# Multilingual Transformer Language Model for Speech Recognition in Low-resource Languages

*Li Miao, Jian Wu, Piyush Behre, Shuangyu Chang, Sarangarajan Parthasarathy*

Microsoft, USA

{limia|jianwu|piyush|shchang|sarangp}@microsoft.com

## Abstract

It is challenging to train and deploy Transformer LMs for hybrid speech recognition 2nd pass re-ranking in low-resource languages due to (1) data scarcity in low-resource languages, (2) expensive computing costs for training and refreshing 100+ monolingual models, and (3) hosting inefficiency considering sparse traffic. In this study, we present a new way to group multiple low-resource locales together and optimize the performance of Multilingual Transformer LMs in ASR. Our Locale-group Multilingual Transformer LMs outperform traditional multilingual LMs along with reducing maintenance costs and operating expenses. Further, for low-resource but high-traffic locales where deploying monolingual models is feasible, we show that fine-tuning our locale-group multilingual LMs produces better monolingual LM candidates than baseline monolingual LMs.

**Index Terms**: Multilingual language model, Transformer language model, speech recognition

## 1. Introduction

Automatic Speech Recognition (ASR) usually involves two passes. The first-pass acoustic models and n-gram language models generate n-best hypotheses from the global search space [1]. In the second pass, Neural Network Language Models (NNLM) are widely used to re-rank the n-best hypotheses [2]. It has been demonstrated that re-ranking using NNLM is effective at reducing WER (Word Error Rate) [3], with Transformer language models producing state-of-the-art results in re-ranking [4].

Today our ASR system supports 100+ locales, but re-ranking is only applied to a few high-resource locales, even though we have proven the higher benefits of re-ranking for low-resource locales like Slovenian. The key challenges today are: (1) the low-resource locales' training data is scarce, which limits our capacity to train the NNLM, (2) it is computationally expensive to train and regularly refresh 100+ monolingual re-ranking models, one for each locale, (3) it is prohibitively expensive and inefficient to host these monolingual models in production, as traffic can be sparse, but each model ends up taking memory and compute to host across all our Speech clusters.

Multilingual Transformer language models [5] [6] [7] provide a very effective general solution to support ASR with pre-trained components and data sources that can be shared across multiple languages. When applied blindly, however, multilingual transformer language models may not always match or beat the monolingual models. We found that grouping multiple akin locales can optimize performance, especially when dealing with low-resource locales. As a result, our Locale-group LMs outperform the general multilingual solutions.

This key insight has helped us tackle the above-listed challenges. Locales with limited resources (scarce data) can benefit from all the data available for their locale group. We would only need to train and maintain a few locale-group Transformer LMs and still can attain locale coverage of 100+ locales for re-ranking, and fewer overall 2nd pass LMs result in hosting and scaling efficiencies across our clusters in production.

In addition, our key finding regarding grouping low-resource locales has been found to work in other related domains, such as improving capitalization and punctuation in recognition outputs, with potential future applications beyond speech.

## 2. Related Work

**Multilingual/Cross-lingual.** The effectiveness of sentence encoders' generative pre-training was first demonstrated for English natural language processing [8] [9] [10]. Multiple approaches have since been proposed to extend it to multilingual/cross-lingual pretraining and show the success in transfer learning, such as mBERT [9], XLM [5], XLM-R [6], Unicoder [11], etc. Large amounts of unlabeled data from multiple languages are used to train these models, with the goal that low-resource languages can benefit from high-resource languages from shared vocabularies and underlying linguistic similarities. mBert trains a BERT model using Wikipedia corpora in 104 languages. XLM introduced a translation language model (TLM) in addition to masked language model (MLM), in which bilingual sentences are concatenated as inputs. To further improve the performance, Unicoder presents three new cross-lingual pre-training tasks, including cross-lingual word recovery, cross-lingual paraphrase classification and cross-lingual masked language model. XLM-R trains exclusively with MLM objective on a huge multilingual dataset at an enormous scale.

Multilingual is also explored in ASR, primarily from an acoustic perspective. In [12], a massive multilingual acoustic model trained with more than 50 languages and more than 16,000 hours of audio is proven to improve recognition performance especially in low-resource languages. XLSR [13], a cross-lingual speech representation learning method is proposed by pre-training a single model from 56,000 hours raw waveform of speech in 53 languages.

Our work mainly focuses on Multilingual Language Model in ASR 2nd pass re-ranking, where a language model score is interpolated with 1st pass LM and AM to select the 1-best recognition candidate.

## 3. Locale-group Transformer LM

Our proposed approach involves two steps. The first is to identify the underlying language group of the low-resource locale using our data-driven method, and the second is to process the low-resource locales' data with shareable Byte Pair Encoding (BPE) tokens [14] and train the large-scale Locale-group Trans-

Figure 1: *Example of lexical similarity scores across 26 European languages.*

|  | ga-IE | bg-BG | mt-MT | tr-TR | el-GR | sv-SE | da-DK | nb-NO | fi-FI | hr-HR | cs-CZ | sk-SK | sl-SI | en-all | de-DE | nl-NL | ca-ES | es-ES | fr-FR | it-IT | pl-PL | pt-PT | ro-RO | lt-LT | lv-LV | et-EE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ga-IE | 1.00 | 0.07 | 0.11 | 0.06 | 0.10 | 0.11 | 0.11 | 0.11 | 0.10 | 0.10 | 0.13 | 0.10 | 0.10 | 0.13 | 0.14 | 0.11 | 0.12 | 0.12 | 0.12 | 0.12 | 0.11 | 0.14 | 0.11 | 0.10 | 0.10 | 0.11 |
| bg-BG | 0.07 | 1.00 | 0.08 | 0.05 | 0.05 | 0.05 | 0.05 | 0.05 | 0.04 | 0.05 | 0.04 | 0.05 | 0.06 | 0.05 | 0.04 | 0.05 | 0.06 | 0.04 | 0.05 | 0.05 | 0.05 | 0.06 | 0.05 | 0.05 | 0.05 | 0.04 |
| mt-MT | 0.11 | 0.08 | 1.00 | 0.10 | 0.14 | 0.17 | 0.17 | 0.17 | 0.15 | 0.16 | 0.18 | 0.15 | 0.14 | 0.18 | 0.20 | 0.17 | 0.18 | 0.17 | 0.18 | 0.20 | 0.17 | 0.18 | 0.17 | 0.14 | 0.16 | 0.16 |
| tr-TR | 0.06 | 0.05 | 0.10 | 1.00 | 0.10 | 0.11 | 0.12 | 0.11 | 0.11 | 0.12 | 0.11 | 0.11 | 0.11 | 0.12 | 0.17 | 0.11 | 0.12 | 0.11 | 0.12 | 0.11 | 0.12 | 0.12 | 0.10 | 0.11 | 0.11 | 0.11 |
| el-GR | 0.10 | 0.05 | 0.14 | 0.10 | 1.00 | 0.10 | 0.12 | 0.12 | 0.09 | 0.11 | 0.09 | 0.10 | 0.09 | 0.12 | 0.08 | 0.11 | 0.14 | 0.10 | 0.11 | 0.10 | 0.11 | 0.14 | 0.09 | 0.12 | 0.10 | 0.10 |
| sv-SE | 0.11 | 0.05 | 0.17 | 0.11 | 0.10 | 1.00 | 0.27 | 0.27 | 0.15 | 0.13 | 0.16 | 0.12 | 0.10 | 0.23 | 0.18 | 0.20 | 0.20 | 0.17 | 0.20 | 0.19 | 0.16 | 0.19 | 0.17 | 0.10 | 0.14 | 0.13 |
| da-DK | 0.11 | 0.05 | 0.17 | 0.12 | 0.12 | 0.27 | 1.00 | 0.39 | 0.16 | 0.14 | 0.18 | 0.14 | 0.11 | 0.24 | 0.21 | 0.23 | 0.21 | 0.18 | 0.22 | 0.20 | 0.19 | 0.21 | 0.19 | 0.11 | 0.16 | 0.15 |
| nb-NO | 0.11 | 0.05 | 0.17 | 0.11 | 0.12 | 0.27 | 0.39 | 1.00 | 0.16 | 0.14 | 0.17 | 0.13 | 0.11 | 0.24 | 0.19 | 0.22 | 0.21 | 0.18 | 0.21 | 0.20 | 0.18 | 0.21 | 0.19 | 0.11 | 0.15 | 0.15 |
| fi-FI | 0.10 | 0.04 | 0.15 | 0.11 | 0.09 | 0.15 | 0.16 | 0.16 | 1.00 | 0.11 | 0.11 | 0.10 | 0.09 | 0.16 | 0.12 | 0.14 | 0.17 | 0.13 | 0.15 | 0.15 | 0.13 | 0.14 | 0.14 | 0.08 | 0.12 | 0.12 |
| hr-HR | 0.10 | 0.05 | 0.16 | 0.12 | 0.11 | 0.13 | 0.14 | 0.14 | 0.11 | 1.00 | 0.17 | 0.18 | 0.31 | 0.17 | 0.16 | 0.14 | 0.17 | 0.14 | 0.15 | 0.16 | 0.15 | 0.19 | 0.13 | 0.16 | 0.16 | 0.13 |
| cs-CZ | 0.13 | 0.04 | 0.18 | 0.11 | 0.09 | 0.16 | 0.18 | 0.17 | 0.11 | 0.17 | 1.00 | 0.35 | 0.16 | 0.28 | 0.17 | 0.19 | 0.22 | 0.19 | 0.23 | 0.21 | 0.20 | 0.22 | 0.20 | 0.11 | 0.15 | 0.12 |
| sk-SK | 0.10 | 0.05 | 0.15 | 0.11 | 0.10 | 0.12 | 0.14 | 0.13 | 0.10 | 0.18 | 0.35 | 1.00 | 0.15 | 0.15 | 0.12 | 0.13 | 0.15 | 0.12 | 0.14 | 0.14 | 0.15 | 0.13 | 0.17 | 0.12 | 0.15 | 0.12 |
| sl-SI | 0.10 | 0.06 | 0.14 | 0.09 | 0.09 | 0.10 | 0.11 | 0.11 | 0.09 | 0.31 | 0.16 | 0.15 | 1.00 | 0.12 | 0.14 | 0.11 | 0.13 | 0.11 | 0.12 | 0.13 | 0.11 | 0.13 | 0.11 | 0.12 | 0.14 | 0.10 |
| en-all | 0.13 | 0.05 | 0.18 | 0.12 | 0.12 | 0.23 | 0.24 | 0.24 | 0.16 | 0.17 | 0.28 | 0.15 | 0.12 | 1.00 | 0.46 | 0.30 | 0.28 | 0.30 | 0.41 | 0.31 | 0.24 | 0.35 | 0.25 | 0.12 | 0.17 | 0.16 |
| de-DE | 0.14 | 0.04 | 0.20 | 0.17 | 0.08 | 0.18 | 0.21 | 0.19 | 0.12 | 0.16 | 0.17 | 0.12 | 0.14 | 0.46 | 1.00 | 0.24 | 0.26 | 0.24 | 0.32 | 0.28 | 0.19 | 0.26 | 0.22 | 0.10 | 0.16 | 0.13 |
| nl-NL | 0.11 | 0.05 | 0.17 | 0.11 | 0.11 | 0.20 | 0.23 | 0.22 | 0.14 | 0.14 | 0.19 | 0.13 | 0.11 | 0.30 | 0.24 | 1.00 | 0.23 | 0.21 | 0.26 | 0.23 | 0.19 | 0.24 | 0.20 | 0.10 | 0.15 | 0.14 |
| ca-ES | 0.12 | 0.06 | 0.18 | 0.12 | 0.14 | 0.20 | 0.21 | 0.21 | 0.17 | 0.17 | 0.22 | 0.15 | 0.13 | 0.28 | 0.26 | 0.23 | 1.00 | 0.37 | 0.29 | 0.30 | 0.22 | 0.34 | 0.23 | 0.13 | 0.17 | 0.17 |
| es-ES | 0.12 | 0.06 | 0.17 | 0.11 | 0.10 | 0.20 | 0.18 | 0.19 | 0.13 | 0.14 | 0.19 | 0.12 | 0.11 | 0.30 | 0.24 | 0.21 | 0.37 | 1.00 | 0.28 | 0.30 | 0.18 | 0.38 | 0.21 | 0.10 | 0.15 | 0.14 |
| fr-FR | 0.12 | 0.05 | 0.18 | 0.12 | 0.11 | 0.20 | 0.22 | 0.21 | 0.15 | 0.15 | 0.23 | 0.14 | 0.12 | 0.41 | 0.32 | 0.26 | 0.29 | 0.28 | 1.00 | 0.32 | 0.21 | 0.32 | 0.23 | 0.11 | 0.16 | 0.15 |
| it-IT | 0.12 | 0.05 | 0.20 | 0.12 | 0.12 | 0.19 | 0.20 | 0.20 | 0.15 | 0.16 | 0.21 | 0.14 | 0.12 | 0.31 | 0.28 | 0.23 | 0.30 | 0.30 | 0.30 | 1.00 | 0.21 | 0.32 | 0.25 | 0.12 | 0.16 | 0.15 |
| pl-PL | 0.11 | 0.05 | 0.17 | 0.11 | 0.10 | 0.16 | 0.19 | 0.18 | 0.13 | 0.15 | 0.20 | 0.15 | 0.13 | 0.24 | 0.19 | 0.19 | 0.21 | 0.18 | 0.21 | 0.21 | 1.00 | 0.20 | 0.19 | 0.10 | 0.15 | 0.13 |
| pt-PT | 0.14 | 0.05 | 0.18 | 0.12 | 0.11 | 0.19 | 0.21 | 0.21 | 0.14 | 0.15 | 0.22 | 0.13 | 0.11 | 0.35 | 0.26 | 0.24 | 0.34 | 0.38 | 0.32 | 0.32 | 0.20 | 1.00 | 0.23 | 0.11 | 0.16 | 0.15 |
| ro-RO | 0.11 | 0.06 | 0.17 | 0.12 | 0.14 | 0.17 | 0.19 | 0.19 | 0.14 | 0.19 | 0.20 | 0.17 | 0.13 | 0.25 | 0.22 | 0.20 | 0.23 | 0.21 | 0.23 | 0.25 | 0.19 | 0.23 | 1.00 | 0.14 | 0.18 | 0.17 |
| lt-LT | 0.10 | 0.05 | 0.14 | 0.10 | 0.09 | 0.10 | 0.11 | 0.11 | 0.08 | 0.13 | 0.10 | 0.12 | 0.11 | 0.12 | 0.10 | 0.10 | 0.13 | 0.10 | 0.11 | 0.12 | 0.10 | 0.11 | 0.14 | 1.00 | 0.16 | 0.11 |
| lv-LV | 0.10 | 0.05 | 0.16 | 0.11 | 0.12 | 0.14 | 0.16 | 0.15 | 0.12 | 0.16 | 0.15 | 0.15 | 0.12 | 0.17 | 0.16 | 0.15 | 0.17 | 0.15 | 0.16 | 0.16 | 0.15 | 0.16 | 0.18 | 0.16 | 1.00 | 0.15 |
| et-EE | 0.11 | 0.04 | 0.16 | 0.11 | 0.10 | 0.13 | 0.15 | 0.15 | 0.12 | 0.13 | 0.12 | 0.12 | 0.10 | 0.16 | 0.13 | 0.14 | 0.17 | 0.14 | 0.15 | 0.15 | 0.13 | 0.15 | 0.17 | 0.11 | 0.15 | 1.00 |

former Language Model. Whenever we lack enough resources and hardware to support individual model development and deployment, we can choose to deploy the group-based multilingual Transformer Language Model, which provides significant Speech Recognition accuracy improvement, maintenance and cost reduction.

### 3.1. Language Group Identification

As one of three organizations selected to potentially partner with the European Parliament in 2020, Microsoft developed a real-time AI-based tool for live transcription and translation of debates. To identify the underlying language groups for 26 European languages, we proposed a two-step data-driven mechanism.

Firstly, we computed bi-lingual similarity score, which can be a measure of the number of phonemes, words, phrases, or the like that are present in both locales. The Figure 1 shows an example of bi-lingual lexical similarity scores for 26 European languages, where higher score indicates closer linguistic relations between the languages. We observed code-switching/loanwords to be common, especially in English. In contrast, per our experiments Bulgarian (bg-BG) does not appear to be close to any other languages based on our collected data. We suspect that some of the data skew was caused by source data filtering.

Secondly, we applied vector-based clustering techniques to categorize similar languages together based on similarity score vectors. As shown in Table 1, this mechanism successfully identifies language family like Balto-Slavic, the group 2, which contains Slovenian, Croatian, Slovak and Czech. Group 3 consists of most high-resource Germanic languages such as English and German, and Latin (Romance) languages like Italian, Spanish and French. In Group 4, Greek is supposed to have its own alphabet - the similarity with other languages is mainly due to code-switching/loanwords.

### 3.2. Shareable BPE Tokens

In our proposed approach, we process all languages with the same shared vocabulary created through Byte Pair Encoding

| Group | Languages |
|---|---|
| 1 | nb-NO, sv-SE, fi-FI, da-DK |
| 2 | sl-SI, hr-HR, cs-CZ, sk-SK |
| 3 | en-all, es-ES, nl-NL, fr-FR, ro-RO, ca-ES, it-IT, pt-PT, pl-PL, de-DE |
| 4 | bg-BG, lv-LV, lt-LT, ga-IE, et-EE, el-GR, mt-MT, tr-TR |

Table 1: *Language groups of 26 European locales*

(BPE) [14]. We provides several BPE format examples in Table 2. This approach can greatly improve the token coverage with limited token set size and standardize the sub-word units across languages that share the same alphabet. For example, with 250K BPE tokens, we achieve almost a 100% coverage of 350M unique words across 26 languages.

| Locale | Word-based Sent | BPE-based Sent |
|---|---|---|
| English | ask consternation to my word list | ask conster@@ nation to my word list |
| Irish | a naoi scoil nach bhfuil seomra acmhainne acu | a naoi scoil nach bhfuil seomra acmha@@ inne acu |
| Estonia | asukoha nimed on sofia ja bulgaaria | asukoha ni@@ med on sofia ja bulgaaria |

Table 2: *Examples of text in BPE format*

### 3.3. Transformer Language Model

#### 3.3.1. Data Balance

We compile a language group training dataset with a balancing mechanism and train the Locale-group Multilingual Transformer Language Model. To ensure balanced data coverage for multiple regions within the same family, we sample sentences with multinomially distributed probabilities $\{q\}_{i=1...N}$, similar as how sentences are sampled in [5],

$$q_i = \frac{p_i^\alpha}{\sum_{j=1}^{N} p_j^\alpha} \quad \text{with} \quad p_i = \frac{n_i}{\sum_{k=1}^{N} n_k} \quad (1)$$

#### 3.3.2. Locale-group Model Training

Using balanced language family data, we train the Locale-group Multi-lingual Transformer Language Model, which is similar to the structure used in [15].

During training, we record the valid loss and perplexity of individual locales, and of the language family. According to our findings, the average loss minimum is within the range of the individual locale's loss minimum, which indicates that the Locale-group Multilingual Transformer Language Model converges for all locales within the identified language group.

#### 3.3.3. General & Masked Fine-tuning

In natural language processing (NLP) and machine translation (MT), fine-tuning a pre-trained language model has become the de facto standard for transfer learning [10] [16]. In our work, our ability to serve efficiently and reduce computation costs will be compromised if we fine-tune the multilingual model towards a target language. However, the possibility of fine-tuning is also worth exploring, as (1) knowledge transfer will be verified if we can achieve better performance with fine-tuned monolingual model than the one trained from scratch, along with a unified training recipe, and (2) certain high-traffic locales will separate themselves from low-resource groups as more data is collected and more traffic is received over time, which allows us to support monolingual model deployment with adequate resources.

In general fine-tuning, we reserve the pre-trained multilingual model parameters, feed the target locale's data in, and train for several more epochs until it converges to a new minimum.

Masked Fine-tuning is designed to force the model to tune to a target locale. For BPE tokens that do not exist in the target locale, we freeze the token embedding updates and set the prediction score to a very small negative number, so the token loss will be close to 0, and avoid paying attention to irrelevant tokens. This process is illustrated in Figure 2.

## 4. Experiment

### 4.1. Dataset

#### 4.1.1. Train

We sampled 5B sentences across 26 European locales from our in-house text data corpus, with an average sentence length of 12. All text data are pre-processed into lexical format with our own text-normalization pipeline. Low-resource locales are up-sampled per Section 2.3.1. The same data is used to train the shareable 250K BPE tokens.

#### 4.1.2. Test

Test audio data used for word-error-rate reductions (WERR) measurement of each locale consists primarily of dictations and
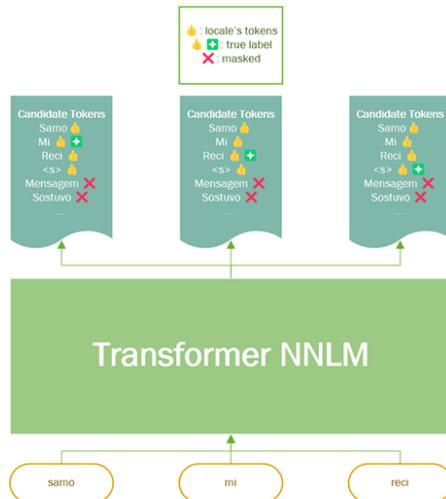


Figure 2: *Mask Fine-tuning Diagram with a Croatian example. "Samo mi reci" means "Just tell me" in English.*

spontaneous conversations. Minimum coverage per locale is 10K sentences.

### 4.2. Model

We experiment with configurations described in Table 4.

We train a baseline 1x1024:512 LSTM LM for each language, where 1 is the number of layers, 1024 is the dimensionality of the LSTM state, 512 is the dimensionality of the embedding and also the output dimensionality of the projection layer.

In addition, we trained one multilingual all-languages-together, four multilingual locale-group, and 26 monolingual Transformer language models with the same shareable 250K BPE token set, same data distribution and similar model configuration. These Transformer language models consist of 12 transformer layers, where each transformer layer contains 4096 feedforward dimensions with 16 heads. The warm-up is set to increase the learning rate gradually to improve the convergence of Transformer LMs [17].

We also applied general fine-tuning and masked fine-tuning as described in section 3.3.3.

### 4.3. Results

In this work, we mainly report word-error-rate reductions (WERR) on several low-resource locales: Croatian (hr-HR), Slovenian (sl-SI), Slovak (sk-SK), Lithuanian (lt-LT), Latvian (lv-LV) and Romanian (ro-RO). Language groups are described in section 3.1.

As shown in Table 3, we observe 3.34% average WERR improvement because of architecture upgrade from LSTM to Transformer, and more parameters. However, it is challenging to deploy those monolingual Transformer LMs as we discussed three limitations in the beginning.

On the other hand, the Locale-group Transformer LM provides a good solution considering the deployment restrictions. Firstly, with one Transformer LM to serve multiple locales in the same group, we can achieve 3.12% average WERR gain compared with the LSTM baseline. Secondly, the locale-group

|        |            | hr-HR | sl-SI | sk-SK | lt-LT | lv-LV | ro-RO | Avg   |
|--------|------------|-------|-------|-------|-------|-------|-------|-------|
| Mono   | LSTM       | 12.15 | 12.13 | 13.21 | 10.12 | 11.78 | 9.28  | 11.45 |
|        | Trans      | 16.66 | 15.13 | 16.12 | 14.37 | 15.06 | 11.28 | 14.79 |
| Multi  | Trans-All  | 13.83 | 13.06 | 15.08 | 10.74 | 11.92 | 9.68  | 12.38 |
|        | Trans-Group| 16.81 | 15.81 | 17.14 | 13.21 | 14.08 | 10.35 | 14.57 |
| Multi+FT | Trans-All | 17.38 | 16.01 | 17.22 | 14.63 | 15.07 | 11.73 | 15.36 |
|        | Trans-Group| 17.6  | 16.21 | 17.42 | 15.21 | 14.97 | 11.48 | 15.50 |
| Multi+MFT | Trans-Group | 18.22 | 17.83 | 17.75 | 15.15 | 16.13 | 11.48 | 16.09 |

Table 3: *WERR of several low-resource locales*

| Config | Description |
|--------|-------------|
| Mono:LSTM | baseline monolingual Long short-term memory (LSTM) LM |
| Mono:Trans | monolingual Transformer LM |
| Multi:Trans-All | all-languages-together Multilingual Transformer LM |
| Multi:Trans-Group | Locale-group Multilingual Transformer LM |
| Multi+FT:Trans-All | fine-tuned all-languages-together Multilingual Transformer LM |
| Multi+FT:Trans-Group | fine-tuned Locale-group Multilingual Transformer LM |
| Multi+MFT:Trans-Group | masked fine-tuned Locale-group Multilingual Transformer LM |

Table 4: *Model Configs*

model in general outperforms the all-data-together model by 2.19% when the implicit language similarity information is included, and limited model capacity can spare more attention to the learning of underlying linguistic patterns in similar languages instead of being distracted by irrelevant signals.

Even more, if we can allocate enough resources to train and deploy dedicated Transformer LM for low-resource locales, we can achieve more with masked fine-tuning. Compared with the monolingual LSTM baseline, our multilingual Locale-group model with masked fine-tuning can provide an additional 4.64% average WERR.

# 5. Discussion

## 5.1. BPE Token Size

We also trained models with 64k shareable BPE tokens to evaluate the impact of BPE token size. WERR on individual locales varies, but compared with 250K, the 64K BPE based Locale-group models generally regress around 2.48%. When we do masked fine-tuning, the gap is reduced to 0.66%. Our hypothesis is that the average text sequence length are elongated by smaller BPE token size, therefore brings in challenges to learn generic sequence patterns under the same number of model pa-

rameters.

## 5.2. Hosting Efficiency

One key motivation for exploring this idea was to improve production efficiency. Our speech service needs to be deployed globally to tens of clusters. Monolingual and general multilingual are two extremes when it comes to hosting costs. Deploying over 100 monolingual models everywhere is quite resource heavy and often wasteful as traffic can be sparse for some locales. On the other hand, general multi-lingual models oversimplify this, at the cost of WERR regression for some locales. Natural data imbalance is also harder to tackle at this scale, and often multi-lingual models can beat monolingual only by increasing model complexity (like adding more Transformer layers). This is also undesirable as this affects individual request latency at serving time. Our approach of locale-group multilingual models finds a more optimal point, needing us to deploy only a few locale-group models across our speech clusters, and better WERR than both monolingual and general multi-lingual models in most cases. We keep the option to deploy monolingual or fine-tuned monolingual models for certain high-traffic locales.

## 5.3. Parameter Tuning

We haven't tuned dropout and hyper parameters extensively for all the models. We believe we will achieve more gains with parameter tuning, and plan to explore this in the future work.

## 5.4. Application in other areas

In our hybrid ASR system, there are many other areas where we use monolingual models or general multi-lingual models. We explored the applicability of this locale-grouping technique to Capitalization models. We achieved similar results with locale-group capitalization models outperforming monolingual and general multi-lingual capitalization models.

# 6. Conclusion

Neural Network Language Model (NNLM) is an essential module in hybrid ASR to deliver the optimal recognition accuracy. In this work, we proposed a general and scalable approach to train and deploy large-scale Locale-group Transformer NNLM to support ASR in low-resource languages, where we observed significant accuracy improvement and reduction in model development and maintenance. Further, our fine-tuning experiments show that this locale-grouping helps create better monolingual models for low-resource languages.

# 7. References

[1] D. Yu and L. Deng, *Automatic Speech Recognition - A Deep Learning Approach*. Springer, October 2014. [Online]. Available: https://www.microsoft.com/en-us/research/publication/automatic-speech-recognition-a-deep-learning-approach/

[2] H. Erdogan, T. Hayashi, J. R. Hershey, T. Hori, C. Hori, W.-N. Hsu, S. Kim, J. L. Roux, Z. Meng, and S. Watanabe, "Multi-channel speech recognition: Lstms all the way through," in *CHiME 2016 - 4th International Workshop on Speech Processing in Everyday Environments*, September 2016. [Online]. Available: https://www.microsoft.com/en-us/research/publication/multi-channel-speech-recognition-lstms-all-the-way-through/

[3] X. Liu, Y. Wang, X. Chen, M. J. F. Gales, and P. C. Woodland, "Efficient lattice rescoring using recurrent neural network language models," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 4908–4912.

[4] K. Irie, A. Zeyer, R. Schlüter, and H. Ney, "Language modeling with deep transformers," *CoRR*, vol. abs/1905.04226, 2019. [Online]. Available: http://arxiv.org/abs/1905.04226

[5] G. Lample and A. Conneau, "Cross-lingual language model pretraining," *CoRR*, vol. abs/1901.07291, 2019. [Online]. Available: http://arxiv.org/abs/1901.07291

[6] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov, "Unsupervised cross-lingual representation learning at scale," *CoRR*, vol. abs/1911.02116, 2019. [Online]. Available: http://arxiv.org/abs/1911.02116

[7] T. Pires, E. Schlinger, and D. Garrette, "How multilingual is multilingual bert?" *CoRR*, vol. abs/1906.01502, 2019. [Online]. Available: http://arxiv.org/abs/1906.01502

[8] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding by generative pre-training," 2018.

[9] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," *CoRR*, vol. abs/1810.04805, 2018. [Online]. Available: http://arxiv.org/abs/1810.04805

[10] J. Howard and S. Ruder, "Fine-tuned language models for text classification," *CoRR*, vol. abs/1801.06146, 2018. [Online]. Available: http://arxiv.org/abs/1801.06146

[11] H. Huang, Y. Liang, N. Duan, M. Gong, L. Shou, D. Jiang, and M. Zhou, "Unicoder: A universal language encoder by pretraining with multiple cross-lingual tasks," 2019.

[12] V. Pratap, A. Sriram, P. Tomasello, A. Hannun, V. Liptchinsky, G. Synnaeve, and R. Collobert, "Massively multilingual asr: 50 languages, 1 model, 1 billion parameters," *arXiv preprint arXiv:2007.03001*, 2020.

[13] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, "Unsupervised cross-lingual representation learning for speech recognition," *CoRR*, vol. abs/2006.13979, 2020. [Online]. Available: https://arxiv.org/abs/2006.13979

[14] R. Sennrich, B. Haddow, and A. Birch, "Neural machine translation of rare words with subword units," *CoRR*, vol. abs/1508.07909, 2015. [Online]. Available: http://arxiv.org/abs/1508.07909

[15] A. Radford and K. Narasimhan, "Improving language understanding by generative pre-training," 2018.

[16] Y. Tang, C. Tran, X. Li, P. Chen, N. Goyal, V. Chaudhary, J. Gu, and A. Fan, "Multilingual translation with extensible multilingual pretraining and finetuning," *CoRR*, vol. abs/2008.00401, 2020. [Online]. Available: https://arxiv.org/abs/2008.00401

[17] M. Popel and O. Bojar, "Training tips for the transformer model," *CoRR*, vol. abs/1804.00247, 2018. [Online]. Available: http://arxiv.org/abs/1804.00247