

# Fine-Tuning mBERT for Icelandic PoS Tagging and Integrating Multilingual NLP Tools

Valgarð Guðni Oddsson

## Abstract

This project explores the fine-tuning of the mBERT multilingual language model for Icelandic Part-of-Speech (PoS) tagging using the MIM-Gold dataset. The fine-tuned model achieves a high accuracy of 97.8%, effectively tagging Icelandic text with detailed grammatical information, including word class, gender, number, case, article, subcategory, mood, and person. To enhance usability, a Vue.js frontend and a FastAPI backend were developed, enabling users to input either Icelandic or English sentences. The system supports multilingual functionality by leveraging AWS Translator, allowing English sentences to be translated into Icelandic for tagging. Additionally, the backend integrates GreynirCorrect to provide grammar and spelling suggestions for Icelandic text and includes a custom-trained classifier for language detection. This comprehensive solution bridges the gap between advanced NLP tools and end-user applications, making Icelandic PoS tagging accessible to a broader audience while demonstrating the potential of fine-tuned multilingual models for low-resource languages.

## 1 Introduction

Natural Language Processing (NLP) is a field of artificial intelligence that focuses on enabling machines to understand and interpret human language. Part-of-Speech (PoS) tagging is one of the fundamental tasks in NLP, involving the assignment of grammatical tags to words in a sentence. In this report, we present the development and evaluation of a PoS tagging system for both English and Icelandic. The system also incorporates additional components such as spelling and grammar suggestions and multilingual translation capabilities.

The primary goal of this work was to create an NLP pipeline that can accurately predict PoS tags for a given sentence, suggest spelling and grammar

improvements, and translate sentences between English and Icelandic. The system is built using a variety of state-of-the-art tools and techniques, including a pre-trained language model, a language classifier, and integration with AWS translation services.

In the following sections, we describe the methodology used to build and fine-tune the system, the results of its performance on a test dataset, the challenges encountered during development, and potential avenues for future work and improvement.

## 2 Methods

### 2.1 Dataset and Preprocessing

The dataset used to finetune the pretrained model was a preprocessed version of the Part-of-Speech tagged Icelandic corpus, the Icelandic Gold Standard; MIM-Gold(Loftsson et al., 2018). The preprocessed version contains one sentence per line, with a "/" between tokens and tags. Each tag contains information about the word, including its word class, gender, case, number, person, and a few more outlined by the file MIM\_GOLD\_DESCRIPTION\_EN\_tagset.pdf found in the project. This was a good dataset as it contained a very large number of sentences, about 58,000, which ensured that the finetuned model had enough data to train and test on.

The dataset was further preprocessed to ensure it could be handled by the mBERT model. The preprocessing involved the following steps:

1. Each line in the dataset was parsed to separate words and their associated PoS tags. The dataset was transformed into a structured format where each sentence was represented as a dictionary containing two lists:
  - sentence: a list of words in the sentence.

080	<ul style="list-style-type: none"> <li>• <b>tags</b>: a list of corresponding PoS tags.</li> </ul>	Fine-tuning was conducted on a system with the following specifications:	123
081	2. The unique PoS tags in the dataset were extracted to create mappings for the fine-tuning process. This produced the dictionary <code>tag2id</code> , which maps each tag to a unique integer ID for training, and the dictionary <code>id2tag</code> , which reverses this mapping.	<ul style="list-style-type: none"> <li>• <b>CPU</b>: Intel i9-12900K</li> </ul>	124
082			
083		<ul style="list-style-type: none"> <li>• <b>GPU</b>: NVIDIA RTX 3090</li> </ul>	125
084			126
085		<ul style="list-style-type: none"> <li>• <b>RAM</b>: 64 GiB</li> </ul>	127
086			
087	3. Sentences were tokenized using the mBERT tokenizer, which splits text into sub-word tokens. To ensure that PoS tags aligned correctly with the tokenized output, each token was associated with the corresponding PoS tag or marked with a -100 for tokens such as padding or sub-word fragments that should be ignored during training. This alignment was achieved by comparing word indices generated during tokenization to the original sentence structure.	The model’s performance was evaluated using the following metrics:	128
088			129
089		<ul style="list-style-type: none"> <li>• <b>Accuracy</b>: The proportion of correctly predicted tags.</li> </ul>	130
090			131
091		<ul style="list-style-type: none"> <li>• <b>Macro Precision</b>: The average precision across all tag classes.</li> </ul>	132
092			133
093		<ul style="list-style-type: none"> <li>• <b>Macro Recall</b>: The average recall across all tag classes.</li> </ul>	134
094			135
095			
096			
097			
098	4. A custom PyTorch class, <code>PosDataset</code> , was implemented to manage the training data. This class leveraged the functions used to carry out the previous steps in order to tokenize each sentence and align its tags, and then returned tokenized inputs and the corresponding label sequence.	Fine-tuning encountered challenges related to CUDA memory limitations. Initially, the batch size for training and evaluation was set to 16. However, during metric computation after each training epoch, the GPU ran out of memory, which was believed to be caused by too large a portion of the dataset being loaded onto the GPU. To test this hypothesis, fine-tuning was done using only 10% of the dataset, which resulted in no CUDA memory issues. However, the model trained on only 10% of the dataset did not perform as desired, so a workaround had to be found.	136
099			137
100			138
101			139
102			140
103			141
104			142
105	These preprocessing steps ensured that the data was properly structured and aligned, enabling the mBERT model to handle them properly during training and evaluation.		143
106			144
107			145
108			146
109	<b>2.2 Model Fine-Tuning</b>	To address the issue:	147
110	The mBERT model was fine-tuned on the MIM-Gold dataset for Part-of-Speech (PoS) tagging in Icelandic. Fine-tuning was performed using the Hugging Face Trainer API (Wolf et al., 2020), with the following configuration:		148
111		1. <b>Reducing the batch size</b> : Training and evaluation batch sizes were reduced, but this alone did not resolve the memory overflow.	149
112			150
113		2. <b>Modifying compute_metrics</b> : The primary cause was identified as the <code>compute_metrics</code> function in the Trainer class, which consumed excessive GPU memory during evaluation. After removing this function, the full dataset could be used for training without memory issues.	151
114			152
115	<ul style="list-style-type: none"> <li>• <b>Evaluation strategy</b>: Per epoch</li> </ul>		153
116	<ul style="list-style-type: none"> <li>• <b>Learning rate</b>: <math>3 \times 10^{-5}</math></li> </ul>		154
117	<ul style="list-style-type: none"> <li>• <b>Batch size</b>: 8 for training, 2 for evaluation</li> </ul>		155
118	<ul style="list-style-type: none"> <li>• <b>Gradient accumulation steps</b>: 2</li> </ul>		156
119	<ul style="list-style-type: none"> <li>• <b>Number of epochs</b>: 3</li> </ul>		157
120	<ul style="list-style-type: none"> <li>• <b>Weight decay</b>: 0.01</li> </ul>	With the resolution of the CUDA memory issue, the training batch size was restored to 8 to maintain stability while optimizing throughput. The fine-tuning process could then be successfully completed, resulting in a model with a high degree of accuracy and robust tagging capabilities.	159
121	<ul style="list-style-type: none"> <li>• <b>Mixed-precision training (fp16)</b>: Enabled</li> </ul>		160
122	<ul style="list-style-type: none"> <li>• <b>Gradient Checkpointing</b>: Enabled</li> </ul>		161
			162
			163
			164

## 2.3 Additional Components

In addition to the fine-tuned mBERT model for Icelandic PoS tagging, the system includes three key components to enhance functionality:

- **Language Classifier:** A custom language classification module was developed to determine whether an input sentence was in Icelandic or English. The model was trained using a pipeline with the Universal Declaration of Human Rights (UDHR) dataset. The training process involved:

1. Generating synthetic sentences by randomly selecting between 3 to 18 words from the UDHR word list for each language, creating a diverse dataset with 5,000 sentences per language.
2. Splitting the dataset into training and testing subsets, with 90% of the data used for training and 10% for testing.
3. Fitting the pipeline to the training data using scikit-learn's `train_test_split` and `shuffle` methods to ensure a robust and randomized training process.

The trained classifier enabled automatic detection of the input language, allowing the system to determine whether translation was required before PoS tagging.

- **Grammar and Spelling Suggestions:** The system integrates the GreynirCorrect library (version 3.4.7) (Team, 2023) for generating grammar and spelling suggestions. This component leverages the `check_single` method to analyze Icelandic sentences and provide corrections. To address compatibility issues, the `islenska` module was downgraded to version 1.0, as recommended by the GreynirCorrect project maintainers. This ensured stable functionality despite outdated documentation and bugs in newer versions of the GreynirCorrect library.
- **Multilingual Translation:** To support English-speaking users, the system incorporates AWS's translation service through the `boto3` Python package (Services, Year). The `translate_text` method from the AWS translate client was used to perform real-time translation of English input sentences into Icelandic, as well as translation of each word in

Icelandic sentences to English. This integration extended the system's utility by allowing users to submit English sentences, which were automatically translated before processing for PoS tagging and grammar suggestions.

These additional components collectively enhanced the usability and versatility of the system, allowing it to handle multilingual input, provide grammatical feedback, and cater to both Icelandic and English speakers.

## 2.4 System Integration

The system was implemented with a FastAPI backend and a Vue.js frontend, enabling simple and user-friendly interaction with the PoS tagging model and supporting features.

- **Backend:** The backend is built using FastAPI and provides a single RESTful endpoint that processes both Icelandic and English sentences. Upon receiving a sentence, the backend performs the following tasks:

1. **Language Prediction:** A custom function uses the trained language classifier to determine whether the input sentence is in Icelandic or English.
2. **Translation:** If the input is identified as English, another function utilizes AWS's `translate_text` method to translate the sentence into Icelandic. After the English sentence has been translated (or if the original sentence was Icelandic), each word in the Icelandic sentence is also translated to provide per-word translations to accompany the PoS information.
3. **PoS Tagging:** A dedicated function generates PoS tag information for the Icelandic sentence using the fine-tuned mBERT model. The output includes grammatical attributes such as word class, gender, number, case, and more.
4. **Grammar and Spelling Suggestions:** The GreynirCorrect library is employed to provide suggestions for improving grammar and spelling in the Icelandic text.

The backend returns a structured JSON object containing:

260	– The PoS-tagged and per-word translated	• <b>Macro Recall:</b> 0.923	306
261	Icelandic sentence.		
262	– Grammar and spelling suggestions.	These metrics highlight the model’s ability to	307
263	– The predicted language of the input.	accurately tag Icelandic sentences across a diverse	308
264		set of grammatical features, including word class,	309
265	This modular architecture ensures flexibility	gender, number, case, mood, and others. While pre-	310
266	and maintainability while accommodating all	cision and recall are slightly lower than accuracy,	311
	system features.	they demonstrate the model’s capability to general-	312
267	• <b>Frontend:</b> The frontend was developed using	ize effectively across all tag categories, even those	313
268	Vue 3, Vite, and Pinia for state management,	with fewer examples.	314
269	with Bootstrap for styling. It serves as a user-		
270	friendly interface that facilitates interaction	<b>3.2 System Functionality</b>	315
271	with the backend. Key features include:	The system was tested end-to-end, with a focus on	316
		validating the integration of its various components.	317
272	– <b>Input Field:</b> A text field allows users	The key results are as follows:	318
273	to input either Icelandic or English sen-		
274	tences.	1. <b>Language Detection and Translation:</b>	319
275	– <b>Integration with Backend:</b> Upon sub-	• The language classifier reliably identified	320
276	mission, the frontend sends the input	Icelandic and English sentences, achiev-	321
277	sentence to the FastAPI backend via an	ing an accuracy of 96.8%.	322
278	HTTP request.	• English inputs were successfully trans-	323
279	– <b>Results Display:</b> The response is parsed	lated into Icelandic using the AWS trans-	324
280	and displayed, showing:	lator; however, only sentence-level trans-	325
281	* The PoS-tagged Icelandic sentence	lations were generally accurate. During	326
282	with detailed grammatical attributes	per-word translation from Icelandic to	327
283	and English translation.	English, the translation often returned	328
284	* Grammar and spelling suggestions	incorrect interpretations due to ambigu-	329
285	for the Icelandic text.	ity and lack of context, or sometimes no	330
286	* The predicted language of the input.	translation at all. This component was	331
287		the least performing out of the system’s	332
288	The frontend simplifies user interaction, mak-	features and has the highest potential for	333
289	ing it intuitive to explore the functionality of	improvement.	334
	the PoS tagger and the additional features.		
290		2. <b>Grammar and Spelling Suggestions:</b>	335
291	This modular integration of the backend and	• The GreynirCorrect library provided	336
292	frontend ensures the system is both flexible and	somewhat accurate and context-aware	337
293	user-friendly, providing comprehensive PoS tag-	grammar and spelling suggestions for	338
294	ging and additional language-related functionali-	Icelandic sentences. Spelling sugges-	339
	ties.	tions were particularly useful, often	340
295	<b>3 Results</b>	catching incorrect spellings and provid-	341
		ing correct alternatives. However, incor-	342
296	<b>3.1 Model Performance</b>	rect grammar was often missed and not	343
297	The fine-tuned mBERT model achieved excellent	corrected, leading to a similar level of	344
298	performance on the PoS tagging task for Icelandic	dissatisfaction as with the per-word trans-	345
299	text. Training was conducted on the full dataset,	lation component.	346
300	comprising 58,412 sentences, while testing was		
301	performed on a held-out 10% subset (5,842 sen-	3. <b>PoS Tagging Output:</b>	347
302	tences). The model achieved the following metrics	• As mentioned previously, the fine-tuned	348
303	during evaluation:	model performed very well and consis-	349
304	• <b>Accuracy:</b> 0.978	tently produced detailed annotations for	350
305	• <b>Macro Precision:</b> 0.936	Icelandic sentences, tagging each word	351
		with the correct attributes.	352

4. **Example Output:** To illustrate the system functionality, consider the following example:

- **Input (English):** "Basketball is one of my favorite sports and I am an active player."
- **Translated to Icelandic:** "Körfubolti er ein af mínum uppáhalds íþróttum og ég er virkur leikmaður."
- **Tagged Output:**
  - Körfubolti (noun, masculine, singular, nominative)
  - er (verb, indicative, active, 3rd person, singular, present)
  - ein (numeral, genitive)
  - af (adverb, governs case)
  - mínum (pronoun, possessive pronoun, feminine, plural, dative)
  - uppáhalds (noun, neutral, singular, genitive)
  - íþróttum (noun, feminine, plural, dative)
  - og (conjunction)
  - ég (pronoun, personal pronoun, 1st person, singular, nominative)
  - er (verb, indicative, active, 1st person, singular, present)
  - virkur (adjective, masculine, singular, nominative, strong declension, positive)
  - leikmaður (noun, masculine, singular, nominative)

These results demonstrate the system's overall effectiveness, with the model performing excellently on PoS tagging and the integrated components working well, though there are areas (such as translation and grammar correction) that could be further improved.

## 4 Discussion

### 4.1 Model Performance Comparison

The fine-tuned mBERT model for Icelandic PoS tagging achieved an accuracy of 97.8%, comparable to results reported for the same task. For instance, a GitHub repository (Jónsson, 2019) also claims 97.8% accuracy for a large PoS tagging model on the MIM-Gold dataset. This similarity in performance demonstrates that the current approach, despite leveraging a general-purpose model

like mBERT, is capable of achieving competitive results in specialized language tasks.

Furthermore, the high accuracy, coupled with the macro precision and macro recall of 93.6% and 92.3% respectively, reflects the model's ability to generalize well across diverse linguistic constructs. This suggests that the dataset used for training and evaluation provided sufficient representation of Icelandic grammatical features, enabling robust tagging across categories.

### 4.2 System Functionality

While the PoS tagging component performed exceptionally well, other system functionalities revealed notable limitations:

1. **Grammar and Spelling Suggestions:** The GreynirCorrect module, while effective at identifying some simple spelling errors, struggled with grammar-related issues. Many grammatical mistakes in test cases were left unaddressed, limiting the tool's utility for users seeking comprehensive language feedback. This lack of performance highlights the need for integrating or developing a more robust grammar correction component for future iterations.
2. **Translation:** The AWS translation service delivered good results for whole-sentence translations, reliably converting English input to Icelandic in most cases. However, word-by-word translation often failed due to the inherent ambiguity of isolated words and the lack of context. Common issues included incorrect or nonsensical translations, as well as instances where no translation was provided at all. This limitation significantly impacts the user experience, especially for English-speaking users relying on accurate word meanings for language learning.

### 4.3 Key Limitations

The primary limitations of the system are as follows:

- **Grammar Suggestion:** The system's inability to consistently identify grammar errors limits its value as a comprehensive language tool, particularly for language learners who might prioritize grammatical accuracy.
- **Word-by-Word Translation:** The failure to account for context during word-by-word



448	translation renders this feature unreliable and	could focus on integrating more advanced grammar	495
449	potentially confusing for users. This is a	correction models and improving context-aware	496
450	critical shortcoming, as accurate translations	word translation to make the system more useful to	497
451	for individual words are essential for aiding	English-speaking users and language learners.	498
452	English-speaking users in understanding Ice-		
453	landic grammar and vocabulary.		
454	<b>4.4 Future Directions</b>	<b>References</b>	499
455	Addressing these limitations offers several direc-	Haukur Páll Jónsson. 2019. <a href="#">Pos</a> . Accessed: 2024-11-	500
456	tions for future improvements:	10.	501
457	<ul style="list-style-type: none"> <li>• <b>Enhanced Grammar Correction:</b> Integrat-</li> </ul>	Hrafn Loftsson, Eiríkur Rögnvaldsson, Sigrún Helgadóttir, Jökull H. Yngvason, Kristján Friðbjörn Sigurðsson, Steinunn Valbjörnsdóttir, Brynhildur Stefánsdóttir, Jón Friðrik Daðason, and Starkaður Barkarson. 2018. <a href="#">Mim-gold: A dataset for icelandic part-of-speech tagging</a> . Accessed: 2024-11-10.	502
458	ing a more sophisticated grammar correction		503
459	model, possibly leveraging fine-tuned trans-		504
460	formers, could dramatically improve this as-		505
461	pect of the system.		506
462	<ul style="list-style-type: none"> <li>• <b>Context-Aware Translation:</b> Developing a</li> </ul>	Amazon Web Services. Year. <a href="#">Amazon translate</a> . Accessed: Date.	508
463	translation model that incorporates sentence-		509
464	level context when translating individual	Greynir Team. 2023. <a href="#">Greynircorrect: Spelling and grammar correction for icelandic</a> . Accessed: 2024-11-10.	510
465	words could significantly enhance this compo-		511
466	nent and provide much better usability for all		512
467	users.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, et al. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . Accessed: YYYY-MM-DD.	513
468	<b>4.5 Strengths and Contributions</b>		514
469	Despite its limitations, the system represents a step		515
470	forward in Icelandic language processing. The in-		516
471	tegration of accurate PoS tagging, grammar correc-		
472	tion, and translation functionalities into a unified		
473	platform is a valuable contribution to the field, and		
474	provides a strong foundation for iterative enhance-		
475	ments and future development.		
476	<b>5 Conclusion</b>		
477	In this report, we presented a comprehensive eval-		
478	uation of a system designed to perform Part-of-		
479	Speech (PoS) tagging, grammar and spelling sug-		
480	gestions, and translation for Icelandic text. The		
481	system leverages a fine-tuned mBERT model for		
482	PoS tagging, achieving an accuracy of 97.8%, com-		
483	parable to similar results on the MIM-Gold dataset.		
484	However, while the PoS tagging componet demon-		
485	strated high performance, the grammar and spelling		
486	suggestions, as well as word-by-word translation re-		
487	vealed some significant limitations. The grammar		
488	correction was under performing, often missing		
489	grammatical issues, while word-by-word transla-		
490	tion failed due to ambiguity and lack of context.		
491	Despite these challenges, the overall system		
492	shows promise, especially in the context of PoS		
493	tagging and full-sentence translation. The work		
494	lays a solid foundation, and future developments		