## The Illiterati: Part-of-Speech Tagging for Magahi and Bhojpuri without Even Knowing the Alphabet

### **Thomas Proisl**

Institute of Cognitive Science Osnabrück University thomas.proisl@uos.de

**Peter Uhrig English and American Studies** FAU Erlangen-Nürnberg peter.uhrig@fau.de

**Philipp Heinrich** 

### **Andreas Blombach**

Comp. Corpus Linguistics FAU Erlangen-Nürnberg philipp.heinrich@fau.de

**Comp.** Corpus Linguistics FAU Erlangen-Nürnberg

**Romance Studies** 

Sefora Mammarella

FAU Erlangen-Nürnberg andreas.blombach@fau.de sefora.mammarella@icloud.com

Natalie Dykes **Computational Corpus Linguistics** FAU Erlangen-Nürnberg natalie.mary.dykes@fau.de

#### Abstract

In this paper, we describe the part-of-speechtagging experiments for Magahi and Bhojpuri that we conducted for our participation in the NSURL 2019 shared tasks 9 and 10 (Lowlevel NLP Tools for (MagahilBhojpuri) Language). We experiment with three different part-of-speech taggers and evaluate the impact of additional resources such as Brown clusters, word embeddings and transfer learning from additional tagged corpora in related languages. In a 10-fold cross-validation on the training data, our best-performing models achieve accuracies of 90.70% for Magahi and 94.08% for Bhojpuri. Accuracy increased to 94.79% for Magahi and dropped to 78.68% for Bhojpuri on the test data.

#### 1 Introduction and Related Work

Magahi and Bhojpuri are two of the three principal languages of the Bihari group (Maithili being the third). There are competing categorizations of the Bihari group within the Indo-Aryan languages (see Grierson, 1903; Cardona, 1974; Jeffers, 1976). While there are few Magahi speakers outside of Southern Bihar, Bhojpuri is spoken in parts of two Indian states, Western Bihar and Eastern Uttar Pradesh, and the Southwest of Nepal. According to the 2011 census, about 51 million people in India stated Bhojpuri as their mother tongue, and about 13 million did so for Magahi. However, these numbers may seriously underestimate the actual number of speakers, since speakers of both languages often name Hindi as their first language - the language used

Besim Kabashi **Computational Corpus Linguistics** FAU Erlangen-Nürnberg besim.kabashi@fau.de

in schools, courts, and other public institutions (Verma, 2003b, p. 547).

Despite these numbers, comparatively few linguistic resources and NLP tools currently exist for both languages, with most of the scarce attention having gone towards Bhojpuri (e.g. Ojha et al., 2015).

It is beyond the scope of this paper and our own expertise to describe both languages in detail (but see, e.g., Verma, 2003b,a). Among the features which appear pertinent to part-of-speech tagging of Magahi and Bhojpuri are SOV order, rich verb morphology, the extensive use of postpositions, and the unusual agreement system of Magahi (where the verb has to agree with subject and object simultaneously).

Table 1 gives an overview of the two datasets of the shared task. While the training set for Bhojpuri is much larger, it also features a more fine-grained tagset.

	Magahi	Bhojpuri
training	61.435	94.692
test	8.205	10.582
tagset size	18	33

Table 1: Sizes of the training and test sets and of the tagsets.

#### 2 **Strategies and Systems**

#### 2.1 Part-of-Speech Taggers

We experiment with three different, freely available part-of-speech taggers:

- SoMeWeTa (Proisl, 2018), a tagger based on the averaged structured perceptron that supports domain adaptation and can incorporate external information sources such as Brown clusters.<sup>1</sup>
- A BiLSTM+CRF sequence tagger by Guillaume Genthial that uses character and word embeddings and supports transfer learning.<sup>2</sup>
- The Stanford Tagger (Toutanova et al., 2003), which is based on a maximum entropy cyclic dependency network.<sup>3</sup>

## 2.2 Additional Resources

In addition to the training data provided by the task organizers, we use the following freely available resources:

- The Hindi UD treebank, which is based on the Hindi Dependency Treebank (HDTB; ca. 352,000 tokens; Bhat et al., 2017; Palmer et al., 2009).<sup>4</sup>
- A POS-tagged Magahi corpus (KMI-Mag; ca. 46,000 tokens) and a corpus of untagged Magahi texts (ca. 2.8 million tokens).<sup>5</sup>
- Wikimedia dumps for Hindi (ca. 34.7 million tokens) and Bihari (ca. 700,000 tokens).<sup>6</sup> We extract the plain text using wikiextractor<sup>7</sup> and tokenize and sentence-split it using the ICU tokenizer via polyglot.<sup>8</sup>
- Brown clusters (Brown et al., 1992) computed from the tokenized Wikimedia dumps and the untagged Magahi corpus (1000 clusters, minimum frequency 5).<sup>9</sup>

- UniversalDependencies/UD\_Hindi-HDTB/
- tree/master

• Pre-trained fastText embeddings for Hindi and Bihari<sup>10</sup>

The additional tagged Magahi corpus (KMI-Mag) is annotated with a tagset consisting of 35 tags which is almost identical to the 33-tag tagset used in the Bhojpuri corpus. KMI-Mag uses three tags that do not occur in the Bhojpuri data (V\_VM\_VF, V\_VM\_VNF and V\_VM\_VNP) and misses one tag that is used for Bhojpuri (RD\_ECH\_B). For our transfer learning experiments targeting Bhojpuri, we simply convert the three verb tags to V\_VM. For targeting Magahi, we map the 35 tags to UD tags.

#### 2.3 Experiments using SoMeWeTa

The distinctive features of SoMeWeTa are its ability to leverage additional resources and its transfer learning or domain adaptation capabilities. Consequently, we focus on these two aspects in our experiments.

For Bhojpuri, we experiment primarily with the Brown clusters computed from the Hindi and Bihari Wikimedia dumps and the untagged additional Magahi corpus (cf. section 2.2). Our crossvalidation experiments show that the Brown clusters have a small positive effect with the best results being obtained by Brown clusters computed from the union of all three additional corpora (cf. Table 2). With KMI-Mag we have a corpus of a closely related language that is annotated with an almost identical tagset (cf. section 2.2). However, pretraining on that and then adapting to Bhojpuri seems to have no noticeable effect.

For Magahi, we experiment with a wide range of transfer learning settings in addition to the different Brown clusters:

- Pretraining on one of KMI-Mag, HDTB or the Bhojpuri dataset (mapped to UD tags).
- Pretraining on all possible combinations of KMI-Mag, HDTB and the Bhojpuri dataset (using the concatenation of these resources).
- Longer pretraining chains where we start with HDTB and adapt to one or two other resources before we make the final adaptation to Magahi.

The best results are obtained by using Brown clusters computed from the Hindi Wikimedia dumps

<sup>&</sup>lt;sup>1</sup>https://github.com/tsproisl/SoMeWeTa
<sup>2</sup>We use the slightly modified version by Riedl
and Padó (2018): https://github.com/riedlma/
sequence\_tagging
<sup>3</sup>https://nlp.stanford.edu/software/

tagger.html <sup>4</sup>https://github.com/

<sup>&</sup>lt;sup>5</sup>https://github.com/kmi-linguistics/ magahi

<sup>&</sup>lt;sup>6</sup>https://dumps.wikimedia.org

<sup>&</sup>lt;sup>7</sup>http://medialab.di.unipi.it/wiki/ Wikipedia\_Extractor

<sup>&</sup>lt;sup>8</sup>http://polyglot-nlp.com/

<sup>&</sup>lt;sup>9</sup>We use the implementation by Liang (2005): https: //github.com/percyliang/brown-cluster

<sup>&</sup>lt;sup>10</sup>https://fasttext.cc/docs/en/ crawl-vectors.html

model	accuracy
No additional resources	91.62 (±0.97)
Hindi Brown clusters	91.79 (±1.00)
Bihari Brown clusters	91.60 (±1.01)
Magahi Brown clusters	91.69 (±0.93)
Hindi+Magahi Brown clusters (hi+mag)	91.99 (±0.83)
<i>Hindi</i> + <i>Bihari</i> + <i>Magahi Brown clusters</i> ( <i>hi</i> + <i>bh</i> + <i>mag</i> )	92.04 (±0.80)
KMI-Mag $\rightarrow$ Bhojpuri, hi+mag	92.03 (±0.90)
$KMI-Mag \rightarrow Bhojpuri, hi+bh+mag$	92.06 (±0.94)

Table 2: Bhojpuri results for SoMeWeTa. We report the mean accuracies and 95% confidence intervals of a 10-fold cross-validation on the training data. The model that we submitted to the shared task is set in italics.

and the untagged additional Magahi corpus. As for Bhojpuri, transfer learning does not seem to have any noticeable effect (cf. Table 3).

# 2.4 Experiments using the BiLSTM-CRF tagger

Neural networks with a BiLSTM-CRF architecture achieve POS-tagging results close to the current state of the art.<sup>11</sup> In our experiments, we focus less on the hyperparameters of the network but rather on the effects of our additional resources. We try out both the Hindi and Bihari fastText embeddings. Since the Bihari embeddings do not perform significantly better than the Hindi embeddings (cf. Table 4) and the Hindi embeddings cover a much larger vocabulary (15.3 million words instead of 8.9 million), we use the Hindi embeddings for our further experiments. In the following, we make use of the tagger's transfer learning abilities and pretrain the models on HDTB or KMI-Mag. The BiLSTM-CRF tagger seems to benefit more from the transfer learning setting than SoMeWeTa and achieves its best results for both languages with a transfer from KMI-Mag. Interestingly, the BiLSTM-CRF outperforms SoMeWeTa only on the Magahi dataset while it performs notably worse on the Bhojpuri dataset.

#### 2.5 Experiments using the Stanford Tagger

The Stanford Log-linear Part-Of-Speech Tagger (Toutanova and Manning 2000; Toutanova et al. 2003) is a mature and stable tagger that still exhibits competitive performance. The system is feature-rich and offers a range of configuration options, the effects of which were initially not fully understood by our research group. It was thus decided to run extensive brute-force hyperparameter tuning making educated guesses about the value ranges of the various parameters. The documentation in the JavaDoc for the MaxentTagger class<sup>12</sup> provides the necessary information. It was decided to set the following parameters with the values or ranges given in Table 5 and Table 6.

Combining all parameters results in 76,800 parameter combinations per language. Although training and testing can be completed in approximately 2 minutes on a modern personal computer, the sheer number of parameter combinations necessitated running the experiments on High-Performance-Computing infrastructure. The setup comprised a central queue of filenames of property files that all involved clients subscribed to.

For Magahi, only two runs with all parameter combinations were performed: one with the top 80% of the training data as actual training data and the bottom 20% as test data and one with the bottom 80% as training data and the top 20% as test data. The values discussed below are the arithmetic mean of the accuracies of those two runs. As the Magahi tagset is Universal-Dependencies-compliant, it was straightforward to identify closed class words by pos tag and to supply the list to the tagger during the training phase.

For Bhojpuri, a full 10-fold cross-validation was carried out for each of the parameter combinations, so the averages discussed below are most likely more reliable than those for Magahi. Since the Bhojpuri tagset was more complicated, we decided to learn the closed class tags automatically based on the default *closedClassTagThreshold* of 40. Thus, a pos tag is only considered a closed class if it is assigned to less than 40 different words.

<sup>&</sup>lt;sup>11</sup>Cf. https://aclweb.org/aclwiki/POS\_ Tagging\_(State\_of\_the\_art)

<sup>&</sup>lt;sup>12</sup>https://nlp.stanford.edu/nlp/javadoc/ javanlp/edu/stanford/nlp/tagger/maxent/ MaxentTagger.html

model	accuracy
No additional resources	88.92 (±1.24)
Hindi Brown cluster	89.07 (±1.24)
Bihari Brown cluster	88.90 (±1.32)
Magahi Brown cluster	89.12 (±1.23)
Hindi+Magahi Brown cluster	89.32 (±1.15)
Hindi+Bihari+Magahi Brown cluster	89.15 (±1.17)
KMI-Mag $\rightarrow$ Magahi, Hindi+Magahi Brown cluster	89.20 (±1.10)
KMI-Mag → Magahi, Hindi+Bihari+Magahi Brown cluster	89.23 (±1.19)
Bhojpuri → Magahi, Hindi+Magahi Brown cluster	89.25 (±1.13)
Bhojpuri $\rightarrow$ Magahi, Hindi+Bihari+Magahi Brown cluster	89.18 (±1.25)
$HDTB \rightarrow Magahi$ , Hindi+Magahi Brown cluster	89.26 (±1.21)
$HDTB \rightarrow Magahi$ , Hindi+Bihari+Magahi Brown cluster	89.17 (±1.18)
HDTB+KMI-Mag → Magahi, Hindi+Magahi Brown cluster	89.22 (±1.12)
HDTB+KMI-Mag $\rightarrow$ Magahi, Hindi+Bihari+Magahi Brown cluster	89.19 (±1.23)
HDTB+Bhojpuri $\rightarrow$ Magahi, Hindi+Magahi Brown cluster	89.23 (±1.13)
HDTB+Bhojpuri $\rightarrow$ Magahi, Hindi+Bihari+Magahi Brown cluster	89.18 (±1.20)
KMI-Mag+Bhojpuri → Magahi, Hindi+Magahi Brown cluster	89.30 (±1.14)
KMI-Mag+Bhojpuri $\rightarrow$ Magahi, Hindi+Bihari+Magahi Brown cluster	89.06 (±1.19)
HDTB+KMI-Mag+Bhojpuri $\rightarrow$ Magahi, Hindi+Magahi Brown cluster	89.21 (±1.17)
HDTB+KMI-Mag+Bhojpuri, Hindi+Bihari+Magahi Brown cluster	89.20 (±1.20)
$HDTB \rightarrow KMI-Mag \rightarrow Magahi$ , Hindi+Magahi Brown cluster	89.24 (±1.20)
$HDTB \rightarrow KMI-Mag \rightarrow Magahi$ , Hindi+Bihari+Magahi Brown cluster	89.22 (±1.18)
$HDTB \rightarrow Bhojpuri \rightarrow Magahi, Hindi+Magahi Brown cluster$	89.27 (±1.14)
$HDTB \rightarrow Bhojpuri \rightarrow Magahi, Hindi+Bihari+Magahi Brown cluster$	89.11 (±1.17)
$HDTB \rightarrow Bhojpuri \rightarrow KMI-Mag \rightarrow Magahi, Hindi+Magahi Brown cluster$	89.22 (±1.11)
$HDTB \rightarrow Bhojpuri \rightarrow KMI-Mag \rightarrow Magahi, Hindi+Bihari+Magahi Brown cluster$	89.20 (±1.19)

Table 3: Magahi results for SoMeWeTa. We report the mean accuracies and 95% confidence intervals of a 10-fold cross-validation on the training data. The model that we submitted to the shared task is set in italics.

model	accuracy
Magahi (Hindi embeddings)	88,97 (±1,14)
Magahi (Bihari embeddings)	89,09 (±1,00)
HDTB → Magahi (Hindi embeddings)	89,85 (±0,99)
KMI-Mag → Magahi (Hindi embeddings)	90,70 (±0,92)
Bhojpuri (Hindi embeddings)	90,78 (±0,55)
Bhojpuri (Bihari embeddings)	90,80 (±0,57)
<i>KMI-Mag → Bhojpuri (Hindi embeddings)</i>	91,23 (±0,68)

Table 4: Results for the BiLSTM-CRF tagger. We report the mean accuracies and 95% confidence intervals of a 10-fold cross-validation on the training data. The models submitted to the shared task are set in italics.

Given that the training dataset is smaller than what is available for more commonly researched languages, we expected that for most thresholds, values below the default values might be more relevant than above, which is why our choice of parameter values is skewed towards smaller numbers.

For both languages, performance decreases abruptly when *rareWordThresh* is set to 1. We exclude this setting for the remainder of the analysis, since it is obviously beneficial for the tagger to treat hapax legomena as rare words. Additionally, performance was insensitive to variation in *veryCommonWordThresh* since this value is ignored by the Tagger in our case. We thus fix the threshold at 250 and use simple linear models without interaction to analyze the influence of all other variables on performance measures:

acc. = 
$$\beta_0 + \beta_1(unicodeshape) + \beta_2(macro)$$
  
+  $\sum_{j=3}^{6} \beta_j \gamma_j + \varepsilon$ 

where  $\beta_i$  are the coefficients,  $\gamma_j$  is one of the integer features (*rareWordThresh*, *curWordMinFeatureThresh*, *minFeatureThresh*, *rareWordMinFeatureThresh*), and  $\varepsilon$  is the residual error.

Accuracy for Bhojpuri reaches around  $\mu \approx$  93.88 with a standard deviation of approximately 0.064 and the linear model yielding an adjusted  $R^2$  of approximately 0.80. For Magahi, overall performance is lower ( $\mu \approx 87.66$ ) and variation is higher ( $\sigma \approx 0.51$ ), but this variation is well-explained by the linear model (adjusted  $R^2 \approx 0.98$ ).

For both languages, the *macro* parameter has the most influence on accuracy. For Bhojpuri, the best *macro* is bidirectional5words (yielding ceteris paribus 0.09 and 0.12 better results compared to generic and left3words, respectively). For Magahi, however, generic

parameter	default value	value/range
closedClassTags	(none)	ADP AUX CCONJ DET NUM PART PRON SCONJ PUNCT
arch - architecture	generic	generic, left3word, bidirectional5words
arch - further unknown-words option	(none)	naacl2003unknowns
arch - unicode shapes for rare words	(none)	unicodeshapes(-2,2), unicodeshapes(-1,1), unicodeshapes(0), (none)
iterations	100	100
learnClosedClassTags	false	false
curWordMinFeatureThresh	2	14
minFeatureThresh	5	15
rareWordMinFeatureThresh	10	110
rareWordThresh	5	18
veryCommonWordThresh	250	100, 150, 200, 250

Table 5: Settings and parameters with ranges for the training of the Stanford PoS Tagger for Magahi.

parameter	default value	value/range
closedClassTags	(none)	(none)
arch - architecture	generic	generic, left3word, bidirectional5words
arch - further unknown-words option	(none)	naacl2003unknowns
arch - unicode shapes for rare words	(none)	unicodeshapes(-2,2), unicodeshapes(-1,1), unicodeshapes(0), (none)
iterations	100	100
learnClosedClassTags	false	true
closedClassTagThreshold	40	40
curWordMinFeatureThresh	2	14
minFeatureThresh	5	15
rareWordMinFeatureThresh	10	110
rareWordThresh	5	18
veryCommonWordThresh	250	100, 150, 200, 250

Table 6: Settings and parameters with ranges for the training of the Stanford PoS Tagger for Bhojpuri.

and left3words give better results (both approximately 1.0 accuracy points better than bidirectional5words). This is surprising, since according to the authors of the Stanford Tagger, "[t]he left3words architectures are faster, but slightly less accurate, than the bidirectional architectures."<sup>13</sup> The only viable explanation that comes to mind is that possibly the Magahi gold standard corpus was annotated with a trigram tagger without sufficient manual correction. This is in line with our observation that in the Magahi data, items that should have been classified as punctuation marks recieved dubious tags, e.g. the grave accent (') was tagged only twice as punctuation, but was categorized as a noun five times, twice as an adposition, once as a verb and once as an auxiliary.

Examining only the respective best-performing *macro*, *rareWordThresh* explains most of the remaining variation, with a significant regression coefficient of about 0.02 for Bhojpuri and 0.07 for Magahi. However, the effect might de-

crease for values higher than the ones tested here (*rareWordThresh*  $\in \{1, ..., 8\}$ ).

unicodeshape has a small effect on performance for Bhojpuri, where (-1, 1) and (-2, 2) yield an increase in performance by about 0.06 compared to (0) and None. This effect cannot be confirmed for Magahi. For both languages, performance decreases in *curWordThresh*, *curWordMin-FeatureThresh*, and *rareWordMinFeatureThresh*, though the effect is negligible and not always significant. In both cases, *minFeatureThresh* does not have a significant influence on accuracy.

#### **3** Results and Error Analysis

#### 3.1 Bhojpuri

The overall results for Bhojpuri are delightful since they are even better than on our training data (see Table 7): Our optimized version of the Stanford tagger scored 95 points macro  $F_1$  (94.78 accuracy), and we thus share first place with our sole competitor (team *NITK-NLP*); SoMeWeta and the BiLSTM tagger are close behind.

We omit the very large confusion matrix  $(33 \times 33$  and predominantly zero off the diagonal)

<sup>&</sup>lt;sup>13</sup>https://nlp.stanford.edu/nlp/javadoc/ javanlp/edu/stanford/nlp/tagger/maxent/ ExtractorFrames.html

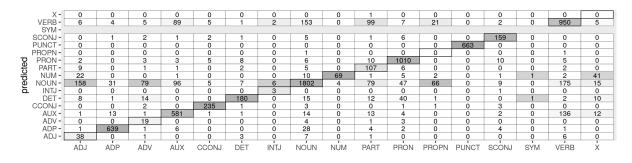


Figure 1: Confusion Matrix for SoMeWeTa predicting Magahi tags on the test data. Absolute numbers are given for all cells; shade represents recall (on the diagonal) and false positive rate, respectively. Actual labels can be found on the abscissa, predicted ones on the ordinate.

rank	submission	$F_1$
1	Stanford	95
1	NITK-NLP_SUB1	95
2	SoMeWeTa	93
3	BiLSTM-CRF	92
4	NITK-NLP_SUB2	89

Table 7: Results for Bhojpuri

and instead provide a quick summary for the Stanford tagger:<sup>14</sup>

- Two tags are not predicted by our tagger at all: RD\_ECH\_B (which appears once in the gold data and was misclassified as N\_NN), and RD\_UNK (classified once as N\_NN and once as V\_VM).
- RP\_INJ appeared five times in the gold standard and was predicted correctly four times. This tag yields the worst recall (apart from the two pathological cases above).
- 30 of the 195 occurrences of RD\_SYM were misclassified (recall 84.6%), mostly as N\_NN (26 cases).
- Further incorrect predictions of N\_NN occur for JJ (11.3% of its occurrences classified as N\_NN, 85.2% recall), RB (7.7%, 89.7% recall), and N\_NNP (6.4%, 92.8% recall).
- Another notable confusion is the pair V\_VM (87.8% recall) and V\_VAUX (86.6% recall);
   V\_VM was predicted as V\_VAUX 64 times, while V\_VAUX was tagged V\_VM 66 times.
   Finally, V\_VM was predicted as N\_NN 85 times.

The results for our other submissions were very much in line with the results discussed here.<sup>15</sup> All in all, the errors made by our submissions are very much what one would expect: Very rare categories are sometimes misclassified, very frequent categories (such as N\_NN) tend to be the go-to label for misclassifications, and similar morphosyntactic categories are confused with each other (V\_VM and V\_AUX, N\_NN and N\_NNP).

#### 3.2 Magahi

With a macro  $F_1$  score of only 77%, our best submissions, SoMeWeTa (78.68 accuracy) and BiLSTM-CRF (78.86 accuracy), rank second in the task of predicting Magahi tags, closely behind the submissions of one of our competing teams (see Table 8). Results are peculiar, since this is a drop of more than ten points compared to our cross-validation on the training data set and far outside our realized confidence intervals (see Table 3).

rank	submission	$F_1$
1	NITK-NLP_SUB2	79
2	SoMeWeTa	77
2	BiLSTM-CRF	77
3	Stanford	74
4	NITK-NLP_SUB1	73

Table 8: Results for Magahi

Figure 1 shows the confusion matrix for SoMeWeTa.<sup>16</sup> Major problems arise for tags ADJ (15.5% recall), ADV (14.8%), PART (32.5%), and PROPN and X (both 0%), since these are quite frequent categories with severe error rates. As with

<sup>&</sup>lt;sup>14</sup>We focus on recall; precision is mostly the same as recall for all frequent labels, and higher for rare ones, since the taggers avoid predicting infrequent labels.

<sup>&</sup>lt;sup>15</sup>One notable exception is that the BiLSTM tagger did non predict the category RD\_ECH at all (another hapax in the gold standard) but did include RD\_ECH\_B (once, incorrectly).

<sup>&</sup>lt;sup>16</sup>Again, results are very similar for our other submissions.

Bhojpuri, the tagger misclassifies them as NOUNS and VERBS, which are the most frequent open classes. Moreover, the tagger frequently mistakes VERB for AUX and vice versa.

#### 4 Conclusion

The results for Bhojpuri are very satisfying. Close to 95% accuracy on a set of 33 tags with approximately 95,000 words of training data is in line with our expectations. It is a bit disappointing, however, that mindless parameter-tuning yields the best results – but the difference may very well not be significant.

The results for Magahi are very disappointing. Since we do not know the language, it is difficult for us to pinpoint the exact reasons for the bad performance, be it an over-generalization of our taggers, a shift in the tag distribution in the test data or an issue with the annotation quality. At least, however, the use of additional resources outperforms mere parameter-tuning.

#### References

- Riyaz Ahmad Bhat, Rajesh Bhatt, Annahita Farudi, Prescott Klassen, Bhuvana Narasimhan, Martha Palmer, Owen Rambow, Dipti Misra Sharma, Ashwini Vaidya, Sri Ramagurumurthy Vishnu, et al. 2017. The Hindi/Urdu treebank project. In Nancy Ide and James Pustejovsky, editors, *Handbook of Linguistic Annotation*, pages 659–697. Springer.
- Peter F. Brown, Vincent J. Della Pietra, Peter V. de Souza, Jennifer C. Lai, and Robert L. Mercer. 1992. Class-based n-gram models of natural language. *Computational Linguistics*, 18(4):467–479.
- George Cardona. 1974. *The Indo-Aryan languages*, 15th edition, volume 9, pages 439–450.
- George Abraham Grierson. 1903. Linguistic survey of India, Vol. V: Indo-Aryan Family, Eastern Group, Pt. II: Specimens of the Bihari and Oriya Languages. Office of the Superintendent of Government Printing, India, Calcutta.
- Robert J. Jeffers. 1976. The position of the Bihārī dialects in Indo-Aryan. *Indo-Iranian Journal*, 18(3):215–225.
- Percy Liang. 2005. Semi-supervised learning for natural language. Master's thesis, Massachusetts Institute of Technology, Department of Electrical Engineering and Computer Science.
- Atul Ku Ojha, Pitambar Behera, Srishti Singh, and Girish N. Jha. 2015. Training & evaluation of POS taggers in Indo-Aryan languages: a case of Hindi,

Odia and Bhojpuri. In *Proceedings of the 7th Language & Technology Conference (LTC 2015)*, pages 524–529.

- Martha Palmer, Rajesh Bhatt, Bhuvana Narasimhan, Owen Rambow, Dipti Misra Sharma, and Fei Xia. 2009. Hindi syntax: Annotating dependency, lexical predicate-argument structure, and phrase structure. In *The 7th International Conference on Natural Language Processing*, pages 14–17.
- Thomas Proisl. 2018. SoMeWeTa: A part-of-speech tagger for German social media and web texts. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, pages 665–670, Miyazaki. European Language Resources Association.
- Office of the Registrar General & Census Commissioner. 2011. 2011 census data, Data on language and mother tongue, Statement 1: Abstract of speakers' strength of languages and mother tongues.
- Martin Riedl and Sebastian Padó. 2018. A named entity recognition shootout for German. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL 2018), Volume 2: Short Papers, pages 120–125, Melbourne.
- Kristina Toutanova, Dan Klein, Christopher D. Manning, and Yoram Singer. 2003. Feature-rich part-ofspeech tagging with a cyclic dependency network. In *Proceedings of HLT-NAACL 2003*, pages 252– 259, Edmonton.
- Kristina Toutanova and Christopher D. Manning. 2000. Enriching the knowledge sources used in a maximum entropy part-of-speech tagger. In *Proceedings* of *EMNLP/VLC-2000*, pages 63–70.
- Manindra K. Verma. 2003a. Bhojpuri. In George Cardona and Dhanesh Jain, editors, *The Indo-Aryan Languages*, pages 566–589. Routledge, London.
- Sheela Verma. 2003b. Magahi. In George Cardona and Dhanesh Jain, editors, *The Indo-Aryan Languages*, pages 547–565. Routledge, London.