



# A Practitioner’s Guide to Building ASR Models for Low-Resource Languages: A Case Study on Scottish Gaelic

Ondřej Klejch<sup>1</sup>, William Lamb<sup>2</sup>, Peter Bell<sup>1</sup>

<sup>1</sup>Centre for Speech Technology Research, University of Edinburgh, United Kingdom

<sup>2</sup>Celtic and Scottish Studies, University of Edinburgh, United Kingdom

{o.klejch,w.lamb,peter.bell}@ed.ac.uk

## Abstract

An effective approach to the development of ASR systems for low-resource languages is to fine-tune an existing multilingual end-to-end model. When the original model has been trained on large quantities of data from many languages, fine-tuning can be effective with limited training data, even when the language in question was not present in the original training data. The fine-tuning approach has been encouraged by the availability of public-domain E2E models and is widely believed to lead to state-of-the-art results. This paper, however, challenges that belief. We show that an approach combining hybrid HMMs with self-supervised models can yield substantially better performance with limited training data. This combination allows better utilisation of all available speech and text data through continued self-supervised pre-training and semi-supervised training. We benchmark our approach on Scottish Gaelic, achieving WER reductions of 32% relative over our best fine-tuned Whisper model.

**Index Terms:** speech recognition, low-resource languages, Scottish Gaelic

## 1. Introduction

Automatic speech recognition (ASR) has been democratised by several public-domain multi-lingual end-to-end (E2E) models. For example, OpenAI’s Whisper [1] is a state-of-the-art multilingual model that supports around 100 languages. Another example, Meta’s MMS [2], supports more than 1000 languages. However, this is only a fraction of all 7000 languages spoken world-wide. Since these models were trained on millions of hours from many languages, they can be fine-tuned even for languages and domains unseen during their training. Popular toolkits, such as HuggingFace Transformers [3], allow laypeople to readily train ASR models for their particular language.

On one hand, E2E models and toolkits are very accessible to beginners. However, this accessibility comes at the cost of flexibility and finer control – these toolkits offer limited options for improving the performance of E2E models beyond simple hyperparameter tuning. Furthermore, fine-tuning of E2E models – which requires transcribed audio data – might not be the optimal choice when having access only to a limited amounts of such data. On the other hand, toolkits for traditional hybrid HMM systems, such as Kaldi [4], have a much steeper learning curve. However, once grasped, they provide practitioners with finer control over every modelling decision, enabling them to train better models with limited training data. They can also leverage unpaired text and audio more straightforwardly than E2E models. Nevertheless, even if both types of systems have access to the same data, hybrid models have been shown to work better even with 1000 hours of transcribed speech [5].

For example, hybrid models can more easily utilize language models trained on all available text corpora for the language in question. These language models can be trained on vast amounts of text crawled from the internet and filtered by a language ID, in datasets like MADLAD-400 [6] or FineWeb [7]. Hybrid models can also work with any language model (LM) type, ranging from n-gram LM, RNN LM [8] to large Transformer LM [9]. Most recently, hybrid models have successfully leveraged large amounts of untranscribed speech by incorporating self-supervised pre-training [10, 11]. They can use large pre-trained self-supervised models (SSL models) as feature extractors similar to multilingual bottleneck features in the past [12]. Furthermore, these SSL models can be adapted to the target language with self-supervised continued pre-training [13] using only audio data.

Hybrid models can also be improved with semi-supervised training [14, 15], which uses a seed ASR model to produce pseudo-labels for untranscribed speech. These pseudo-labels can then be used for standard supervised training. Semi-supervised training can significantly improve a poor acoustic model if there is access to a good language model [16]: for example a language model trained on the aforementioned web-scale text corpora. Furthermore, it has been shown that self-supervised pre-training and semi-supervised training are complementary [17] and can be used to improve the performance of ASR systems for low-resource languages [18]. E2E models can also be improved with semi-supervised training [19], but since they cannot so easily leverage knowledge from text, the performance of semi-supervised training with E2E models is likely to be bounded by the quality of the seed model. Still, it is possible to train a good hybrid model first and then use it to transcribe vast amounts of data to train a powerful E2E model [20].

This paper demonstrates that hybrid models can continue to outperform end-to-end models in low-resource language settings. Specifically, we present an approach that combines hybrid models with features extracted from self-supervised models, subword n-gram and RNN language models. Furthermore, we find that replacing grapheme acoustic units with byte pair encoding (BPE) [21] subword units yields better performance in code-switched scenarios. We benchmark our approach on Gaelic, a Celtic language with 69,700 speakers in Scotland [22]. Despite the low number of speakers, Gaelic is somewhat atypical of low-resource languages because it has a long textual history with a standardised orthography, and a century of media broadcasts. These text and speech resources were previously used to build Gaelic ASR models with the standard hybrid recipes [23, 24]. However, with our carefully tuned approach leveraging recent advancements, we achieve 54% better performance than the previous models. Furthermore, our approach is also 32% better than our best fine-tuned Whisper model.

## 2. Method

To train the best possible Gaelic ASR system with a limited amount of manually transcribed data, it is important to leverage all available text corpora, which might be more accessible than transcribed speech. Hybrid models [25] can effectively combine acoustic models trained on transcribed data with language models trained on all available text. Therefore, we train an acoustic model on the limited transcribed speech and train a language model on all available Gaelic text. Since the amount of Gaelic text is limited compared to well-resourced languages, we experiment with using subword language models, which have been shown to work better than word language models with limited text data [26]. We also train RNN language models for use in lattice rescoring [8] to better utilise the text data.

The performance can be boosted further by using large pre-trained models. In particular, we use hidden representations extracted with very large multilingual pre-trained models [10, 11] instead of MFCC features. These large pre-trained models have seen very limited Gaelic data during pre-training. Therefore, we run continued self-supervised pre-training on our training data [13] to improve their performance on Gaelic data.

Languages like Gaelic with a strong focus on oral traditions may contain substantial quantities of speech data in media archives, and untranscribed speech might even be more readily available than text. Therefore, it is important to be able to leverage this untranscribed speech: in this paper, we do it in two ways. We use our initial ASR model to generate automatic transcriptions that are then manually corrected by professional transcribers, akin to active learning methods [27]. We also experiment with semi-supervised training [14], by directly using automatic transcriptions as pseudo-labels for training.

## 3. Data

We trained our models on several datasets: see Table 1. We used the same training data presented in [24] at the Celtic Language Technology Workshop, which we label “CLTW train”. This dataset consists of 103 hours of teaching videos, traditional narratives, and audio books, which were automatically aligned with normalised transcripts. Upon analysing this data, we noticed that the transcripts were not normalised properly. Therefore, to prevent any possible performance degradation, we applied our own normalisation script, which ensured that Gaelic accents are used consistently and which mapped words containing non-Gaelic letters to spoken noise.<sup>1</sup>

We also used 40 hours of manually transcribed speech from the historical BBC Radio nan Gàidheal programme, Prògram Choinnich, which ran for 30 years. Finally, since code-switching between Gaelic and English is very common – as in many marginalised languages – we also used 145 hours of accurately transcribed English broadcast data from the MGB dataset [28]. In the later stages of our experiments we used an additional 68 hours of Prògram Choinnich data that was automatically transcribed with our initial ASR models and manually corrected by professional transcribers. We also used 184 hours of untranscribed news data for semi-supervised training.

We used all available text data for training the language models. This included manual transcripts, books and web data. The most important datasets were the CLTW language model

<sup>1</sup>We also noticed that apostrophes, which have a special role in Gaelic, were removed from the training transcripts. It was not possible to fix this with normalisation, so we decided to ignore substitution errors with leading or trailing apostrophes during evaluation on this data.

Table 1: Summary of training datasets

dataset	language	# hours
CLTW train [24]	Gaelic	103
Prògram Choinnich	Gaelic	40
MGB	English	145
Extra Prògram Choinnich	Gaelic	68
News (untranscribed)	Gaelic	184

Table 2: Summary of testing datasets

dataset	# hours
CLTW test [24]	0.9
News	2.6
BBC	2.5
PC	1.7

data (18M words), the Gaelic portion of MADLAD-400 [6] (86M words) and English text from MGB [28] (645M words).

We used four datasets for evaluating our model: see Table 2. First, we used the test set from [24], “CLTW test”, which comprises 0.9 hours of speech. Similar to the training data we found many problems with transcriptions of this dataset. Therefore, our professional transcribers manually corrected the transcriptions. The WER between original and corrected transcriptions was 16.3%. Due to its small size, we decided to use this dataset as a validation set during our experiments. Since our ultimate goal is to provide subtitles for Gaelic broadcasts, we created two test sets from Gaelic BBC Alba programmes. The first, “News”, contains 7 episodes of an evening news programme broadcasted from 16th October 2023 to 22nd October 2023. The second, “BBC”, contains all Gaelic programmes broadcasted on 18th October 2023. It consists of several genres ranging from kids cartoons, news, talk shows to drama. Finally, we also held out 2 episodes of Prògram Choinnich (“PC”) as a test set.

## 4. Experiments

### 4.1. Whisper Fine-Tuning

We fine-tuned Whisper [1] with the HuggingFace Transformers library [3] as a baseline approach for training Gaelic ASR. Due to the limited amounts of transcribed training data (288 hours) and computation constraints, we decided to fine-tune Whisper-Turbo [1] with LoRA [29] using 8-bit quantisation of the frozen weights. We trained the models on 4 NVIDIA GeForce RTX 2080 Titan GPUs for 20k steps, which corresponds to 11.3 epochs. We used the default AdamW optimizer with 50 warm-up steps and a linear learning rate schedule with a peak learning rate 0.001. The effective batch size was 32. We optimised the LoRA rank  $r = \{64, 128, 256, 512\}$ .

### 4.2. Baseline Hybrid Models

As a hybrid baseline, we used an improved version of the Gaelic ASR model from [24], which is available online via the Tarsgrìobhadair API. This model uses phoneme pronunciation lexicons and CNN-TDNN neural network architecture on top of MFCC features and i-vectors trained with LF-MMI [30].

We trained our own baseline hybrid models and all other subsequent acoustic models with the Kaldi toolkit [4]. Our

baseline model was a standard TDNN-F model [4] trained with LF-MMI [30]. It used MFCC features, but it did not use i-vectors. Even though a pronunciation lexicon exists for Gaelic,<sup>2</sup> we decided to use grapheme units to allow easier integration of subword n-gram language models. We trained all n-gram language models with the SRILM toolkit [31]. We experimented with a word 3-gram LM and a subword 4-gram LM trained on the training transcripts. The subwords were produced by a BPE tokenizer [21] with 10k tokens trained on the training transcripts. Since the amount of training transcripts is limited, we also used other text sources as explained in Section 3.

### 4.3. SSL Features

We used pre-trained SSL models as feature extractors to enhance the performance of our Gaelic ASR model. We replaced traditional MFCC features with features extracted with SSL models as in [18]. In particular, we used two SSL models: XLS-R 300M [10] and XEUS [11]. We extracted features from the 18th layer of both models. Since neither of these models had been trained on Gaelic, we experimented with continued self-supervised pre-training of XLS-R 300M on our training data. We used fairseq [32] with the default XLS-R 300M pre-training configuration for additional 40k steps. This took 2 days on 8 NVIDIA GeForce RTX 2080 Titan GPUs. To be able to train on these GPUs with 12 GB VRAM, we reduced the maximum duration of each utterance to 5s.

Subsequently, we trained a standard TDNN-F hybrid model [4] on top of these SSL features. Since both SSL models output 50 features per second, we adapted the Kaldi training recipe to use a frame subsampling factor of 1 instead of the traditionally used 3. We also removed the LDA initial layer, because estimating the transform is very slow for SSL features.

### 4.4. BPE Acoustic Units

In our initial experiments, we observed that our model did not perform well on code-switched utterances. We hypothesized that this was due to the fact that graphemes are pronounced very differently in English and Gaelic. Whilst this could be fixed by using phone pronunciation dictionaries instead of grapheme dictionaries, that solution would make using subword LMs complicated because it would be necessary to infer a pronunciation for each subword unit. Instead, we decided to replace grapheme units with contextual graphemes as in [33, 34]: we used BPE tokens [21] as a replacement for graphemes, experimenting with various sizes of BPE inventory ( $\{500, 1000, 2000, 5000\}$ ). We used word-position independent BPE units, reducing the n-gram order of the denominator graph used for LF-MMI [30] to 2 to prevent memory explosion. Note that we still used ‘bi-phones’ and we clustered them with tree-based clustering.

### 4.5. RNN LM

Our initial model made a lot of errors that we attributed to the weak language model. Therefore, we applied RNN-LM rescoring [8]. We trained RNN-LMs with Kaldi on all available text training data mixed with English MGB text data for 40 epochs. We trained three RNN-LMs with increasing size of the embedding and the LSTM hidden cells, respectively  $\{(512, 128), (1024, 256), (2048, 512)\}$ .

<sup>2</sup><https://www.faclair.com/index.aspx>

### 4.6. Using Untranscribed Data

To further improve the performance of our system, we investigated three approaches to increase the amount of training data. First, we used noise augmentation [35] to make the models more noise robust, and to increase the amount of training data by using 3 noisy copies of each utterance. Second, we used an earlier iteration of our ASR system to produce automatic transcriptions of 68 hours of the Prògram Choinnich data. Human transcribers then manually corrected these automatic transcriptions. The use of ASR made the manual transcription process 50% faster compared to transcribing the recordings from scratch. Third, we performed semi-supervised training [15] by using our ASR system to produce automatic transcription of 300 hours of An Là data. We decoded this data in 30 s chunks and we segmented the data based on these first-pass transcripts. We only used segments that were 5-30 s long for training, resulting in an additional 184 hours of training data. Note that we applied noise augmentation to both Prògram Choinnich and An Là data and we pooled this data with the original noise augmented training data.

## 5. Results

We used the Tar-sgrìobhadair API and Whisper as baselines for our experiments. The Tar-sgrìobhadair API achieved an average WER of 28.0%. Looking at fine-tuned Whisper-Turbo models, we see that the performance improves with the number of fine-tuned parameters. The best average WER of 22.0% was achieved with a LoRA rank 512. This suggests that we might achieve even better results when fine-tuning the whole Whisper-Turbo model. We also fine-tuned Whisper with the additional manually and automatically transcribed data, reducing the average WER to 19.0%. Our grapheme TDNN-F baseline model performed best with the web subword LM, with an average of WER 28.4%, which is close to the Tar-sgrìobhadair baseline.

Subsequently, we explored the performance of models using SSL features instead of MFCC features. We can see from the results that using SSL features makes a substantial difference, reducing the average WER by 33.5% to 37.7% relative. XEUS appears to be a better base model than XLS-R 300M, but continued pre-training of XLS-R 300M, called XLS-R 300M CP, can achieve better performance than XEUS with an average WER of 17.7%. Whilst we believe that XEUS would also benefit from continued pre-training, unfortunately, code for this has not been integrated into ESPnet yet; therefore, we decided to use XLS-R 300M CP for the remainder of the experiments.

We next explored replacing graphemes with BPE units. We can see that models using BPE units achieve 6.7% – 9.6% lower WER than the model using grapheme units. The best models use 1000 or 2000 BPE units and achieve an average WER of 16.0%. Therefore we decided to use 1000 BPE units.

After that, we rescored lattices with RNN language models. Among these language models, the largest one achieves the best performance with an average WER of 15.0%, a 6.3% relative improvement compared to the first pass results. We hypothesise that we could achieve even better performance by rescoring with a Transformer LM [9] or a large language model fine-tuned on Gaelic. In all the subsequent experiments, we report the WER obtained after rescoring with the largest RNN-LM.

Finally, we increased the amount of training data by noise augmentation, adding the additional manually transcribed data, and by using automatically transcribed data. Furthermore, we make additional use of all the untranscribed data, we decided

Table 3: Results

<i>Baselines</i>					CLTW	News	BBC	PC	Avg.
Tar-sgrìobhadair API (an improved model from [24])					20.2	23.7	38.7	29.2	28.0
Whisper Turbo, LoRA 64					25.7	23.1	31.5	27.4	26.9
Whisper Turbo, LoRA 128					21.9	21.0	28.4	22.8	23.5
Whisper Turbo, LoRA 256					21.3	19.8	27.3	22.4	22.7
Whisper Turbo, LoRA 512					20.8	19.2	26.9	21.2	22.0
+ 68 hours of manually transcribed radio talk shows					19.5	17.1	25.8	19.3	20.4
+ 184 hours of automatically transcribed news					<b>19.4</b>	<b>15.0</b>	<b>23.5</b>	<b>18.1</b>	<b>19.0</b>
<i>Grapheme AM Baseline</i>									
AM Unit	AM Feature	LM Unit	LM Data	RNN LM					
grapheme	MFCC	300k words	train		26.8	25.3	44.1	24.2	30.1
grapheme	MFCC	10k BPEs	train		27.1	25.2	44.4	23.7	30.1
grapheme	MFCC	10k BPEs	web		<b>26.5</b>	<b>21.7</b>	<b>42.8</b>	<b>22.6</b>	<b>28.4</b>
<i>SSL Features</i>									
grapheme	XLS-R 300M	10k BPEs	web		18.1	15.8	25.3	16.3	18.9
grapheme	XLS-R 300M CP	10k BPEs	web		<b>16.0</b>	15.9	24.2	<b>14.7</b>	<b>17.7</b>
grapheme	XEUS	10k BPEs	web		17.3	<b>13.9</b>	<b>23.7</b>	17.8	18.2
<i>BPE Acoustic Units</i>									
0.5k BPEs	XLS-R 300M CP	10k BPEs	web		14.9	13.9	22.1	14.3	16.3
1k BPEs	XLS-R 300M CP	10k BPEs	web		14.6	<b>13.7</b>	<b>21.6</b>	<b>14.0</b>	<b>16.0</b>
2k BPEs	XLS-R 300M CP	10k BPEs	web		<b>14.4</b>	14.0	21.7	14.1	<b>16.0</b>
5k BPEs	XLS-R 300M CP	10k BPEs	web		15.0	14.7	22.2	14.2	16.5
<i>RNN LM Rescoring</i>									
1k BPEs	XLS-R 300M CP	10k BPEs	web	RNN LM 512	14.5	13.1	21.1	13.9	15.6
1k BPEs	XLS-R 300M CP	10k BPEs	web	RNN LM 1024	13.9	12.7	20.5	13.6	15.2
1k BPEs	XLS-R 300M CP	10k BPEs	web	RNN LM 2048	<b>13.8</b>	<b>12.4</b>	<b>20.3</b>	<b>13.4</b>	<b>15.0</b>
<i>Using Untranscribed Data</i>									
+ noise augmentation					13.6	11.7	18.9	13.4	14.4
+ 68 hours of manually transcribed radio talk shows					13.1	11.2	18.0	12.5	13.7
+ continual pre-training on all data					<b>11.8</b>	10.8	19.7	<b>10.3</b>	13.2
+ 184 hours of automatically transcribed news					12.0	<b>10.4</b>	<b>17.7</b>	10.9	<b>12.8</b>

to continue pre-training XLS-R 300M on all our transcribed and untranscribed data for 100k iterations. Combining all these techniques yielded the final average WER of 12.8%, which is 32% relative better than our best fine-tuned Whisper model. Looking at individual test sets, we see that WER on CLTW, News and PC ranges from 10.4% to 12.0%. However, the WER on the BBC test set is much higher at 17.7%. This is due to a large number of children’s programs in the dataset, which have high WER because of deletion errors caused by background music and children’s voices.

Overall, these results show a promising trend, suggesting that we could achieve further gains by using larger amounts of untranscribed data. We hope that we might further reduce the WER by doing continued self-supervised pre-training and lattice-based semi-supervised training with LF-MMI [15] on a much larger and more diverse corpus.

Since languages like Gaelic might have significantly more speech data available than text data, we also experimented with unsupervised language model adaptation [36]. We interpolated our web n-gram LM with an n-gram LM trained on automatic transcripts. Our initial experiments showed that the adapted LM significantly degraded accuracy, therefore we left this for future experiments.

## 6. Conclusions

In this paper we showed that optimized hybrid models can achieve very good results on low-resource languages such as Scottish Gaelic. Our model outperformed the previous best model deployed in the Tar-sgrìobhadair API by 54%. Furthermore, it outperformed our best fine-tuned Whisper-Turbo model by 32% relative. To get the best possible performance, we leveraged all available text data to train n-gram and RNN language models, used BPE acoustic units instead of graphemes, used SSL features instead of MFCCs, and performed noise augmentation, active and semi-supervised learning to increase the amount of training data. We believe that this approach is applicable to all low-resource languages that have similar amounts of speech and text data available as Scottish Gaelic.

In the future, we will try fine-tuning other end-to-end models [2, 37] for Scottish Gaelic. To get more transcribed data for training these end-to-end models, we plan to transcribe large quantities of untranscribed data with our best hybrid model as in [20]. Furthermore, based on the promising RNN LM results, we will attempt to fine-tune large language models on Gaelic, similar to recent work on Basque [38]. Such a model could be used for rescoring, generative error correction [39] or directly for speech recognition as in LLaMA-Omni [40].

## 7. Acknowledgements

This work was supported by the Scottish Government (Grant name: ‘Ecosystem for Interactive Speech Technologies’). We thank BBC Alba for providing the data and Cailean Gordan, Alison Diack and Fearchar MacIllFhinnein for transcribing training and testing data.

## 8. References

- [1] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” in *ICML*, 2023.
- [2] V. Pratap, A. Tjandra, B. Shi *et al.*, “Scaling speech technology to 1,000+ languages,” *JMLR*, vol. 25, no. 97, pp. 1–52, 2024.
- [3] T. Wolf, L. Debut, V. Sanh *et al.*, “Transformers: State-of-the-art natural language processing,” *EMNLP*, 2020.
- [4] D. Povey, G. Cheng, Y. Wang, K. Li, H. Xu, M. Yarmohammadi, and S. Khudanpur, “Semi-Orthogonal Low-Rank Matrix Factorization for Deep Neural Networks,” in *Interspeech*, 2018.
- [5] A. Rouhe, T. Grósz, and M. Kurimo, “Principled comparisons for end-to-end speech recognition: Attention vs hybrid at the 1000-hour scale,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 32, pp. 623–638, 2023.
- [6] S. Kudugunta, I. Caswell, B. Zhang, X. Garcia, D. Xin, A. Kusupati, R. Stella, A. Bapna, and O. Firat, “Madlad-400: A multilingual and document-level large audited dataset,” *NeurIPS*, 2024.
- [7] G. Penedo, H. Kydlíček, L. Ben allal, A. Lozhkov, M. Mitchell, C. A. Raffel, L. Von Werra, and T. Wolf, “The fineweb datasets: Decanting the web for the finest text data at scale,” in *NeurIPS*, 2024.
- [8] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, “Recurrent neural network based language model,” in *Interspeech*, 2010.
- [9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017.
- [10] A. Babu, C. Wang, A. Tjandra *et al.*, “XLS-R: Self-supervised Cross-lingual Speech Representation Learning at Scale,” in *Interspeech*, 2022.
- [11] W. Chen, W. Zhang, Y. Peng, X. Li, J. Tian, J. Shi, X. Chang, S. Maiti, K. Livescu, and S. Watanabe, “Towards robust speech representation learning for thousands of languages,” in *EMNLP*, 2024.
- [12] K. Veselý, M. Karafiát, F. Grézl, M. Janda, and E. Egorova, “The language-independent bottleneck features,” in *SLT*, 2012.
- [13] J.-H. Lee, C.-W. Lee, J.-S. Choi, J.-H. Chang, W. K. Seong, and J. Lee, “CTRL: Continual Representation Learning to Transfer Information of Pre-trained for WAV2VEC 2.0,” *Interspeech*, 2022.
- [14] L. Lamel, J.-L. Gauvain, and G. Adda, “Lightly supervised and unsupervised acoustic model training,” *Computer Speech and Language*, vol. 16, no. 1, pp. 115–229, 2002.
- [15] V. Manohar, H. Hadian, D. Povey, and S. Khudanpur, “Semi-supervised training of acoustic models using lattice-free MMI,” in *ICASSP*, 2018.
- [16] E. Wallington, B. Kershenbaum, O. Klejch, and P. Bell, “On the learning dynamics of semi-supervised training for ASR,” in *Interspeech*, 2021.
- [17] Q. Xu, A. Baevski, T. Likhomanenko, P. Tomasello, A. Conneau, R. Collobert, G. Synnaeve, and M. Auli, “Self-training and pre-training are complementary for speech recognition,” in *ICASSP*, 2021.
- [18] L.-M. Lam-Yee-Mui, L. O. Yang, and O. Klejch, “Comparing self-supervised pre-training and semi-supervised training for speech recognition in languages with weak language models,” in *Interspeech*, 2023.
- [19] D. S. Park, Y. Zhang, Y. Jia, W. Han, C.-C. Chiu, B. Li, Y. Wu, and Q. V. Le, “Improved noisy student training for automatic speech recognition,” in *Interspeech*, 2020.
- [20] J. Silovsky, L. Deng, A. Argueta, T. Arvizo, R. Hsiao, S. Kuznetsov, Y.-C. Lin, X. Xiao, and Y. Zhang, “Cross-lingual knowledge transfer and iterative pseudo-labeling for low-resource speech recognition with transducers,” *arXiv preprint arXiv:2305.13652*, 2023.
- [21] R. Sennrich, B. Haddow, and A. Birch, “Neural machine translation of rare words with subword units,” in *ACL*, 2016.
- [22] National Records of Scotland, “Scotland’s census 2022,” <https://www.scotlandscensus.gov.uk>, 2022, accessed: 2024-08-14.
- [23] R. Rasipuram, P. Bell, and M. M. Doss, “Grapheme and multilingual posterior features for under-resourced speech recognition: a study on Scottish Gaelic,” in *ICASSP*, 2013.
- [24] L. Evans, W. Lamb, M. Sinclair, and B. Alex, “Developing automatic speech recognition for Scottish Gaelic,” in *Celtic Language Technology Workshop*, 2022.
- [25] H. A. Bourlard and N. Morgan, *Connectionist speech recognition: a hybrid approach*. Springer Science & Business Media, 2012, vol. 247.
- [26] P. Smit, S. Virpioja, and M. Kurimo, “Advances in subword-based HMM-DNN speech recognition across languages,” *Computer Speech & Language*, vol. 66, p. 101158, 2021.
- [27] T. Drugman, J. Pytkkönen, and R. Kneser, “Active and semi-supervised learning in ASR: Benefits on the acoustic and language models,” in *Interspeech*, 2016.
- [28] P. Bell, M. J. Gales, T. Hain *et al.*, “The MGB challenge: Evaluating multi-genre broadcast media recognition,” in *ASRU*, 2015.
- [29] E. J. Hu, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen *et al.*, “LoRA: Low-Rank Adaptation of Large Language Models,” in *ICLR*, 2022.
- [30] D. Povey, V. Peddinti, D. Galvez, P. Ghahremani, V. Manohar, X. Na, Y. Wang, and S. Khudanpur, “Purely sequence-trained neural networks for ASR based on lattice-free MMI,” in *Interspeech*, 2016.
- [31] A. Stolcke *et al.*, “SRILM - an extensible language modeling toolkit,” in *Interspeech*, 2002.
- [32] M. Ott, S. Edunov, A. Baevski, A. Fan, S. Gross, N. Ng, D. Grangier, and M. Auli, “fairseq: A fast, extensible toolkit for sequence modeling,” in *NAACL-HLT: Demonstrations*, 2019.
- [33] D. Le, X. Zhang, W. Zheng, C. Fügen, G. Zweig, and M. L. Seltzer, “From senones to chenones: Tied context-dependent graphemes for hybrid speech recognition,” in *ASRU*, 2019.
- [34] F. Zhang, Y. Wang, X. Zhang, C. Liu, Y. Saraf, and G. Zweig, “Faster, simpler and more accurate hybrid ASR systems using wordpieces,” in *Interspeech*, 2020.
- [35] T. Ko, V. Peddinti, D. Povey, M. L. Seltzer, and S. Khudanpur, “A study on data augmentation of reverberant speech for robust speech recognition,” in *ICASSP*, 2017.
- [36] M. Bacchiani and B. Roark, “Unsupervised language model adaptation,” in *ICASSP*, 2003.
- [37] K. C. Puvvada, P. Zelasko, H. Huang *et al.*, “Less is More: Accurate Speech Recognition & Translation without Web-Scale Data,” in *Interspeech*, 2024.
- [38] J. Etxaniz, O. Sainz, N. Miguel, I. Aldabe, G. Rigau, E. Agirre, A. Ormazabal, M. Artetxe, and A. Soroa, “Latxa: An open language model and evaluation suite for Basque,” in *ACL*, 2024.
- [39] C. Chen, Y. Hu, C.-H. H. Yang, S. M. Siniscalchi, P.-Y. Chen, and E.-S. Chng, “Hyporadise: An open baseline for generative speech recognition with large language models,” *NeurIPS*, 2024.
- [40] Q. Fang, S. Guo, Y. Zhou, Z. Ma, S. Zhang, and Y. Feng, “LLaMA-omni: Seamless speech interaction with large language models,” in *ICLR*, 2025.