

## Abstract

When parsing morphologically-rich languages with neural models, it is beneficial to model input at the character level, and it has been claimed that this is because character-level models learn morphology. We test these claims by comparing character-level models to an oracle with access to explicit morphological analysis on twelve languages with varying morphological typologies. Our results highlight many strengths of character-level models, but also show that they are poor at disambiguating some words, particularly in the face of case syncretism. We then demonstrate that explicitly modeling morphological case improves our best model, showing that character-level models can benefit from targeted forms of explicit morphological modeling.

## Introduction

Modeling language input at the character level BIBREF0 , BIBREF1 is effective for many NLP tasks, and often produces better results than modeling at the word level. For parsing, ballesteros-dyer-smith:2015:EMNLP have shown that character-level input modeling is highly effective on morphologically-rich languages, and the three best systems on the 45 languages of the CoNLL 2017 shared task on universal dependency parsing all use character-level models BIBREF2 , BIBREF3 , BIBREF4 , BIBREF5 , showing that they are effective across many typologies.

The effectiveness of character-level models in morphologically-rich languages has raised a question and indeed debate about explicit modeling of morphology in NLP. BIBREF0 propose that “prior information regarding morphology ... among others, should be incorporated” into character-level models, while BIBREF6 counter that it is “unnecessary to consider these prior information” when modeling characters.

Whether we need to explicitly model morphology is a question whose answer has a real cost: as [ballesteros-dyer-smith:2015:EMNLP](#) note, morphological annotation is expensive, and this expense could be reinvested elsewhere if the predictive aspects of morphology are learnable from strings.

Do character-level models learn morphology? We view this as an empirical claim requiring empirical evidence. The claim has been tested implicitly by comparing character-level models to word lookup models [BIBREF7](#) , [BIBREF8](#) . In this paper, we test it explicitly, asking how character-level models compare with an oracle model with access to morphological annotations. This extends experiments showing that character-aware language models in Czech and Russian benefit substantially from oracle morphology [BIBREF9](#) , but here we focus on dependency parsing (§ "Dependency parsing model" )—a task that benefits substantially from morphological knowledge—and we experiment with twelve languages using a variety of techniques to probe our models.

Our summary finding is that character-level models lag the oracle in nearly all languages (§ "Experiments" ). The difference is small, but suggests that there is value in modeling morphology. When we tease apart the results by part of speech and dependency type, we trace the difference back to the character-level model's inability to disambiguate words even when encoded with arbitrary context (§ "Analysis" ). Specifically, it struggles with case syncretism, in which noun case—and thus syntactic function—is ambiguous. We show that the oracle relies on morphological case, and that a character-level model provided only with morphological case rivals the oracle, even when case is provided by another predictive model (§ "Characters and case syncretism" ). Finally, we show that the crucial morphological features vary by language (§ "Understanding head selection" ).

## Dependency parsing model

We use a neural graph-based dependency parser combining elements of two recent models [BIBREF10](#) ,

BIBREF11 . Let  $w = w_1, \dots, w_{|w|}$  be an input sentence of length  $|w|$  and let  $w_0$  denote an artificial Root token. We represent the  $i$ th input token  $w_i$  by concatenating its word representation (§ "Computing word representations" ),  $\mathbf{e}(w_i)$  and part-of-speech (POS) representation,  $\mathbf{p}_i$  . Using a semicolon  $(;)$  to denote vector concatenation, we have:

$$\mathbf{x}_i = [\mathbf{e}(w_i); \mathbf{p}_i] \quad (\text{Eq. 2})$$

We call  $\mathbf{x}_i$  the embedding of  $w_i$  since it depends on context-independent word and POS representations. We obtain a context-sensitive encoding  $\mathbf{h}_i$  with a bidirectional LSTM (bi-LSTM), which concatenates the hidden states of a forward and backward LSTM at position  $i$  . Using  $\mathbf{h}_i^f$  and  $\mathbf{h}_i^b$  respectively to denote these hidden states, we have:

$$\mathbf{h}_i = [\mathbf{h}_i^f; \mathbf{h}_i^b] \quad (\text{Eq. 3})$$

We use  $\mathbf{h}_i$  as the final input representation of  $w_i$  .

## Head prediction

For each word  $w_i$  , we compute a distribution over all other word positions  $j \in \{0, \dots, |w| \}$  / $i$  denoting the probability that  $w_j$  is the headword of  $w_i$  .

$$P_{\text{head}}(w_j \mid w_i, w) = \frac{\exp(a(\mathbf{h}_i, \mathbf{h}_j))}{\sum_{j' \neq i} \exp(a(\mathbf{h}_i, \mathbf{h}_{j'}))} \quad (\text{Eq. 5})$$

Here,  $a$  is a neural network that computes an association between  $w_i$  and  $w_j$  using model parameters  $\mathbf{U}_a, \mathbf{W}_a$  and  $\mathbf{v}_a$  .

$$a(\mathbf{h}_i, \mathbf{h}_j) = \mathbf{v}_a \tanh(\mathbf{U}_a \mathbf{h}_i + \mathbf{W}_a \mathbf{h}_j)$$

(Eq. 6)

## Label prediction

Given a head prediction for word  $w_i$ , we predict its syntactic label  $\ell_k \in L$  using a similar network.

$$P_{\{label\}}(\ell_k \mid w_i, w_j, w) = \frac{\exp(f(\mathbf{h}_i, \mathbf{h}_j)[k])}{\sum_{k'=1}^{|L|} \exp(f(\mathbf{h}_i, \mathbf{h}_{\{j\}})[k'])} \quad (\text{Eq. 8})$$

where  $L$  is the set of output labels and  $f$  is a function that computes label score using model parameters  $\mathbf{U}_{\ell}$ ,  $\mathbf{W}_{\ell}$ , and  $\mathbf{V}_{\ell}$ :

$$f(\mathbf{h}_i, \mathbf{h}_j) = \mathbf{V}_{\ell} \tanh(\mathbf{U}_{\ell} \mathbf{h}_i + \mathbf{W}_{\ell} \mathbf{h}_j) \quad (\text{Eq. 9})$$

The model is trained to minimize the summed cross-entropy losses of both head and label prediction. At test time, we use the Chu-Liu-Edmonds algorithm to ensure well-formed, possibly non-projective trees.

## Computing word representations

We consider several ways to compute the word representation  $\mathbf{e}(w_i)$  in Eq. 2:

Every word type has its own learned vector representation.

Characters are composed using a bi-LSTM BIBREF0 , and the final states of the forward and backward LSTMs are concatenated to yield the word representation.

Characters are composed using a convolutional neural network BIBREF1 .

Character trigrams are composed using a bi-LSTM, an approach that we previously found to be effective across typologies BIBREF9 .

We treat the morphemes of a morphological annotation as a sequence and compose them using a bi-LSTM. We only use universal inflectional features defined in the UD annotation guidelines. For example, the morphological annotation of “chases” is  $\langle \text{chase, person=3rd, num-SG, tense=Pres} \rangle$  .

For the remainder of the paper, we use the name of model as shorthand for the dependency parser that uses that model as input (Eq. 2 ).

We experiment on twelve languages with varying morphological typologies (Table 1 ) in the Universal Dependencies (UD) treebanks version 2.0 BIBREF14 . Note that while Arabic and Hebrew follow a root & pattern typology, their datasets are unvocalized, which might reduce the observed effects of this typology. Following common practice, we remove language-specific dependency relations and multiword token annotations. We use gold sentence segmentation, tokenization, universal POS (UPOS), and morphological (XFEATS) annotations provided in UD.

Our Chainer BIBREF15 implementation encodes words (Eq. 3 ) in two-layer bi-LSTMs with 200 hidden units, and uses 100 hidden units for head and label predictions (output of Eqs. 4 and 6). We set batch size to 16 for char-cnn and 32 for other models following a grid search. We apply dropout to the

embeddings (Eq. 2 ) and the input of the head prediction. We use Adam optimizer with initial learning rate 0.001 and clip gradients to 5, and train all models for 50 epochs with early stopping. For the word model, we limit our vocabulary to the 20K most frequent words, replacing less frequent words with an unknown word token. The char-lstm, trigram-lstm, and oracle models use a one-layer bi-LSTM with 200 hidden units to compose subwords. For char-cnn, we use the small model setup of kim2015.

Table 2 presents test results for every model on every language, establishing three results. First, they support previous findings that character-level models outperform word-based models—indeed, the char-lstm model outperforms the word model on LAS for all languages except Hindi and Urdu for which the results are identical. Second, they establish strong baselines for the character-level models: the char-lstm generally obtains the best parsing accuracy, closely followed by char-cnn. Third, they demonstrate that character-level models rarely match the accuracy of an oracle model with access to explicit morphology. This reinforces a finding of BIBREF9 : character-level models are effective tools, but they do not learn everything about morphology, and they seem to be closer to oracle accuracy in agglutinative rather than in fusional languages.

In character-level models, orthographically similar words share many parameters, so we would expect these models to produce good representations of OOV words that are morphological variants of training words. Does this effect explain why they are better than word-level models?

Table 3 shows how the character model improves over the word model for both non-OOV and OOV words. On the agglutinative languages Finnish and Turkish, where the OOV rates are 23% and 24% respectively, we see the highest LAS improvements, and we see especially large improvements in accuracy of OOV words. However, the effects are more mixed in other languages, even with relatively high OOV rates. In particular, languages with rich morphology like Czech, Russian, and (unvocalised) Arabic see more improvement than languages with moderately rich morphology and high OOV rates like

Portuguese or Spanish. This pattern suggests that parameter sharing between pairs of observed training words can also improve parsing performance. For example, if “dog” and “dogs” are observed in the training data, they will share activations in their context and on their common prefix.

Let's turn to our main question: what do character-level models learn about morphology? To answer it, we compare the oracle model to char-lstm, our best character-level model.

In the oracle, morphological annotations disambiguate some words that the char-lstm must disambiguate from context. Consider these Russian sentences from [baerman-brown-corbett-2005](#):

Maša čitaet pis'mo

Masha reads letter

`Masha reads a letter.'

Na stole ležit pis'mo

on table lies letter

`There's a letter on the table.' Pis'mo (“letter”) acts as the subject in ( [UID28](#) ), and as object in ( [UID28](#) ).

This knowledge is available to the oracle via morphological case: in ( [UID28](#) ), the case of pis'mo is nominative and in ( [UID28](#) ) it is accusative. Could this explain why the oracle outperforms the character model?

To test this, we look at accuracy for word types that are empirically ambiguous—those that have more

than one morphological analysis in the training data. Note that by this definition, some ambiguous words will be seen as unambiguous, since they were seen with only one analysis. To make the comparison as fair as possible, we consider only words that were observed in the training data. Figure 1 compares the improvement of the oracle on ambiguous and seen unambiguous words, and as expected we find that handling of ambiguous words improves with the oracle in almost all languages. The only exception is Turkish, which has the least training data.

Now we turn to a more fine-grained analysis conditioned on the annotated part-of-speech (POS) of the dependent. We focus on four languages where the oracle strongly outperforms the best character-level model on the development set: Finnish, Czech, German, and Russian. We consider five POS categories that are frequent in all languages and consistently annotated for morphology in our data: adjective (ADJ), noun (NOUN), pronoun (PRON), proper noun (PROPN), and verb (VERB).

Table 4 shows that the three noun categories—ADJ, PRON, and PROPN—benefit substantially from oracle morphology, especially for the three fusional languages: Czech, German, and Russian.

We analyze results by the dependency type of the dependent, focusing on types that interact with morphology: root, nominal subjects (nsubj), objects (obj), indirect objects (iobj), nominal modifiers (nmod), adjectival modifier (amod), obliques (obl), and (syntactic) case markings (case).

Figure 2 shows the differences in the confusion matrices of the char-lstm and oracle for those words on which both models correctly predict the head. The differences on Finnish are small, which we expect from the similar overall LAS of both models. But for the fusional languages, a pattern emerges: the char-lstm consistently underperforms the oracle on nominal subject, object, and indirect object dependencies—labels closely associated with noun categories. From inspection, it appears to frequently mislabel objects as nominal subjects when the dependent noun is morphologically ambiguous. For



example, in the sentence of Figure 3 , Gelände (“terrain”) is an object, but the char-lstm incorrectly predicts that it is a nominal subject. In the training data, Gelände is ambiguous: it can be accusative, nominative, or dative.

In German, the char-lstm frequently confuses objects and indirect objects. By inspection, we found 21 mislabeled cases, where 20 of them would likely be correct if the model had access to morphological case (usually dative). In Czech and Russian, the results are more varied: indirect objects are frequently mislabeled as objects, obliques, nominal modifiers, and nominal subjects. We note that indirect objects are relatively rare in these data, which may partly explain their frequent mislabeling.

So far, we've seen that for our three fusional languages—German, Czech, and Russian—the oracle strongly outperforms a character model on nouns with ambiguous morphological analyses, particularly on core dependencies: nominal subjects, objects and indirect objects. Since the nominative, accusative, and dative morphological cases are strongly (though not perfectly) correlated with these dependencies, it is easy to see why the morphologically-aware oracle is able to predict them so well. We hypothesized that these cases are more challenging for the character model because these languages feature a high degree of syncretism—functionally distinct words that have the same form—and in particular case syncretism. For example, referring back to examples ( UID28 ) and ( UID28 ), the character model must disambiguate *pis'mo* from its context, whereas the oracle can directly disambiguate it from a feature of the word itself.

To understand this, we first designed an experiment to see whether the char-lstm could successfully disambiguate noun case, using a method similar to BIBREF8 . We train a neural classifier that takes as input a word representation from the trained parser and predicts a morphological feature of that word—for example that its case is nominative (Case=Nom). The classifier is a feedforward neural network with one hidden layer, followed by a ReLU non-linearity. We consider two representations of each word: its

embedding (  $\mathbf{x}_i$  ; Eq. 2 ) and its encoding (  $\mathbf{h}_i$  ; Eq. 3 ). To understand the importance of case, we consider it alongside number and gender features as well as whole feature bundles.

Table 5 shows the results of morphological feature classification on Czech; we found very similar results in German and Russian (Appendix "Results on morphological tagging" ). The oracle embeddings have almost perfect accuracy—and this is just what we expect, since the representation only needs to preserve information from its input. The char-lstm embeddings perform well on number and gender, but less well on case. This results suggest that the character-level models still struggle to learn case when given only the input text. Comparing the char-lstm with a baseline model which predicts the most frequent feature for each type in the training data, we observe that both of them show similar trends even though character models slightly outperforms the baseline model.

The classification results from the encoding are particularly interesting: the oracle still performs very well on morphological case, but less well on other features, even though they appear in the input. In the character model, the accuracy in morphological prediction also degrades in the encoding—except for case, where accuracy on case improves by 12%.

These results make intuitive sense: representations learn to preserve information from their input that is useful for subsequent predictions. In our parsing model, morphological case is very useful for predicting dependency labels, and since it is present in the oracle's input, it is passed almost completely intact through each representation layer. The character model, which must disambiguate case from context, draws as much additional information as it can from surrounding words through the LSTM encoder. But other features, and particularly whole feature bundles, are presumably less useful for parsing, so neither model preserves them with the same fidelity.

Our analysis indicates that case is important for parsing, so it is natural to ask: Can we improve the neural model by explicitly modeling case? To answer this question, we ran a set of experiments, considering two ways to augment the char-lstm with case information: multitask learning BIBREF16 and a pipeline model in which we augment the char-lstm model with either predicted or gold case. For example, we use  $\langle p, i, z, z, a, \text{Nom} \rangle$  to represent pizza with nominative case. For MTL, we follow the setup of BIBREF17 and BIBREF18 . We increase the biLSTMs layers from two to four and use the first two layers to predict morphological case, leaving out the other two layers specific only for parser. For the pipeline model, we train a morphological tagger to predict morphological case (Appendix "Morphological tagger" ). This tagger does not share parameters with the parser.

Table 6 summarizes the results on Czech, German, and Russian. We find augmenting the char-lstm model with either oracle or predicted case improve its accuracy, although the effect is different across languages. The improvements from predicted case results are interesting, since in non-neural parsers, predicted case usually harms accuracy BIBREF19 . However, we note that our taggers use gold POS, which might help. The MTL models achieve similar or slightly better performance than the character-only models, suggesting that supplying case in this way is beneficial. Curiously, the MTL parser is worse than the pipeline parser, but the MTL case tagger is better than the pipeline case tagger (Table 7 ). This indicates that the MTL model must learn to encode case in the model's representation, but must not learn to effectively use it for parsing. Finally, we observe that augmenting the char-lstm with either gold or predicted case improves the parsing performance for all languages, and indeed closes the performance gap with the full oracle, which has access to all morphological features. This is especially interesting, because it shows using carefully targeted linguistic analyses can improve accuracy as much as wholesale linguistic analysis.

The previous experiments condition their analysis on the dependent, but dependency is a relationship between dependents and heads. We also want to understand the importance of morphological features to

the head. Which morphological features of the head are important to the oracle?

To see which morphological features the oracle depends on when making predictions, we augmented our model with a gated attention mechanism following kuncoro-EtAl:2017:EACLLong. Our new model attends to the morphological features of candidate head  $w_j$  when computing its association with dependent  $w_i$  (Eq. 5), and morpheme representations are then scaled by their attention weights to produce a final representation.

Let  $f_{i1}, \dots, f_{ik}$  be the  $k$  morphological features of  $w_i$ , and denote by  $\mathbf{f}_{i1}, \dots, \mathbf{f}_{ik}$  their corresponding feature embeddings. As in § "Dependency parsing model",  $\mathbf{h}_i$  and  $\mathbf{h}_j$  are the encodings of  $w_i$  and  $w_j$ , respectively. The morphological representation  $\mathbf{m}_j$  of  $w_j$  is:

$$\mathbf{m}_j = [\mathbf{f}_{j1}, \dots, \mathbf{f}_{jk}]^{\top} \mathbf{k} \quad (\text{Eq. 43})$$

where  $\mathbf{k}$  is a vector of attention weights:

$$\mathbf{k} = \text{softmax}([\mathbf{f}_{j1}, \dots, \mathbf{f}_{jk}]^{\top} \mathbf{V} \mathbf{h}_i) \quad (\text{Eq. 44})$$

The intuition is that dependent  $w_i$  can choose which morphological features of  $w_j$  are most important when deciding whether  $w_j$  is its head. Note that this model is asymmetric: a word only attends to the morphological features of its (single) parent, and not its (many) children, which may have different functions.

We combine the morphological representation with the word's encoding via a sigmoid gating mechanism.

$$\mathbf{z}_j = \mathbf{g} \odot \mathbf{h}_j + (1 - \mathbf{g}) \odot \mathbf{m}_j \\ \mathbf{g} = \sigma(\mathbf{W}_1 \mathbf{h}_j + \mathbf{W}_2 \mathbf{m}_j) \quad (\text{Eq. 46})$$

where  $\odot$  denotes element-wise multiplication. The gating mechanism allows the model to choose between the computed word representation and the weighted morphological representations, since for some dependencies, morphological features of the head might not be important. In the final model, we replace Eq. 5 and Eq. 6 with the following:

$$P_{\text{head}}(\mathbf{w}_j | \mathbf{w}_i, \mathbf{w}) = \frac{\exp(a(\mathbf{h}_i, \mathbf{z}_j))}{\sum_{j'=0}^N \exp(a(\mathbf{h}_i, \mathbf{z}_{j'}))} \\ a(\mathbf{h}_i, \mathbf{z}_j) = \mathbf{v}_a \tanh(\mathbf{U}_a \mathbf{h}_i + \mathbf{W}_a \mathbf{z}_j) \quad (\text{Eq. 47})$$

The modified label prediction is:

$$P_{\text{label}}(\ell_k | \mathbf{w}_i, \mathbf{w}_j, \mathbf{w}) = \frac{\exp(f(\mathbf{h}_i, \mathbf{z}_j)[k])}{\sum_{k'=0}^{|L|} \exp(f(\mathbf{h}_i, \mathbf{z}_{j'})[k'])} \quad (\text{Eq. 48})$$

where  $f$  is again a function to compute label score:

$$f(\mathbf{h}_i, \mathbf{z}_j) = \mathbf{V}_\ell \tanh(\mathbf{U}_\ell \mathbf{h}_i + \mathbf{W}_\ell \mathbf{z}_j) \quad (\text{Eq. 49})$$

We trained our augmented model (oracle-attn) on Finnish, German, Czech, and Russian. Its accuracy is very similar to the oracle model (Table 8 ), so we obtain a more interpretable model with no change to our main results.

Next, we look at the learned attention vectors to understand which morphological features are important, focusing on the core arguments: nominal subjects, objects, and indirect objects. Since our model knows the case of each dependent, this enables us to understand what features it seeks in potential heads for each case. For simplicity, we only report results for words where both head and label predictions are correct.

Figure 4 shows how attention is distributed across multiple features of the head word. In Czech and Russian, we observe that the model attends to Gender and Number when the noun is in nominative case. This makes intuitive sense since these features often signal subject-verb agreement. As we saw in earlier experiments, these are features for which a character model can learn reliably good representations. For most other dependencies (and all dependencies in German), Lemma is the most important feature, suggesting a strong reliance on lexical semantics of nouns and verbs. However, we also notice that the model sometimes attends to features like Aspect, Polarity, and VerbForm—since these features are present only on verbs, we suspect that the model may simply use them as convenient signals that a word is verb, and thus a likely head for a given noun.

Character-level models are effective because they can represent OOV words and orthographic regularities of words that are consistent with morphology. But they depend on context to disambiguate words, and for some words this context is insufficient. Case syncretism is a specific example that our analysis identified, but the main results in Table 2 hint at the possibility that different phenomena are at play in different languages.

While our results show that prior knowledge of morphology is important, they also show that it can be used in a targeted way: our character-level models improved markedly when we augmented them only with case. This suggests a pragmatic reality in the middle of the wide spectrum between pure machine learning from raw text input and linguistically-intensive modeling: our new models don't need all prior

linguistic knowledge, but they clearly benefit from some knowledge in addition to raw input. While we used a data-driven analysis to identify case syncretism as a problem for neural parsers, this result is consistent with previous linguistically-informed analyses BIBREF20 , BIBREF19 . We conclude that neural models can still benefit from linguistic analyses that target specific phenomena where annotation is likely to be useful.

Clara Vania is supported by the Indonesian Endowment Fund for Education (LPDP), the Centre for Doctoral Training in Data Science, funded by the UK EPSRC (grant EP/L016427/1), and the University of Edinburgh. We would like to thank Yonatan Belinkov for the helpful discussion regarding morphological tagging experiments. We thank Sameer Bansal, Marco Damonte, Denis Emelin, Federico Fancellu, SORCHA Gilroy, Jonathan Mallinson, Joana Ribeiro, Naomi Saphra, Ida Szubert, Sabine Weber, and the anonymous reviewers for helpful discussion of this work and comments on previous drafts of the paper.

We adapt the parser's encoder architecture for our morphological tagger. Following notation in Section "Dependency parsing model" , each word  $w_i$  is represented by its context-sensitive encoding,  $\text{vec}(w_i)$  (Eq. 3 ). The encodings are then fed into a feed-forward neural network with two hidden layers—each has a ReLU non-linearity—and an output layer mapping the to the morphological tags, followed by a softmax. We set the size of the hidden layer to 100 and use dropout probability 0.2. We use Adam optimizer with initial learning rate 0.001 and clip gradients to 5. We train each model for 20 epochs with early stopping. The model is trained to minimize the cross-entropy loss.

Since we do not have additional data with the same annotations, we use the same UD dataset to train our tagger. To prevent overfitting, we only use the first 75% of training data for training. After training the taggers, we predict the case for the training, development, and test sets and use them for dependency parsing.

Table 9 and 10 present morphological tagging results for German and Russian. We found that German and Russian have similar pattern to Czech (Table 5 ), where morphological case seems to be preserved in the encoder because they are useful for dependency parsing. In these three fusional languages, contextual information helps character-level model to predict the correct case. However, its performance still behind the oracle.

We observe a slightly different pattern on Finnish results (Table 11 ). The character embeddings achieves almost similar performance as the oracle embeddings. This results highlights the differences in morphological process between Finnish and the other fusional languages. We observe that performance of the encoder representations are slightly worse than the embeddings.