

# ***Tabula* nearly *rasa*: Probing the linguistic knowledge of character-level neural language models trained on unsegmented text**

**Michael Hahn\***

Stanford University  
mhahn2@stanford.edu

**Marco Baroni**

Facebook AI Research  
UPF Linguistics Department  
Catalan Institution for Research  
and Advanced Studies  
mbaroni@gmail.com

## **Abstract**

Recurrent neural networks (RNNs) have reached striking performance in many natural language processing tasks. This has renewed interest in whether these generic sequence processing devices are inducing genuine linguistic knowledge. Nearly all current analytical studies, however, initialize the RNNs with a vocabulary of known words, and feed them tokenized input during training. We present a multi-lingual study of the linguistic knowledge encoded in RNNs trained as character-level language models, on input data with word boundaries removed. These networks face a tougher and more cognitively realistic task, having to discover any useful linguistic unit from scratch based on input statistics. The results show that our “near *tabula rasa*” RNNs are mostly able to solve morphological, syntactic and semantic tasks that intuitively presuppose word-level knowledge, and indeed they learned, to some extent, to track word boundaries. Our study opens the door to speculations about the necessity of an explicit, rigid word lexicon in language learning and usage.

## **1 Introduction**

Recurrent neural networks (RNNs, Elman, 1990), in particular in their Long-Short-Term-Memory variant (LSTMs, Hochreiter and Schmidhuber, 1997), are widely used in natural language processing. RNNs, often pre-trained on the simple *language modeling* objective of predicting the next symbol in natural text, are often a crucial component of state-of-the-art architectures for machine translation, natural language inference and text categorization (Goldberg, 2017).

RNNs are very general devices for sequence processing, hardly assuming any prior linguistic

knowledge. Moreover, the simple prediction task they are trained on in language modeling is well-attuned to the core role prediction plays in cognition (e.g., Bar, 2007; Clark, 2016). RNNs have thus long attracted researchers interested in language acquisition and processing. Their recent successes in large-scale tasks has rekindled this interest (e.g., Frank et al., 2013; Lau et al., 2017; Kirov and Cotterell, 2018; Linzen et al., 2018; McCoy et al., 2018; Pater, 2018).

The standard pre-processing pipeline of modern RNNs assumes that the input has been tokenized into word units that are pre-stored in the RNN vocabulary (Goldberg, 2017). This is a reasonable practical approach, but it makes simulations less interesting from a linguistic point of view. First, discovering words (or other primitive constituents of linguistic structure) is one of the major challenges a learner faces, and by pre-encoding them in the RNN we are facilitating its task in an unnatural way (not even the staunchest nativists would take specific word dictionaries to be part of our genetic code). Second, assuming a unique tokenization into a finite number of discrete word units is in any case problematic. The very notion of what counts as a word in languages with a rich morphology is far from clear (e.g., Dixon and Aikhenvald, 2002; Bickel and Zúñiga, 2017), and, universally, lexical knowledge is probably organized into a not-necessarily-consistent hierarchy of units at different levels: morphemes, words, compounds, constructions, etc. (e.g., Goldberg, 2005). Indeed, it has been suggested that the notion of word cannot even be meaningfully defined cross-linguistically (Haspelmath, 2011).

Motivated by these considerations, we study here RNNs that are trained without any notion of word in their input or in their architecture. We train our RNNs as *character-level neural language models* (CNLMs, Mikolov et al., 2011; Sutskever et al., 2011; Graves, 2014) by removing white-

---

Work partially done while interning at Facebook AI Research.

space from their input, so that, like children learning a language, they don't have access to explicit cues to wordhood.<sup>1</sup> This setup is almost as *tabula rasa* as it gets. By using unsegmented orthographic input (and assuming that, in the alphabetic writing systems we work with, there is a reasonable correspondence between letters and phonetic segments), we are only postulating that the learner figured out how to map the continuous speech stream to a sequence of phonological units, an ability children already possess few months after birth (e.g., Maye et al., 2002; Kuhl, 2004). We believe that focusing on language modeling of an unsegmented phoneme sequence, abstracting away from other complexities of a fully realistic child language acquisition setup, is particularly instructive, in order to study which linguistic structures naturally emerge.

We evaluate our character-level networks on a bank of linguistic tests in German, Italian and English. We focus on these languages due to resource availability and ease of benchmark construction. Also, well-studied synthetic languages with a clear, orthographically-driven notion of word might be a better starting point to test non-word-centric models, compared to agglutinative or polysynthetic languages, where the very notion of what counts as a word is problematic.

Our tasks require models to develop the latent ability to parse characters into word-like items associated to morphological, syntactic and broadly semantic features. The RNNs pass most of the tests, suggesting that they are in some way able to construct and manipulate the right lexical objects. In a final experiment, we look more directly into *how* the models are handling word-like units. We find, confirming an earlier observation by Kementchedjhieva and Lopez (2018), that the RNNs specialized some cells to the task of detecting word boundaries (or, more generally, salient linguistic boundaries, in a sense to be further discussed below). Taken together, our results suggest that character-level RNNs capture forms of linguistic knowledge that are traditionally thought to be word-based, without being exposed to an explicit segmentation of their input and, more importantly, without possessing an explicit word lexicon. We will discuss the implications of these

---

<sup>1</sup>We do not erase punctuation marks, reasoning that they have a similar function to prosodic cues in spoken language.

findings in the conclusion.<sup>2</sup>

## 2 Related work

**On the primacy of words** Several linguistic studies suggest that words, at least as delimited by whitespace in some writing systems, are neither necessary nor sufficient units of linguistic analysis. Haspelmath (2011) claims that there is no cross-linguistically valid definition of the notion of word (see also Schiering et al., 2010, who address specifically the notion of prosodic word). Others have stressed the difficulty of characterizing words in polysynthetic languages (Bickel and Zúñiga, 2017). Children are only rarely exposed to words in isolation during learning (Tomasello, 2003),<sup>3</sup> and it is likely that the units that adult speakers end up storing in their lexicon are of variable size, both smaller and larger than conventional words (e.g., Jackendoff, 2002; Goldberg, 2005). From a more applied perspective, Schütze (2017) recently defended tokenization-free approaches to NLP, proposing a general non-symbolic approach to text representation. We hope our results will contribute to the theoretical debate on word primacy, suggesting, through computational simulations, that word priors are not crucial to language learning and processing.

**Character-based neural language models** received attention in the last decade because of their greater generality compared to word-level models. Early studies (Mikolov et al., 2011; Sutskever et al., 2011; Graves, 2014) established that CNLMs might not be as good at language modeling as their word-based counterparts, but lag only slightly behind. This is particularly encouraging in light of the fact that character-level sentence prediction involves a much larger search space than prediction at the word level, as a character-level model must make a prediction after each character, rather than after each word. Sutskever et al. (2011) and Graves (2014) ran qualitative analyses showing that CNLMs capture some basic linguistic properties of their input. The latter, who used LSTM cells, also showed,

---

<sup>2</sup>Our input data, test sets and pre-trained models are available at <https://github.com/m-hahn/tabula-rasa-rnns>.

<sup>3</sup>Single-word utterances are not uncommon in child-directed language, but they are still rather the exception than the rule, and many important words, such as determiners, never occur in isolation (Christiansen et al., 2005).

qualitatively, that CNLMs are sensitive to hierarchical structure. In particular, they balance parentheses correctly when generating text.

Most recent work in the area has focused on *character-aware* architectures combining character- and word-level information to develop state-of-the-art language models that are also effective in morphologically rich languages (e.g., Bojanowski et al., 2016; Kim et al., 2016; Gerz et al., 2018). For example, Kim and colleagues perform prediction at the word level, but use a character-based convolutional network to generate word representations. Other work focuses on splitting words into morphemes, using character-level RNNs and an explicit segmentation objective (e.g., Kann et al., 2016). These latter lines of work are only distantly related to our interest in probing what a purely character-level network trained on running text has implicitly learned about linguistic structure. There is also extensive work on segmentation of the linguistic signal that does not rely on neural methods, and is not directly relevant here, (e.g., Brent and Cartwright, 1996; Goldwater et al., 2009; Kamper et al., 2016, and references there).

**Probing linguistic knowledge of neural language models** is currently a popular research topic (Li et al., 2016; Linzen et al., 2016; Shi et al., 2016; Adi et al., 2017; Belinkov et al., 2017; Kádár et al., 2017; Hupkes et al., 2018; Conneau et al., 2018; Ettinger et al., 2018; Linzen et al., 2018). Among studies focusing on character-level models, Elman (1990) already reported a proof-of-concept experiment on implicit learning of word segmentation. Christiansen et al. (1998) trained a RNN on phoneme-level language modeling of transcribed child-directed speech with tokens marking utterance boundaries, and found that the network learned to segment the input by predicting the utterance boundary symbol also at word edges. More recently, Sennrich (2017) explored the grammatical properties of character- and subword-unit-level models that are used as components of a machine translation system. He concluded that current character-based decoders generalize better to unseen words, but capture less grammatical knowledge than subword units. Still, his character-based systems lagged only marginally behind the subword architectures on grammatical tasks such as handling agreement and negation. Radford et al. (2017) focused on

CNLMs deployed in the domain of sentiment analysis, where they found the network to specialize a unit for sentiment tracking. We will discuss below how our CNLMs also show single-unit specialization, but for boundary tracking. Godin et al. (2018) investigated the rules implicitly used by supervised character-aware neural morphological segmentation methods, finding linguistically sensible patterns. Alishahi et al. (2017) probed the linguistic knowledge induced by a neural network that receives unsegmented acoustic input. Focusing on phonology, they found that the lower layers of the model process finer-grained information, whereas higher layers are sensitive to more abstract patterns. Kementchedjhieva and Lopez (2018) recently probed the linguistic knowledge of an English CNLM trained with whitespace in the input. Their results are aligned with ours. The model is sensitive to lexical and morphological structure, and it captures morphosyntactic categories as well as constraints on possible morpheme combinations. Intriguingly, the model tracks word/morpheme boundaries through a single specialized unit, suggesting that such boundaries are salient (at least when marked by whitespace, as in their experiments) and informative enough that it is worthwhile for the network to devote a special mechanism to process them. We replicated this finding for our networks trained on whitespace-free text, as discussed in Section 4.4 below, where we discuss it in the context of our other results.

### 3 Experimental setup

We extracted plain text from full English, German and Italian Wikipedia dumps with WikiExtractor.<sup>4</sup> We randomly selected test and validation sections consisting of 50,000 paragraphs each, and used the remainder for training. The training sets contained 16M (German), 9M (Italian), and 41M (English) paragraphs, corresponding to 819M, 463M and 2,333M words, respectively. Paragraph order was shuffled for training, without attempting to split by sentences. All characters were lower-cased. For benchmark construction and word-based model training, we tokenized and tagged the corpora with TreeTagger (Schmid,

---

<sup>4</sup><https://github.com/attardi/wikiextractor>

1999).<sup>5</sup> We used as vocabularies the most frequent characters from each corpus, setting thresholds so as to ensure that all characters representing phonemes were included, resulting in vocabularies of sizes 60 (English), 73 (German), and 59 (Italian). We further constructed *word-level neural language models* (WordNLMs); their vocabulary included the most frequent 50,000 words per corpus.

We trained RNN and LSTM CNLMs; we will refer to them simply as *RNN* and *LSTM*, respectively. The “vanilla” RNN will serve as a baseline to ascertain if/when the longer-range information-tracking abilities afforded to the LSTM by its gating mechanisms are necessary. Our WordNLMs are always LSTMs. For each model/language, we applied random hyperparameter search. We terminated training after 72 hours.<sup>6</sup> None of the models had overfitted, as measured by performance on the validation set.<sup>7</sup>

Language modeling performance on the test partitions is shown in Table 1. Recall that we removed whitespace, which is both easy to predict, and aids prediction of other characters. Consequently, the fact that our character-level models are below the state of the art is expected.<sup>8</sup> For example, the best model of Merity et al. (2018) achieved 1.23 English BPC on a Wikipedia-derived dataset. On EuroParl data, Cotterell et al. (2018) report 0.85 for English, 0.90 for German, and 0.82 for Italian. Still, our English BPC is comparable to that reported by Graves (2014) for his static character-level LSTM trained on space-delimited Wikipedia data, suggesting that we are achieving reasonable performance. The perplexity of the word-level model might not be comparable to that of highly-optimized state-of-the-art architectures, but it is at the expected level for a well-

<sup>5</sup><http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

<sup>6</sup>This was due to resource availability. The reasonable language-modeling results in Table 1 suggest that no model is seriously underfit, but the weaker overall RNN results in particular should be interpreted in the light of the following qualification: models are compared *given equal amount of training, but possibly at different convergence stages*.

<sup>7</sup>Hyperparameter details are in the appendix. Chosen architectures (layers/embedding size/hidden size): LSTM: En. 3/200/1024, Ge. 2/100/1024, It. 2/200/1024; RNN: En. 2/200/2048, Ge. 2/50/2048, It. same; WordNLM: En. 2/1024/1024, Ge. 2/200/1024, It. same.

<sup>8</sup>Training our models with whitespace, without further hyperparameter tuning, resulted in BPCs of 1.32 (English), 1.28 (German), and 1.24 (Italian).

	<i>LSTM</i>	<i>RNN</i>	<i>WordNLM</i>
English	1.62	2.08	48.99
German	1.51	1.83	37.96
Italian	1.47	1.97	42.02

Table 1: Performance of language models. For CNLMs, we report bits-per-character (BPC). For WordNLMs, we report perplexity.

tuned vanilla LSTM language model. For example, Gulordava et al. (2018) report 51.9 and 44.9 perplexities respectively in English and Italian for their best LSTMs trained on Wikipedia data with same vocabulary size as ours.

## 4 Experiments

### 4.1 Discovering morphological categories

Words belong to part-of-speech categories, such as nouns and verbs. Moreover, they typically carry inflectional features such as number. We start by probing whether CNLMs capture such properties. We use here the popular method of “diagnostic classifiers” (Hupkes et al., 2018). That is, we treat the hidden activations produced by a CNLM whose weights were fixed after language model training as input features for a shallow (logistic) classifier of the property of interest (e.g., plural vs. singular). If the classifier is successful, this means that the representations provided by the model are encoding the relevant information. The classifier is deliberately shallow and trained on a small set of examples, as we want to test whether the properties of interest are robustly encoded in the representations produced by the CNLMs, and amenable to a simple linear readout (Fusi et al., 2016). In our case, we want to probe word-level properties in models trained at the character level. To do this, we let the model read each target word character-by-character, and we treat the state of its hidden layer after processing the last character in the word as the model’s implicit representation of the word, on which we train the diagnostic classifier. The experiments focus on German and Italian, as it’s harder to design reliable test sets for the impoverished English morphological system.

**Word classes (nouns vs. verbs)** For both German and Italian, we sampled 500 verbs and 500 nouns from the Wikipedia training sets, requiring that they are unambiguously tagged in the corpus by TreeTagger. Verbal and nominal forms are often cued by suffixes. We removed this con-

found by selecting examples with the same ending across the two categories (*-en* in German: *Westen* ‘west’,<sup>9</sup> *stehen* ‘to stand’; and *-re* in Italian: *autore* ‘author’, *dire* ‘to say’). We randomly selected 20 training examples (10 nouns and 10 verbs), and tested on the remaining items. We repeated the experiment 100 times to account for random train-test split variation.

While we controlled for suffixes as described above, it could still be the case that other substrings reliably cue verbs or nouns. We thus considered a baseline trained on word-internal information only, namely a character-level LSTM autoencoder trained on the Wikipedia datasets to reconstruct words in isolation.<sup>10</sup> The hidden state of the LSTM autoencoder should capture discriminating orthographic features, but, by design, will have no access to broader contexts. We further considered word embeddings from the output layer of the WordNLM. Unlike CNLMs, the WordNLM cannot make educated guesses about words that are not in its training vocabulary. These OOV words are by construction less frequent, and thus likely to be in general more difficult. To get a sense of both “best-case-scenario” and more realistic WordNLM performance, we report its accuracy both excluding and including OOV items (WordNLM<sub>subs.</sub> and WordNLM in Table 2, respectively). In the latter case, we let the model make a random guess for OOV items. The percentage of OOV items over the entire dataset, balanced for nouns and verbs, was 92.3% for German and 69.4% for Italian. Note that none of the words were OOV for the CNLM, as they all were taken from the Wikipedia training set.

Results are in Table 2. All language models outperform the autoencoders, showing that they learned categories based on broader distributional evidence, not just typical strings cuing nouns and verbs. Moreover, the LSTM CNLM outperforms the RNN, probably because it can track broader contexts. Not surprisingly, the word-based model fares better on in-vocabulary words, but the gap,

<sup>9</sup>German nouns are capitalized; this cue is unavailable to the CNLM as we lower-case the input.

<sup>10</sup>The autoencoder is implemented as a standard LSTM sequence-to-sequence model (Sutskever et al., 2014). For each language, autoencoder hyperparameters were chosen using random search, as for the language models; details are in supplementary material to be made available upon publication. For both German and Italian models, the following parameters were chosen: 2 layers, 100 embedding dimensions, 1024 hidden dimensions.

	<i>German</i>	<i>Italian</i>
Random	50.0	50.0
Autoencoder	65.1 ( $\pm$ 0.22)	82.8 ( $\pm$ 0.26)
LSTM	89.0 ( $\pm$ 0.14)	95.0 ( $\pm$ 0.10)
RNN	82.0 ( $\pm$ 0.64)	91.9 ( $\pm$ 0.24)
WordNLM	53.5 ( $\pm$ 0.18)	62.5 ( $\pm$ 0.26)
WordNLM <sub>subs.</sub>	97.4 ( $\pm$ 0.05)	96.0 ( $\pm$ 0.06)

Table 2: Accuracy of diagnostic classifier on predicting word class, with standard errors across 100 random train-test splits. ‘subs.’ marks in-vocabulary subset evaluation, not comparable to the other results.

especially in Italian, is rather narrow, and there is a strong negative impact of OOV words (as expected, given that WordNLM is at random on them).

**Number** We turn next to number, a more granular morphological feature. We study German, as it possesses a rich system of nominal classes forming plural through different morphological processes. We train a diagnostic number classifier on a subset of these classes, and test on the others, in order to probe the abstract number generalization capabilities of the tested models. If a model generalizes correctly, it means that the CNLM is sensitive to number as an abstract feature, independently of its surface expression.

We extracted plural nouns from the Wiktionary and the German UD treebank (McDonald et al., 2013; Brants et al., 2002). We selected nouns with plurals in *-n*, *-s*, or *-e* to train the classifier (e.g., *Geschichte(n)* ‘story(-ies)’, *Radio(s)* ‘radio(s)’, *Pferd(e)* ‘horse(s)’, respectively). We tested on plurals formed with *-r* (e.g., *Lieder* for singular *Lied* ‘song’), or through vowel change (*Umlaut*, e.g., *Äpfel* from singular *Apfel* ‘apple’). Certain nouns form plurals through concurrent suffixing and Umlaut. We grouped these together with nouns using the same suffix, reserving the Umlaut group for nouns *only* undergoing vowel change (e.g., *Saft/Säfte* ‘juice(s)’ would be an instance of *-e* suffixation). The diagnostic classifier was trained on 15 singulars and plurals randomly selected from each training class. As plural suffixes make words longer, we sampled singulars and plurals from a single distribution over lengths, to ensure that their lengths were approximately matched. Moreover, since in uncontrolled samples from our training classes a final *-e* vowel would constitute a strong surface cue to plurality, we balanced the distribution of this property

	train classes	test classes	
	<i>-n/-s/-e</i>	<i>-r</i>	<i>Umlaut</i>
Random	50.0	50.0	50.0
Autoencoder	61.4 ( $\pm 0.9$ )	50.7 ( $\pm 0.8$ )	51.9 ( $\pm 0.4$ )
LSTM	71.5 ( $\pm 0.8$ )	78.8 ( $\pm 0.6$ )	60.8 ( $\pm 0.6$ )
RNN	65.4 ( $\pm 0.9$ )	59.8 ( $\pm 1.0$ )	56.7 ( $\pm 0.7$ )
WordNLM	77.3 ( $\pm 0.7$ )	77.1 ( $\pm 0.5$ )	74.2 ( $\pm 0.6$ )
WordNLM <sub>subs.</sub>	97.1 ( $\pm 0.3$ )	90.7 ( $\pm 0.1$ )	97.5 ( $\pm 0.1$ )

Table 3: German number classification accuracy, with standard errors computed from 200 random train-test splits. ‘subs.’ marks in-vocabulary subset evaluation, not comparable to the other results.

across singulars and plurals in the samples. For the test set, we selected all plurals in *-r* (127) or Umlaut (38), with their respective singulars. We also used all remaining plurals ending in *-n* (1467), *-s* (98) and *-e* (832) as in-domain test data. To control for the impact of training sample selection, we report accuracies averaged over 200 random train-test splits and standard errors over these splits. For WordNLM OOV, there were 45.0% OOVs in the training classes, 49.1% among the *-r* forms, and 52.1% for Umlaut.

Results are in Table 3. The classifier based on word embeddings is the most successful. It outperforms in most cases the best CNLM even in the more cogent OOV-inclusive evaluation. This confirms the common observation that word embeddings reliably encode number (Mikolov et al., 2013b). Again, the LSTM-based CNLM is better than the RNN, but both significantly outperform the autoencoder. The latter is near-random on new class prediction, confirming that we properly controlled for orthographic confounds.

We observe a considerable drop in the LSTM CNLM performance between generalization to *-r* and Umlaut. On the one hand, the fact that performance is still clearly above chance (and autoencoder) in the latter condition shows that the LSTM CNLM has a somewhat abstract notion of number not tied to specific orthographic exponents. On the other, the *-r* vs. Umlaut difference suggests that the generalization is not completely abstract, as it works more reliably when the target is a new suffixation pattern, albeit one that is distinct from those seen in training, than when it is a purely non-concatenative process.

## 4.2 Capturing syntactic dependencies

Words encapsulate linguistic information into units that are then put into relation by syntactic

rules. A long tradition in linguistics has even claimed that syntax is blind to sub-word-level processes (e.g., Chomsky, 1970; Di Sciullo and Williams, 1987; Bresnan and Mchombo, 1995; Williams, 2007). Can our CNLMs, despite the lack of an explicit word lexicon, capture relational syntactic phenomena, such as agreement and case assignment? We investigate this by testing them on syntactic dependencies between non-adjacent words. We adopt the “grammaticality judgment” paradigm of Linzen et al. (2016). We create minimal sets of grammatical and ungrammatical phrases illustrating the phenomenon of interest, and let the language model assign a likelihood to all items in the set. The language model is said to “prefer” the grammatical variant if it assigns a higher likelihood to it than to its ungrammatical counterparts. We must stress two methodological points. First, since a character-level language model assigns a probability to each character of a phrase, and the phrase likelihood is the product of these values (all between 0 and 1), minimal sets must be controlled for character length. This makes existing benchmarks unusable. Second, the “distance” of a relation is defined differently for a character-level model, and it is not straightforward to quantify. Consider the German phrase in (1) below. For a word model, two items separate the article from the noun. For a (space-less) character model, 8 characters intervene until the noun onset, but the span to consider will typically be longer. For example, *Baum* could be the beginning of the feminine noun *Baumwolle* ‘cotton’, which would change the agreement requirements on the article. So, until the model finds evidence that it fully parsed the head noun, it cannot reliably check agreement. This will typically require parsing at least the full noun and the first character following it. We again focus on German and Italian, as their richer inflectional morphology simplifies the task of constructing balanced minimal sets.

### 4.2.1 German

**Article-noun gender agreement** Each German noun belongs to one of three genders (masculine, feminine, neuter), morphologically marked on the article. As the article and the noun can be separated by adjectives and adverbs, we can probe knowledge of lexical gender together with long-distance agreement. We create stimuli of the form

- (1) {der, die, das} sehr rote Baum  
 the very red tree

where the correct nominative singular article (*der*, in this case) matches the gender of the noun. We then run the CNLM on the three versions of this phrase (removing whitespace) and record the probabilities it assigns to them. If the model assigns the highest probability to the version with the right article, we count it as a hit for the model. To avoid phrase segmentation ambiguities (as in the *Baum/Baumwolle* example above), we present phrases surrounded by full stops.

To build the test set, we select all 4,581 nominative singular nouns from the German UD treebank: 49.3% feminine, 26.4% masculine, 24.3% neuter. WordNLM OOV noun ratios are: 40.0% for masculine, 36.2% for feminine, 41.5% for neuter. We construct four conditions varying the number of adverbs and adjectives between article and noun. We first consider stimuli where no material intervenes. In the second condition, an adjective with the correct case ending, randomly selected from the training corpus, is added. Crucially, the ending of the adjective does not reveal the gender of the noun. We only used adjectives occurring at least 100 times, and not ending in *-r*.<sup>11</sup> We obtained a pool of 9,742 adjectives to sample from, also used in subsequent experiments. 74.9% of these were OOV for the WordNLM. In the third and fourth conditions, one (*sehr*) or two adverbs (*sehr extrem*) intervene between article and adjective. These do not cue gender either. We obtained 2,290 (m.), 2,261 (f.), and 1,111 (n.) stimuli, respectively. To control for surface co-occurrence statistics in the input, we constructed an n-gram baseline picking the article most frequently occurring before the phrase in the training data, breaking ties randomly. OOVs were excluded from WordNLM evaluation, resulting in an easier test for this rival model. However, here and in the next two tasks, CNLM performance on this reduced set was only slightly better, and we do not report it here. We report accuracy averaged over nouns belonging to each of the three genders. By design, the random baseline accuracy is 33%.

Results are presented in Figure 1 (left).

<sup>11</sup>Adjectives ending in *-r* often reflect lemmatization problems, as TreeTagger occasionally failed to remove the inflectional suffix *-r* when lemmatizing. We needed to extract lemmas, as we constructed the appropriate inflected forms on their basis.

WordNLM performs best, followed by the LSTM CNLM. The n-gram baseline performs similarly to the CNLM when there is no intervening material, which is expected, as a noun will often be preceded by its article in the corpus. However, its accuracy drops to chance level (0.33) in the presence of an adjective, whereas the CNLM is still able to track agreement. The RNN variant is much worse. It is outperformed by the n-gram model in the adjacent condition, and it drops to random accuracy as more material intervenes. We emphasized at the outset of this section that CNLMs must track agreement across much wider spans than word-based models. The LSTM variant ability to preserve information for longer might play a crucial role here.

**Article-noun case agreement** We selected the two determiners *dem* and *des*, which unambiguously indicate dative and genitive case, respectively, for masculine and neuter nouns:

- (2) a. {dem, des} sehr roten Baum  
 the very red tree (*dative*)  
 b. {dem, des} sehr roten Baums  
 the very red tree (*genitive*)

We selected all noun lemmas of the appropriate genders from the German UD treebank, and extracted morphological paradigms from Wiktionary to obtain case-marked forms, retaining only nouns unambiguously marking the two cases (4,509 nouns). We created four conditions, varying the amount of intervening material, as in the gender agreement experiment (4,509 stimuli per condition). For 81.3% of the nouns, at least one of the two forms was OOV for the WordNLM, and we tested the latter on the full-coverage subset. Random baseline accuracy is 50%.

Results are in Figure 1 (center). Again, WordNLM has the best performance, but the LSTM CNLM is competitive as more elements intervene. Accuracy stays well above 80% even with three intervening words. The n-gram model performs well if there is no intervening material (again reflecting the obvious fact that article-noun sequences are frequent in the corpus), and at chance otherwise. The RNN CNLM accuracy is above chance with one and two intervening elements, but drops considerably with distance.

**Prepositional case subcategorization** German verbs and prepositions lexically specify their ob-

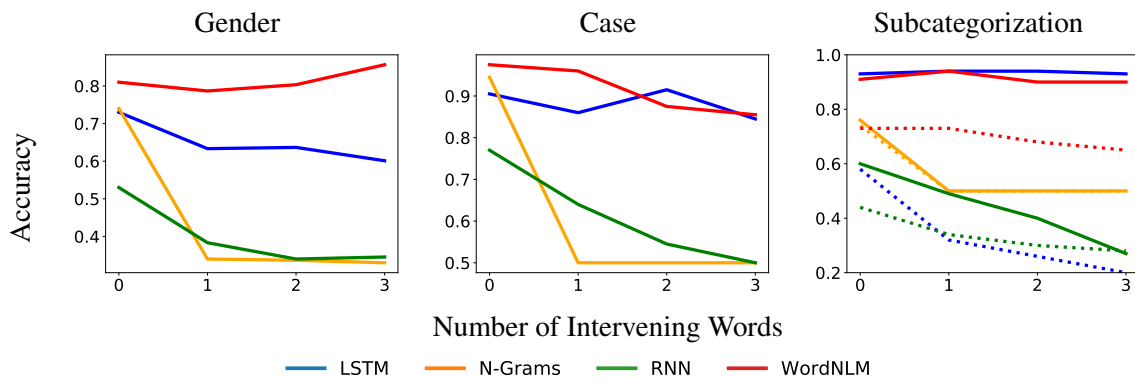


Figure 1: Accuracy in the German syntax tasks, in function of number of intervening words.

ject’s case. We study the preposition *mit* ‘with’, which selects a dative object. We focus on *mit*, as it unambiguously requires a dative object, and it is extremely frequent in the Wikipedia corpus we are using. To build the test set, we select objects whose head noun is a nominalized adjective, with regular, overtly marked case inflection. We use the same adjective pool as in the preceding experiments. We then select all sentences containing a *mit* prepositional phrase in the German Universal Dependencies treebank, subject to the constraints that (1) the head of the noun phrase governed by the preposition is not a pronoun (replacing such items with a nominal object often results in ungrammaticality), and (2) the governed noun phrase is continuous, in the sense that it is not interrupted by words that do not belong to it.<sup>12</sup> We obtained 1,629 such sentences. For each sentence, we remove the prepositional phrase and replace it by a phrase of the form

- (3) *mit der sehr {rote, roten}*  
with the very red one

where only the *-en* (dative) version of the adjective is compatible with the case requirement of the preposition (and the intervening material does not disambiguate case). We construct three conditions by varying the presence and number of adverbs (*sehr* ‘very’, *sehr extrem* ‘very extremely’, *sehr extrem unglaublich* ‘very extremely incredibly’). Note that here the correct form is longer than the wrong one. As the overall likelihood is the product of character probabilities ranging be-

tween 0 and 1, if this introduces a length bias, the latter will work against the character models. Note also that we embed test phrases into full sentences (e.g., *Die Figur hat mit der roten gespielt und meistens gewonnen*. ‘The figure played with the red one and mostly won’). We do this because this will disambiguate the final element of the phrase as a noun (not an adjective), and exclude the reading in which *mit* is a particle not governing the noun phrase of interest (Dudenredaktion, 2019).<sup>13</sup> When running the WordNLM, we excluded OOV adjectives as in the previous experiments, but did not apply further OOV filtering to the sentence frames. For the n-gram baseline, we only counted occurrences of the prepositional phrase, omitting the sentential contexts. Random baseline accuracy is 50%.

We also created control stimuli where all words up to and including the preposition are removed (the example sentence above becomes: *der roten gespielt und meistens gewonnen*). If a model’s accuracy is lower on these control stimuli than on the full ones, its performance cannot be simply explained by the different unigram probabilities of the two adjective forms.

Results are shown in Figure 1 (right). Only the n-gram baseline fails to outperform control accuracy (dotted). Surprisingly, the LSTM CNLM slightly outperforms the WordNLM, even though the latter is evaluated on the easier full-lexical-coverage stimulus subset. Neither model shows accuracy decay as the number of adverbs increases. As before, the n-gram model drops to chance as adverbs intervene, while the RNN

<sup>12</sup>The main source of noun phrase discontinuity in the German UD corpus is extraposition, a common phenomenon where part of the noun phrase is separated from the rest by the verb.

<sup>13</sup>An example of this unintended reading of *mit* is: *Ich war mit der erste, der hier war*. ‘I was one of the first who arrived here.’ In this context, dative *ersten* would be ungrammatical.



CNLM starts with low accuracy that progressively decays below chance.

#### 4.2.2 Italian

**Article-noun gender agreement** Similar to German, Italian articles agree with the noun in gender; however, Italian has a relatively extended paradigm of masculine and feminine nouns differing only in the final vowel (*-o* and *-a*, respectively), allowing us to test agreement in fully controlled paradigms such as the following:

- (4) a. {il, la} congeniale candidato  
the congenial candidate (m.)  
b. {il, la} congeniale candidata  
the congenial candidate (f.)

The intervening adjective, ending in *-e*, does not cue gender. We constructed the stimuli with words appearing at least 100 times in the training corpus. We required moreover the *-a* and *-o* forms of a noun to be reasonably balanced in frequency (neither form is more than twice as frequent as the other), or both rather frequent (appear at least 500 times). As the prenominal adjectives are somewhat marked, we only considered *-e* adjectives that occur prenominally with at least 10 distinct nouns in the training corpus. Here and below, stimuli were manually checked removing nonsensical adjective-noun (below, adverb-adjective) combinations. Finally, adjective-noun combinations that occurred in the training corpus were excluded, so that an n-gram baseline would perform at chance level. We obtained 15,005 stimulus pairs in total. 35.8% of them contained an adjective or noun that was OOV for the WordNLM. Again, we report this model’s results on its full-coverage subset, where the CNLM performance is only slightly above the one reported.

Results are shown on the first line of Table 4. WordNLM shows the strongest performance, closely followed by the LSTM CNLM. The RNN CNLM performs strongly above chance (50%), but again lags behind the LSTM.

**Article-adjective gender agreement** We next consider agreement between articles and adjectives with an intervening adverb:

- (5) a. il meno {alieno, aliena}  
the (m.) less alien one  
b. la meno {alieno, aliena}  
the (f.) less alien one

	CNLM		WordNLM
	LSTM	RNN	
Noun Gender	93.1	79.2	97.4
Adj. Gender	99.5	98.9	99.5
Adj. Number	99.0	84.5	100.0

Table 4: Italian agreement results. Random baseline accuracy is 50% in all three experiments.

where we used the adverbs *più* ‘more’, *meno* ‘less’, *tanto* ‘so much’. We considered only adjectives that occurred 1K times in the training corpus (as adjectives ending in *-al-o* are very common). We excluded all cases in which the adverb-adjective combination occurred in the training corpus, obtaining 88 stimulus pairs. Due to the restriction to common adjectives, there were no WordNLM OOVs. Results are shown on the second line of Table 4; all three models perform almost perfectly. Possibly, the task is made easy by the use of extremely common adverbs and adjectives.

**Article-adjective number agreement** Finally, we constructed a version of the last test that probed number agreement. For feminine forms, it is possible to compare same-length phrases such as:

- (6) a. la meno {aliena, aliene}  
the (s.) less alien one(s)  
b. le meno {aliena, aliene}  
the (p.) less alien one(s)

Stimulus selection was as in the last experiment, but we used a 500-occurrences threshold for adjectives, as feminine plurals are less common, obtaining 99 pairs. Again, no adverb-adjective combination was attested. There were no OOV items for the WordNLM. Results are shown on the third line of Table 4; the LSTMs perform almost perfectly, and the RNN is strongly above chance.

#### 4.3 Semantics-driven sentence completion

We probe whether CNLMs are capable of tracking the shallow form of word-level semantics required in a fill-the-gap test. We turn now to English, as for this language we can use the Microsoft Research Sentence Completion task (Zweig and Burges, 2011). The challenge consists of sentences with a gap, with 5 possible choices to fill it. Language models can be directly applied to the task, by calculating the likelihood of sentence variants with all possible completions, and selecting

the one with the highest likelihood.

The creators of the benchmark took multiple precautions to insure that success on the task implies some command of semantics. The multiple choices were controlled for frequency, and the annotators were encouraged to choose confounders whose elimination required “semantic knowledge and logical inference” (Zweig and Burges, 2011). For example, the right choice in “*Was she his [client|musings|discomfiture|choice|opportunity], his friend, or his mistress?*” depends on the cue that the missing word is coordinated with *friend* and *mistress*, and the latter are animate entities.

The task domain (Sherlock Holmes novels) is very different from the Wikipedia data-set we originally trained our models on. For a fairer comparison with previous work, we re-trained our models on the corpus provided with the benchmark, consisting of 41 Million words from 19th century English novels (we removed whitespace from this corpus as well).

Results are in Table 5. We confirm the importance of in-domain training, as the models trained on Wikipedia perform poorly (but still above chance level, which is at 20%). With in-domain training, the LSTM CNLM outperforms many earlier word-level neural models, and is only slightly below our WordNLM. The RNN is not successful even when trained in-domain, contrasting with the *word*-based vanilla RNN from the literature, whose performance, while still below LSTMs, is much stronger. Once more, this suggests that capturing word-level generalizations with a word-lexicon-less character model requires the long-span processing abilities of an LSTM.

#### 4.4 Boundary tracking in CNLMs

The good performance of CNLMs on most tasks above suggests that, although they lack a hard-coded word vocabulary and they were trained on unsegmented input, there is enough pressure from the language modeling task for them to learn to track word-like items, and associate them with various morphological, syntactic and semantic properties. In this section, we take a direct look at *how* CNLMs might be segmenting their input. Kementchedjhieva and Lopez (2018) found a *single* unit in their English CNLM that seems, qualitatively, to be tracking morpheme/word boundaries. Since they trained the model with whitespace, the main function of this unit could simply

<i>Our models (wiki/in-domain)</i>			
LSTM		34.1/59.0	
RNN		24.3/24.0	
WordNLM		37.1/63.3	
<i>From the literature</i>			
KN5	40.0	Skipgram	48.0
Word RNN	45.0	Skipgram + RNNs	58.9
Word LSTM	56.0	PMI	61.4
LdTreeLSTM	60.7	Context-Embed	65.1

Table 5: Results on MSR Sentence Completion. For our models (top), we show accuracies for Wikipedia (left) and in-domain (right) training. We compare with language models from prior work (left): Kneser-Ney 5-gram model (Mikolov, 2012), Word RNN (Zweig et al., 2012), Word LSTM and LdTreeLSTM (Zhang et al., 2016). We further report models incorporating distributional encodings of semantics (right): Skipgram(+RNNs) from Mikolov et al. (2013a), the PMI-based model of Woods (2016), and the Context-Embedding-based approach of Melamud et al. (2016).

be to predict the very frequent whitespace character. We conjecture instead (like them) that the ability to segment the input into meaningful items is so important when processing language that CNLMs will specialize units for boundary tracking even when trained without whitespace.

To look for “boundary units”, we created a random set of 10,000 positions from the training set, balanced between those corresponding to a word-final character and those occurring word-initially or word-medially. We then computed, for each hidden unit, the Pearson correlation between its activations and a binary variable that takes value 1 in word-final position and 0 elsewhere. For each language and model (LSTM or RNN), we found very few units with a high correlation score, suggesting that the models have indeed specialized units for boundary tracking. We further study the units with the highest correlations, that are, for the LSTMs, 0.58 (English), 0.69 (German), and 0.57 (Italian). For the RNNs, the highest correlations are 0.40 (English), 0.46 (German and Italian).<sup>14</sup>

**Examples** We looked at the behaviour of the selected LSTM units qualitatively by extracting random sets of 40-character strings from the development partition of each language (left-aligned with word onsets) and plotting the corresponding boundary unit activations. Figure 2 reports illus-

<sup>14</sup>In an early version of this analysis, we arbitrarily imposed a minimum 0.70 correlation threshold, missing the presence of these units. We thank the reviewer who encouraged us to look further into the matter.

trative examples. In all languages, most peaks in activation mark word boundaries. However, other interesting patterns emerge. In English, we see how the unit reasonably treats *co-* and *pro-* in *co-produced* as separate elements, and it also posits a weaker boundary after the prefix *pro-*. As it proceeds left-to-right, with no information on what follows, the network posits a boundary after *but* in *Buttrich*. In the German example, we observe how the complex word *Hauptaufgabe* (‘main task’) is segmented into the morphemes *haupt*, *auf* and *gabe*. Similarly, in the final *transformati-* fragment, we observe a weak boundary after the prefix *trans*. In the pronoun *deren* ‘whose’, the case suffix *-n* is separated. In Italian, *in seguito a* is a lexicalized multi-word sequence meaning ‘following’ (literally: ‘in continuation to’). The boundary unit does not spike inside it. Similarly, the fixed expression *Sommo Pontefice* (referring to the Pope) does not trigger inner boundary unit activation spikes. On the other hand, we notice peaks after *di* and *mi* in *dimissioni*. Again, in left-to-right processing, the unit has a tendency to immediately posit boundaries when frequent function words are encountered.

**Detecting word boundaries** To gain a more quantitative understanding of how well the boundary unit is tracking word boundaries, we trained a single-parameter diagnostic classifier on the activation of the unit (the classifier simply sets an optimal threshold on the unit activation to separate word boundaries from word-internal positions). We ran two experiments. In the first, following standard practice, we trained and tested the classifier on uncontrolled running text. We used 1k characters for training, 1M for testing, both taken from the left-out Wikipedia test partitions. We will report F1 performance on this task.

We also considered a more cogent evaluation regime, in which we split training and test data so that the number of boundary and non-boundary conditions are balanced, and there is no overlap between training and test words. Specifically, we randomly selected positions from the test partitions of the Wikipedia corpus, such that half of these were the last character of a token, and the other half were not. We sampled the test data points subject to the constraint that the word (in the case of a boundary position) or word prefix (in the case of a word-internal position) ending at the selected character does not overlap with the train-

ing set. This ensures that a classifier cannot succeed by looking for encodings reflecting specific words. For each datapoint, we fed a substring of the 40 preceding characters to the CNLM. We collected 1,000 such points for training, and tested on 1M additional datapoints. In this case, we will report classification accuracy as figure of merit. For reference, in both experiments we also trained diagnostic classifiers on the *full* hidden layer of the LSTMs.

Looking at the F1 results on uncontrolled running text (Table 6), we observe first that the LSTM-based full-hidden-layer classifier has strong performance in all 3 languages, confirming that the LSTM model encodes boundary information. Moreover, in all languages, a large proportion of this performance is already accounted for by the single-parameter classifier using boundary unit activations. This confirms that tracking boundaries is important enough for the network to devote a specialized unit to this task. Full-layer RNN results are below LSTM-level but still strong. There is however a stronger drop from full-layer to single-unit classification. This is in line with the fact that, as reported above, the candidate RNN boundary units have lower boundary correlations than the LSTM ones.

Results for the balanced classifiers tested on new-word generalization are shown in Table 7 (because of the different nature of the experiments, these are not directly comparable to the F1 results in Table 6). Again, we observe a strong performance of the LSTM-based full-hidden-layer classifier across the board. The LSTM single-parameter classifier using boundary unit activations is also strong, even outperforming the full classifier in German. Moreover, in this more cogent setup, the single-unit LSTM classifier is at least competitive with the full-layer RNN classifier in all languages. The weaker results of RNNs in the word-centric tasks of the previous sections might in part be due to their poorer overall ability to track word boundaries, as specifically suggested by this stricter evaluation setup.

**Error analysis** As a final way to characterize the function and behaviour of the boundary units, we inspected the most frequent under- and over-segmentation errors made by the classifier based on the single boundary units, in the more difficult balanced task. We discuss German here, as it is the language where the classifier reaches highest ac-

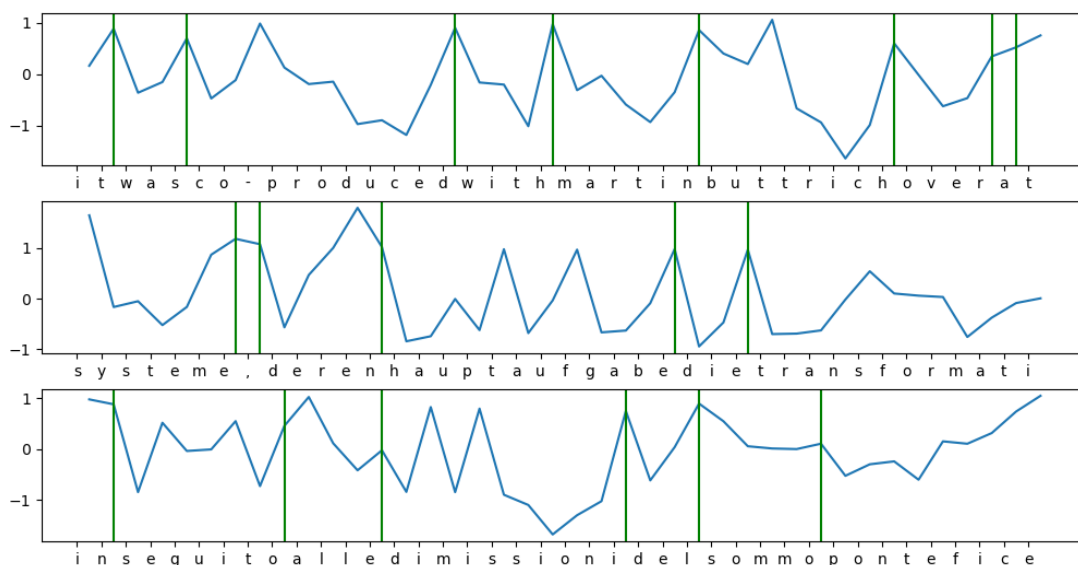


Figure 2: Examples of the LSTM CNLM boundary unit activation profile, with ground-truth word boundaries marked in green. English: *It was co-produced with Martin Buttrich over at...* German: *Systeme, deren Hauptaufgabe die transformati(-on)* ‘systems, whose main task is the transformation...’. Italian: *in seguito alle dimissioni del Sommo Pontefice* ‘following the resignation of the Supreme Pontiff...’.

	<i>LSTM single</i>	<i>LSTM full</i>	<i>RNN single</i>	<i>RNN full</i>
English	87.7	93.0	65.6	90.5
German	86.6	91.9	70.4	85.0
Italian	85.6	92.2	71.3	91.5

Table 6: F1 of single-unit and full-hidden-state word-boundary diagnostic classifiers, trained and tested on uncontrolled running text.

	<i>LSTM single</i>	<i>LSTM full</i>	<i>RNN single</i>	<i>RNN full</i>
English	77.5	90.0	65.9	76.8
German	80.8	79.7	67.0	75.8
Italian	75.5	82.9	71.4	75.9

Table 7: Accuracy of single-unit and full-hidden-state word-boundary diagnostic classifiers, trained and tested on balanced data requiring new-word generalization. Chance accuracy is at 50%.

curacy, and its tendency to have long, morphologically complex words makes it particularly interesting. However, similar patterns were also detected in Italian and, to a lesser extent, English (in the latter, there are fewer and less interpretable common oversegmentations, probably because words are on average shorter and morphology more limited).

Considering first the 30 most common undersegmentations, the large majority (24/30) are com-

mon sequences of grammatical terms or very frequent items that can sometimes be reasonably re-analyzed as single function words or adverbs (e.g., *bis zu*, ‘up to’ (lit. ‘until to’), *je nach* ‘depending on’ (lit. ‘per after’), *bis heute* ‘to date’ (lit. ‘until today’)). 3 cases are multi-word city names (*Los Angeles*). The final 3 cases interestingly involve *Bau* ‘building’ followed by *von* ‘of’ or genitive determiners *der/des*. In its eventive reading, this noun requires a patient licensed by either a preposition or the genitive determiner, e.g., *Bau der Mauer* ‘building of the wall’ (lit. ‘building the-GEN wall’). Apparently the model decided to absorb the case assigner into the form of the noun.

We looked next at the 30 most common oversegmentations, that is, at the substrings that were wrongly segmented out of the largest number of distinct words. We limited the analysis to those containing at least 3 characters, because shorter strings were ambiguous and hard to interpret. Among the top oversegmentations, 6 are prefixes that can also occur in isolation as prepositions or verb particles (*auf* ‘on’, *nach* ‘after’, etc.). 7 are content words that form many compounds (e.g., *haupt* ‘main’, occurring in *Hauptstadt* ‘capital’, *Hauptbahnhof* ‘main station’ etc.; *Land* ‘land’, occurring in *Deutschland* ‘Germany’, *Landkreis* ‘district’, etc.). Another 7 items can be classified as suffixes (e.g., *-lich* as in *südlich* ‘southern’, *wissenschaftlich* ‘scientific’), although their segmen-

tation is not always canonical (e.g., *-chaft* instead of the expected *-schaft* in *Wissenschaft* ‘science’). 4 very common function words are often wrongly segmented out of longer words (e.g., *sie* ‘she’ from *sieben* ‘seven’). The *kom* and *kon* cases are interesting, as the model segments them as stems (or stem fragments) in forms of the verbs *kommen* ‘to come’ and *können* ‘to be able to’, respectively (e.g., *kommt* and *konnte*), but it also treats them as pseudo-affixes elsewhere (*komponist* ‘composer’, *kontakt* ‘contact’). The remaining 3 oversegmentations, *rie*, *run* and *ter* don’t have any clear interpretation.

To conclude, the boundary unit, even when analyzed through the lens of a classifier that was optimized on word-level segmentation, is actually tracking salient linguistic boundaries at different levels. While in many cases these boundaries naturally coincide with words (hence the high classifier performance), the CNLM is also sensitive to frequent morphemes and compound elements, as well as to different types of multi-word expressions. This is in line with a view of wordhood as a useful but “soft”, emergent property, rather than a rigid primitive of linguistic processing.

## 5 Discussion

We probed the linguistic information induced by a character-level LSTM language model trained on unsegmented text. The model was found to possess implicit knowledge about a range of intuitively word-mediated phenomena, such as sensitivity to lexical categories and syntactic and shallow-semantics dependencies. A model initialized with a word vocabulary and fed tokenized input was in general superior, but the performance of the word-less model did not lag much behind, suggesting that word priors are helpful but not strictly required. A character-level RNN was less consistent than the LSTM, suggesting that the latter’s ability to track information across longer time spans is important to make the correct generalizations. The character-level models consistently outperformed n-gram controls, confirming they are tapping into more abstract patterns than local co-occurrence statistics.

As a first step towards understanding *how* character-level models handle supra-character phenomena, we searched and found specialized boundary-tracking units in them. These units are not only and not always sensitive to word bound-

aries, but also respond to other salient items, such as morphemes and multi-word expressions, in accordance with an “emergent” and flexible view of the basic constituents of language (Schiering et al., 2010).

Our results are preliminary in many ways. Our tests are relatively simple. We did not attempt, for example, to model long-distance agreement in presence of distractors, a challenging task even for humans (Gulordava et al., 2018). The results on number classification in German suggest that the models might not be capturing linguistic generalizations of the correct degree of abstractness, settling for shallower heuristics. Still, as a whole, our work suggests that a large corpus, combined with the weak priors encoded in an LSTM, might suffice to learn generalizations about word-mediated linguistic processes without a hard-coded word lexicon or explicit wordhood cues.

Nearly all contemporary linguistics recognizes a central role to the lexicon (see, e.g., Sag et al., 2003; Goldberg, 2005; Radford, 2006; Bresnan et al., 2016; Ježek, 2016, for very different perspectives). Linguistic formalisms assume that the lexicon is essentially a dictionary of words, possibly complemented by other units, not unlike the list of words and associated embeddings in a standard word-based NLM. Intriguingly, our CNLMs captured a range of lexical phenomena *without* anything resembling a word dictionary. Any information a CNLM might acquire about units larger than characters must be stored in its recurrent weights. This suggests a radically different and possibly more neurally plausible view of the lexicon as implicitly encoded in a distributed memory, that we intend to characterize more precisely and test in future work (similar ideas are being explored in a more applied NLP perspective, e.g., Gillick et al., 2016; Lee et al., 2017; Cherry et al., 2018).

Concerning the model input, we would like to study whether the CNLM successes crucially depend on the huge amount of training data it receives. Are word priors more important when learning from smaller corpora? In terms of comparison with human learning, the Wikipedia text we fed our CNLMs is far from what children acquiring a language would hear. Future work should explore character/phoneme-level learning from child-directed speech corpora. Still, by feeding our networks “grown-up” prose, we are ar-

guably making the job of identifying basic constituents harder than it might be when processing the simpler utterances of early child-directed speech (Tomasello, 2003).

As discussed, a rigid word notion is problematic both cross-linguistically (cf. polysynthetic and agglutinative languages) and within single linguistic systems (cf. the view that the lexicon hosts units at different levels of the linguistic hierarchy, from morphemes to large syntactic constructions, e.g., Jackendoff, 1997; Croft and Cruse, 2004; Goldberg, 2005). This study provided a necessary initial check that word-free models can account for phenomena traditionally seen as word-based. Future work should test whether such models can also account for grammatical patterns that are harder to capture in word-based formalisms, exploring both a typologically wider range of languages and a broader set of grammatical tests.

## Acknowledgments

We would like to thank Piotr Bojanowski, Alex Cristia, Kristina Gulordava, Urvashi Khandelwal, Germán Kruszewski, Sebastian Riedel, Hinrich Schütze and the anonymous reviewers for feedback and advice.

## References

- Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. 2017. Fine-grained analysis of sentence embeddings using auxiliary prediction tasks. In *Proceedings of ICLR Conference Track*, Toulon, France. Published online: <https://openreview.net/group?id=ICLR.cc/2017/conference>.
- Afra Alishahi, Marie Barking, and Grzegorz Chrupała. 2017. Encoding of phonology in a recurrent neural model of grounded speech. In *Proceedings of CoNLL*, pages 368–378, Vancouver, Canada.
- Moshe Bar. 2007. The proactive brain: using analogies and associations to generate predictions. *Trends in Cognitive Science*, 11(7):280–289.
- Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. What do neural machine translation models learn about morphology? In *Proceedings of ACL*, pages 861–872, Vancouver, Canada.
- Balthasar Bickel and Fernando Zúñiga. 2017. The ‘word’ in polysynthetic languages: Phonological and syntactic challenges. In Michael Fortescue, Marianne Mithun, and Nicholas Evans, editors, *Oxford Handbook of Polysynthesis*, pages 158–186. Oxford University Press, Oxford, UK.
- Piotr Bojanowski, Armand Joulin, and Tomas Mikolov. 2016. Alternative structures for character-level RNNs. In *Proceedings of ICLR Workshop Track*, San Juan, Puerto Rico. Published online: <https://openreview.net/group?id=ICLR.cc/2016/workshop>.
- Sabine Brants, Stefanie Dipper, Silvia Hansen, Wolfgang Lezius, and George Smith. 2002. The TIGER treebank. In *Proceedings of the workshop on treebanks and linguistic theories*, volume 168.
- Michael Brent and Timothy Cartwright. 1996. Distributional regularity and phonotactic constraints are useful for segmentation. *Cognition*, 61:93–125.
- Joan Bresnan, Ash Asudeh, Ida Toivonen, and Stephen Wechsler. 2016. *Lexical-Functional Syntax, 2nd ed.* Blackwell, Malden, MA.
- Joan Bresnan and Sam Mchombo. 1995. The lexical integrity principle: Evidence from Bantu. *Natural Language and Linguistic Theory*, pages 181–254.
- Colin Cherry, George Foster, Ankur Bapna, Orhan Firat, and Wolfgang Macherey. 2018. Revisiting character-based neural machine translation with capacity and compression. *arXiv preprint arXiv:1808.09943*.
- Noam Chomsky. 1970. Remarks on nominalization. In Roderick Jacobs and Peter Rosenbaum, editors, *Readings in English Transformational Grammar*, pages 184–221. Ginn, Waltham, MA.
- Morten Christiansen, Christopher Conway, and Suzanne Curtin. 2005. Multiple-cue integration in language acquisition: A connectionist model of speech segmentation and rule-like behavior. In James Minett and William Wang, editors,

- Language Acquisition, Change and Emergence: Essays in Evolutionary Linguistics*, pages 205–249. City University of Hong Kong Press, Hong Kong.
- Morten Christiansen, Allen Joseh, and Mark Seidenberg. 1998. Learning to segment speech using multiple cues: A connectionist model. *Language and Cognitive Processes*, 13(2/3):221–268.
- Andy Clark. 2016. *Surfing Uncertainty*. Oxford University Press, Oxford, UK.
- Alexis Conneau, Germán Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single \$&!#\* vector: Probing sentence embeddings for linguistic properties. In *Proceedings ACL*, pages 2126–2136, Melbourne, Australia.
- Ryan Cotterell, Sebastian J Mielke, Jason Eisner, and Brian Roark. 2018. Are all languages equally hard to language-model? In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, volume 2, pages 536–541.
- William Croft and Alan Cruse. 2004. *Cognitive Linguistics*. Cambridge University Press, Cambridge, UK.
- Anna-Maria Di Sciullo and Edwin Williams. 1987. *On the Definition of Word*. MIT Press, Cambridge, MA.
- Robert Dixon and Alexandra Aikhenvald, editors. 2002. *Word: A cross-linguistic typology*. Cambridge University Press, Cambridge, UK.
- Dudenredaktion. 2019. mit (Adverb). In *Duden online*. <https://www.duden.de/node/152710/revision/152746>, retrieved June 3, 2019.
- Jeffrey Elman. 1990. Finding structure in time. *Cognitive Science*, 14:179–211.
- Allyson Ettinger, Ahmed Elgohary, Colin Phillips, and Philip Resnik. 2018. Assessing composition in sentence vector representations. In *Proceedings of COLING*, pages 1790–1801, Santa Fe, NM.
- Robert Frank, Donald Mathis, and William Badecker. 2013. The acquisition of anaphora by simple recurrent networks. *Language Acquisition*, 20(3):181–227.
- Stefano Fusi, Earl Miller, and Mattia Rigotti. 2016. Why neurons mix: High dimensionality for higher cognition. *Current Opinion in Neurobiology*, 37:66–74.
- Daniela Gerz, Ivan Vulić, Edoardo Maria Ponti, Jason Naradowsky, Roi Reichart, and Anna Korhonen. 2018. Language modeling for morphologically rich languages: Character-aware modeling for word-level prediction. *Transactions of the Association for Computational Linguistics*, 6:451–465.
- Dan Gillick, Cliff Brunk, Oriol Vinyals, and Amarnag Subramanya. 2016. Multilingual language processing from bytes. In *Proceedings of NAACL-HLT*, pages 1296–1306.
- Frédéric Godin, Kris Demuynck, Joni Dambre, Wesley De Neve, and Thomas Demeester. 2018. Explaining character-aware neural networks for word-level prediction: Do they discover linguistic rules? In *Proceedings of EMNLP*, Brussels, Belgium. In press.
- Adele Goldberg. 2005. *Constructions at work: The nature of generalization in language*. Oxford University Press, Oxford, UK.
- Yoav Goldberg. 2017. *Neural Network Methods for Natural Language Processing*. Morgan & Claypool, San Francisco, CA.
- Sharon Goldwater, Thomas L. Griffiths, and Mark Johnson. 2009. A Bayesian framework for word segmentation: Exploring the effects of context. *Cognition*, 112(1):21–54.
- Alex Graves. 2014. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850v5.
- Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. Colorless green recurrent networks dream hierarchically. In *Proceedings of NAACL*, pages 1195–1205, New Orleans, LA.
- Martin Haspelmath. 2011. The indeterminacy of word segmentation and the nature of morphology and syntax. *Folia Linguistica*, 45(1):31–80.

- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- Dieuwke Hupkes, Sara Veldhoen, and Willem Zuidema. 2018. Visualisation and ‘diagnostic classifiers’ reveal how recurrent and recursive neural networks process hierarchical structure. *Journal of Artificial Intelligence Research*, 61:907–926.
- Ray Jackendoff. 1997. Twistin’ the night away. *Language*, 73:534–559.
- Ray Jackendoff. 2002. *Foundations of Language: Brain, Meaning, Grammar, Evolution*. Oxford University Press, Oxford, UK.
- Elisabetta Ježek. 2016. *The Lexicon: An Introduction*. Oxford University Press, Oxford, UK.
- Àkos Kádàr, Grzegorz Chrupała, and Afra Alishahi. 2017. Representation of linguistic form and function in recurrent neural networks. *Computational Linguistics*, 43(4):761–780.
- Herman Kamper, Aren Jansen, and Sharon Goldwater. 2016. Unsupervised word segmentation and lexicon discovery using acoustic word embeddings. *IEEE Transactions on Audio, Speech and Language Processing*, 24(4):669–679.
- Katharina Kann, Ryan Cotterell, and Hinrich Schütze. 2016. Neural morphological analysis: Encoding-decoding canonical segments. In *Proceedings of EMNLP*, pages 961–967, Austin, Texas.
- Yova Kementchedjheva and Adam Lopez. 2018. ‘Indicatements’ that character language models learn English morpho-syntactic units and regularities. In *Proceedings of the EMNLP BlackboxNLP Workshop*, pages 145–153, Brussels, Belgium.
- Yoon Kim, Yacine Jernite, David Sontag, and Alexander Rush. 2016. Character-aware neural language models. In *Proceedings of AAAI*, pages 2741–2749, Phoenix, AZ.
- Christo Kirov and Ryan Cotterell. 2018. Recurrent neural networks in linguistic theory: Revisiting Pinker and Prince (1988) and the past tense debate. *Transactions of the Association for Computational Linguistics*. In press.
- Patricia Kuhl. 2004. Early language acquisition: Cracking the speech code. *Nature Reviews Neuroscience*, 5(11):831–843.
- Jey Han Lau, Alexander Clark, and Shalom Lapin. 2017. Grammaticality, acceptability, and probability: A probabilistic view of linguistic knowledge. *Cognitive Science*, 41(5):1202–1241.
- Jason Lee, Kyunghyun Cho, and Thomas Hofmann. 2017. Fully character-level neural machine translation without explicit segmentation. *Transactions of the Association for Computational Linguistics*, 5:365–378.
- Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. Visualizing and understanding neural models in NLP. In *Proceedings of NAACL*, pages 681–691, San Diego, CA.
- Tal Linzen, Grzegorz Chrupała, and Afra Alishahi, editors. 2018. *Proceedings of the EMNLP BlackboxNLP Workshop*. ACL, Brussels, Belgium.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of LSTMs to learn syntax-sensitive dependencies. *Transactions of the Association for Computational Linguistics*, 4:521–535.
- Jessica Maye, Janet Werker, and LouAnn Gerken. 2002. Infant sensitivity to distributional information can affect phonetic discrimination. *Cognition*, 82(3):B101–B111.
- Thomas McCoy, Robert Frank, and Tal Linzen. 2018. Revisiting the poverty of the stimulus: Hierarchical generalization without a hierarchical bias in recurrent neural networks. In *Proceedings of CogSci*, pages 2093–2098, Madison, WI.
- Ryan McDonald, Joakim Nivre, Yvonne Quirnbach-Brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith Hall, Slav Petrov, Hao Zhang, Oscar Täckström, et al. 2013. Universal dependency annotation for multilingual parsing. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 92–97.



- Oren Melamud, Jacob Goldberger, and Ido Dagan. 2016. context2vec: Learning generic context embedding with bidirectional lstm. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 51–61.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2018. An analysis of neural language modeling at multiple scales. *arXiv preprint arXiv:1803.08240*.
- Tomas Mikolov. 2012. *Statistical language models based on neural networks*. Dissertation, Brno University of Technology.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781.
- Tomas Mikolov, Ilya Sutskever, Anoop Deoras, Hai-Son Le, Stefan Kombrink, and Jan Cernocký. 2011. Subword language modeling with neural networks. <http://www.fit.vutbr.cz/~imikolov/rnnlm/>.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013b. Linguistic regularities in continuous space word representations. In *Proceedings of NAACL*, pages 746–751, Atlanta, Georgia.
- Joe Pater. 2018. Generative linguistics and neural networks at 60: Foundation, friction, and fusion. *Language*. In press.
- Alec Radford, Rafal Józefowicz, and Ilya Sutskever. 2017. Learning to generate reviews and discovering sentiment. *CoRR*, abs/1704.01444.
- Andrew Radford. 2006. Minimalist syntax revisited. <http://www.public.asu.edu/~gelderren/Radford2009.pdf>.
- Ivan Sag, Thomas Wasow, and Emily Bender. 2003. *Syntactic Theory: A Formal Introduction*. CSLI, Stanford, CA.
- René Schiering, Balthasar Bickel, and Kristine Hildebrandt. 2010. The prosodic word is not universal, but emergent. *Journal of Linguistics*, 46(3):657–709.
- Helmut Schmid. 1999. Improvements in part-of-speech tagging with an application to german. In *Natural language processing using very large corpora*, pages 13–25. Springer.
- Hinrich Schütze. 2017. Nonsymbolic text representation. In *Proceedings of EACL*, pages 785–796, Valencia, Spain.
- Rico Sennrich. 2017. How grammatical is character-level neural machine translation? assessing MT quality with contrastive translation pairs. In *Proceedings of EACL (Short Papers)*, pages 376–382, Valencia, Spain.
- Xing Shi, Inkit Padhi, and Kevin Knight. 2016. Does string-based neural MT learn source syntax? In *Proceedings of EMNLP*, pages 1526–1534, Austin, Texas.
- Ilya Sutskever, James Martens, and Geoffrey Hinton. 2011. Generating text with recurrent neural networks. In *Proceedings of ICML*, pages 1017–1024, Bellevue, WA.
- 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.
- Michael Tomasello. 2003. *Constructing a Language: A Usage-Based Theory of Language Acquisition*. Harvard University Press, Cambridge, MA.
- Edwin Williams. 2007. Dumping lexicalism. In Gillian Ramchand and Charles Reiss, editors, *The Oxford Handbook of Linguistic Interfaces*. Oxford University Press, Oxford, UK.
- Aubrie Woods. 2016. Exploiting linguistic features for sentence completion. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 438–442.
- Xingxing Zhang, Liang Lu, and Mirella Lapata. 2016. Top-down tree long short-term memory networks. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 310–320.
- Geoffrey Zweig and Christopher Burges. 2011. The Microsoft Research sentence completion

challenge. Technical Report MSR-TR-2011-129, Microsoft Research.

Geoffrey Zweig, John C Platt, Christopher Meek, Christopher JC Burges, Ainur Yessenalina, and Qiang Liu. 2012. Computational approaches to sentence completion. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 601–610. Association for Computational Linguistics.

	LSTM			RNN			WordNLM		
	En.	Ge.	It.	En.	Ge.	It.	En.	Ge.	It.
Batch Size	128	512	128	256	256	256	128	128	128
Embedding Size	200	100	200	200	50	50	1024	200	200
Dimension	1024	1024	1024	2048	2048	2048	1024	1024	1024
Layers	3	2	2	2	2	2	2	2	2
Learning Rate	3.6	2.0	3.2	0.01	0.1	0.1	1.1	0.9	1.2
Decay	0.95	1.0	0.98	0.9	0.95	0.95	1.0	1.0	0.98
BPTT Length	80	50	80	50	30	30	50	50	50
Hidden Dropout	0.01	0.0	0.0	0.05	0.0	0.0	0.15	0.15	0.05
Embedding Dropout	0.0	0.01	0.0	0.01	0.0	0.0	0.0	0.1	0.0
Input Dropout	0.001	0.0	0.0	0.001	0.01	0.01	0.01	0.001	0.01
Nonlinearity	–	–	–	ReLU	tanh	tanh	–	–	–

Table 8: Chosen hyperparameters