

Phone-ing it in: Towards Flexible, Multi-Modal Language Model Training using Phonetic Representations of Data

Anonymous ACL submission

Abstract

Multi-modal techniques offer significant untapped potential to unlock improved NLP technology for local languages. However, many advances in language model pre-training are focused on text, a fact that only increases systematic inequalities in the performance of NLP tasks across the world’s languages. In this work, we propose a multi-modal approach to train language models using whatever text and/or audio data might be available in a language. Initial experiments using Swahili and Kinyarwanda data suggest the viability of the approach for downstream Named Entity Recognition (NER) tasks, with models pre-trained on phone data showing an improvement of up to 6% F1-score above models that are trained from scratch.¹

1 Introduction

Pre-trained language models are increasingly applied in ways that are agnostic to targeted downstream tasks (Brown et al., 2020). This usage has led to a proliferation of large language models trained on enormous amounts of data. For example, the recent Megatron-Turing NLG 530B model was trained on the Pile, which includes 800GB+ of text (Gao et al., 2021), and other large language models utilize large portions of the 200TB+ common crawl data.² These large data sets include impressive amounts of text, but all languages are not represented equally (or at all) in that text. The reality is that only a negligible fraction of the 7000+ currently spoken languages (Eberhard et al., 2021) have sufficient text corpora to train state-of-the-art language models. This data scarcity results in systematic inequalities in the performance of NLP tasks across the world’s languages (Blasi et al., 2021).

Local language communities that are working to develop and preserve their languages are producing

diverse sets of data beyond pure text. The Bloom software project,³ for example, is being used by local language communities to create and translate "shell" or "template" books into many languages (426 languages at the time this paper is being written). However, Bloom allows users to do more than just translate text. Users are also recording audio tracks and sign language videos, which has resulted in 1600+ oral translations. Other examples showing the multi-modal nature of data in local languages include: (i) the creation of ChoCo: a multimodal corpus of the Choctaw language (Brixey and Artstein, 2021); (ii) SIL International’s 15+ year effort to document endangered Austronesian languages via text, audio, and video (Quakenbush, 2007); (iii) the grassroots Masakhane effort catalyzing the creation and use of diverse sets of African language data (v et al., 2020); and (iv) work with the Me’phaa language of western Mexico that is producing digital recordings (video and audio) along with vocabulary, grammar and texts (Marlett and Weathers, 2018). These diverse data sources are effectively unusable by traditional text-based NLP techniques. In the light of data scarcity on these languages, they offer significant untapped potential to unlock improved NLP technology, if text data can be leveraged along with audio, image and video data. Furthermore, flexible multi-modal technology such as this will make it easier to include diverse people and communities such as those described above within the NLP technology development process - audio-based technology reducing the need for literacy, for example.

In this paper, we propose a multi-modal approach to train both language models and models for downstream NLP tasks using whatever text and/or audio data might be available in a language (or even in a related language). Our method utilizes recent advances in phone recognition and text/grapheme-to-phone transliteration to convert

¹Preprocessing and training code will be released after publication.

²<https://commoncrawl.org/>

³<https://bloomlibrary.org/>

001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038

039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078

079 input audio and text into a common phonetic rep- 126
080 resentation (the IPA phone inventory). We then 127
081 pre-train character-based language models in this 128
082 phone-space. Finally, we fine-tune models for 129
083 downstream tasks by mapping text-based training 130
084 data into the phonetic representation. Thus, in ad- 131
085 dition to flexibility in pre-training, our method pro- 132
086 vides a way to reuse labeled text data for common 133
087 NLP tasks, like Named Entity Recognition or Sen- 134
088 timent Analysis, in the context of audio inputs.

089 We demonstrate our phonetic approach by train- 135
090 ing Named Entity Recognition (NER) models for 136
091 Swahili [swh]⁴ using various combinations of 137
092 Swahili text data, Swahili audio data, Kinyarwanda 138
093 [kin] text data, and Kinyarwanda audio data. These 139
094 two languages both originate from from the same 140
095 language family, Bantu, and are spoken by millions 141
096 of people in Eastern Africa, but are both considered 142
097 low-resource languages. Kinyarwanda in particu- 143
098 lar, though spoken by nearly 10 million people, 144
099 has very little text data available in that language, 145
100 with fewer than 3,000 articles on the Kinyarwanda- 146
101 language Wikipedia, and Swahili comparatively 147
102 ahead but still poorly resourced at approximately 148
103 68,000 articles, far less than many European lan- 149
104 guages.⁵ On the other hand, Kinyarwanda is 150
105 uniquely placed as a language to leverage speech- 151
106 based technologies, due to well-organized efforts⁶ 152
107 to collect voice data for that language. It is in fact 153
108 one of the largest subsets available on the Com- 154
109 mon Voice Dataset (Ardila et al., 2019), with 1,183 155
110 hours of voice clips collected and validated. Choos- 156
111 ing these two languages allowed us to test the use 157
112 of the technique on legitimately low-resourced lan- 158
113 guages that could benefit from improved NLP tech- 159
114 nology, and which as part of the same family of 160
115 languages might be similar enough in vocabulary, 161
116 grammar, sound systems and so on, to benefit from 162
117 cross-lingual training. 163

118 We find that simple NER models, which just 164
119 look for the presence or absence of entities, can 165
120 be trained on small amounts of data (around 2000 166
121 samples) in the phonetic representation. Models 167
122 trained for complicated NER tasks in the phonetic 168
123 representation, which look for entities and their 169
124 locations within a sequence, are improved (by up 170
125 to 6+% in F1 score) through pre-training a phonetic 171

⁴Language codes formatted according to ISO 639-3 stan-
dard: <https://iso639-3.sil.org/>

⁵https://meta.wikimedia.org/wiki/List_of_Wikipedias

⁶<https://foundation.mozilla.org/en/blog/how-rwanda-making-voice-tech-more-open/>

126 language model using a combination of text and 127
128 audio data. We see this improvement when fine- 129
130 tuning either a Swahili or Kinyarwanda language 131
132 model for downstream Swahili tasks, which implies 133
134 that one could make use of text and audio data in 134

2 Related Work 135

136 There have been a series of attempts to utilize pho- 137
138 netic representations of language to improve or 138
139 extend automatic speech recognition (ASR) mod- 139
140 els. Some of these jointly model text and audio 140
141 data using sequences of phonemes combined with 141
142 sequences of text characters. Sundararaman et al. 142
143 (2021), for example, uses a joint transformer archi- 143
144 tecture that encodes sequences of phonemes and 144
145 sequences of text simultaneously. However, this 145
146 joint model is utilized to learn representations that 146
147 are more robust to transcription errors. The archi- 147
148 tecture still requires text inputs (from ASR tran- 148
149 scriptions) and generates outputs in both text and 149
150 phoneme representations. In contrast, our approach 150
151 allows for text input, audio input, or text plus audio 151
152 input to language models. 151

152 Baevski et al. (2021) transforms unlabeled text 152
153 (i.e., not aligned with corresponding audio files) 153
154 into phonemes in a scheme to train speech recogni- 154
155 tion models without any labeled data. This scheme 155
156 involves a generator model trained jointly with a 156
157 discriminator model. The generator model converts 157
158 audio, segmented into phonetic units into predicted 158
159 phonemes, and the discriminator model attempts 159
160 to discriminate between these predicted phonemes 160
161 and the phonemes transliterated from unlabeled 161
162 text. Although both text and audio are utilized in 162
163 this work, they are not input to the same model 163
164 and the primary output of the training scheme is a 164
165 model that creates good phonetic speech represen- 165
166 tations from input audio. 166

167 Outside of speech recognition focused 167
168 work, Shen et al. (2020) (and other researchers 168
169 cited therein) attempt to "fuse" audio and text at the 169
170 word level for emotion recognition. They introduce 170
171 another architecture that internally represents both 171
172 audio and text. However, the so-called WISE 172
173 framework relies on speech recognition to generate 173
174 the text corresponding to audio frames in real-time. 174
175 The current work explicitly avoids reliance 175

on speech recognition. The 2021 Multimodal Sentiment Analysis (MuSe) challenge continues this vein of research integrating audio, video, text, and physiology data in an emotion recognition task (Stappen et al., 2021). Contributions to this challenge, such as Vlasenko et al. (2021), introduce a variety of ways to "fuse" audio and text inputs. However, these contributions are squarely focused on emotion/sentiment analysis and do not propose methods for flexible, phonetic language models.

Lakhotia et al. (2021) introduced functionality for "textless" NLP. They explored possibility of creating a dialogue system from only audio inputs (i.e., without text). As part of this system, language models are directly trained on audio units without any text. This advances the state-of-the-art with regard to self-supervised speech methods, but it does not provide the flexibility in audio and/or text language modeling introduced here.

3 Methodology

Our approach is inspired by the fact that many languages are primarily oral, with writing systems that represent spoken sounds. We convert both text and audio into single common representation of sounds, or "phones," represented using the International Phonetic Alphabet, or IPA. Then, we perform both language model pre-training and the training of models for downstream tasks in this phonetic representation. Well-tested architectures, such as BERT-style transformer models (Vaswani et al., 2017), are thus flexibly extended to either speech or audio data.

Regarding the conversion process of text and audio data, we leverage recent advances to transliterate this data into corresponding sounds represented by IPA phonetic symbols. This transliteration is possible for speech/audio data using tools such as the Allosaurus universal phone recognizer, which can be applied without additional training to any language (Li et al., 2020), though it can benefit from fine-tuning (Siminyu et al., 2021). To convert text data to phonemes we can use tools such as the Epitran grapheme-to-phoneme converter (Mortensen et al., 2018), which is specifically designed to provide precise phonetic transliterations in low-resource scenarios.

Fig. 1 shows how downstream models for certain NLP tasks, like Named Entity Recognition (NER), are performed in the phonetic representation. La-

beled data sets for NLP tasks need to be mapped or encoded into the phonetic representation to train downstream models. However, once this mapping is accomplished, models trained in the phonetic representation can perform tasks with audio input that are typically restricted to processing text input.

3.1 Phonetic Language Modeling

One complication arising from direct speech-to-phone transcription is the loss of word boundaries in the transcription. This is expected, as natural speech does not typically include long pauses between word utterances. This does, however, result in merging text data sets containing clear word boundaries with speech data sets containing no clear word boundaries.

Borrowing from techniques used on languages that do not indicate word boundaries by the use of whitespace, we address the problem by removing all whitespace from our data sets after phone transliteration. We train a *character-based* language models over the resulting data. Character-based models such as CharFormer (Tay et al., 2021) or ByT5 (Xue et al., 2021) have shown promise in recent years for language modeling, even if this approach is known to have some trade offs related to shorter context windows.

3.2 Potential Information Losses

The transliteration of text and audio data into phonetic representations presents several other challenges related to potential loss of information or injection of noise:

1. *Loss of intonation or "suprasegmental" effects, extending across segments:* In some languages, meaning may be encoded through intonation or *across* sounds. Particularly for tonal languages such as Mandarin Chinese [cmn], this loss can represent a significant informational loss particularly for homophones with different tones, as seen in (Amrhein and Sennrich, 2020). While IPA symbols can represent these intricacies, it adds complexity
2. *Phone/phoneme differences:* As noted in (Li et al., 2020), speech sounds which are physically different (different *phones*), may be *perceived* as the same (one *phoneme*) by speakers of one language, but these same sounds could perhaps be distinguished by speakers of another language. For example, the Spanish

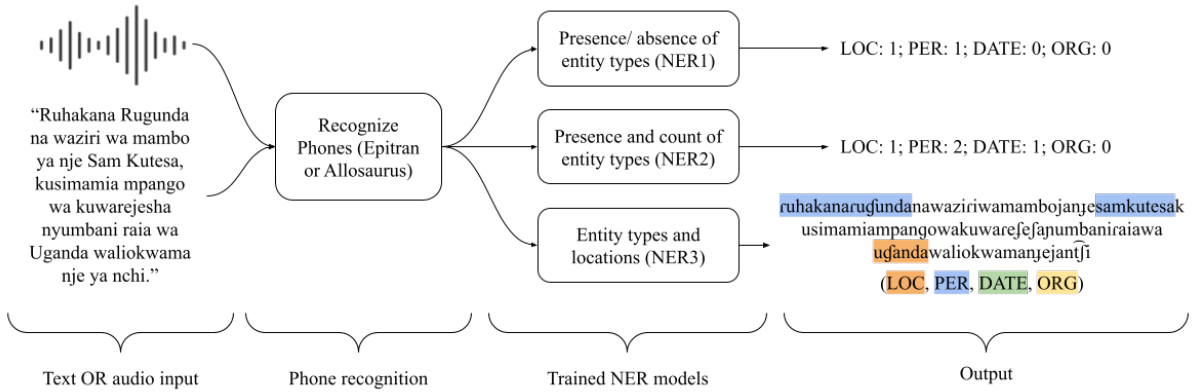


Figure 1: Our approach: input from either modality can be converted by phone recognition, e.g. Epitran for text, Allosaurus for speech. Then we test on several downstream tasks which we designate NER1, NER2, NER3.

words *anos*, and *años* contain phones (n and ñ) which sound "the same" to English speakers, but are semantically different to Spanish speakers. In other words, in English, both *phones* map to the same *phoneme* perceptually. As the Allosaurus phone recognizer recognizes the actual phones/sounds, not their perceived phonemes, it would transcribe these two phones to different representations even for English speech. This can be mitigated to an extent by customizing the output of Allosaurus on a per-language basis, see Sec. 4.3.

3. *Simple errors in phone recognition*: As noted in (Siminyu et al., 2021), even the best-trained Allosaurus models, fine-tuned on language-specific data, have a non-trivial Phone Error Rate (PER).

An important question, therefore, is whether these added sources of noise/information losses are outweighed by the potential benefits in terms of flexibility. Does working in a phonetic representation cause a prohibitive amount of information loss? We constructed our experiments and data sets in order to answer this question.

4 Experiments

In order to evaluate the quality of learned phonetic representations, we transliterate several text and audio data sets in the Swahili [swh] language. We pre-train phonetic language models on various combinations of these data sets and evaluate downstream performance on NER tasks. See Fig. 2 for a detailed overview of these various combinations.

We refer to these combinations as denoted by downstream tasks (SNER for Swahili NER), and

pre-training language ((**K** for **K**inyarwanda, **S** for **S**wahili) as well as data modality (**T** for text, **A** for audio). By way of example, the SNER+ST2 model results from pre-training using 2 swh text datasets (ST2) and fine-tuning on the swh NER (SNER) task, whereas the SNER+SAT model results from pre-training using swh audio and text data (SAT).

Kinyarwanda [kin] data is used in our experiments as a language related to the target language (swh) with existing text and audio resources that, in some ways, surpasses those available in the target language. Thus, we pre-train some models on kin data while fine-tuning for the downstream NER task using swh data.

Three different formulations of the NER task, from more simple (NER1) to more complicated/granular (NER3), are used (see Fig. 2) to help determine the applicability of our methods to less challenging (NER1) to more challenging (NER3) tasks. The NER1 task tries to determine the presence or absence of certain kinds of entities within an input. For our task we use PER, ORG, DATE, and LOC entities. The NER2 task additionally requires models to predict the correct numbers of these entities within an input. Finally, the NER3 task requires models to determine entities at the correct locations with an input sequence of phones.

For all of these tasks, we first convert text data to phones using Epitran and audio data to phones using Allosaurus. Then, we pre-train on various combinations of data, before fine-tuning on NER.

4.1 Data Sources

For swh pre-training data we use: (i) the "Language Modeling Data for Swahili" dataset (Shikali and Refuoe, 2019) hosted on Hugging Face (which

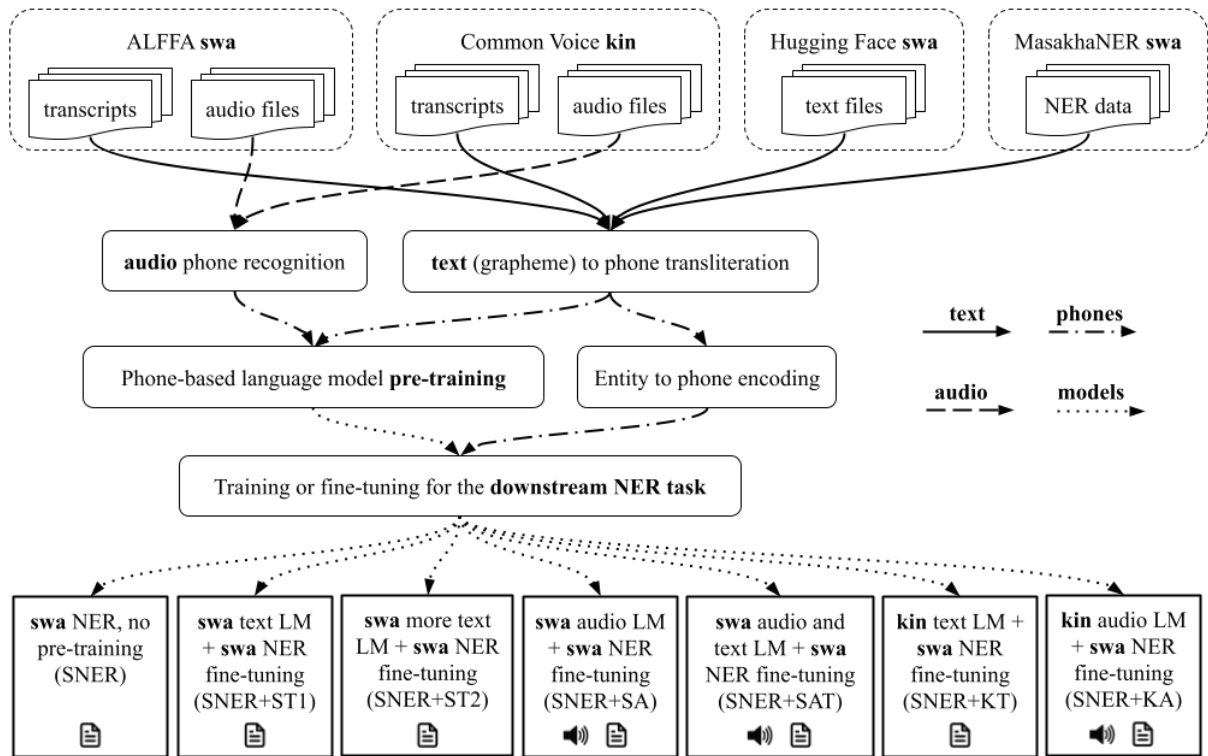


Figure 2: Training scenarios: we pre-train on various combinations of phonemized datasets, evaluating on the downstream NER task. SNER-ST denotes "Swahili Text (ST) pre-training, Swahili NER (SNER) fine-tuning", SNER-SAT denotes Swahili NER with Swahili Audio and Text (SAT) pre-training, SNER-KA uses Kinyarwanda Audio (KA), etc.

we refer to as the "HF Swahili" data set); and (ii) the ALFFA speech dataset (Gelas et al., 2012). For ALFFA data we process both the audio files (using Allosaurus) and the original "gold" text transcriptions (using Epitran).

For Kinyarwanda pre-training data, we use the Common Voice (CV) Kinyarwanda 6.1 subset (Ardila et al., 2019). Again, we utilize both the audio files and transcriptions. Due to the large size of the CV 6.1 Kinyarwanda subset, we processed only about 80% of the audio files.

For fine-tuning the downstream NER task, we use the MasakhaNER data set (Adelani et al., 2021). As with other text-based data sets, we transform the NER sample with Epitran to map the samples into the phonetic representation.

4.2 Entity to Phone Encoding

For the downstream NER tasks we map or encode the NER annotations into the phonetic representation. We thus edited the labels (PER, ORG, DATE, and LOC) to convert them from word-level labels to phone-level labels as shown in Fig. 3. Unlike (Kuru et al., 2016), we leave in the B- and I- prefixes.

Our fork of the MasakhaNER data set, which im-

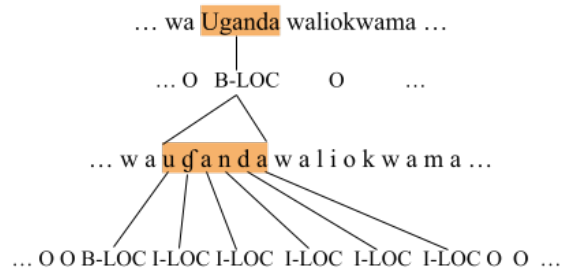


Figure 3: Adaptation of word-level NER annotations to character-level annotations.

plements our phonetic representations of the labels, is published on Github.⁷

4.3 Phone Inventory Considerations

As mentioned already, we use Allosaurus for phone recognition with audio inputs. In order to ensure consistency with Epitran, we took advantage of Allosaurus's inventory customization feature, giving it the phone inventories specified by the same language in Epitran. The inventory used throughout this work (for swa) is the swa-Latn inventory

⁷<https://anonymous.4open.science/r/masakhane-ner-5CC1/README.md>

377 from Epitran.⁸ When this inventory is supplied as
 378 input, Allosaurus will only output symbols from
 379 the inventory. We followed similar practice when
 380 transliterating Kinyarwanda data.

381 We compare the output of Epitran and Al-
 382 losaurus on the ALFFA dataset. Following
 383 the practice of (Li et al., 2020), we used the
 384 `editdistance`⁹ library to calculate the Phone
 385 Error Rate (PER). Having no ground truth phone
 386 annotations, we instead take Epitran’s outputs as
 387 "ground truth" for the sake of comparison. The
 388 mean PER between the outputs is 23.7%. This
 389 result is consistent with Siminyu et al. (2021),
 390 which finds PERs as high as 72.8% when testing
 391 on on the Bukusu (bxx), Saamia (lsm) and East
 392 Tusom languages (an endangered subdialect of the
 393 Tungkulic language family). However, by train-
 394 ing the phone recognizer on even minimal amounts
 395 of data in these languages, PERs were improved
 396 significantly.

397 A spreadsheet with detailed results for 10k sam-
 398 ples from ALFFA can be found online.¹⁰

399 4.4 Model Architecture and Training

400 All models use the SHIBA implementation of CA-
 401 NINE (Tanner and Hagiwara, 2021). SHIBA was
 402 designed for use on the Japanese [jpn] language,
 403 which does not include spaces between its charac-
 404 ters (similar to our phonetic representations without
 405 word boundaries). We used the default hyperpa-
 406 rameter settings for SHIBA pre-training and fine-
 407 tuning, because we are primarily concerned with
 408 the relative impact of various combinations of pre-
 409 training data on the downstream NER tasks. We
 410 use the Hugging Face library (Wolf et al., 2020) to
 411 train all models.

412 Because of the small size of the NER data
 413 set used during fine-tuning, we enabled Hugging
 414 Face’s early stopping callback for all downstream
 415 training runs. We stopped these runs if they did not
 416 improve training loss after 20 evaluations. Nonethe-
 417 less, we found after a number of trials that the
 418 models quickly overfit using this setting. We also
 419 experimented with modifying this on several tri-
 420 als to stop based on the evaluation loss instead,
 421 but this change did not significantly influence the
 422 evaluation results.

423 Following the example of Adelanı et al. (2021),
 424 we do not run downstream model trainings once,

⁸<https://bit.ly/30f8YCI>

⁹<https://github.com/roy-ht/editdistance>

¹⁰<https://bit.ly/3F0is3t>

Model	F1 NER1	F1 NER2
SNER	0.829	0.753
SNER+ST1	0.827	0.770
SNER+ST2	0.824	0.747
SNER+SA	0.817	0.751
SNER+SAT	0.818	0.763
SNER+KT	0.823	0.771
SNER+KA	0.846	0.763

Table 1: Mean results for presence/absence of entity types (NER1) and presence and *count* of entity types (NER2). Average of at least three trials per experiment, calculated with the scikit-learn library. (Pedregosa et al., 2011)

but multiple times. We also pre-trained each pho-
 netic language model multiple times with different
 random seeds. We report averages of these multiple
 trials in the following.

5 Results and Discussion

Table 1 presents the F1 scores for our training sce-
 narios in the downstream NER1 and NER2 tasks.
 The models that utilize pre-training on the kin
 audio and text data give the best results. However,
 pre-training does not appear to dramatically influ-
 ence the level. F1 scores in the range of 74-85%
 suggests the minimum viability of these phonetic
 models for simple NLP tasks.

Table 2 presents the F1 scores for our various
 training scenarios in the downstream NER3 task,
 which should be the most challenging for our pho-
 netic models. The influence of pre-training is more
 noticeable for this task. Further, the models pre-
 trained on the kin audio and text data have the
 best performance. This is likely due to the fact that
 the kin data is both large and higher quality (in
 terms of sound quality) as compared to the ALFFA
 Swahili data. This benefit of this data size and
 quality appears to outweigh any degradation due to
 the pre-training occurring in a different (although
 related) language.

The importance (or relative impact) of pre-
 training phonetic language models increases with
 the complexity of the NER task. Fig. 4 shows
 the maximum percentage improvement due to pre-
 training for each of our NER tasks. This suggests
 that simple NLP tasks with a small number of out-
 put classes are much easier to port to phonetic rep-
 resentations, even without pre-training, while more
 complicated NLP tasks may require a more sig-
 nificant amount of text and/or audio data for pre-

Model	F1	F1 (strict)
SNER	0.357	0.161
SNER+ST1	0.401	0.213
SNER+ST2	0.394	0.166
SNER+SA	0.363	0.163
SNER+SAT	0.405	0.203
SNER+KT	0.408	0.217
SNER+KA	0.397	0.197

Table 2: Prediction of entity types and precise locations (NER3) Average of at least three trials per experiment, scores calculated with seqeval library. (Nakayama, 2018)

training. We expect this trend to carry through to tasks like sentiment analysis, which could be formulated as a simple classification task with NEG, NEU, and POS sentiment labels or a more complicated aspect based sentiment analysis task.

6 Conclusions and Further Work

The proposed method for multi-modal training using phonetic representations of data has minimum viability for simple NER tasks. For more complicated NER tasks, pre-training phonetic language models boosts downstream model performance by up to 6% in F1 scores. This pre-training can be performed in the target language or in a related language using text and/or audio data. Thus, the method provides flexibility in the data needed to train language models, while also allowing for audio and/or text inputs to models trained on downstream NLP tasks.

We anticipate exploring various extensions to and validations of this method in the future. Specifically, we would like to explore methods that might mitigate performance degradation due to a lack of word boundaries in our method. Subword tokenization techniques, such as Byte-Pair Encodings (BPE) (Sennrich et al., 2016; Gage, 1994), or character-based word segmentation techniques might help in detecting and exploiting repeating patterns within the phonetic representation.

We would also like to validate our methods on a variety of other data sets and tasks. We selected the MasakhaNER dataset for evaluation because we specifically wished to evaluate results on actual low-resource languages supported by both Allosaurus and Epitran. While there are still, we argue, detectable improvements in downstream results with our method, further work would benefit

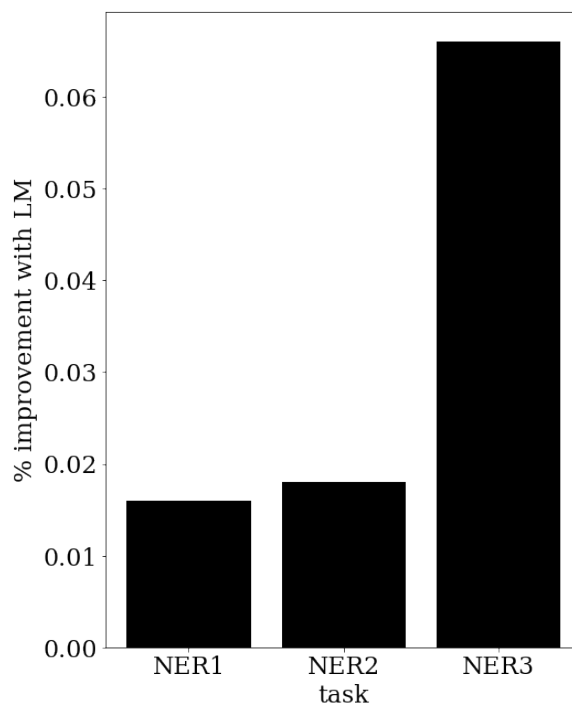


Figure 4: The max percentage improvement with fine-tuning for each kind of NER task that was explored. presence/absence of entity types (NER1), presence and count of entity types (NER2), and prediction of entity types and precise locations (NER3)

from additional evaluations on other data sets or tasks. In particular, the Swahili News Classification corpus (David, 2020) corpus may provide a useful evaluation.

Finally, it has been shown by Siminyu et al. (2021) that it is possible to improve phone recognition with even small amounts (approximately 100 sentences) of annotation. It may be possible to improve phonetic language modeling results by performing this fine-tuning in the target language.

References

- D. Adelani, Jade Z. Abbott, Graham Neubig, Daniel D’souza, Julia Kreutzer, Constantine Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, Stephen Mayhew, Israel Abebe Azime, Shamsuddeen Hassan Muhammad, Chris C. Emezue, Joyce Nakatumba-Nabende, Perez Ogayo, Anuoluwapo Aremu, Catherine Gitau, Derguene Mbaye, J. Alabi, Seid Muhie Yimam, Tajuddeen R. Gwadabe, Ignatius Ezeani, Rubungo Andre Niyongabo, Jonathan Mukiiibi, Verrah A Otiende, Iroro Orife, Davis David, Samba Ngom, Tosin P. Adewumi, Paul Rayson, Mofetoluwa Adeyemi, Gerald Muriuki, Emmanuel Anebi, Chiamaka Ijeoma Chukwuneke, Nkiruka Bridget Odu, Eric Peter Wairagala, S. Ajiboye Oyerinde, Clemencia Siro, Tobius Saul Bateesa,

523	Temilola Oloyede, Yvonne Wambui, Victor Akinode,	Fagbohunbe, Solomon Oluwole Akinola, Sham-	579
524	Deborah Nabagereka, Maurice Katusiime, Ayodele	suddee Hassan Muhammad, Salomon Kabongo, Sa-	580
525	Awokoya, Mouhamadane Mboup, Dibora Gebrey-	lomey Osei, and others. 2020. Participatory research	581
526	ohannes, Henok Tilaye, Kelechi Nwaike, Degaga	for low-resourced machine translation: A case study	582
527	Wolde, Abdoulaye N Faye, Blessing Sibanda, Ore-	in african languages. <i>Findings of EMNLP</i> .	583
528	vaoghene Ahia, Bonaventure F. P. Dossou, Kelechi		
529	Ogueji, Thierno Ibrahima Diop, Abdoulaye Diallo,	Philip Gage. 1994. A new algorithm for data compres-	584
530	Adeiwale Akinfaderin, Tendai Munyaradzi Maren-	sion. <i>The C Users Journal archive</i> , 12:23–38.	585
531	gereke, and Salomey Osei. 2021. MasakhaNER:		
532	Named Entity Recognition for African Languages.	Leo Gao, Stella Rose Biderman, Sid Black, Laurence	586
533	<i>Transactions of the Association for Computational</i>	Golding, Travis Hoppe, Charles Foster, Jason Phang,	587
534	<i>Linguistics</i> , 9:1116–1131.	Horace He, Anish Thite, Noa Nabeshima, Shawn	588
		Presser, and Connor Leahy. 2021. The pile: An	589
535	Chantal Amrhein and Rico Sennrich. 2020. On Roman-	800gb dataset of diverse text for language modeling.	590
536	ization for model transfer between scripts in neural	<i>ArXiv</i> , abs/2101.00027.	591
537	machine translation . In <i>Findings of the Association</i>		
538	<i>for Computational Linguistics: EMNLP 2020</i> , pages	Hadrien Gelas, Laurent Besacier, and Francois Pelle-	592
539	2461–2469, Online. Association for Computational	grino. 2012. Developments of Swahili resources for	593
540	Linguistics.	an automatic speech recognition system . In <i>SLTU</i>	594
		<i>- Workshop on Spoken Language Technologies for</i>	595
541	Rosana Ardila, Megan Branson, Kelly Davis, Michael	<i>Under-Resourced Languages</i> , Cape-Town, Afrique	596
542	Henretty, Michael Kohler, Josh Meyer, Reuben	Du Sud.	597
543	Morais, Lindsay Saunders, Francis M. Tyers, and		
544	Gregor Weber. 2019. Common voice: A massively-	Onur Kuru, Ozan Arkan Can, and Deniz Yuret. 2016.	598
545	multilingual speech corpus . <i>CoRR</i> , abs/1912.06670.	CharNER: Character-level named entity recognition .	599
		In <i>Proceedings of COLING 2016, the 26th Inter-</i>	600
546	Alexei Baevski, Wei-Ning Hsu, Alexis Conneau, and	<i>national Conference on Computational Linguistics:</i>	601
547	Michael Auli. 2021. Unsupervised speech recogni-	<i>Technical Papers</i> , pages 911–921, Osaka, Japan. The	602
548	tion. <i>ArXiv</i> , abs/2105.11084.	COLING 2016 Organizing Committee.	603
549	Damián E. Blasi, Antonios Anastasopoulos, and Gra-	Kushal Lakhotia, Evgeny Kharitonov, Wei-Ning	604
550	ham Neubig. 2021. Systematic inequalities in lan-	Hsu, Yossi Adi, Adam Polyak, Benjamin Bolte,	605
551	guage technology performance across the world’s	Tu Nguyen, Jade Copet, Alexei Baevski, Adel Ben	606
552	languages. <i>ArXiv</i> , abs/2110.06733.	Mohamed, and Emmanuel Dupoux. 2021. Gener-	607
		ative spoken language modeling from raw audio.	608
553	Jacqueline Brixey and Ron Artstein. 2021. Choco: a	<i>ArXiv</i> , abs/2102.01192.	609
554	multimodal corpus of the choctaw language. <i>Lan-</i>		
555	<i>guage Resources and Evaluation</i> , 55:241–257.	Xinjian Li, Siddharth Dalmia, Juncheng Li, Matthew	610
		Lee, Patrick Littell, Jiali Yao, Antonios Anastasopou-	611
556	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	los, David R Mortensen, Graham Neubig, Alan W	612
557	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	Black, and Metze Florian. 2020. Universal phone	613
558	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	recognition with a multilingual allophone system.	614
559	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	In <i>ICASSP 2020-2020 IEEE International Confer-</i>	615
560	Gretchen Krueger, Tom Henighan, Rewon Child,	<i>ence on Acoustics, Speech and Signal Processing</i>	616
561	Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens	<i>(ICASSP)</i> , pages 8249–8253. IEEE.	617
562	Winter, Chris Hesse, Mark Chen, Eric Sigler, Mate-		
563	usz Litwin, Scott Gray, Benjamin Chess, Jack	Stephen A. Marlett and Mark L. Weathers. 2018. The	618
564	Clark, Christopher Berner, Sam McCandlish, Alec	sounds of me’phaa (tlapanec): A new assessment.	619
565	Radford, Ilya Sutskever, and Dario Amodei. 2020.	<i>SIL-Mexico Electronic Working Papers</i> , 25.	620
566	Language models are few-shot learners . In <i>Ad-</i>		
567	<i>vances in Neural Information Processing Systems</i> ,	David R. Mortensen, Siddharth Dalmia, and Patrick	621
568	volume 33, pages 1877–1901. Curran Associates,	Littell. 2018. Epitran: Precision G2P for many lan-	622
569	Inc.	guages. In <i>Proceedings of the Eleventh International</i>	623
		<i>Conference on Language Resources and Evaluation</i>	624
570	Davis David. 2020. Swahili : News classification	<i>(LREC 2018)</i> , Paris, France. European Language Re-	625
571	dataset . The news version contains both train and	sources Association (ELRA).	626
572	test sets.		
		Hiroki Nakayama. 2018. sequeval: A python framework	627
573	David M. Eberhard, Gary F. Simons, and Charles D.	for sequence labeling evaluation . Software available	628
574	Fennig. 2021. <i>Ethnologue: Languages of the World</i> ,	from https://github.com/chakki-works/sequeval .	629
575	twenty-fourth edition. SIL International, Dallas,		
576	Texas.	Fabian Pedregosa, Gaël Varoquaux, Alexandre Gram-	630
		fort, Vincent Michel, Bertrand Thirion, Olivier Grisel,	631
577	∀, Wilhelmina Nekoto, Vukosi Marivate, Tshinondiwa	Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vin-	632
578	Matsila, Timi Fasubaa, Tajudeen Kolawole, Taiwo	cent Dubourg, et al. 2011. Scikit-learn: Machine	633
		learning in python. <i>Journal of machine learning re-</i>	634
		<i>search</i> , 12(Oct):2825–2830.	635

636	J. S. Quakenbush. 2007. Chapter 4. sil international and endangered austronesian languages. In <i>LD&C Special Publication No. 1: Documenting and Revitalizing Austronesian Languages</i> .	In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	691 692 693 694
640	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2021. Byt5: Towards a token-free future with pre-trained byte-to-byte models . <i>CoRR</i> , abs/2105.13626.	695 696 697 698 699
647	Guanghu Shen, Riwei Lai, Rui Chen, Yu Zhang, Kejia Zhang, Qilong Han, and Hongtao Song. 2020. Wise: Word-level interaction-based multimodal fusion for speech emotion recognition. In <i>INTERSPEECH</i> .		
651	Shivachi Casper Shikali and Mokhosi Refuoe. 2019. Language modeling data for Swahili . Type: dataset.		
653	Kathleen Siminyu, Xinjian Li, Antonios Anastasopoulos, David Mortensen, Michael R. Marlo, and Graham Neubig. 2021. Phoneme recognition through fine tuning of phonetic representations: a case study on luhya language varieties .		
658	Lukas Stappen, Alice Baird, Lea Schumann, and Björn W. Schuller. 2021. The multimodal sentiment analysis in car reviews (muse-car) dataset: Collection, insights and improvements. <i>ArXiv</i> , abs/2101.06053.		
662	Mukuntha Narayanan Sundararaman, Ayush Kumar, and Jithendra Vepa. 2021. Phoneme-bert: Joint language modelling of phoneme sequence and asr transcript. <i>ArXiv</i> , abs/2102.00804.		
666	Joshua Tanner and Masato Hagiwara. 2021. SHIBA: Japanese CANINE model . Publication Title: GitHub repository.		
669	Yi Tay, Vinh Tran, Sebastian Ruder, Jai Gupta, Hyung Won Chung, Dara Bahri, Zhen Qin, Simon Baumgartner, Cong Yu, and Donald Metzler. 2021. Charformer: Fast character transformers via gradient-based subword tokenization. <i>ArXiv</i> , abs/2106.12672.		
674	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . <i>CoRR</i> , abs/1706.03762.		
678	Bogdan Vlasenko, RaviShankar Prasad, and Mathew Magimai.-Doss. 2021. Fusion of acoustic and linguistic information using supervised autoencoder for improved emotion recognition. <i>Proceedings of the 2nd on Multimodal Sentiment Analysis Challenge</i> .		
683	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing .		