

Character Eyes: Seeing Language through Character-Level Taggers

Yuval Pinter

School of Interactive Computing
Georgia Institute of Technology
uyp@gatech.edu

Marc Marone *

Microsoft
mmarone6@gatech.edu

Jacob Eisenstein

Facebook AI Research
jacobeisenstein@fb.com

Abstract

Character-level models have been used extensively in recent years in NLP tasks as both supplements and replacements for closed-vocabulary token-level word representations. In one popular architecture, character-level LSTMs are used to feed token representations into a sequence tagger predicting token-level annotations such as part-of-speech (POS) tags.

In this work, we examine the behavior of POS taggers across languages from the perspective of individual hidden units within the character LSTM. We aggregate the behavior of these units into language-level metrics which quantify the challenges that taggers face on languages with different morphological properties, and identify links between synthesis and affixation preference and emergent behavior of the hidden tagger layer. In a comparative experiment, we show how modifying the balance between forward and backward hidden units affects model arrangement and performance in these types of languages.

1 Introduction

Subword vector representations are now a standard part of neural architectures for natural language processing (e.g., [Bojanowski et al., 2017](#); [Peters et al., 2018](#)). In particular, character representations have been shown to handle out-of-vocabulary words in supervised tagging tasks ([Ling et al., 2015](#); [Lample et al., 2016](#)). These advantages generalize across multiple languages, where morphological formation may differ greatly but the character composition of words remains a relatively reliable primitive ([Plank et al., 2016](#)).

While the advantages of character-level models are readily apparent, existing evaluation methods fail to explain the mechanism by which these models encode linguistic knowledge about morphology and orthography. Different languages exhibit

character-word correspondence in very different patterns, and yet the bi-directional LSTM appears to be, or is assumed to be, capable of capturing them all. In large multilingual settings, it is not uncommon to tune hyperparameters on a handful of languages, and apply them to the rest (e.g., [Pinter et al., 2017](#)).

In this work, we challenge this implicit generalization. We train character-based sequence taggers on a large selection of languages exhibiting various strategies for word formation, and subject the resulting models to a novel analysis of the behavior of individual units in the character-level Bi-LSTM hidden layer. This reveals differences in the ability of the Bi-LSTM architecture to identify parts-of-speech, based on typological properties: hidden layers trained on agglutinative languages find more regularities on the character level than in fusional languages; languages that are suffix-heavy give a stronger signal to the backward-facing hidden units, and vice versa for prefix-heavy languages. In short, character-level recurrent networks function differently depending on how each language expresses morphosyntactic properties in characters.

These empirical results motivate a novel Bi-LSTM architecture, in which the number of hidden units is unbalanced across the forward and backward directions. We find empirical correspondence between the analytical findings above and performance of such unbalanced Bi-LSTM models, allowing us to translate the typological properties of a language into concrete recommendations for model selection.¹

2 Related Work

Several recent papers attempt to explain neural network performance by investigating hidden state activation patterns on auxiliary or downstream

*Work done while at Georgia Institute of Technology.

¹<https://github.com/ruyimarone/character-eyes>

tasks. On the word level, [Linzen et al. \(2016\)](#) trained LSTM language models, evaluated their performance on grammatical agreement detection, and analyzed activation patterns within specific hidden units. We build on this analysis strategy as we aggregate (character-) sequence activation patterns across all hidden units in a model into quantitative measures.

Substantial prior work exists on the character level as well ([Karpathy et al., 2015](#); [Vania and Lopez, 2017](#); [Kementchedjhieva and Lopez, 2018](#); [Gerz et al., 2018](#)). [Smith et al. \(2018\)](#) examined the character component in multilingual parsing models empirically, comparing it to the contribution of POS embeddings and pre-trained embeddings. [Chaudhary et al. \(2018\)](#) leveraged cross-lingual character-level correspondence to train NER models for low-resource languages. [Godin et al. \(2018\)](#), compared CNN and LSTM character models on a type-level prediction task on three languages, using the post-network softmax values to see which models identify useful character sequences. Unlike their analysis, we examine a more applied token-level task (POS tagging), and focus on the hidden states within the LSTM model in order to analyze its raw view of word composition.

Our analysis assumes a characterization of unit roles, where each hidden unit is observed to have some specific function. Findings from [Linzen et al. \(2016\)](#) and others suggest that a single hidden unit can learn to track complex syntactic rules. [Radford et al. \(2017\)](#) found that a character-level language model can implicitly assign a single unit to track sentiment, without being directly supervised. [Kementchedjhieva and Lopez \(2018\)](#) also examined individual units in a character model and found complex behavior by inspecting activation patterns by hand. Most recently, [Dalvi et al. \(2019\)](#) performed post-hoc tuning of neurons trained in language model and machine translation components, and examined their ability to predict grammatical functions. Like them, we perform an aggregative analysis of individual units to reach measurable quantities of models at a whole, but apply our method to taggers trained directly on supervised grammatical tasks, and focus on cross-lingual variation as the main object of investigation.

| Language | Affix [†] | Morph synth [‡] | POS Accuracy % | |
|------------|--------------------|--------------------------|----------------|-------|
| | | | Dev | Test |
| Arabic | S | int | 96.11 | 95.93 |
| Bulgarian | S | fus | 97.91 | 97.80 |
| Coptic | p | agg | 92.54 | 92.51 |
| Danish | S | fus | 95.59 | 95.46 |
| Greek | S | fus | 96.13 | 96.46 |
| English | S | fus | 93.65 | 93.30 |
| Spanish | S | fus | 95.75 | 95.00 |
| Basque | = | agg | 92.99 | 92.43 |
| Persian | s | fus | 96.07 | 96.10 |
| Irish | = | fus | | 89.35 |
| Hebrew | s | int | 95.71 | 94.60 |
| Hindi | S | fus | 95.03 | 94.91 |
| Hungarian | S | agg | 94.14 | 92.00 |
| Indonesian | S | iso | 92.55 | 92.68 |
| Italian | S | fus | 96.82 | 96.95 |
| Latvian | s | fus | 94.70 | 93.09 |
| Russian | S | fus | 95.29 | 95.25 |
| Swedish | S | fus | 95.80 | 95.73 |
| Tamil | S | agg | 86.46 | 87.58 |
| Thai | ∅ | fus | 91.37 | |
| Turkish | S | agg | 92.08 | 92.48 |
| Ukrainian | S | fus | 95.68 | 95.26 |
| Vietnamese | ∅ | iso | 88.51 | 86.58 |
| Chinese | S | iso | 93.05 | 93.11 |

Table 1: Attributes and tagging accuracy by language (Irish and Thai do not have both dev and test sets). [†]Affixation: S/s is strongly/weakly suffixing; P/p is strongly/weakly prefixing; = is equally prefixing/suffixing; ∅ is little affixation. [‡]Morphological synthesis: agglutinative, fusional, introflexive, isolating.

3 Tagging Task

We train a set of LSTM tagging models, following the setup of [Ling et al. \(2015\)](#). A word representation trained from a character-level LSTM submodule is fed into a word-level bidirectional LSTM, with each word’s hidden state subsequently fed into a two-layer perceptron producing tag scores, which are then softmaxed to produce a tagging distribution. For languages with additional morphosyntactic attribute tagging, we follow the architecture in [Pinter et al. \(2017\)](#) where the same word-level Bi-LSTM states are used to predict each attribute’s value using its own perceptron+softmax scaffolding. In order to produce character models which would be as informative as possible to our subsequent analysis, we do not include word-level embeddings, pre-trained or otherwise, in our setup.

3.1 Language Selection

As our goal is to examine the relationship between character-level modeling and linguistic properties, we drove language selection based on two morphological properties deemed relevant to the archi-

tectural effects examined. All 24 datasets were obtained from Universal Dependencies (UD) version 2.3 (Nivre et al., 2018), and linguistic properties were found in the World Atlas of Language Structures (Bickel and Nichols, 2013; Dryer, 2013). The selected languages and their properties are presented in Table 1. We note that eleven of the 24 languages selected are not Indo-European.

Affixation. To evaluate the role of forward and backward units in a bidirectional model, we selected all languages available in UD which are not classified as either weakly or strongly suffixing in inflectional morphology (the vast majority of UD languages). This includes a single prefixing language (Coptic), two equally suffixing and prefixing languages (Basque and Irish), and two languages with little affixation (Thai and Vietnamese).

Morphological Synthesis. Linguistically functional features vary between being expressed as distinct tokens (isolating languages), detectable unique character substrings (agglutinative), fused together but still distinguishable from the stem (fusional), and non-linearly represented within the word form (introflexive). This property has previously been found to affect performance in character-level models (Pinter et al., 2017; Gerz et al., 2018; Chaudhary et al., 2018), and thus we select representatives of each group, including most available non-fusional languages.

3.2 Technical Setup

Most of our selected languages have only a single UD 2.3 treebank. For languages with multiple treebanks we selected the largest, except in the cases of Spanish and Indonesian, where we selected the GSD treebanks. The Irish IDT treebank has only a train and test split, so we used the test set for early stopping. The Thai PUD treebank only provided a single dataset with 1000 instances, which we shuffled and partitioned into a 850/150 split. Tokens were normalized to remove noisy data: tokens containing ‘http’ were replaced with ‘URL’ and tokens containing ‘@’ were replaced with ‘EMAIL’. This was most relevant (293 replacements) for the English treebank, which contained many long URLs.

Hyperparameters. For the initial bidirectional character-level LSTM, we used a total hidden state size of 128 (64 units in each direction). The char-

acter embedding size is set to 256, initialized using the method of Glorot and Bengio (2010). The word-level bidirectional LSTM has two layers and a hidden state size of 128, with 50% dropout applied in the style of Gal and Ghahramani (2016). Each attribute-prediction MLP has a single hidden layer that is the same size as the tagset size for that attribute, and includes a tanh nonlinearity. Models were trained for up to 80 epochs, and we select the model with the highest POS tagging accuracy on the dev set. Training used SGD with 0.9 momentum, and all models were implemented using DyNet 2.0 (Neubig et al., 2017).

3.3 Results

In our initial setup, we represent words using a concatenation of the final states from a bidirectional character-level LSTM with 64 forward and backward hidden units each. The results for POS tagging, presented in Table 1, are on par with similar models (Plank et al., 2016, for example) despite not including a word-level type embedding component. We attribute this success to our large character embedding size of 256, corroborating findings reported by Smith et al. (2018).

4 Analysis

We next analyze the models trained on the tagging task in an attempt to see how their character-level hidden states encode different manifestations of linguistic information. We suggest that individual hidden units in the character-level sequence model attune to track patterns in the words which would indicate their linguistic roles (POS and morphological properties), and so patterns in character-role regularity across typologically different languages would manifest themselves in an observable form at the individual unit activation level. This motivates us to devise metrics which would characterize languages through aggregation of individual unit behaviour.

4.1 Metrics

For each language, we run the character-level BiLSTM from the trained tagger on POS-unambiguous word types occurring frequently in the training set, grouped into their parts of speech.² This filtering was done in order to focus

²We used 8 as our frequency threshold, and define unambiguous forms as ones tagged at least 60% of the time with a single POS.

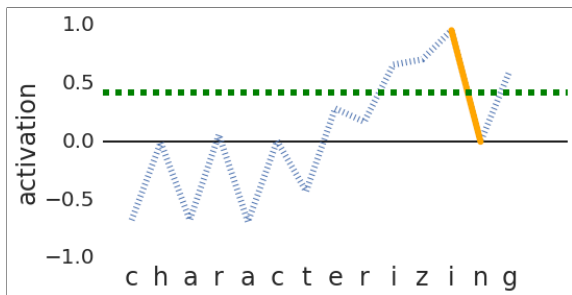


Figure 1: Activations of the English model’s unit 42 (forward) on the word *characterizing*. $b_{\text{avg}|\cdot|}$ is 0.42, and b_{mad} is 0.96 (the drop from the second i to n).

on the more consistent generalizations found by the taggers during training, as our goal is to qualify properties of languages.³ On each word w , we observe each hidden unit h_i ’s activation level (output) on each character h_i^c . We obtain a **base measure** $b(w, i)$ based on the activation pattern. For example, an *average absolute* base measure is defined as the average of absolute value activations:

$$b_{\text{avg}|\cdot|}(w, i) = \frac{1}{|w|} \sum_{c=1}^{|w|} |h_i^c|.$$

The *max absolute diff* base measure is defined as:

$$b_{\text{mad}}(w, i) = \max_{c=1}^{|w|-1} |h_i^{c+1} - h_i^c|.$$

Figure 1 demonstrates these two metrics for a sample (word, unit) pair, showing how the former captures the general level of activation the word caused on the unit, while the latter captures the local character pattern deemed most important by it. We intentionally did not consider metrics based on the final activation values, the direct signals used by the later layers in the model, as these bear no insight into the effect of a word’s composition on the learned model.

Next, we derive a language-level metric for each hidden unit, based on the principle of Mutual Information (MI). The base metric’s range ($[0, 1]$ for $b_{\text{avg}|\cdot|}$, $[0, 2]$ for b_{mad}) is divided into B bins of equal size, and base activations from each word are summed across each of the T POS tag categories⁴, then normalized to produce a joint probability distribution. The mutual information is com-

³This consideration also motivated our choice of UD data, which is tokenized to separate syntactic fusion such as Hebrew and Arabic function words, or Spanish *del*.

⁴We omit the following ‘character-simple’ part-of-speech tags: INTJ, NUM, PROP, PUNCT, SYM, X.

puted as:

$$\sum_{t=1}^T \sum_{b=1}^B P(t, b) [\ln P(t, b) - \ln P(t) - \ln P(b)],$$

and we call the resulting number the POS-Discrimination Index, or **PDI**. Intuitively, a higher PDI implies that the unit activates differently on words of different parts of speech, i.e. it is a better discriminator for the task.

At this point a language produces a set of d_h PDI scores, one for each unit. We sort them from high to low, and define two language-level metrics: The **mass** is the sum of PDI values for all units, $\mathcal{M}(\mathcal{L}) := \sum_{i=1}^{d_h} \text{PDI}(\mathcal{L}, i)$, intuitively meant to quantify the degree of success the model has in assigning hidden units to discriminate POS in this language. The **head forwardness** is the proportion of forward-directional units (under the sorted ordering) before the point at which half of the mass accumulates (in a random setup, this number would tend to 0.5):

$$\frac{\left| \left\{ k : \sum_{i=1}^k \text{PDI}(\mathcal{L}, i) \leq \frac{\mathcal{M}(\mathcal{L})}{2} \wedge h_k \text{ is forward} \right\} \right|}{\left| \left\{ k : \sum_{i=1}^k \text{PDI}(\mathcal{L}, i) \leq \frac{\mathcal{M}(\mathcal{L})}{2} \right\} \right|}$$

This metric aims to quantify the relative importance of forward and backward units in discriminating POS for \mathcal{L} . We use only the top units for the metric as a de-noising heuristic, under the assumption that all units end up with some minimal amount of mass even without performing a function.

4.2 PDI Patterns

The PDI patterns on the $b_{\text{avg}|\cdot|}$ base measure with $B = 16$ bins on all 24 languages are presented in Table 2. We see that agglutinative languages, where we can expect a better discrimination signal to emerge from the consistently-formed morphemes, cluster mostly at the top of the PDI mass scale, suggesting more individual character-level units extract these signals successfully. Intoflexive languages, where character sequences seldom correspond to useful indications of POS or morphosyntactic attributes, cluster towards the bottom.

We present the full unit-level PDI value distributions for Coptic, a prefixing agglutinative language, and English, a suffixing fusional language,

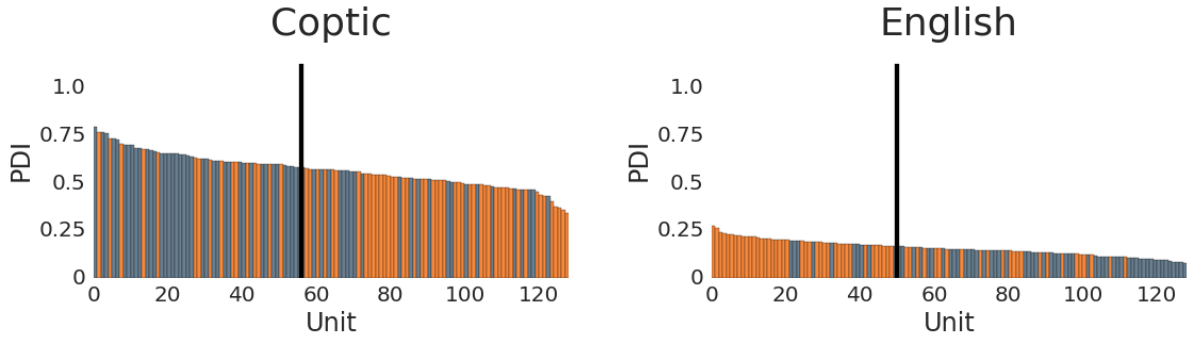


Figure 2: Distribution of PDI values ($b_{\text{avg}|\cdot}$) across hidden units in Coptic and English, shown in ordered PDI values from largest to smallest, with blue (orange) bars indicating forward (backward) units. The black line demarcates the median point of mass accumulation.

| Language | Mass | Mass median index | % of forward units until median |
|------------------|------|-------------------|---------------------------------|
| Tamil | 71.0 | 55 | 49.1 |
| Irish | 62.0 | 56 | 42.9 |
| Coptic | 58.1 | 56 | 71.4 |
| Hungarian | 47.9 | 55 | 50.9 |
| Greek | 31.2 | 55 | 45.5 |
| Turkish | 30.1 | 54 | 57.4 |
| Russian | 25.9 | 54 | 40.7 |
| Thai | 25.9 | 55 | 47.3 |
| Ukrainian | 25.0 | 54 | 37.0 |
| Vietnamese | 24.2 | 55 | 36.4 |
| Chinese | 23.8 | 47 | 42.6 |
| Danish | 21.7 | 54 | 44.4 |
| Swedish | 20.8 | 53 | 34.0 |
| Basque | 20.6 | 51 | 64.7 |
| Indonesian | 20.3 | 45 | 71.1 |
| Latvian | 17.0 | 52 | 42.3 |
| Spanish | 16.1 | 45 | 33.3 |
| English | 16.0 | 50 | 20.0 |
| Bulgarian | 15.6 | 52 | 46.2 |
| Italian | 14.1 | 48 | 56.2 |
| <i>Arabic</i> | 12.6 | 46 | 58.7 |
| <i>Hebrew</i> | 11.4 | 51 | 74.5 |
| Persian | 10.3 | 50 | 46.0 |
| Hindi | 8.4 | 51 | 41.2 |

Table 2: PDI statistics for UD 2.3 models, $b_{\text{avg}|\cdot}$ metric, sorted by the mass metric (sum of PDIs). Agglutinative languages in **bold**, inflexive in *italics*.

in Figure 2 (trends for b_{mad} are similar). Consistent with other agglutinative languages, Coptic’s cumulative mass is very large ($\mathcal{M}(\text{cop}) = 58.1$), suggesting the predictive qualities of the sequence-based LSTM allows good discrimination from the character signal, as one might expect from an agglutinative language. Conversely, $\mathcal{M}(\text{eng}) = 16$, demonstrating the difficulty presented by fusional languages. The accumulation of 71% forward (80% backward) units in the head of the Coptic (English) value ranking suggests an interesting relationship between affixation and LSTM direction: LSTM units are likely to hone in on POS-indicative signals, which often occur as affixes, in the beginning of their run, causing activation values to rise (in absolute value) and stay large throughout the subsequent traversal of the stem. Unfortunately, since no other prefixing languages are available in UD, we were not able to pursue this hypothesis further.

4.3 Asymmetric Directionality

Based on these observations, we conduct a directionality balance study, where we vary the number of hidden units in the forward and backwards dimensions. In addition to the models analyzed above, which use 64 forward and 64 backward units (denoted hereafter 64/64), we trained models with imbalanced directionality (128/0, 96/32, 32/96, 0/128). We test the hypothesis that imbalanced models affect languages differently based on their linguistic properties and statistical metrics. We note that these settings do not maintain parameter set size: intra-direction transition operations are quadratic in that direction’s hidden layer size, and so this adds a possible advantage in favor of direction-imbalanced models.

| Language Type | 128/0 | 96/32 | 64/64 (base) | 32/96 | 0/128 |
|------------------------------------|--------------|--------------|-----------------|--------------|--------------|
| Inflectional Affixation Categories | | | | | |
| S. suffix | +0.22 | +0.07 | 94.50 | -0.06 | -0.02 |
| W. suffix | +0.26 | +0.12 | 95.46 | -0.07 | -0.01 |
| Equal p/s | +0.61 | +0.32 | 90.99 | -0.07 | +0.06 |
| Little aff. | -0.06 | -0.21 | 89.59 | -0.16 | -0.22 |
| W. prefix | +0.52 | +0.22 | 92.91 | +0.40 | +0.33 |
| Morphological Synthesis Categories | | | | | |
| Introflex. | +0.17 | +0.05 | 95.87 | -0.06 | +0.01 |
| Fusional | +0.22 | +0.07 | 94.95 | +0.01 | +0.06 |
| Agglutina. | +0.59 | +0.27 | 91.58 | -0.16 | -0.15 |
| Isolating | -0.14 | -0.13 | 91.15 | -0.15 | -0.13 |
| Overall | +0.25 | +0.08 | 93.85 | -0.05 | -0.01 |

Table 3: Imbalanced models’ mean POS accuracy on UD development data (differences between three averaged random runs in all models; **boldfaced** when significant at $p < 0.05$ using a paired two-tailed t -test).

The results for this study are presented in Table 3 as averages for the language categories listed in Table 1 (the full, raw results are available in Table 4).

One trend which emerges is the preference of **agglutinative** languages for imbalanced models, whereas the other languages are little affected by this change. This could be explained by the increase in inter-unit interaction in the larger direction of an imbalanced model – contiguous character sequences consistently code reliable linguistic features in these languages. A second finding is the slight bias of suffixing languages towards more forward units and of the prefixing language to more backward units, indicating that hidden LSTM units are better in detecting formations close to their final state. Coupled with the findings regarding PDI mass distribution in the different directional units in § 4.2, we suggest that a subtle relation exists between morphological information and model directionality: units which end their run on the affix are more important for detecting the POS signal, so having more of them helps the model. We also note the stability of isolating and little-affixing languages to directionality balance, possibly owing to the relatively small significance of contiguous character sequences in detecting word role. Lastly, we point out that the compromise *sesquidirectional* models 96/32 and 32/96 did not tend to stand out significantly on our tested language categories, suggesting there is no substantial middle-ground between the two popular techniques of unidirectional and bidirectional

| Language | 128/0 | 96/32 | 64/64 | 32/96 | 0/128 |
|------------|-------|-------|-------|-------|-------|
| Arabic | 96.29 | 96.08 | 96.06 | 96.09 | 96.16 |
| Bulgarian | 97.95 | 97.86 | 97.84 | 97.74 | 97.71 |
| Coptic | 93.10 | 92.80 | 92.58 | 92.98 | 92.91 |
| Danish | 95.93 | 95.68 | 95.61 | 95.60 | 95.70 |
| Greek | 96.19 | 96.07 | 96.01 | 96.00 | 95.93 |
| English | 93.86 | 93.74 | 93.65 | 93.80 | 93.87 |
| Spanish | 95.74 | 95.63 | 95.64 | 95.64 | 95.77 |
| Basque | 93.52 | 93.13 | 92.89 | 92.59 | 92.90 |
| Persian | 96.31 | 96.20 | 96.11 | 96.02 | 96.20 |
| Irish | 89.54 | 89.35 | 88.95 | 89.11 | 89.07 |
| Hebrew | 95.76 | 95.72 | 95.60 | 95.50 | 95.57 |
| Hindi | 95.35 | 95.22 | 95.12 | 95.11 | 95.25 |
| Hungarian | 94.25 | 94.29 | 94.20 | 93.97 | 94.00 |
| Indonesian | 92.42 | 92.34 | 92.49 | 92.53 | 92.55 |
| Italian | 97.00 | 96.78 | 96.87 | 96.88 | 97.01 |
| Latvian | 95.10 | 94.84 | 94.69 | 94.58 | 94.61 |
| Russian | 95.51 | 95.39 | 95.32 | 95.31 | 95.36 |
| Swedish | 95.93 | 95.69 | 95.64 | 95.52 | 95.85 |
| Tamil | 87.54 | 87.28 | 86.88 | 86.28 | 85.99 |
| Thai | 91.52 | 91.27 | 91.38 | 91.47 | 91.32 |
| Turkish | 93.14 | 92.45 | 92.06 | 92.03 | 92.09 |
| Ukrainian | 95.72 | 95.76 | 95.63 | 95.68 | 95.66 |
| Vietnamese | 87.98 | 87.92 | 88.23 | 87.83 | 87.85 |
| Chinese | 93.01 | 93.17 | 93.12 | 93.03 | 93.04 |

Table 4: Full scores for the directionality balance experiment, each point averaged over three random seed runs.

LSTMs.

5 Conclusion

While character-level Bi-LSTM models compute meaningful word representations across many languages, the way they do it depends on each language’s typological properties. These observations can guide model selection: for example, in agglutinative languages we observe a strong preference for a single direction of analysis, motivating the use of unidirectional character-level LSTMs for at least this type of language. In future work, we plan to introduce further control into our metrics by incorporating dataset attributes such as tag distribution and number of instances, as well as learning-related properties like convergence rate and effect of initialization.

Acknowledgments

We would like to thank Sebastian Mielke, anonymous reviewers from NAACL and BlackboxNLP, and the members of the Computational Linguistics lab at Georgia Tech for their valuable notes. YP is a Bloomberg Data Science PhD Fellow. MM’s work was funded by the Georgia Tech Undergraduate Research Opportunities Program.

References

- Balthasar Bickel and Johanna Nichols. 2013. [Fusion of selected inflectional formatives](#). In Matthew S. Dryer and Martin Haspelmath, editors, *The World Atlas of Language Structures Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association of Computational Linguistics*, 5(1):135–146.
- Aditi Chaudhary, Chunting Zhou, Lori Levin, Graham Neubig, David R Mortensen, and Jaime Carbonell. 2018. Adapting word embeddings to new languages with morphological and phonological subword representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3285–3295.
- Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, D Anthony Bau, and James Glass. 2019. [What is one grain of sand in the desert? analyzing individual neurons in deep nlp models](#). In *Proceedings of the 32th AAAI Conference on Artificial Intelligence*.
- Matthew S. Dryer. 2013. [Prefixing vs. suffixing in inflectional morphology](#). In Matthew S. Dryer and Martin Haspelmath, editors, *The World Atlas of Language Structures Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- Yarin Gal and Zoubin Ghahramani. 2016. [A theoretically grounded application of dropout in recurrent neural networks](#). In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 1019–1027. Curran Associates, Inc.
- 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 of Computational Linguistics*, 6:451–465.
- Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256.
- 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 the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3275–3284.
- Andrej Karpathy, Justin Johnson, and Li Fei-Fei. 2015. Visualizing and understanding recurrent networks. *arXiv preprint arXiv:1506.02078*.
- Yova Kementchedjheva and Adam Lopez. 2018. ‘indications’ that character language models learn english morpho-syntactic units and regularities. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 145–153.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. Neural architectures for named entity recognition. In *Proceedings of NAACL-HLT*, pages 260–270.
- Wang Ling, Tiago Luís, Luís Marujo, Ramón Fernandez Astudillo, Silvio Amir, Chris Dyer, Alan W Black, and Isabel Trancoso. 2015. Finding function in form: Compositional character models for open vocabulary word representation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. 2016. Assessing the ability of lstms to learn syntax-sensitive dependencies. *Transactions of the Association of Computational Linguistics*, 4(1):521–535.
- Graham Neubig, Chris Dyer, Yoav Goldberg, Austin Matthews, Waleed Ammar, Antonios Anastasopoulos, Miguel Ballesteros, David Chiang, Daniel Clothiaux, Trevor Cohn, Kevin Duh, Manaal Faruqui, Cynthia Gan, Dan Garrette, Yangfeng Ji, Lingpeng Kong, Adhiguna Kuncoro, Gaurav Kumar, Chaitanya Malaviya, Paul Michel, Yusuke Oda, Matthew Richardson, Naomi Saphra, Swabha Swayamdipta, and Pengcheng Yin. 2017. Dynet: The dynamic neural network toolkit. *arXiv preprint arXiv:1701.03980*.
- Joakim Nivre et al. 2018. [Universal dependencies 2.3](#). LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University.
- Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, volume 1, pages 2227–2237.
- Yuval Pinter, Robert Guthrie, and Jacob Eisenstein. 2017. Mimicking Word Embeddings using Subword RNNs. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 102–112.
- Barbara Plank, Anders Søgaard, and Yoav Goldberg. 2016. [Multilingual part-of-speech tagging with bidirectional long short-term memory models and auxiliary loss](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 412–418. Association for Computational Linguistics.

Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. 2017. Learning to generate reviews and discovering sentiment. *arXiv preprint arXiv:1704.01444*.

Aaron Smith, Miryam de Lhoneux, Sara Stymne, and Joakim Nivre. 2018. An investigation of the interactions between pre-trained word embeddings, character models and pos tags in dependency parsing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2711–2720.

Clara Vania and Adam Lopez. 2017. From characters to words to in between: Do we capture morphology? In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 2016–2027.