests that adverbs are most cleanly separated from other parts of speech along the information axis (consistent with Figure 3), which may indicate that morphemes such as affixes have important effects in embedding space. A detailed investigation is beyond the scope of this paper and may require analyses through alternative heuristics such as pseudomorphology and lexical neighborhood density (Needle et al., 2020).

Many extant words near the midpoint of the information axis are, or may be, compound words; the boundary between derivative and compound words is thought to be fuzzy because many derivational suffixes developed from words are frequently used in compounding (Trost, 1992). Both derivative and compound words populate other spaces of the extant language region, but conflicting definitions hamper straightforward statistical analysis.

Morphological traits such as adjectival suffixes *-ness*, *-ism*, and *-able*, or the adverbial suffix *-ly* correlate to clear embedding mappings, but the boundaries for morphological classes are not distinct. Garble ending in “s” occupies a closer region to extant language than most other garble, arguably due to the semantic associations of ending in “s” (e.g. regarding pluralization) derived from the suffix *-s*. Note, morphological heuristics like affixation apply to lexicalized words but not garble. Pseudowords ending in “s” share that region of garble ending in “s”, however, such seemingly plural pseudowords tend closer to extant language, reflecting the notion that word form similarity increases with semantic similarity (Dautriche et al., 2017). Given the fuzziness of morphology and the opaqueness of English spelling (Needle et al., 2020), pseudowords ending in “s” may or may not be due to affixation.

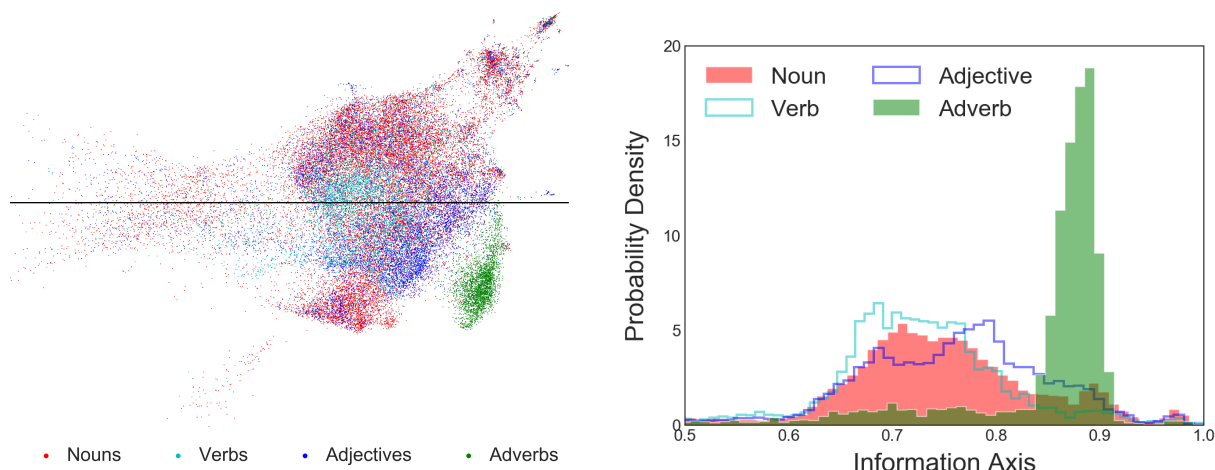


Figure 3: *Left panel*: UMAP projection of CharacterBERT embeddings for extant words split by part-of-speech into nouns (red), verbs (cyan), adjectives (blue), and adverbs (green). *Right panel*: Probability density of extant words, split by part-of-speech, as a function of minmax-normalized position along the information axis shown in Figure 1.

### 5.3 Concreteness/Abstractness

The internal positioning of different parts-of-speech within the extant language space of our low-dimensional UMAP projection suggests that the representations also capture notions of concreteness (e.g nouns) and abstractness (e.g adverbs) which we explore by projecting concreteness scores from the (Brybaert et al., 2014) study. We calculate the center of extant UMAP coordinates with no weighting and with weighting by minmax-normalized concreteness and used those points to define a *concreteness axis*, which demonstrates that concreteness varies in a direction roughly orthogonal to our information axis (see Figure 4). The bootstrap-resampled angle distribution between information and concreteness axes is  $86.6 \pm 1.2$  degrees.

Thus, the information axis and word concreteness capture two crucial and largely distinct aspects of the many latent features underlying CharacterBERT representations. This finding is particularly relevant in light of recent work showing not only that word concreteness is a psychologically rich dimension that shapes semantic processing (Brybaert et al., 2016; Guilbeault et al., 2020), but also that word concreteness is surprisingly effective at enriching the predictive capacities of word embedding models, such as for the purpose of automated metaphor detection (Srinivasa Desikan et al., 2020). We leave a detailed investigation of this finding, including its relation to the visual information (Brybaert et al., 2016) carried by concrete and abstract words, to future work.

### 5.4 Markov Chain Model

We also create a language model using the Prediction by Partial Matching (PPM) variable order Markov model (VOMM) to estimate the probability of each of these character  $n$ -grams (Begleiter et al., 2004). The model calculates the *logpdf* for each character  $n$ -gram in which more commonly occurring character  $n$ -grams have a lower score, and less commonly occurring character  $n$ -grams receive a higher score. The model is trained on extant words, then used to score all of the extant, pseudowords and garble character  $n$ -grams. We use this score to capture the likelihood of character  $n$ -grams in our character sequence space (Figure 5).

These Markov model values correlate with our information axis measure. In particular, the Spearman correlation coefficient between information axis and Markov chain information content is 0.4 (highly significant) for randomly generated  $n$ -grams, and 0.007 (not significant) for extant words. Thus, for random character  $n$ -grams, our information axis measure is correlated with statistical properties of the character  $n$ -grams from the Markov model (see the left panel of Figure 5). However, our information axis measure more clearly separates extant and garble  $n$ -grams, indicating that it incorporates information beyond purely statistical properties of  $n$ -gram classes (see the right panel of Figure 5). This suggests that the CharacterBERT model learns information beyond character-level statistical information, even for  $n$ -grams that never explicitly appear in the training data.

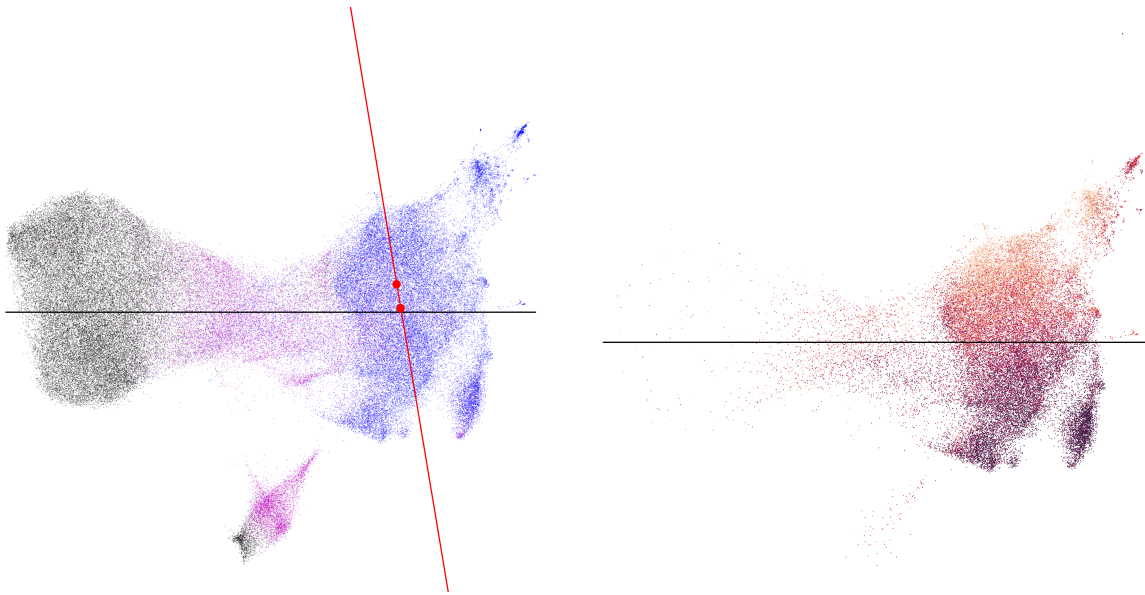


Figure 4: *Left panel*: UMAP projection of CharacterBERT embeddings for extant words (blue), pseudowords (magenta), and randomly generated character  $n$ -grams (black). The solid black line shows the information axis that we define in this work, and the red line shows the axis that captures variability in word concreteness, computed by connecting the unweighted average UMAP position for extant words with that weighted by minmax-normalized concreteness (red dots). *Right panel*: UMAP of only extant words, colored by minmax-normalized concreteness, with lighter colors indicating more concrete words.

## 6 Discussion and Conclusion

Using the CharacterBERT model, we embedded a large corpus of character level  $n$ -grams outside of extant language to study how the primitive information they contain relates to the semantic information carried by extant language. The key findings of this paper are:

1. Extant words and randomly generated character  $n$ -grams are separated along a particular axis in our UMAP projection of CharacterBERT embedding space (Figures 1–2);
2. Pseudowords lie between extant and randomly generated  $n$ -grams along this axis, but there is no distinct boundary between these classes of  $n$ -grams (Figures 1–2);
3. The structure of CharacterBERT embeddings of extant language, including structure based on part-of-speech and morphology, is correlated with the information axis (Figure 3);
4. Word concreteness varies along a dimension that is roughly orthogonal to the information axis in our UMAP projection (Figure 4);
5. Separation between extant and randomly generated  $n$ -grams captured by CharacterBERT

is correlated with and more coherent than that based purely on the statistical properties of  $n$ -grams (Figure 5).

These findings suggest that character-based Transformer models are largely able to explore the relation between extant words and randomly generated character strings. In particular, character-level models capture complex structure in the space of words, pseudowords, and randomly generated  $n$ -grams. These findings are consistent with work suggesting that character-level and morpheme-aware representations are rich in meaning, even compared to word or sub-word models (Al-Rfou et al., 2019; El Boukkouri et al., 2020; Ma et al., 2020; Hofmann et al., 2020, 2021).

Our study is limited to extant words in English and randomly generated character  $n$ -grams using the English alphabet. Given the unique impact of a specific language and alphabet on representation spaces, there is motivation to see whether the relationships we identify generalise to other languages and alphabets. Finally, we reiterate that our analysis was limited to the last embedding layer of the CharacterBERT model; future work may focus on weights in earlier layers, including attention mechanisms explored by other BERTology studies (Clark et al., 2019; Jawahar et al., 2019). By only analysing the final embedding layer, we study



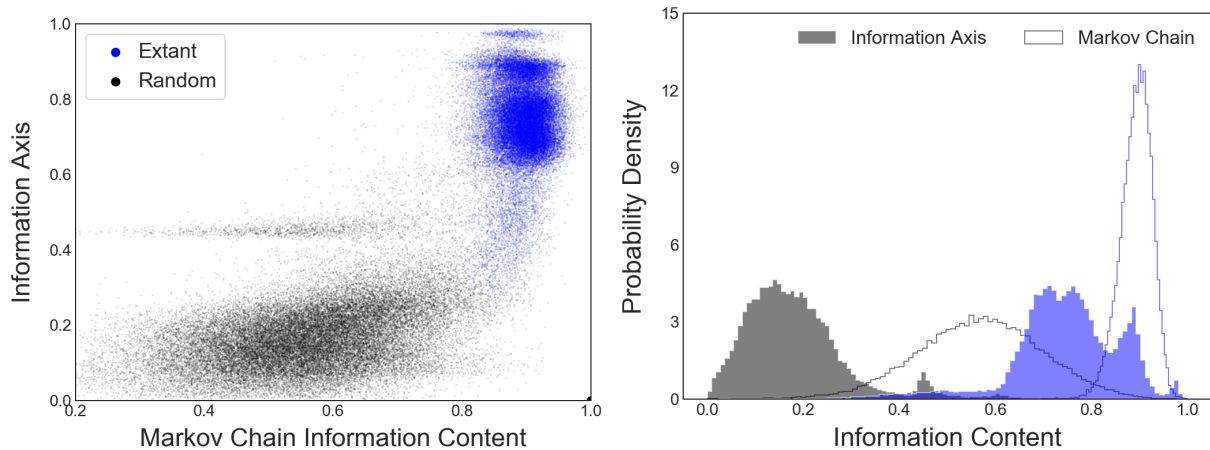


Figure 5: *Left panel*: Minmax-normalized position along the information axis shown in Figure 1 vs. minmax-normalized information content from our Markov Chain model, for extant words (blue) and randomly generated character  $n$ -grams (black). *Right panel*: Probability density of minmax-normalized information content measures from our UMAP projection (filled histograms) and Markov Chain model (unfilled histograms).

the ‘psychology’ of such character-level models; in analogy, much may be gained by studying the ‘neuroscience’ of such models encoded in their attention weights (Wang, 2020).

Our study also has important practical implications for the widespread use of pseudowords as an experimental tool in psycholinguistic research. Pseudowords are frequently used as stimuli to observe the psychological and neurocognitive processes underlying the interpretation of novel words (Price et al., 1996; Stark and McClelland, 2000; Keuleers and Brysbaert, 2010; Lupyan and Casasanto, 2015; Davis et al., 2019b). However, the lion’s share of this research treats all pseudoword stimuli as equivalent in their novelty, based on *prima facie* human judgments. By contrast, our method shows that not all pseudowords are created equal. Due to various features of character sequences, including morphological structure, some pseudowords encode disproportionately more information according to character-aware language models, and are therefore represented as significantly more similar to extant words, whereas other pseudowords are recognized by these models as random character sequences. This variation is especially striking given that the algorithms used to generate pseudowords are highly constrained and designed to produce morphologically coherent words (Keuleers and Brysbaert, 2010); that some pseudowords are evaluated as random by CharacterBERT reveals not only asymmetries in the coherence of pseudowords that may be of psychological relevance, but also assumptions and limitations

in terms of which morphological units CharacterBERT and related models recognize as signatures of extant words. Our study thus provides a quantitative method for evaluating pseudoword plausibility, without relying on variable human judgments, while also revealing insights into key differences between how humans and contemporary language models evaluate the plausibility of pseudowords.

To allow for further explorations and replicability, we release all of our data and code on GitHub<sup>3</sup>. Our findings reveal new avenues for future work using character-aware embeddings of extant, pseudoword, and garble  $n$ -grams, including analyses of nonsense poetry like Lewis Carroll’s “Jabberwocky” or of the innovative idiosyncrasies of rap lyricists and graffiti artists. The embeddings we study may also complement philological studies (especially if dynamic analyses are employed), as well as research into novel category formation (Lupyan and Casasanto, 2015; Guilbeault et al., 2021). Also, language acquisition studies of the distinction between language and noise may benefit from character-level embeddings beyond the realm of extant language (Yang, 2006; Carey, 2000). By investigating a broadened embedding space to include randomly generated  $n$ -grams, we found new structures of meaning through the context of meaningless information; further studies may extend our garble-based approach across different media and modes to contribute to more general understandings of human meaning.

<sup>3</sup><https://github.com/comp-syn/garble>

## References

- Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. 2019. Character-level language modeling with deeper self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3159–3166.
- Ron Begleiter, Ran El-Yaniv, and Golan Yona. 2004. On prediction using variable order markov models. *Journal of Artificial Intelligence Research*, 22:385–421.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates.
- Marc Brysbaert, Michaël Stevens, Paweł Mandera, and Emmanuel Keuleers. 2016. How many words do we know? practical estimates of vocabulary size dependent on word definition, the degree of language input and the participant’s age. *Frontiers in psychology*, 7:1116.
- Marc Brysbaert, Amy Beth Warriner, and Victor Kuperman. 2014. Concreteness ratings for 40 thousand generally known english word lemmas. *Behavior research methods*, 46(3):904–911.
- Susan Carey. 2000. The origin of concepts. *Journal of Cognition and Development*, 1(1):37–41.
- Shamil Chollampatt and Hwee Tou Ng. 2018. A multi-layer convolutional encoder-decoder neural network for grammatical error correction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. 2019. What does bert look at? an analysis of bert’s attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286.
- Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning*, 20(3):273–297.
- Isabelle Dautriche, Kyle Mahowald, Edward Gibson, and Steven T Piantadosi. 2017. Wordform similarity increases with semantic similarity: An analysis of 100 languages. *Cognitive science*, 41(8):2149–2169.
- B. L. Davis and P. F. MacNeilage. 1995. The articulatory basis of babbling. *Journal of Speech & Hearing Research*, 38(6):1199–1211.
- Charles Davis, Hannah Morrow, and Gary Lupyan. 2019a. What does a horgous look like? nonsense words elicit meaningful drawings. *Cognitive Science*, 43.
- Charles P Davis, Hannah M Morrow, and Gary Lupyan. 2019b. What does a horgous look like? nonsense words elicit meaningful drawings. *Cognitive Science*, 43(10):e12791.
- David Deutsch. 2011. *The beginning of infinity: Explanations that transform the world*. Penguin UK.
- David Deutsch and Chiara Marletto. 2015. Constructor theory of information. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 471(2174):20140540.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*.
- Lisa Susan Ede. 1975. *The nonsense literature of Edward Lear and Lewis Carroll*. The Ohio State University.
- Hicham El Boukkouri, Olivier Ferret, Thomas Lavergne, Hiroshi Noji, Pierre Zweigenbaum, and Jun’ichi Tsujii. 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.
- Guy Emerson. 2020. What are the goals of distributional semantics? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7436–7453.
- Kawin Ethayarajh. 2019. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and gpt-2 embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65.
- John R Firth. 1957. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*.
- Manuela Friedrich and Angela D. Friederici. 2005. Phonotactic Knowledge and Lexical-Semantic Processing in One-year-olds: Brain Responses to Words and Nonsense Words in Picture Contexts. *Journal of Cognitive Neuroscience*, 17(11):1785–1802.
- Yoav Goldberg. 2019. Assessing bert’s syntactic abilities. *arXiv preprint arXiv:1901.05287*.
- Douglas Guilbeault, Andrea Baronchelli, and Damon Centola. 2021. Experimental evidence for scale-induced category convergence across populations. *Nature communications*, 12(1):1–7.

- Douglas Guilbeault, Ethan O Nadler, Mark Chu, Donald Ruggiero Lo Sardo, Aabir Abubaker Kar, and Bhargav Srinivasa Desikan. 2020. Color associations in abstract semantic domains. *Cognition*, 201:104306.
- V Hofmann, J Pierrehumbert, and H Schütze. 2020. Dagobert: generating derivational morphology with a pretrained language model. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (and forerunners)(EMNLP)*. ACL Anthology.
- Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. 2021. Superbizarre is not superb: Derivational morphology improves bert’s interpretation of complex words. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3594–3608.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does bert learn about the structure of language? In *ACL 2019-57th Annual Meeting of the Association for Computational Linguistics*.
- Daniel Jurafsky and James H Martin. 2021. *Speech and language processing* 3rd edition.
- Emmanuel Keuleers and Marc Brysbaert. 2010. Wuggy: A multilingual pseudoword generator. *Behavior research methods*, 42(3):627–633.
- Yoon Kim, Yacine Jernite, David Sontag, and Alexander M Rush. 2016. Character-aware neural language models. In *Thirtieth AAAI conference on artificial intelligence*.
- Kamran Kowsari, Kiana Jafari Meimandi, Mojtaba Heidarysafa, Sanjana Mendu, Laura Barnes, and Donald Brown. 2019. Text classification algorithms: A survey. *Information*, 10(4):150.
- Thomas K Landauer and Susan T Dumais. 1997. A solution to plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological review*, 104(2):211.
- Jean-Jacques Lecercle. 2012. *Philosophy of nonsense: The intuitions of Victorian nonsense literature*. Routledge.
- Gary Lupyan and Daniel Casasanto. 2015. Meaningless words promote meaningful categorization. *Language and Cognition*, 7(2):167–193.
- Wentao Ma, Yiming Cui, Chenglei Si, Ting Liu, Shijin Wang, and Guoping Hu. 2020. Charbert: Character-aware pre-trained language model. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 39–50.
- David Marr and Ellen Hildreth. 1980. Theory of edge detection. *Proceedings of the Royal Society of London. B. Biological Sciences*, 207(1167):187–217.
- Frank J Massey Jr. 1951. The kolmogorov-smirnov test for goodness of fit. *Journal of the American statistical Association*, 46(253):68–78.
- Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. 2018. Umap: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29).
- Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *Interspeech*, volume 2, pages 1045–1048. Makuhari.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Bhaskar Mitra, Nick Craswell, et al. 2018. *An introduction to neural information retrieval*. Now Foundations and Trends.
- J Needle, J Pierrehumbert, and Jennifer B Hay. 2020. Phonological and morphological effects in the acceptability of pseudowords.
- Winfried Noth. 1990. *Handbook of semiotics*. Indiana University Press.
- Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of NAACL-HLT*, pages 2227–2237.
- Cathy J Price, RJS Wise, and RSJ Frackowiak. 1996. Demonstrating the implicit processing of visually presented words and pseudowords. *Cerebral cortex*, 6(1):62–70.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in bertology: What we know about how bert works. *Transactions of the Association for Computational Linguistics*, 8:842–866.
- Claude E Shannon. 1951. The redundancy of english. In *Cybernetics; Transactions of the 7th Conference, New York: Josiah Macy, Jr. Foundation*, pages 248–272.
- Claude Elwood Shannon. 2001. A mathematical theory of communication. *ACM SIGMOBILE mobile computing and communications review*, 5(1):3–55.
- Bhargav Srinivasa Desikan, Tasker Hull, Ethan Nadler, Douglas Guilbeault, Aabir Abubakar Kar, Mark Chu, and Donald Ruggiero Lo Sardo. 2020. comp-syn: Perceptually grounded word embeddings with color. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1744–1751.
- Craig EL Stark and James L McClelland. 2000. Repetition priming of words, pseudowords, and nonwords. *Journal of experimental psychology: Learning, memory, and cognition*, 26(4):945.

- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.
- Timothy D Sweeny, Nicole Wurnitsch, Alison Gopnik, and David Whitney. 2015. Ensemble perception of size in 4–5-year-old children. *Developmental science*, 18(4):556–568.
- Harald Trost. 1992. Computational Morphology. In *Morphology and Computation*. The MIT Press.
- Xin Wang. 2020. The curious case of developmental bertology: On sparsity, transfer learning, generalization and the brain. *arXiv preprint arXiv:2007.03774*.
- Ziang Xie, Anand Avati, Naveen Arivazhagan, Dan Jurafsky, and Andrew Y Ng. 2016. Neural language correction with character-based attention. *arXiv preprint arXiv:1603.09727*.
- Charles Yang. 2006. *The infinite gift: How children learn and unlearn the languages of the world*. Simon and Schuster.

## A Robustness to Alternative Dimensionality Reduction Techniques

Our main analyses use the UMAP algorithm to project garble, pseudoword, and extant word CharacterBERT embeddings into an interpretable, low-dimensional space. Here, we demonstrate that our key results are not sensitive to this choice of dimensionality reduction technique by recreating our findings using t-SNE, a popular alternative to UMAP. Figure 6 shows the extant, pseudoword, and garble embeddings resulting from the `scikit-learn` t-SNE algorithm (run with  $n_{\text{components}} = 2$  and perplexity = 10). The qualitative structure is unchanged relative to the UMAP embedding shown in Figure 1: garble and extant  $n$ -grams are separated along a new information axis that captures roughly the same amount of variance as our original UMAP information axis, and pseudowords embeddings connect these two clusters. Furthermore, some particular aspects of the UMAP structure are preserved, including a distinct cluster of garble and pseudoword  $n$ -grams ending in “s” near the bottom of Figure 6. In general, the separation among t-SNE  $n$ -gram clusters is somewhat less distinct compared to the UMAP case, which we attribute to UMAP’s better preservation of global structure (McInnes et al., 2018).

The results of this t-SNE projection are also quantitatively consistent with our main findings. In particular, the UMAP information axis summary statistics presented in Table 1 become  $0.70 \pm 0.15$ ,  $0.54 \pm 0.13$ , and  $0.26 \pm 0.12$  for extant, pseudoword, and randomly generated  $n$ -grams, respectively; these results are all consistent with our UMAP results at the  $1\sigma$  level. Similarly, KS two-sample tests between the extant, pseudoword, and garble information axis distributions all remain highly significant ( $p \ll 0.001$ ), and their ordering is consistent with our UMAP results: the t-SNE information axis KS statistic scores are 0.86 for (extant, random), 0.76 for (pseudoword, random), and 0.50 for (extant, pseudoword). Relative to our fiducial UMAP results, the slightly larger scatter for the information axis summary statistics and the slightly weaker KS statistic scores are consistent with the increases in scatter orthogonal to the information axis in the t-SNE projection (Figure 6) relative to the UMAP projection (Figure 1). Thus, our main results are not sensitive to the dimensionality reduction method employed.

## B Global Structure of Garble and Pseudoword Embeddings

Here, we qualitatively explore the main features of pseudoword and randomly generated  $n$ -grams’ structure in our UMAP projection, deferring a more detailed exposition to future work. Figure 7 highlights several distinct groups of randomly generated (black) and pseudoword (magenta)  $n$ -grams that we describe in detail below.

Beginning with randomly generated  $n$ -grams, we first note that there is a significant correlation between their string length and information axis score, such that randomly generated  $n$ -grams with low information axis scores tend to contain more characters, and vice versa. Indeed, the high-information tail of the garble distribution shown in Figure 2 has a power-law exponent that is quantitatively consistent with the low-length tail of the underlying string length distribution.<sup>4</sup> Figure 7 highlights two notable exceptions to this rule: a cluster of randomly generated  $n$ -grams with strings that tend to be short and often contain repeated characters, and a garble cluster in which strings tend to end in “s.” We refer to the remaining randomly generated  $n$ -grams as “typical garble.” To illustrate, we provide ten examples of  $n$ -grams in each category:

- Typical garble: kiwbckodaffzhjxkvpfh, ijhtsfjsu, ojcfere, fsgnwy, qiqa, nevm, uzp, tgj, bv, w;
- Short repeated garble: cureuul, fbxoon, gallm, alln, ffod, ido, obb, tek, aa, hq;
- -s garble: dddgvasbbzaeuius, wdycrynylhyos, nkeccmosls, ilvtubdts, eoubazos, ptfjqs, hslxls, xwkss, gehs, jgs.

A particularly interesting feature of short repeated garble is that it encodes considerably more information along our axis than the pseudowords in our sample. This is striking because the pseudowords were generated using an algorithm designed to generate morphologically plausible words, whereas the garble is generated purely randomly at the character level. In this way, our garble embeddings provide novel insights into the string patterns that CharacterBERT identifies as information rich and predictive of word plausibility (in terms of proximity to extant words in embedding

<sup>4</sup>We remind the reader that our randomly generated string length distribution is matched to that of our extant sample.

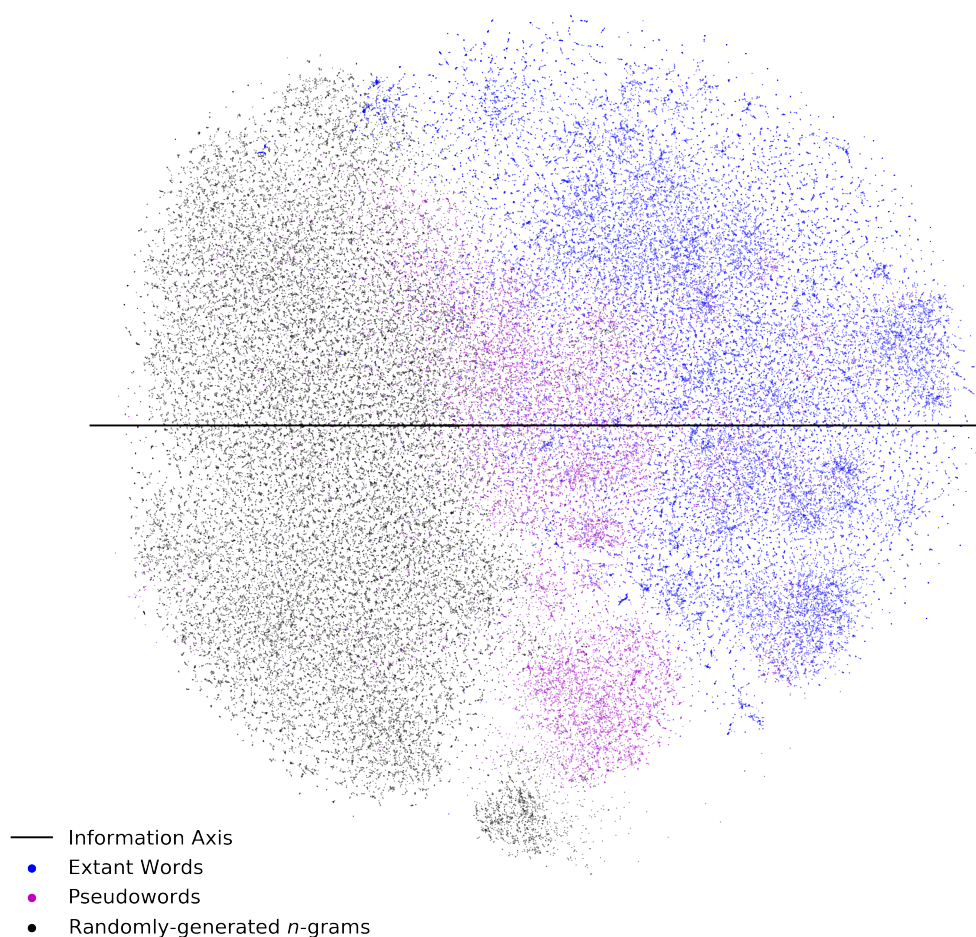


Figure 6: t-SNE projection of CharacterBERT embeddings for extant words (blue), pseudowords (magenta), and randomly generated character  $n$ -grams (black). The solid black line shows the information axis that we define in this work. The bottom-most cluster of random and pseudoword character  $n$ -grams is comprised of character  $n$ -grams ending in “s.”

space); specifically, it reveals that CharacterBERT identifies repeated characters in the same string as information rich, even though these repeated character sequences often lack morphological hallmarks of extant words.

These sequences of repeated sounds share similarities with early-childhood vocalizations (“babbling”), as well as stylistic features of child-oriented speech in the context of early word learning, for example, words such as *mama* and *dada*. The role of simple character repetitions in child development is often studied from the phonetic standpoint as a mechanism for a child to become proficient in the diversity of sounds appearing in a language (Davis and MacNeilage, 1995). However, the results from CharacterBERT suggest that repetitive sequences are especially rich in their information at the character-level, which may confer additional syntactic or lexical benefits as children learn to differentiate random sounds from linguisti-

cally meaningful units. We leave further investigation of this to future work.

Pseudoword embeddings also display a clear cluster in which strings tend to end in “s.” In addition, there is a distinct pseudoword group near extant adverbs (see Figure 3) in which strings tend to end in “ly.” We refer to the remaining pseudoword  $n$ -grams as “typical pseudowords.” To illustrate, we provide ten examples of  $n$ -grams in each category:

- Typical pseudowords: hypnostementer, eatmendownwald, eninardister, unalgion, conquing, ambooked, runton, ditity, etbarn;
- -s pseudowords: sacrembelcones, irstuphories, unnessnells, herepairs, finihips, littoes, warposs, quards, prects, gicass;
- -ly pseudowords: queepecturusly, unbornordardly, remechlocally, expotputtly, musteetly, confully, popubly, ectoily, artfaly, mously.

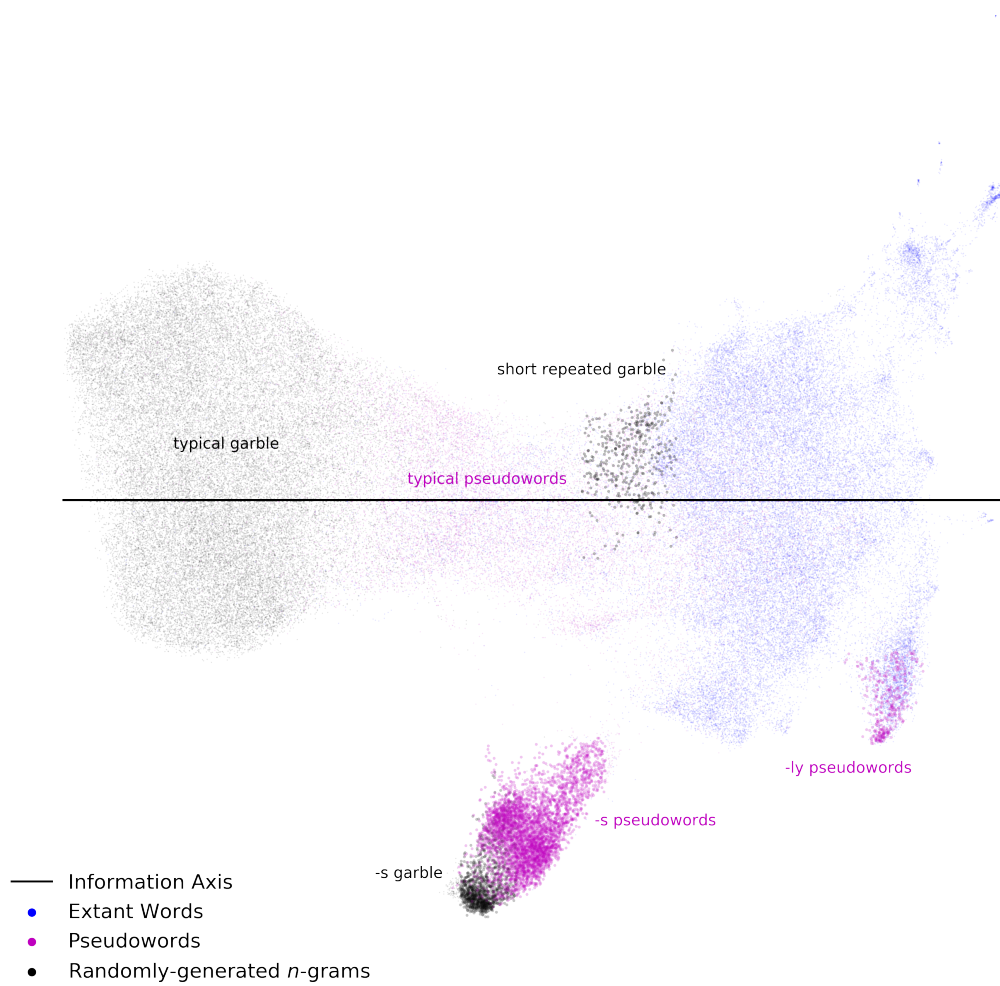


Figure 7: UMAP projection of CharacterBERT embeddings for extant words (blue), pseudowords (magenta), and randomly generated character  $n$ -grams (black). The solid black line shows the information axis that we define in this work. We discuss the highlighted clusters of pseudoword and garble  $n$ -grams in Appendix B.