

Evaluating Ensemble-Based Classifiers for Khasi Named Entity Recognition with Data Imbalance Constraints

Ransly Hoojon¹, Amitabha Nath^{2*} and Saralin Lyngdoh³

^{1,2} Information Technology Department, North Eastern Hill University, Umshing
Mawkynroh, Shillong, 793022, Meghalaya, India

³ Linguistic Department, North Eastern Hill University, Umshing
Mawkynroh, Shillong, 793022, Meghalaya, India

Emails: ranslyh@gmail.com¹, amitabha.me@gmail.com², saralyngdoh@gmail.com³

* Corresponding Author

Abstract

Named Entity Recognition (NER), a core task of NLP, has been the subject of much research in resource-rich languages. However, little has been attempted with respect to low-resource languages such as Khasi. Based on our previous efforts, the existing works on Khasi NER applied transformer models such as RoBERTa, XLM-RoBERTa, and BERT-base-multilingual-cased using transfer learning (TL) and achieved an F1 score of 91.25% with RoBERTa. In this paper, ensemble methods are investigated to improve performance even further. We used a novel Khasi named entity recognition dataset annotated in span-based and converted to IOB2 formats for token-level predictions for both TL and ensemble methods. Ensemble methods, including majority voting, weighted voting, and stacking with logistic regression and XGBoost, were employed. Two settings were tested, one where two models are implemented (RoBERTa and XLM-RoBERTa) and the other with three models, with the addition of BERT-base-multilingual-cased, where these models act as learners. The ensemble of three models with XGBoost scored the highest F1 score of 94.06%, which is an improvement of +2.81% over the previous best. It also achieved significant improvements on rare entities, such as WORK_OF_ART, LANGUAGE, ORDINAL, and EVENT, although the dataset was skewed toward the PERSON class. The results show how ensemble methods

34 can improve NER performance on imbalanced datasets and low-resource languages such as
35 Khasi.

36 **Keywords:** Named Entity Recognition, Low-Resource Language, Khasi Language,
37 Transformer-based Models, Ensemble Learning

38 **1.0 Introduction**

39 Named entity recognition (NER) whose task is to classify named entities (NEs) into
40 predefined categories (Sharma *et al.*, 2022) such as names of persons, organizations, and
41 locations, among others (Li *et al.*, 2022), from textual data serves as a one of the important
42 aspect in the NLP pipeline processing, which plays a vital role in NLP applications including
43 question answering (Diefenbach *et al.*, 2018; Rogers *et al.*, 2023; Lewis *et al.*, 2019),
44 information retrieval system (Gupta and Dixit, 2023; Eswaraiah and Syed, 2023), linking of
45 entities (Al-Moslmi *et al.*, 2020; Santana *et al.*, 2023; Mulwad *et al.*, 2023), and search
46 engines (Riaz, 2018; Ssemwogerere *et al.*, 2023; Orekhov *et al.*, 2023). In resource-rich
47 languages, NER benefits from abundant annotated datasets and relatively well-distributed
48 classes, leading to strong results. Whereas, for a low-resource language like Khasi spoken by
49 around one million people in North East India, the situation is markedly different. Besides the
50 rise in digitizing textual data, unique challenges remain for entity extraction, such as low
51 quality digitization, inconsistent spellings, limitations of OCR and the lack of textual data
52 covering diverse topics, the challenges are further compounded with the unavailability of an
53 annotated NER dataset where annotation has to be done manually from scratch, which
54 introduces the problem of class imbalance, where frequent entities like GPE (9,436 instances)
55 overshadow less common ones like WORK OF ART (15 instances), hindering the
56 effectiveness of standard NER methods. Ambiguity and multi-word entities are another
57 challenge when working with Khasi NER, e.g., “Ka Phan Nonglait,” which can be a location
58 or a person’s name. Thus, contextual understanding of the language is needed to be able to

59 disambiguate and produce the correct label. An additional issue is the presence of nested
60 entities, an entity inside another entity, e.g., “U Soso Tham Auditorium” where “U Soso
61 Tham” is a person and “U Soso Tham Auditorium” can be labelled as FAC. In our previous
62 studies, we utilized pre-trained models like RoBERTa, XLM-RoBERTa, and BERT-base-
63 multilingual-cased for Khasi NER through transfer learning (TL) fine-tuned based on the
64 annotated Khasi NER dataset consisting a total of 119,932 annotations (Hoojon, 2025),
65 achieving some encouraging results of 91.25%, 91.17%, and 90.39%, respectively, however,
66 these initiatives have drawbacks which fail to positively identify rare entity types, a limitation
67 that is amplified by unbalanced datasets. This goes to show that relying solely on single
68 models may not be a suitable approach, indicating the need for more creative approaches
69 tailored to the specific requirements of low-resource situations. Thus, in this paper, we
70 investigate a different method using ensemble techniques as a way to enhance the
71 performance of the Khasi NER. Using our previously fine-tuned models, we combined them
72 together, with the intention of leveraging their complementary advantages to boost predictive
73 accuracy and address class imbalances. Our contributions in this paper include the following:

- 74 1. Showcasing the effectiveness of ensemble methods in improving NER outcomes for
75 low-resource languages.
- 76 2. Substantially enhancing the identification of rare entities despite class imbalance.
- 77 3. Presenting a generalizable framework that can guide NER advancements for other
78 under-resourced languages.

79 **2.0 Related Work**

80 The use of multiple models to take advantage of their strengths and produce better predictive
81 performance has been a founding principle of machine learning and ensemble learning in
82 particular, going back many years in applications such as classification, time series
83 forecasting, and, more recently, NER. Ensemble methods, which combine models using

84 methods such as stacking and voting, overcome the shortcomings and limitations of a single
85 classifier, such as overfitting, bias, and sensitivity to data scarcities. This section summarizes
86 some of the more recent developments in ensemble learning, including their datasets,
87 methods, results, and conclusions, to provide background for the proposed study.

88 Ensemble methodologies utilizing XGBoost are extensively examined within diverse
89 application domains owing to their capability to capture complex nonlinear relationships,
90 thereby significantly augmenting predictive accuracy. XGBoost is helpful in feature selection
91 and gaining insight into an understanding of relationships in the data, as it can help find those
92 features in the dataset that are most significant (Banga *et al.*, 2021). XGBoost has also been
93 used in stock market prediction (Naik and Mohan, 2021) and sales forecasting (Shilong *et al.*,
94 2021). XGBoost is frequently applied to binary and multi-class classification tasks, and has
95 been utilized, for example, for real-time accident detection (Parsa *et al.*, 2020), customer
96 churn prediction (Tang *et al.*, 2020), credit risk analysis (Li *et al.*, 2022; Liu *et al.*, 2022;
97 Wang *et al.*, 2022), and medical diagnosis (Ogunleye and Wang, 2019). The class imbalance
98 issue is dealt with through the modeling process via loss function weighting among other
99 hyperparameters already available in XGBoost (Wang *et al.*, 2020). Among others, XGBoost
100 has been applied to class imbalance issues in medicine by (Létinier *et al.*, 2021) for rare
101 adverse drug reactions and to multiple imbalanced datasets (Zuech *et al.*, 2021; Le *et al.*,
102 2022; Mishra *et al.*, 2021; Mushava and Murray, 2022).

103 The paper (Pavlyshenko, 2018) explores the use of multi-level stacking for time series
104 predictions as well as logistic regression and supports the benefits of multi-layered structures
105 for complicated predictions. This research utilizes three datasets published in Kaggle: the
106 “Rossmann Store Sales” contains 1,115,000 records, “Grupo Bimbo Inventory Demand”
107 covers 800,000 stores with sales data, and “Bosch Production Line Performance” has
108 1,183,747 records with imbalanced labels regarding manufacturing failure. In the case of time

109 series, first-level predictions obtained from first-level models such as GradientBoosting or
110 ExtraTrees on the validation sets are incorporated as features into second-level models of the
111 Lasso-type linear regression or Random Forests. Temporal integrity is maintained by not
112 performing regular cross-validation but rather splitting based on time. In the case of logistic
113 regression, stacking combines XGBoost classifiers with a second-level linear or Bayesian
114 regression while using undersampling and one-hot encoding to tackle the class imbalance.
115 “Grupo Bimbo” solution implements a three-level stack where the first and second levels are
116 XGBoost and ExtraTrees, and the third layer is a weighted average. Stacking achieves an
117 accuracy of 5-10% better than the accuracy of individual models, while the three-level Grupo
118 Bimbo model reached a top score in Kaggle competitions with a log loss of 0.45. In the case
119 of “Bosch”, stacked XGBoost achieves an MCC of 0.35, which is higher than single
120 XGBoost’s MCC of 0.30. It is shown that stacking is capable of representing variations of
121 parameters of the models and exploiting heterogeneity of the first-level predictions to be more
122 robust, especially in the case of imbalanced data.

123 To increase accuracy in classification proposed a probability-weighted voting ensemble
124 (Rojarath and Songpan, 2020). The benchmarks employed in this work are the CPU dataset,
125 with 8,192 instances and 21 features, Hepatitis, with 155 instances and 19 features,
126 Ionosphere with 351 instances and 34 features, Sonar with 208 instances and 60 features, and
127 Vote with 435 instances and 16 features; these datasets include binary and multiclass
128 problems and are from UCI. They combine six base classifiers: Naïve Bayes, Bayesian
129 Networks, Decision Trees, Multilayer Perceptrons, SVM, k-NN with six ensemble
130 techniques: AdaBoost, Bagging, Stacking, Voting, Random Forests, Random Subspace. The
131 new method proposes to weight classifiers according to the accuracy obtained in the training
132 set, computing a weighted average of probabilities to predict a class. Among the experiments,
133 the 3PW-Ensemble, which fuses three out of three classifiers, yields the maximum accuracy

134 rates of 85.36% over Random Forests 74.99% and other ensembles 4PW: 84.87%.
135 Particularly for imbalanced datasets such as Hepatitis, precision and recall have significant
136 improvements (80% F1 vs. 70% of Random Forests). The study finds probability-weighted
137 voting to be beneficial to improve accuracy when votes are cast according to the true class
138 conditional distributions and finds that it is a flexible technique useful for many classification
139 problems. Performance, but, starts to deteriorate in cases of highly noisy data, indicating that
140 robust to noise weighting schemes are required.

141 The study (Won and Martins, 2018) examines name recognition in historical text. The two
142 historical corpora utilized are the Mary Hamilton Papers with 1,200 place entities and the
143 Samuel Hartlib Papers with 400 place entities. Five different NER systems, Stanford NER,
144 NER-Tagger, Edinburgh Geoparser, spaCy, and Polyglot, are tested and pooled using an
145 ensemble NER system that uses votes per tool. The ensemble 2-3 votes obtain F1 scores of
146 70% in all scenarios, and the highest 73.3% F1 on Hartlib letters. Stanford NER is the best for
147 Hartlib (F1 approximately 68%), while Polyglot performs the best for Hamilton (F1
148 approximately 65%). The ensemble's higher recall of 75% outperforms individual tools.
149 (Ullah *et al.*, 2024) addresses NER in a low-resource language, Urdu. UNER, which
150 originally contains 50,692 tokens and includes person, location, organization and
151 miscellaneous entities, has been augmented as UNER-II, which contains a total of 160,132
152 tokens and a total of 114,912 labeled entities. To supplement the dataset, CWEA produces
153 synthetic entities (47,600 persons, 30,940 locations, 30,900 organization). Four transformer
154 models, BERT-multilingual, RoBERTa-Urdu-small, BERT-base-cased, BERT-large-cased, as
155 well as handwritten RNN and BiLSTM baselines, are all fine-tuned with a 70:10:20 train:
156 validation: test split. BERT-multilingual reports a macro F1 of 98.2% on UNER-II, which is
157 higher than RoBERTa-Urdu-small 88.4%, BERT-base-cased 90.8%, BERT-large-cased
158 91.6%, RNN 97%, and BiLSTM 96%. When evaluated on the original UNER it achieves

159 84.6% which is better than the previous RNN-based approach which scored 77%.
160 Performance is greatly improved by CWEA as it diversifies the dataset and achieves error
161 rates of 1-2% for its common entities. The types of errors encountered indicate that this can be
162 improved by using better tokenization techniques, such as Word Piece, and more
163 comprehensive gazetteers. The study highlights BERT-multilingual's adaptability to low-
164 resource languages, with potential applications to Pashto and Persian.

165 In the medical domain, (Jin *et al.*, 2025) tackles NER and relation extraction (RE) problems.
166 The four datasets employed are: BioRED (biomedical texts with 26,956 entities and 5,151
167 relations), RDD (rare disease texts with 4,256 entities and 1,957 relations), ADE (11,070
168 entities on adverse drug effects), and DIANN (1,656 mentions of disability annotations).
169 PubMed-T5, COMCARE combines the PubMedBERT and PubMed T5 embedding with
170 BiLSTM layers and a CRF to maintain sequence consistency and employs threshold 0.6
171 collaborative decision fusion which integrates two matrices of emissions. For RE, it uses
172 BERT, PubMedBERT, and PubMed-T5, along with semantic chunking using 512 tokens with
173 a stride of 128 tokens to address long-range dependencies while leveraging NER outputs to
174 inform RE extraction. COMCARE's F1 scores are: 93.76% NER and 68.73% RE on
175 BioRED, 77.86% NER and 86.79% RE on RDD, 82.48% NER on ADE, and 99.36% NER on
176 DIANN; these scores exceed those of baselines such as BERT-BiLSTM-CRF (92.51% NER
177 on BioRED) and LLMs such as GPT-4 (0% on DIANN). (Munthe, 2024) applies stacking to
178 finance forecasting. The analysis employs 335 news records related to the trends of the
179 Indonesian stock market Jakarta Composite Index which are classified as "Up" with 714
180 records and "Down" with 335 records, resulting in a dataset that has class imbalance. SMOTE
181 creates synthetic "Down" samples in order to re-balance the data. It involves tokenization,
182 stemming and TF-IDF vectorization. Predictions are then made by base models XGBoost,
183 Random Forest and passed to meta models Logistic Regression and LSTM. Two stacking

184 techniques are explored, using either XGBoost/Random Forest with Logistic Regression or
185 LSTM. Stacking with Logistic Regression obtains 86% accuracy, precision 85-86%, recall
186 84-87% and F1 85-86% outperforming Random Forest 85%, XGBoost 82%, Logistic
187 Regression 81% and stacking with LSTM achieved 83% accuracy. A ROC-AUC score of
188 0.91% indicates a strong discrimination. SMOTE mitigates majority class bias to enhance the
189 prediction of the minority class. Relative to previous research, such as SVM performance at
190 84% and CNN-LSTM performance at 85%, this approach performs better given the simplicity
191 and interpretability of Logistic Regression as the meta-model.

192 **3.0 Methodology**

193 The datasets are collected from various online sources, such as newspapers and books, with
194 an effort to cover different ranges of topics, which are then cleaned, standardized for
195 diacritics, and corrected for spelling errors, which are prevalent when downloading online
196 content, and then the tasks of manual annotations are performed across 19 different entity
197 types with a total of 30,776 annotated entities. Our annotators annotated the test set in flat
198 format by labeling spans of entities with the following inputs: “text”: “Shillong”, “start”: 0,
199 “end”: 8, “label”: “GPE”, which we converted to IOB2 format later. The conversion from flat
200 to IOB2 is done by using B to denote the first token of a multi-token entity combined with the
201 type of that entity, for example, B-PERSON for the beginning of a person’s name. Tokens
202 that are inside an entity are labeled with I, for example, I-PERSON, and tokens that are not
203 part of any entity are tagged with O, which stands for outside. This was done for the reasons
204 of better interoperability, less ambiguous definitions of multi-token entities, and better metrics
205 reporting on the level of single tokens. Table 1 shows the IOB2 tagging example for a Khasi
206 sentence “U John na Shillong” and Table 2 provides an overview of the distribution of entities
207 in the Khasi NER Test Dataset.

Table 1: IOB2 Tagging

Token	IOB Tag
U	B-PERSON
John	I-PERSON
na	O
Shillong	I-GPE

208

Table 2: Entity Distribution with Descriptions and Examples

Entity	Count	Description	Example
GPE	9436	Countries, cities, states	India, Shillong, Delhi
PERSON	6046	People or fictional characters	U Soso Tham, John Smith
ORG	3675	Companies or institutions	Microsoft, FBI
NORP	3195	Nationalities or groups	Republican Party
LOC	2134	Non-GPE locations	Nile River
DATE	1849	Calendar dates	06/09/23
CARDINAL	1631	Numerals	2, Three, ar, lai
MONEY	623	Monetary amounts	\$20, ₹500, 10 tyngka
LAW	521	Legal documents	Roe v. Wade
TIME	388	Time periods	4 hours, 5 kynta, 10 baje
DAY	303	Days of the week	Monday, Sngi Nyngkong
PRODUCT	239	Physical goods	iPhone, Pizza
EVENT	191	Events or occurrences	FIFA World Cup
QUANTITY	171	Measurements	50kg, 1km
PERCENT	153	Percentages	80%
FAC	126	Facilities	Logan Airport
ORDINAL	82	Ordered numbers	9th, Ninth
LANGUAGE	48	Languages	English, Khasi
WORK_OF_ART	15	Titles of works	Mona Lisa

209

210 3.1 Base Models

211 In our ensemble experiment, we utilized the previously fine-tuned models, including
 212 RoBERTa-base, XLM-RoBERTa-base, and BERT-base-multilingual-cased, by combining
 213 them with the default Spacy NER layer (DSNER) for entity prediction, which was trained on
 214 the Khasi NER dataset featuring a training sample size of 10,756 examples that encompass a
 215 total of 119,932 entities (Hoojon, 2025). These transformer models were chosen to investigate

216 whether there is a notable difference in performance for Khasi NER between monolingual and
217 multilingual approaches.

218 **3.2 Ensemble Learning**

219 In this work, we experimented with three ensemble learning techniques: Majority Voting,
220 Weighted Voting, and Stacking using Logistic Regression and XGBoost, each aimed at
221 utilizing model diversity while tackling the sparse and imbalanced dataset. These approaches
222 are outlined below, adhering to the standard principles of ensemble learning (Ganaie *et al.*,
223 2022).

224 **3.2.1 MAJORITY VOTING**

225 Majority Voting selects the label with the most votes from the models for each token. For a
226 token i , let $M = \{m_1, m_2, \dots, m_k\}$ denote the set of K models ($K = 2$ for the two-model
227 ensemble, $K = 3$ for the three-model ensemble), and $y_m^{(i)} \in L$ be the predicted label by model
228 m , where $L = \{\text{B-PERSON, I-PERSON, O, B-GPE, . . .}\}$ is the label set. The number of votes
229 for label $l \in L$ is:

$$230 \quad V_l^{(i)} = \sum_{m \in M} 1(y_m^{(i)} = l) \quad (1)$$

231 The final prediction is:

$$232 \quad \hat{y}^{(i)} = \arg \max_{l \in L} V_l^{(i)} \quad (2)$$

233

234 In case of ties when multiple labels have the same number of votes, we default to the
235 RoBERTa prediction, as it is the model with the overall best performance. This method
236 ensures robustness by favoring consensus, which is important for low-resource languages
237 where individual models may overfit (Rahimi *et al.*, 2019; Chen *et al.*, 2022).

238 **3.2.2 Weighted Voting**

239 In weighted voting, the validation performance of each model is assigned a weight, where
 240 models with higher accuracy are given priority. We use F1 scores from a held-out Khasi
 241 validation set: $W_{\text{RoBERTa}} = 0.9125$, $W_{\text{XLM-RoBERTa}} = 0.9117$, and $W_{\text{BERT}} = 0.9039$ (for the three-
 242 model ensemble). The weighted score for label $l \in L$ is:

$$243 \quad W_l^{(i)} = \sum_{m \in M: y_m^{(i)} = l} w_m \quad (3)$$

244 The final prediction is:

$$245 \quad \hat{y}^{(i)} = \arg \max_{l \in L} W_l^{(i)} \quad (4)$$

246 While this approach leverages the fine-tuned performance of all three models, it also takes
 247 advantage by incorporating cross-lingual insights from XLM-RoBERTa and BERT-base-
 248 multilingual-cased models. Weighted Voting is particularly effective as it balances model
 249 strengths without the need for additional training (Lima *et al.*, 2023).

250 3.2.3 Stacking with Logistic Regression

251 Stacking trains, a meta-classifier to combine model predictions (Li, 2025), learning which
 252 model is most reliable for specific tokens. For token i , the feature vector is
 253 $(x^{(i)} = [y_{m_1}^{(i)}, y_{m_2}^{(i)}, \dots, y_{m_K}^{(i)}])$, where predictions are encoded numerically using a Label Encoder
 254 (e.g., O \rightarrow 0, B-PERSON \rightarrow 1). A logistic regression (LR) meta-classifier, parameterized by
 255 θ , predicts the probability of each label:

$$256 \quad P(y^{(i)} = l | x^{(i)}, \theta) = \frac{\exp(\theta_l \cdot \phi(x^{(i)}))}{\sum_{l' \in L} \exp(\theta_{l'} \cdot \phi(x^{(i)}))} \quad (5)$$

257 where,

- 258 • $\phi(x^{(i)})$ is the encoded feature vector.
- 259 • θ_l represents the parameters (weights) of the Logistic Regression model for label l ,
- 260 learned during training and
- 261 • $\exp(\theta_l \cdot \phi(x^{(i)}))$ computes an unnormalized score (logit) for label l .

262 The final prediction is:

$$263 \quad \hat{y}^{(i)} = \arg \max_{l \in L} P(y^{(i)} = l | x^{(i)}, \theta) \quad (6)$$

264 The meta-classifier uses cross-entropy loss during training, which measures the difference
265 between predicted probabilities and true labels by adjusting θ to maximize its potential to
266 handle complex patterns.

267 **3.2.4 XGBoost as a Meta-Classifier**

268 XGBoost, which stands for Extreme Gradient Boosting, is used as the meta-classifier in our
269 NER system to combine predictions of the base models, which include RoBERTa, XLM-
270 RoBERTa and BERT-multilingual-cased, in a stacking ensemble. XGBoost is a scalable and
271 efficient implementation of gradient boosting, which is good at picking up complex patterns
272 and is well suited for combining the outputs of base models in NER (Bentéjac *et al.*, 2021;
273 Bartz-Beielstein *et al.*, 2023; Biau and Cadre, 2021).

274 **3.2.4.1 Xgboost**

275 XGBoost develops a collection of decision trees, wherein the forecast for an input feature
276 vector \mathbf{x}_i (comprising encoded predictions from the base model for a token) is derived from
277 the aggregated outputs across K trees. In the context of multi-class classification of IOB tags,
278 such as “B-PERSON” or “O,” the score for a given class c is calculated as follows:

$$279 \quad \hat{y}_{i,c} = \sum_{k=1}^K f_{k,c}(\mathbf{x}_i) \quad (7)$$

280 where $f_{k,c}(\mathbf{x}_i)$ is the score from the k^{th} tree for class c , and $\hat{y}_{i,c}$ is the raw score for instance i .

281 Probabilities are obtained via the softmax function:

$$282 \quad P(y_i = c | \mathbf{x}_i) = \frac{\exp(\hat{y}_{i,c})}{\sum_{j=1}^C \exp(\hat{y}_{i,j})} \quad (8)$$

283 where C is the number of IOB tags. The predicted class is:

$$284 \quad \hat{c}_i = \arg \max_c P(y_i = c | \mathbf{x}_i) \quad (9)$$

285 Here, \mathbf{x}_i comprises encoded IOB predictions from base models, and classes represent IOB tags
286 for 19 entity types plus “O”.

287 3.2.4.2 Training Objective

288 XGBoost optimizes an objective function balancing accuracy and complexity:

$$289 \text{Obj} = \sum_{i=1}^n \ell(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k) \quad (10)$$

290 where $\ell(y_i, \hat{y}_i)$ is the multi-class log loss:

$$291 \ell(y_i, \hat{y}_i) = - \sum_{c=1}^C \mathbb{1}(y_i = c) \log(P(y_i = c | \mathbf{x}_i)) \quad (11)$$

292 and $\Omega(f_k) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$ regularizes the model, with T as the number of leaves, w_j as the
293 leaf weights, γ as the leaf penalty, and λ as the L2 regularization parameter. Leaf weights are
294 computed as:

$$295 w_j = - \frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda} \quad (12)$$

296 where I_j is the set of instances in leaf j , $g_i = \frac{\partial \ell}{\partial \hat{y}_i}$ is the gradient, and $h_i = \frac{\partial^2 \ell}{\partial \hat{y}_i^2}$ is the Hessian
297 (Biau and Cadre, 2021).

298 3.3 Algorithms

299 Algorithm 1: Data Preparation

300 **Input:** Dataset source with texts and entity annotations

301 1: Initialize $D \leftarrow []$, $T \leftarrow []$

302 2: **for** each pair in the parsed dataset **do**

303 3: **if** the pair has valid text and entities, then

304 4: Append pair to D , entities to T

305 5: **else**

306 6: raise error: “Invalid dataset structure.”

```
307 7: end if  
308 8: end for  
309 9: return D, T
```

310 Algorithm 1 shows how a dataset and entity annotations are preprocessed and validated to
311 serve as the basis of the ensemble learning pipeline. It populates two lists: D, which holds the
312 valid text-entity pairs, and T, which contains the corresponding entity labels, iterating through
313 the data set while appending valid pairs to D and the respective entities to T, raising an error
314 in case of invalid entries so as to ensure data integrity. This algorithm is key in the preparation
315 of structured data for model predictions and ensemble techniques.

316 **Algorithm 2: Model Loading and Prediction Generation**

```
317 Input: Dataset D, configurations for base models  
318 1: Initialize dictionaries M (models), P (predictions)  
319 2: for each model name, path in configurations do  
320 3: M[name] ← load model from path  
321 4: end for  
322 5: for each model name, model in M do  
323 6: P[name] ← model predictions on texts from D  
324 7: end for  
325 8: return M, P
```

326 Algorithm 2 loads predefined base models and generates predictions on dataset D. It
327 bootstraps dictionaries M and P for the models and their predictions, respectively, loading
328 each model and using it for analyzing the texts in D to arrive at predictions. This process is
329 critical in providing the initial predictions needed for majority voting, weighted voting, and
330 stacking.

331 **Algorithm 3: Stacking Meta-Classifer Training**

332 **Input:** Dataset D, true labels T, models M, predictions P

333 1: Initialize lists X (stacking features), Y (true labels), L (label set)

334 2: **for** each example, index in D **do**

335 3: Extract text and true spans; process with reference model to get tokens

336 4: Convert true spans to labels \rightarrow Y; add unique labels to L

337 5: **for** each model name in M **do**

338 6: Convert P[name][index] to labels or probabilities \rightarrow X; update L

339 7: **end for**

340 8: **end for**

341 9: Fit label encoder on L, encode X and Y

342 10: Train meta-classifier (logistic regression or XGBoost) on encoded X, Y

343 11: **return** meta-classifier, label encoder

344 Algorithm 3 implements the stacking ensemble approach by transforming the base model
345 predictions into features and training a meta-classifier, logistic regression, or XGBoost, which
346 combines the base predictions into a single output. It initializes lists X for stacking features,
347 the list of true labels Y, and the list of unique labels L, by iterating through each example in D
348 to produce the true spans and transforming the predictions of the basemodel contained in P
349 into features. A label encoder is fitted to encode X and Y, enabling the training of the meta-
350 classifier on these features.

351 **Algorithm 4: Chunking and Ensemble Prediction**

352 **Input:** Dataset D, true labels T, models M, predictions P, meta-classifier, label encoder,

353 predefined weights $L = (l_1, l_2, l_3)$, test dataset T_b

354 1: Initialize empty list C for labeled chunks; position $\leftarrow 0$

355 2: Split D into sentences; compute sentence embeddings S
356 3: **while** position < length of D **do**
357 4: Determine chunk end by semantic similarity in S and entity spans in T
358 5: Extract chunk text $D[\text{position}:\text{end}]$; tokenize using reference model
359 6: **For** each model, collect predictions from P and convert to token-level labels
360 7: Apply ensemble methods:
361 8: - Majority Voting: Select label with the most votes per token
362 9: - Weighted Voting: Select label with the highest weighted sum using L
363 10: - Stacking: Use a meta-classifier to predict labels from encoded predictions
364 11: Append (chunk text, final labels) to C
365 12: Update position to chunk end
366 13: **end while**
367 14: Repeat steps 2–12 for T_b to generate test predictions C_t
368 15: **return** C, C_t

369 Algorithm 4 employs the ensemble methods majority voting, weighted voting, and stacking
370 by partitioning the dataset D into semantically related segments and assigning final labels to
371 every token in each segment, and processing the test dataset T_b in a similar way. It maintains a
372 list C of labeled chunks, computes sentence embeddings for segmentation, and for each
373 chunk, it tokenizes the chunk text and collects predictions from P . The algorithm uses
374 majority voting to compute the most common label, weighted voting where the weights are
375 given to compute a weighted sum, and stacking, where it applies the meta-classifier to make
376 the predictions. Labeled chunks are stored in C , and the process is repeated for T_b to produce
377 test predictions C_t .

378 **Algorithm 5: Evaluation and Output**

379 **Input:** Chunks C , true labels T , test chunks C_t , test labels T_t

```
380 1: Validate that C and T correspond to the same text spans; if Ct and Tt are provided, validate
381 similarly
382 2: Initialize an empty dictionary M for evaluation metrics
383 3: if T not in IOB2 format then
384 4:   Convert T to IOB2 format by mapping true labels to chunk text spans
385 5: end if
386 6: for each chunk, true label in C, T do
387 7:   Compute micro-averaged precision, recall, and F1 score for predicted vs. true IOB2
388 labels
389 8:   Optionally compute per-entity-type scores (e.g., for “Person”, “Organization”)
390 9:   Store scores in M (e.g., M[“precision”], M[“Person_f1”])
391 10: end for
392 11: if Ct and Tt are provided then
393 12:   if Tt not in IOB2 format then
394 13:     Convert Tt to IOB2 format
395 14:   end if
396 15:   for each chunk, true label in Ct, Tt do
397 16:     Compute micro-averaged precision, recall, and F1 score
398 17:     Store scores in M (e.g., M[“test_precision”])
399 18:   end for
400 19: end if
401 20: return M (precision, recall, F1 scores for training and test data)
```

402 Here, Algorithm 5 constitutes the last component of the ensemble workflow, which aims to
403 perform a quantitative analysis of the achieved predicted chunks against the ground-truth
404 labels in order to validate the strength of the entity recognition and chunking pipeline. The

405 algorithm takes as input two main parameters, C , the predicted chunks that result from the
406 ensemble techniques of Algorithm 4 (i.e., majority, weighted voting, stacking), and T , the true
407 labels from the data preparation of Algorithm 1, which represent the ground-truth entity
408 annotations. It optionally takes C_t (test predictions) and T_t (test true labels) for performance
409 evaluation on unseen data. The algorithm starts by checking if the predicted chunks in C align
410 with the text spans in T (and chunks in C_t and spans in T_t , if applicable). A new empty
411 dictionary M is initialized as a place to store metrics for evaluation. If T (or T_t) is not in IOB2
412 format, it is first converted to IOB2 by tagging spans of text covered by chunks with the
413 corresponding IOB2 tags, as it matches the predicted in C , which are also IOB2. The core
414 evaluation iterates over each chunk and true labels, averaging precision, recall and F1 at the
415 level of individual characters for the purpose of describing the accuracy of predicted IOB2
416 labels in comparison to true labels using the micro-averaging method. Per-entity type scores
417 (for example “Person” or “Organization”) can optionally be computed as well. These metrics
418 are recorded in M with keys like $M[\text{“precision”}]$ or $M[\text{“Person_f1”}]$. When C_t and T_t are
419 available for test data, a parallel computation of the same metrics is carried out and stored
420 separately (for example, $M[\text{“test_precision”}]$). The algorithm produced an output M ,
421 encapsulating a comprehensive set of performance metrics. The following formulas were used
422 to calculate precision, recall, and F1 score (Jurafsky and Martin, 2020):

$$423 \quad \text{Precision} = \frac{TP}{TP + FP} \quad (13)$$

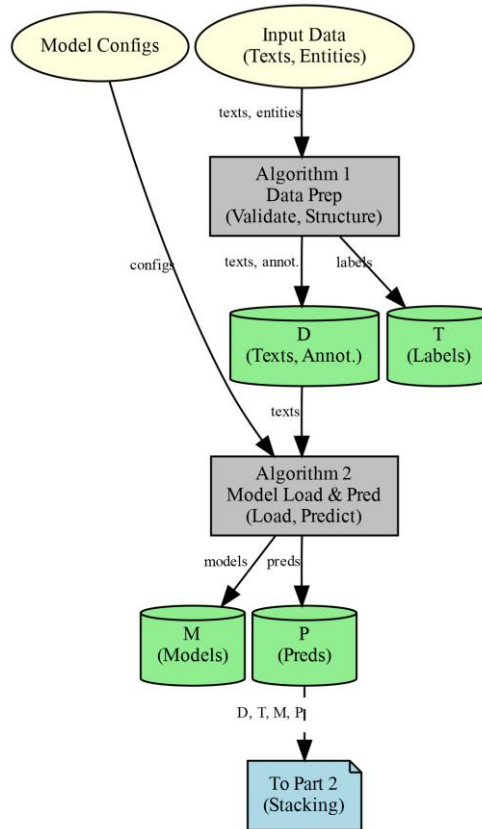
$$424 \quad \text{Recall} = \frac{TP}{TP + FN} \quad (14)$$

$$425 \quad F_1 = 2 \cdot \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (15)$$

426 **4.0 Experimental Setup**

427 To improve our predictions and make them more reliable, we used ensemble methods in two
428 settings: one with two models (RoBERTa and XLM-RoBERTa-base) and another with three

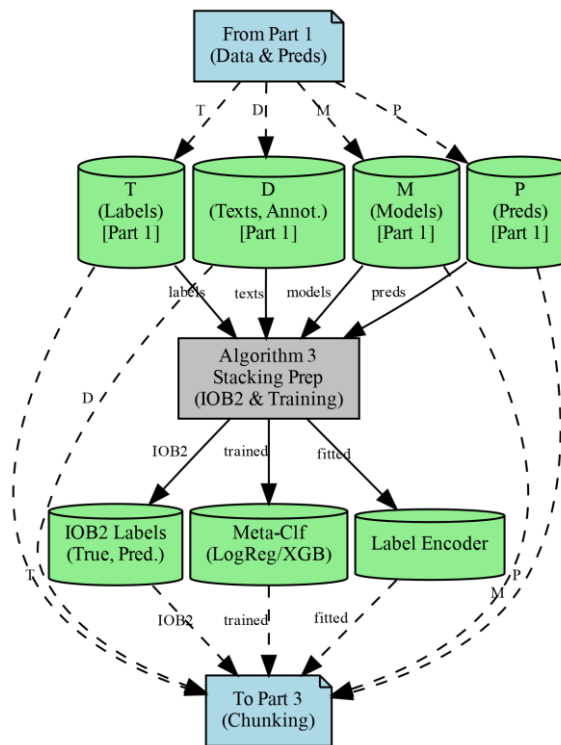
429 models (RoBERTa, XLM-RoBERTa-base, and BERT-base-multilingual-cased). We tested
430 these on a dataset with 2,675 examples from 19 different categories, totaling 30,766 labeled
431 entities. These pre-trained models were fine-tuned for the Khasi NER task. The goal of the
432 experiment was to evaluate the different models both separately and in combination by means
433 of several ensemble techniques: majority voting, weighted voting, and stacking via meta-
434 classifiers such as logistic regression (LR) and XGBoost. For each model, the input text was
435 processed by passing it through the corresponding tokenizers of the models, and the
436 predictions were made at the span level, which were then converted to predictions at the token
437 level by using character offsets. For the ensemble methods, we used three different strategies
438 for aggregating predictions of the models. In cases of ties, the selection would default to the
439 label suggested by RoBERTa, which is the best-performing model. In weighted voting, the
440 predictions were aggregated according to the accuracy of each model, with weights equal to
441 0.9125, 0.9117, and 0.9039 for RoBERTa, XLM-RoBERTa, and BERT-multilingual,
442 respectively. The stacking meta-classifiers, which were trained on the predictions of the base
443 models, provide predictions for the final token label. Evaluation is done on standard NER
444 metrics precision, recall, and F1 score calculated at the token level. Figure 1 through 4 shows
445 the system setup for the ensemble method.



446

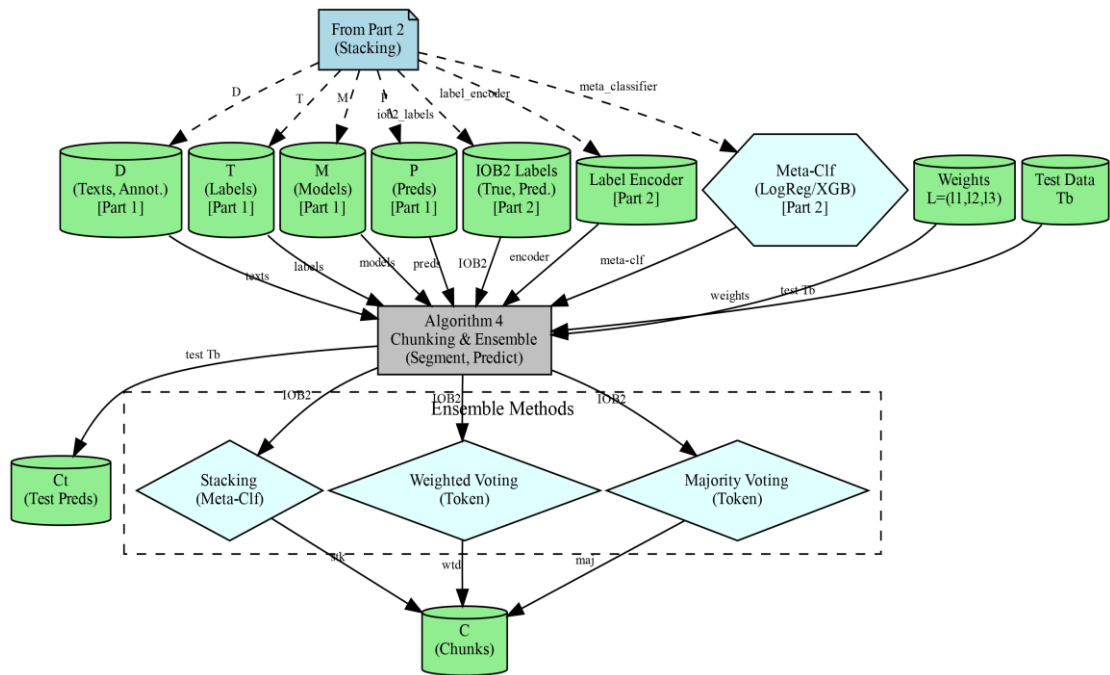
447

Figure 1: Text-entity preprocessing and validation



459

Figure 2: Base model loading and prediction generation pipeline

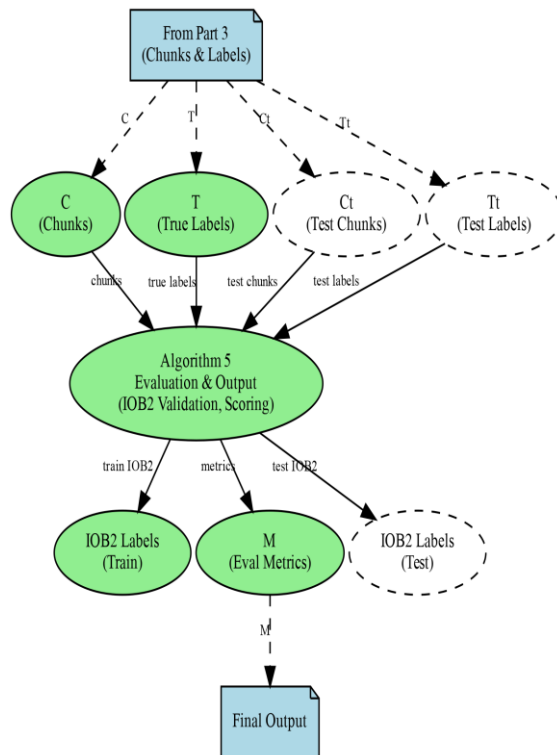


460

461

Figure 3: Ensemble labeling with segmentation and voting methods

462



463

464

Figure 4: Evaluation of ensemble predictions using precision, recall, and F1

465

466

Table 3: Model descriptions

Abbreviation	Description
RoBERTa	Robustly optimized BERT pretraining approach
XLM-R	Cross-lingual RoBERTa
BERT-M	Multilingual BERT
Maj-2	Majority Voting with 2 models
Wgt-2	Weighted Voting with 2 models
Stk-2 (LR)	Stacking Ensemble with 2 models (Logistic Regression)
Stk-2 (XGBoost)	Stacking Ensemble with 2 models (XGBoost)
Maj-3	Majority Voting with 3 models
Wgt-3	Weighted Voting with 3 models
Stk-3 (LR)	Stacking Ensemble with 3 models (Logistic Regression)
Stk-3 (XGBoost)	Stacking Ensemble with 3 models (XGBoost)

468 Referring to Table 3, we evaluated the performance of the ensemble techniques on Khasi
469 NER dataset by comparing with the baseline models. From the results, it is evident that 2-
470 model (RoBERTa and XLM-RoBERTa) and 3-model (RoBERTa, XLM-RoBERTa, and
471 BERT-multilingual-cased) ensembles for Named Entity Recognition (NER) outperform
472 individual baseline models considerably, but also indicate that their effectiveness is dependent
473 on the specific ensemble approach, as well as on the entity type. The baseline models,
474 RoBERTa with F1-score 91.25%, XLM-RoBERTa with F1-score 91.17%, and BERT-
475 multilingual-cased with F1-score 90.39%, have already excellent results, with RoBERTa
476 reaching the best performance overall. In the 2-model ensembles, majority voting (Maj-2) and
477 weighted voting (Wgt-2) show the same performance with an F1-score of 93.29% (+2.04%
478 respect to RoBERTa), and XGBoost stacking (Stk-2 XGBoost) outperforms both with an F1-
479 score of 93.65% (+2.40%). Logistic regression stacking (Stk-2), on the other hand, performs
480 very poorly with F1-score of 54.19% (-37.06%) and proves to be an inappropriate meta-
481 learner. The 3-model ensembles perform even better, with majority and weighted voting
482 (Maj-3/Wgt-3) obtaining an F1-score of 93.54% (+2.29%) and XGBoost stacking (Stk-3
483 XGBoost) reaching the highest F1-score of 94.06% (+2.81%) as an attestation of XGBoost's

484 ability to exploit complementary predictions provided by other models. Stk-3 with logistic
485 regression, still remains ineffective with F1-score of 62.87% (-28.38%). For high frequency
486 entities like PERSON with F1-score of 96.44% to 96.86%, GPE with F1-score of 94.19% to
487 94.64% and DATE with F1-score of 97.47% to 98.49% the per-entity analysis shows good
488 performance, with Stk-3 XGBoost prevailing in all cases. Rare entities such as
489 WORK_OF_ART with F1-score of 34.29% to 38.36%, PRODUCT with F1-score of 62.47%
490 to 68.81% and FAC with F1-score of 66.99% to 77.79% all have low F1 scores due to
491 support, which ranged from 57 to 355, where Stk-3 XGBoost only performed slightly better.
492 Logistic regression fails across multiple entities (e.g., F1-score of 0.0% for FAC,
493 PRODUCT), highlighting its limitations. The 3-model ensembles marginally outperform than
494 the 2-model ones, especially for entities like PERCENT with F1-score of 95.45% vs. 94.01%
495 and FAC with F1-score of 77.79% vs. 73.94%, suggesting that BERT-multilingual-cased
496 makes complementary predictions even if it is a weaker baseline. Micro-averaged F1-scores
497 (93.26%-94.02%) reflect the prevalence of entities with high support, whereas macro-
498 averaged F1-scores (84.48% -86.69%) suggest that Stk-3 XGBoost performs better at
499 classifying rare entities. In view of these results, we can state that the XGBoost stacking,
500 especially in the case of 3 models, is the best approach to improve NER accuracy, but
501 majority as well as weighted voting are also simpler, yet successful, alternatives. Table 4
502 compares the performance of the various ensemble approaches with the three base models in
503 terms of overall precision, recall, and F1 scores.

Table 4: Overall Performance Comparison of Individual and Ensemble Models

Method	Models Included	Precision	Recall	F1	Others vs. RoBERTa
RoBERTa	RoBERTa	91.12	91.39	91.25	---
XLM-R	XLM-RoBERTa	90.55	91.79	91.17	-0.08
BERT-M	BERT-multilingual-cased	90.24	90.54	90.39	-0.86
Maj-2	RoBERTa, XLM-RoBERTa	92.88	93.7	93.29	+2.04

Wgt-2	RoBERTa, XLM-RoBERTa	92.88	93.7	93.29	+2.04
Stk-2(LR)	RoBERTa, XLM-RoBERTa	54.74	53.65	54.19	-37.06
Stk-2 (XGBoost)	RoBERTa, XLM-RoBERTa	93.58	93.71	93.65	+2.4
Maj-3	RoBERTa, XLM-R, BERT-M	93.36	93.72	93.54	+2.29
Wgt-3	RoBERTa, XLM-R, BERT-M	93.36	93.72	93.54	+2.29
Stk-3 (LR)	RoBERTa, XLM-R, BERT-M	63.74	62.02	62.87	-28.38
Stk-3 (XGBoost)	RoBERTa, XLM-R, BERT-M	94.2	93.92	94.06	+2.81

504

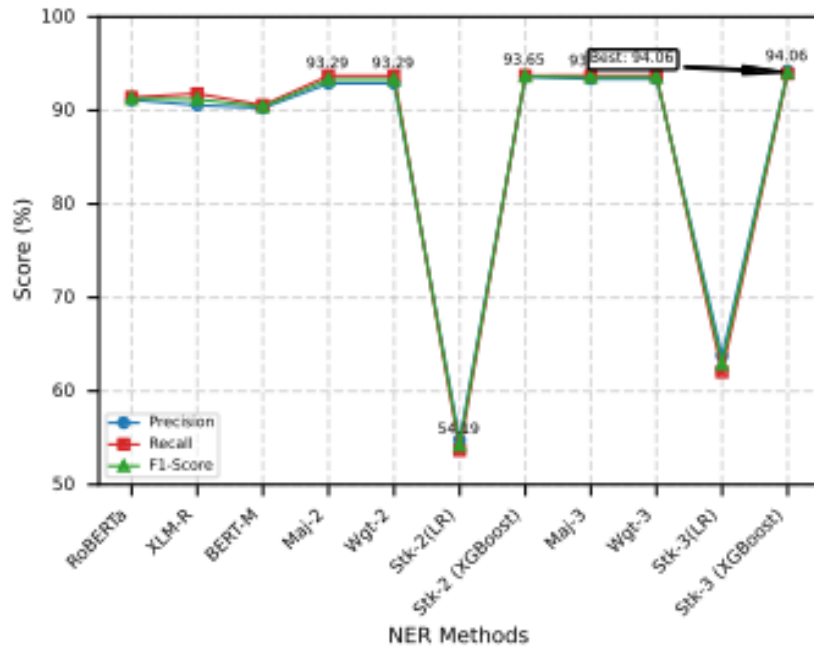
505 For prediction on new sentences, we intentionally selected rare entities to test the performance
506 of the various ensemble methods:

507 Sentence 1 features the rare entity *Ka Akor Kaba Tam Bynta-III* (name of a book), a
508 WORK_OF_ART. Only the 2-model and 3-model Stacking with XGBoost accurately
509 classified this entity. No other ensemble methods, were able to identify it. For instance, the
510 phrase “*la pyllait paidbah ia ka kot [Ka Akor Kaba Tam Bynta-III] da u Kongsan*
511 *Rangbah...*” shows a true positive only when stacked with XGBoost. This performance
512 demonstrates the higher capability of the XGBoost model, when compared to simpler
513 ensemble strategies, to understand the contextual subtleties required for identifying
514 complicated multi-token named entities.

515 Sentence 2 also investigates the distribution of the ordinal word *ba ar* (meaning second),
516 which is labeled with the ORDINAL category, e.g., ‘*...la tip kum ka bneng [ba ar]...*’ and
517 ‘*.ka bneng [ba ar] sha pyrthei...*’. All of the ensemble methods correctly classified them
518 except for 2-model and 3-model (stk-2(LR) and stk-3(LR)).

519 In Sentence 3, we tested with the entity LANGUAGE, e.g., ‘*...ka Jaka Hikai Ktien [Khasi] ne*
520 *ka Khasi Learning Centre.....ka Kim kam nang ban kren [Khasi]...*’. All ensemble methods,
521 2-model and 3-model variants of majority, weighted voting, as well as XGBoost stacking,
522 were able to successfully classify every instance of the positive class. The performance is
523 relatively consistent in this case, which indicates that the ensemble models are strong when

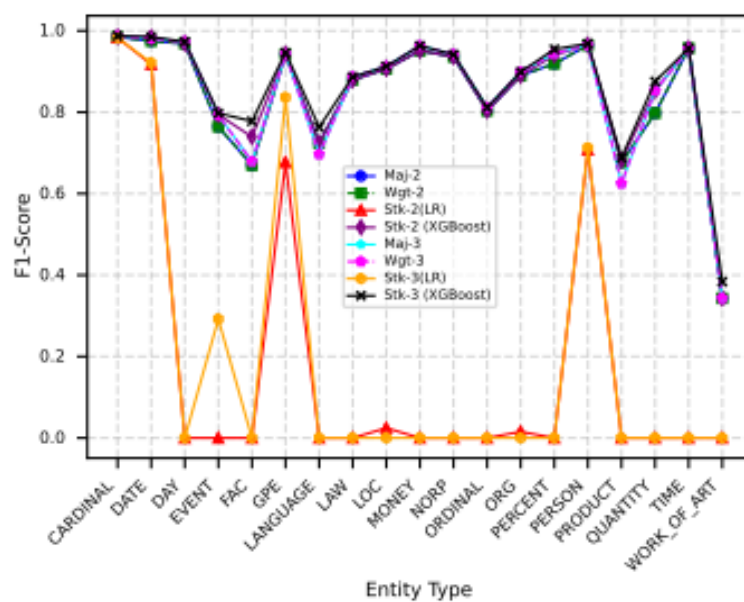
524 rare items occur in very indicative contexts, such as language names with the verb
 525 “*ktien/kren*” which means “*language/speak*”, with the exception of logistic stacking (stk-
 526 2(LR) and stk-3(LR)), which failed to identify the entity correctly. Plots of various model
 527 performances over all ensemble settings and per-entity predictions are shown in Figure 5 and
 528 Figure 6.



529

530

Figure 5: Evaluation metrics across ensemble and baseline models



531

532

Figure 6: Entity-wise precision, recall, and F1 of ensemble models

533 **6.0 Conclusion**

534 In conclusion, this paper aims to investigate whether ensemble learning methods can help to
535 improve NER system performance on the Khasi language. Our ensemble learning approach is
536 based on our previous work on individual transformer models, RoBERTa, XLM-RoBERTa,
537 and BERT-multilingual-cased trained on Khasi NER datasets. The experiment showed
538 improved performance over individual models; RoBERTa, the strongest baseline model with
539 an F1-score of 91.25%, was outperformed by XGBoost as a meta-learner using the three base
540 learners, with an improvement of +2.81%. This experiment also showed an improvement in
541 detecting rare and under-represented entities such as WORK_OF_ART, LANGUAGE,
542 ORDINAL, and EVENT, even though the dataset is skewed towards other classes such as
543 GPE, PERSON, ORG, NORP, LOC, DATE, and CARDINAL entities. Findings also show
544 that logistic regression as a meta-learner performs very poorly, and other limitations, such as
545 an imbalanced dataset and overdependence on other high-resource languages as foundational
546 models. Our future goal is to increase the dataset size, covering a wider range of topic
547 domains with well-balanced classes, adopting data augmentation, and exploring more
548 sophisticated meta-learners or methods that are best suited to the language. Overall, our
549 experimental approach shows that ensemble learning is a promising and scalable strategy for
550 improving NER performance in low-resource languages such as Khasi.

551

552

553 **REFERENCES**

554

555 Sharma, A., Chakraborty, S., Kumar, S. *et al.* (2022). Named entity recognition in natural
556 language processing: a systematic review. In: *Proceedings of the Second Doctoral Symposium*
557 *on Computational Intelligence*. Singapore: Springer, pp. 817–828.
558 https://doi.org/10.1007/978-981-16-3647-9_59

559

560 Li, J., Sun, A., Han, J. and Li, C. (2022). A survey on deep learning for named entity
561 recognition. *IEEE Transactions on Knowledge and Data Engineering*, 34(1), pp. 50–70.
562 <https://doi.org/10.1109/tkde.2020.2981314>

563

564 Diefenbach, D., Lopez, V., Singh, K. and Maret, P. (2018). Core techniques of question
565 answering systems over knowledge bases: a survey. *Knowledge and Information Systems*,
566 55(3), pp. 529–569. <https://doi.org/10.1007/s10115-017-1100-y>
567

568 Rogers, A., Gardner, M. and Augenstein, I. (2023). QA dataset explosion: a taxonomy of NLP
569 resources for question answering and reading comprehension. *ACM Computing Surveys*,
570 55(10), pp. 1–45. <https://doi.org/10.1145/3560260>
571

572 Lewis, P., Oğuz, B., Rinott, R., Riedel, S. and Schwenk, H. (2019). MLQA: evaluating cross-
573 lingual extractive question answering. *arXiv*. Available at: <https://arxiv.org/abs/1910.07475>
574

575 Gupta, V. and Dixit, A. (2023). Recent query reformulation approaches for information
576 retrieval system – a survey. *Recent Advances in Computer Science and Communications*,
577 16(1), pp. 94–107. <https://doi.org/10.2174/2666255815666220404091920>
578

579 Eswaraiah, P. and Syed, H. (2023). An efficient ontology model with query execution for
580 accurate document content extraction. *Indonesian Journal of Electrical Engineering and*
581 *Computer Science*, 29(2), pp. 981–989. <https://doi.org/10.11591/ijeecs.v29.i2.pp981-989>
582

583 Al-Moslmi, T., Ocaña, M.G., Opdahl, A.L. and Veres, C. (2020). Named entity extraction for
584 knowledge graphs: a literature overview. *IEEE Access*, 8, pp. 32862–32881.
585 <https://doi.org/10.1109/ACCESS.2020.2973928>
586

587 Santana, B., Campos, R., Amorim, E., Jorge, A., Silvano, P. and Nunes, S. (2023). A survey
588 on narrative extraction from textual data. *Artificial Intelligence Review*, 56(8), pp. 8393–8435.
589 <https://doi.org/10.1007/s10462-022-10338-7>
590

591 Mulwad, V., Finin, T., Kumar, V.S., Williams, J.W., Dixit, S. and Joshi, A. (2023). A
592 practical entity linking system for tables in scientific literature. *arXiv*.
593 <https://doi.org/10.48550/arXiv.2306.10044>
594

595 Riaz, K. (2018). Improving search via named entity recognition in morphologically rich
596 languages: a case study in Urdu. PhD thesis, University of Minnesota. Available at:
597 <https://hdl.handle.net/11299/195403>

598 Ssemwogerere, R., Sajo, A.A., Mutwalibi, N. and Mzee, A.K. (2023). A survey about the
599 application of artificial intelligence in search engines. In: *Handbook of Research on AI*
600 *Methods and Applications in Computer Engineering*. Hershey, PA: IGI Global, pp. 161–178.
601 <https://doi.org/10.4018/978-1-6684-6937-8.ch008>
602

603 Orekhov, S., Godlevsky, M., Malyhon, H. and Goncharenko, T. (2023). A new method of
604 search engine optimization based on semantic kernel idea. In: *Advances in Artificial Systems*
605 *for Medicine and Education VI*. Cham: Springer, pp. 67–77. https://doi.org/10.1007/978-3-031-24468-1_7
606

607

608 Hoojon, R. (2025). Optimizing low-resource Khasi NER through transfer learning. In:
609 *Proceedings of the International Conference on North East India AI Summit (NEIAIS-2025)*.
610 Manipur, India
611

612 Banga, A., Ahuja, R. and Sharma, S.C. (2021). Performance analysis of regression algorithms
613 and feature selection techniques to predict PM2.5 in smart cities. *International Journal of*
614 *System Assurance Engineering and Management*, pp. 1–14
615

616 Naik, N. and Mohan, B.R. (2021). Novel stock crisis prediction technique: a study on Indian
617 stock market. *IEEE Access*, 9, pp. 86230–86242.
618 <https://doi.org/10.1109/ACCESS.2021.3088999>
619

620 Shilong, Z. *et al.* (2021). Machine learning model for sales forecasting by using XGBoost. In:
621 *IEEE International Conference on Consumer Electronics and Computer Engineering*. pp.
622 480–483
623

624 Parsa, A.B., Movahedi, A., Taghipour, H., Derrible, S. and Mohammadian, A.K. (2020).
625 Toward safer highways: application of XGBoost and SHAP for real-time accident detection.
626 *Accident Analysis & Prevention*, 136, p. 105405. <https://doi.org/10.1016/j.aap.2019.105405>
627

628 Tang, Q., Xia, G., Zhang, X. and Long, F. (2020). A customer churn prediction model based
629 on XGBoost and MLP. In: *International Conference on Computer Engineering and*
630 *Application*. pp. 608–612. <https://doi.org/10.1109/ICCEA50009.2020.00133>
631

632 Li, Y., Stasinakis, C. and Yeo, W.M. (2022). A hybrid XGBoost-MLP model for credit risk
633 assessment. *Forecasting*, 4(1), pp. 184–207
634

635 Liu, J., Zhang, S. and Fan, H. (2022). A two-stage hybrid credit risk prediction model. *Expert*
636 *Systems with Applications*, 195, p. 116624. <https://doi.org/10.1016/j.eswa.2022.116624>
637

638 Wang, K., Li, M., Cheng, J., Zhou, X. and Li, G. (2022). Research on personal credit risk
639 evaluation based on XGBoost. *Procedia Computer Science*, 199, pp. 1128–1135.
640 <https://doi.org/10.1016/j.procs.2022.01.143>
641

642 Ogunleye, A. and Wang, Q.G. (2019). XGBoost model for chronic kidney disease diagnosis.
643 *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 17(6), pp. 2131–
644 2140
645

646 Wang, C., Deng, C. and Wang, S. (2020). Imbalance-XGBoost: leveraging weighted and
647 focal losses. *Pattern Recognition Letters*, 136, pp. 190–197.
648 <https://doi.org/10.1016/j.patrec.2020.05.035>
649

650 Létinier, L. *et al.* (2021). Artificial intelligence for unstructured healthcare data. *Clinical*
651 *Pharmacology & Therapeutics*, 110(2), pp. 392–400. <https://doi.org/10.1002/cpt.2266>
652

653 Zuech, R., Hancock, J. and Khoshgoftaar, T.M. (2021). Detecting web attacks using random
654 undersampling. *Journal of Big Data*, 8(1), p. 75
655

656 Le, T.T.H., Oktian, Y.E. and Kim, H. (2022). XGBoost for imbalanced multiclass
657 classification-based intrusion detection systems. *Sustainability*, 14(14), p. 8707
658

659 Mishra, M., Patnaik, B., Bansal, R.C., Naidoo, R., Naik, B. and Nayak, J. (2021). DTCDWT-
660 SMOTE-XGBoost-based islanding detection. *IEEE Systems Journal*, 16(2), pp. 2008–2019
661

662 Mushava, J. and Murray, M. (2022). A novel XGBoost extension for credit scoring. *Expert*
663 *Systems with Applications*, 202, p. 117233. <https://doi.org/10.1016/j.eswa.2022.117233>
664

665 Pavlyshenko, B. (2018). Using stacking approaches for machine learning models. In: *IEEE*
666 *International Conference on Data Stream Mining and Processing*. pp. 255–258.
667 <https://doi.org/10.1109/DSMP.2018.8478522>
668

669 Rojarath, A. and Songpan, W. (2020). Probability-weighted voting ensemble learning.
670 *Journal of Advances in Information Technology*, 11(4), pp. 217–227.
671 <https://doi.org/10.12720/jait.11.4.217-227>
672

673 Won, M. and Martins, B. (2018). Ensemble named entity recognition. *Frontiers in Digital*
674 *Humanities*, 5. <https://doi.org/10.3389/fdigh.2018.00002>
675

676 Ullah, F., Gelbukh, A., Zamir, M., Riverón, E. and Sidorov, G. (2024). Enhancement of
677 named entity recognition in low-resource languages. *Computers*, 13(10), p. 258.
678 <https://doi.org/10.3390/computers13100258>
679

680 Jin, M., Choi, S.M. and Kim, G.W. (2025). COMCARE: a collaborative ensemble framework.
681 *Electronics*, 14, p. 328. <https://doi.org/10.3390/electronics14020328>
682

683 Munthe, I. (2024). Implementation of stacking technique combining machine learning and
684 deep learning algorithms. *Journal of Applied Data Sciences*, 5, pp. 2079–2091.
685 <https://doi.org/10.47738/jads.v5i4.421>
686

687 Ganaie, M.A., Hu, M., Malik, A.K., Tanveer, M. and Suganthan, P.N. (2022). Ensemble deep
688 learning: a review. *Engineering Applications of Artificial Intelligence*, 115, p. 105151.
689 <https://doi.org/10.1016/j.engappai.2022.105151>
690

691 Rahimi, A., Li, Y. and Cohn, T. (2019). Massively multilingual transfer for NER. *arXiv*.
692 Available at: <https://arxiv.org/abs/1902.00193>
693

694 Chen, Y., Zhong, R., Zha, S., Karypis, G. and He, H. (2022). Meta-learning via language
695 model in-context tuning. In: *Proceedings of the 60th Annual Meeting of the Association for*
696 *Computational Linguistics*. pp. 719–730. <https://doi.org/10.18653/v1/2022.acl-long.53>.

697 Lima, K.A. *et al.* (2023). A novel data and model centric artificial intelligence approach for
698 Bengali NER. *PLOS ONE*, 18(9), p. e0287818. <https://doi.org/10.1371/journal.pone.0287818>
699

700 Li, Y. (2025). Enhanced logistic regression using stacking algorithm. *Highlights in Science,*
701 *Engineering and Technology*, 136, pp. 1–11. <https://doi.org/10.54097/xmphgz15>
702

703 Bentéjac, C., Csörgő, A. and Martínez-Muñoz, G. (2021). A comparative analysis of gradient
704 boosting algorithms. *Artificial Intelligence Review*, 54(3), pp. 1937–1967.
705 <https://doi.org/10.1007/s10462-020-09896-5>
706

707 Bartz-Beielstein, T., Chandrasekaran, S. and Rehbach, F. (2023). Case study II: tuning of
708 gradient boosting (XGBoost). In: *Hyperparameter tuning for machine and deep learning with*
709 *R*. Singapore: Springer, pp. 221–234. https://doi.org/10.1007/978-981-19-5170-1_9
710

711 Biau, G. and Cadre, B. (2021). Optimization by gradient boosting. In: *Advances in*
712 *Contemporary Statistics and Econometrics*. Cham: Springer, pp. 23–44.
713 https://doi.org/10.1007/978-3-030-73249-3_2
714
715 Jurafsky, D. and Martin, J.H. (2020). *Speech and language processing*. 3rd ed. Stanford:
716 Stanford University