
Google USM: Scaling Automatic Speech Recognition Beyond 100 Languages

Yu Zhang Wei Han James Qin Yongqiang Wang Ankur Bapna Zhehuai Chen
Nanxin Chen Bo Li Vera Axelrod Gary Wang Zhong Meng Ke Hu
Andrew Rosenberg Rohit Prabhavalkar Daniel S. Park Parisa Haghani
Jason Riesa Ginger Perng Hagen Soltau Trevor Strohman
Bhuvana Ramabhadran Tara Sainath Pedro Moreno Chung-Cheng Chiu
Johan Schalkwyk Françoise Beaufays Yonghui Wu^{*†}

Abstract

We introduce the Universal Speech Model (USM), a single large model that performs automatic speech recognition (ASR) across 100+ languages. This is achieved by pre-training the encoder of the model on a large unlabeled multilingual dataset of 12 million (M) hours spanning over 300 languages, and fine-tuning on a smaller labeled dataset. We use multilingual pre-training with random-projection quantization and speech-text modality matching to achieve state-of-the-art performance on downstream multilingual ASR and speech-to-text translation tasks. We also demonstrate that despite using a labeled training set 1/7-th the size of that used for the Whisper model [1], our model exhibits comparable or better performance on both in-domain and out-of-domain speech recognition tasks across many languages.

1 Introduction

Recent advances in self-supervised learning have ushered in a new era for speech recognition. Whereas previous works focused mostly on improving the quality of monolingual models for mainstream languages, recent studies have increasingly turned to “universal” models [1–4]. These may take the form of a single model that performs well on multiple tasks [1, 2], or one that covers multiple domains [2, 3], or one that supports multiple languages [1, 5]. In this work, we explore the frontiers of language expansion. Our long-term goal is to train a universal ASR model that covers all the spoken languages in the world.

A fundamental challenge in scaling speech technologies to many languages is obtaining enough data to train high-quality models. With conventional supervised training approaches, audio data needs to be manually transcribed, which is lengthy and expensive, or collected from existing transcribed sources which are hard to find for tail languages. While transcribed speech may be scarce in many

*All authors are affiliated with Google Inc.

†Contact author at ngyuzh@google.com.

languages, untranscribed speech and text data are practically unlimited. Recent developments in semi-supervised algorithms for speech recognition makes it possible to leverage such data for pre-training and produce high-quality speech models with a limited amount of transcribed data [3, 6].

Moreover, recent studies have shown that a single large model can utilize large data sets more effectively than smaller models [1, 4]. This all points to a promising direction where large amounts of unpaired multilingual speech and text data and smaller amounts of transcribed data can contribute to training a single large universal ASR model.

1.1 Our approach

We produce large “Universal Speech Models” (USMs) through a training pipeline that utilizes three types of datasets:

- **Unpaired Audio:**
 - **YT-NTL-U:** A large unlabeled multilingual dataset consisting of 12M hours of YouTube-based audio covering over 300 languages.
 - **Pub-U:** 429k hours of unlabeled speech in 51 languages based on public datasets.
- **Unpaired Text:**
 - **Web-NTL:** A large multilingual text-only corpus with 28B sentences spanning over 1140 languages.
- **Paired ASR Data:** We utilize two corpora of paired audio-text data with O(10k) hours of audio for supervised training.
 - **YT-SUP+:** 90k hours of labeled multilingual data covering 73 language and 100k hours of en-US pseudo-labeled data generated by noisy student training (NST) [7, 8] from YT-NTL-U.
 - **Pub-S:** 10k hours of labeled multi-domain en-US public data and 10k labeled multilingual public data covering 102 languages.

2B-parameter Conformer [9] models are built using these datasets through the following steps:

1. **Unsupervised Pre-training:** BEST-RQ (**B**ERT-based **S**peech pre-**T**raining with **R**andom-projection **Q**uantizer) [10] is used to pre-train the encoder of the model with YT-NTL-U.
2. **MOST (Multi-Objective Supervised pre-Training):** The model can optionally be further prepared by a multi-objective supervised pre-training pipeline that utilizes all three kinds of datasets: YT-NTL-U, Pub-U, Web-NTL and Pub-S. Here, a weighted sum of the BEST-RQ masked language model loss [11], along with the text-injection losses (including the supervised ASR loss and modality matching losses) [12, 13] is optimized during training.
3. **Supervised ASR Training:** We produce generic ASR models trained with connectionist temporal classification (CTC) [14] and Listen, Attend, and Spell (LAS) [15] transducers for downstream tasks.

Two types of models are produced through this pipeline—pre-trained models that can be fine-tuned on downstream tasks, and generic ASR models for which we assume no downstream fine-tuning occurs. The generic ASR models are trained with chunk-wise attention, which we introduce later in this report.

Table 1: USM models prepared in this work. The generic ASR models are trained on a large "upstream" ASR corpus and not finetuned further, while the pre-trained models are fine-tuned on downstream tasks.

Model	BEST-RQ	MOST	Model-Type	Decoder	Upstream ASR Dataset	Chunk-wise Attention
USM		N	Pre-trained	Downstream Dependent	-	N
USM-M	YT-NTL-U	Y	Pre-trained	Downstream Dependent	-	N
USM-LAS		N	Generic ASR	LAS	YT-SUP+	Y
USM-CTC		N	Generic ASR	CTC	YT-SUP+	Y

We denote the pre-trained models USM and USM-M, where the appendix **-M** indicates that **MOST** has been utilized for the preparation of the model. The USM and USM-M models can be further fine-tuned on the downstream task of choice with an appropriate transducer unit, which can be a CTC, LAS or RNN transducer (RNN-T) unit. We evaluate our USM models on two types of benchmarks:

- **Automatic Speech Recognition (ASR):** We use YouTube data to train USMs for YouTube (e.g., closed captions). We evaluate the USMs on two public benchmarks, SpeechStew [2] and FLEURS [16]. We also report results on the long-form test set CORAAL [17] for which only the evaluation set is available.
- **Automatic Speech Translation (AST):** We test AST performance on CoVoST 2 [18].

As indicated in Table 1, the generic ASR models are trained with YT-SUP+ and not fine-tuned on domain-specific datasets for downstream ASR tasks. We, however, explore the possibility of attaching additional “adapter” units [19] to both generic and pre-trained ASR models and training adapter weights while keeping the rest of the model frozen.

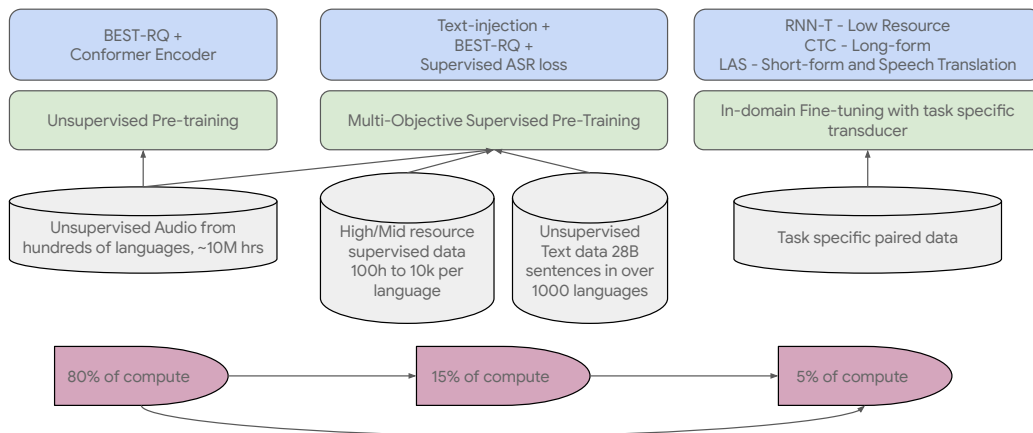


Figure 1: An overview of our approach. Training is split into three stages. (i) The first stage trains a conformer backbone on a large unlabeled speech dataset, optimizing for the BEST-RQ objective. (ii) We continue training this speech representation learning model while optimizing for multiple objectives, the BEST-RQ objective on unlabeled speech, the modality matching, supervised ASR and duration modeling losses on paired speech and transcript data and the text reconstruction objective with an RNN-T decoder on unlabeled text. (iii) The third stage fine-tunes this pre-trained encoder on the ASR or AST tasks.

The overall training pipeline of our models is summarized in Fig. 1. In our design, once a large amount of compute is expended in the pre-training stages, the downstream application can be conveniently fine-tuned from a model trained from stage-1 or stage-2 with a task-specific transducer. Our experimental results demonstrate that this pipelined training framework enables us to build both generic multilingual ASR systems and domain specific models with state-of-the-art performance.

We next present the key findings of our research, provide an overall view of the report, and review related work.

1.2 Key Findings

SoTA results for downstream multilingual speech tasks: Our USM models achieve state-of-the-art performance for multilingual ASR and AST for multiple datasets in multiple domains. This includes SpeechStew (mono-lingual ASR) [2], CORAAL (African American Vernacular English (AAVE) ASR) [17], FLEURS (multi-lingual ASR) [16], YT (multilingual long-form ASR), and CoVoST (AST from English to multiple languages). We depict our model’s performance in the first panel of Fig. 2. We also build an ASR model for YouTube captioning – i.e., the transcription of speech in YouTube videos, that achieves < 30% WER on 73 languages. With only 90k hours of supervised data, this model performs better than Whisper [1], a strong general ASR system trained on more than

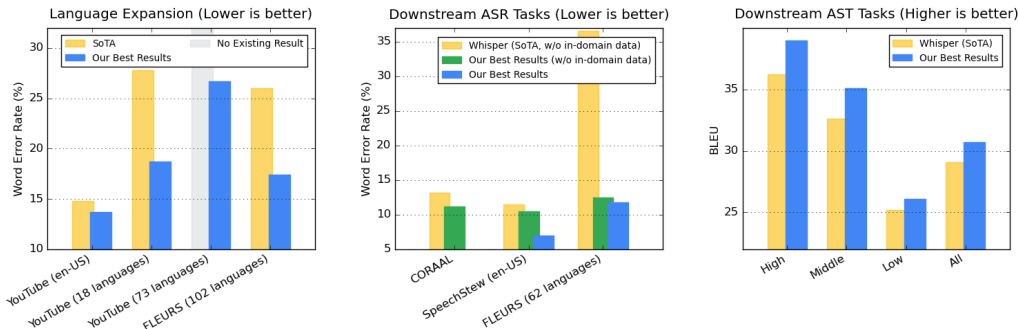


Figure 2: **(Left)**[†] WERs (%) Our language expansion effort to support more languages on YouTube (73 languages) and extending to 100+ languages on the public dataset (FLEURS). Lower is better. To the best of our knowledge, no published model can successfully decode all 73 languages from our YouTube set, thus we only list our results. **(Middle)**[†] Our results on ASR benchmarks, with or without in-domain data. Lower is better. **(Right)** SoTA results on public speech translation tasks. Results presented are presented as high/middle/low resources languages defined in [20]. Higher is better.

400k hours of transcribed data (we select 18 languages that Whisper can successfully decode with lower than 40% WER). The second panel of Fig. 2 demonstrates that our YouTube captions model generalizes well to unseen domains.

BEST-RQ is a scalable speech representation learner: We find that BEST-RQ pre-training can effectively scale to the very large data regime with a 2B parameter Conformer-based backbone, comparing favorably against Wav2Vec 2.0 [6] and W2v-BERT [21] in this setting.

MOST (BEST-RQ + text-injection) is a scalable speech and text representation learner: We demonstrate that MOST is an effective method for utilizing large scale text data for improving quality on downstream speech tasks, as demonstrated by quality gains exhibited for the FLEURS and CoVoST 2 tasks. Fig. 2 depicts USM’s performance, establishing a new state-of-the-art on the FLEURS benchmark across 102 languages for ASR and on CoVoST 2 across 21 languages on AST.

Representations from MOST (BEST-RQ + text-injection) can quickly adapt to new domains: We find that it is possible to obtain powerful downstream ASR/AST models by attaching and training light-weight residual adapter modules, which only add 2% of additional parameters, while keeping the rest of the model frozen.

Chunk-wise attention for robust long-form speech recognition: We introduce chunk-wise attention, an effective, scalable method for extending the performance of ASR models trained on shorter utterances to very long speech inputs. We find that the USM-CTC/LAS models trained with chunk-wise attention is able to produce high-quality transcripts for very long utterances in the YouTube evaluation sets.

1.3 Outline

The outline of this report is as follows:

Methods: We review the architecture and the methods used in the paper. We provide brief summaries of the Conformer [9], BEST-RQ [10], text-injection [12, 13] used for MOST, and Noisy Student Training (NST) [7, 8]. We also introduce chunk-wise attention for scalable training on long utterances.

Data: We describe the four types of datasets used to train our models: the unlabeled multilingual speech dataset YT-NTL-U, the multilingual text corpus Web-NTL, labeled datasets, and pseudo-labeled datasets.

Key Results: We present the performance of our USM models on downstream ASR and AST tasks. We demonstrate that USM establishes new states-of-the-art on several speech understanding benchmarks.

Analysis and Ablations: We present analysis of the effects of the key components of our work and compare their performance against existing methods.

1.4 Related Work

There is extensive literature on pre-training [6, 12, 22–33] and self-training [8, 34–44] for ASR. Large speech models trained on large datasets have been studied previously in both monolingual [3] and multilingual contexts [1, 4]. Large multi-modal speech models have been explored in [13, 20, 45–54]. Various unsupervised pre-training methods for speech models have been proposed and applied in [6, 10, 21].

Our work is an extension of a host of recent research efforts [3, 10, 13, 53, 55] that have studied semi-supervised learning for ASR in the context of deep-learning. Large speech models ($> 1B$) were first studied in [3]; we expand upon this approach to train multilingual speech models in this work. We improve the methods used in [3] by employing a more scalable self-supervised learning algorithm (BEST-RQ) and additionally applying multi-modal pre-training (text-injection) to prepare the models. We introduce an improvement to BEST-RQ [10] by utilizing a multi-softmax loss. We also incorporate Multi-Objective Supervised Training (BEST-RQ with text-injection) to improve the quality of speech representations learnt during pre-training, by utilizing transcribed data and unlabeled text. Long-form ASR has been studied in [1, 56, 57]; we propose chunk-wise attention as an alternative solution to chunk-based decoding.

In this paper, we propose a scalable self-supervised training framework for multilingual ASR which extends to hundreds of languages. In particular:

- We demonstrate that USMs pre-trained on 300 languages can successfully adapt to both ASR and AST tasks in new languages with a small amount of supervised data.
- We build a generic ASR model on 73 languages by fine-tuning pre-trained models on 90k hours of supervised data. We show that the generic ASR models can carry out inference efficiently on TPUs and can reliably transcribe hours-long audio on YouTube Caption ASR benchmarks.
- We conduct a systematic study on the effects of pre-training, noisy student training, text injection, and model size for multilingual ASR.

2 Methods

2.1 Model Architecture: Conformer

We use the convolution-augmented transformer [9], or Conformer, with relative attention [58] as an encoder model. For downstream speech tasks such as ASR or AST, the features produced by the Conformer are either used as an input to a connectionist temporal classification (CTC) [14], RNN transducer (RNN-T) [59] or a Listen, Attend, and Spell (LAS) [15] unit after additional projection. As will be discussed further, BEST-RQ pre-training is exclusively applied to the encoder, while other forms of training (e.g., T5 [60]) train the entire task network as a whole.

For our experiments, we consider two models with 600M and 2B parameters respectively. While the main results presented have been obtained using the 2B model, the 600M model is utilized for ablation studies and observing model scaling behavior. Some features of the models are listed in Table 2.

Table 2: Conformer model parameters.

Model	# Params (B)	# Layers	Dimension	Att. Heads	Conv. Kernel Size
Conformer-0.6	0.6	24	1024	8	5
Conformer-2B	2.0	32	1536	16	5

2.2 Pre-training: BEST-RQ

We select BEST-RQ [10] as the method to pre-train our networks with speech audio. BEST-RQ provides a simple framework with a small number of hyperparameters for unsupervised training on

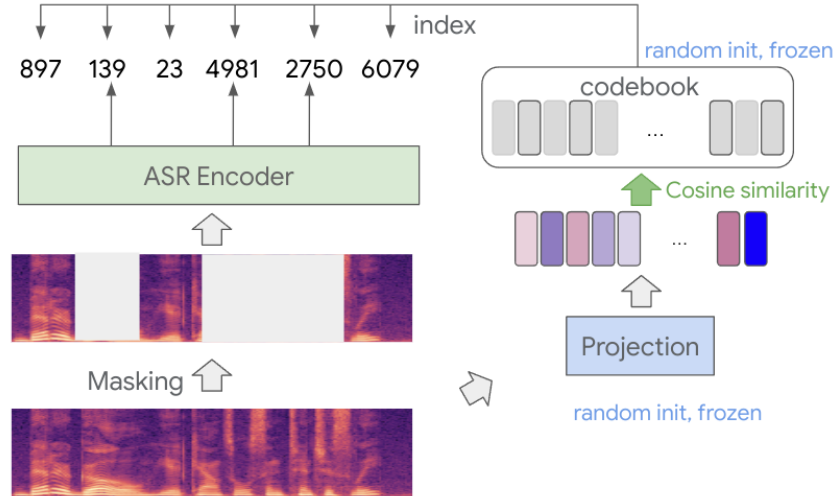


Figure 3: BEST-RQ based pre-training with conformer encoder.

large-scale unlabeled audio data. We discuss the comparative advantage of BEST-RQ against other pre-training methods in section 5.3.

BEST-RQ employs a BERT-style training task for the audio input that attempts to predict masked speech features. To make the task compatible with BERT-style training, the original speech features corresponding to the masked frames are quantized, and the task requires predicting the quantized label of these features. For a given number of quantization targets c , random “codebook” vectors v_0, \dots, v_{c-1} are chosen in an embedding space. The discrete label of the speech feature is obtained by first projecting the feature into the embedding space by a randomly initialized, frozen projection matrix and then finding the closest codebook vector. The index of this codebook vector is identified as the label of the speech feature. Cosine similarity is used as the distance measure for determining the code.

We note that while w2v-BERT [21] pre-training has proven to be an effective method for unsupervised pre-training, it requires an additional quantization module which introduces more complexity. As we increase the model size and language coverage, the learnt codebook module proves costly to tune and can impede progress of model development. Meanwhile, the BEST-RQ algorithm does not require such a module, making it a more scalable method for pre-training.

2.2.1 Multi-softmax

Instead of utilizing a single codebook [10], we use multiple codebooks to improve BEST-RQ training in this study. More precisely, we use N softmax layers to produce N probability predictions from the output of the encoder to compare against N independent quantization targets obtained from the masked speech features. We train the network with equal weights for each softmax layer. The use of multiple codebooks improves the stability and convergence of the model.

2.3 Self-training: Noisy Student Training

We utilize noisy student training (NST) [7, 8] to generate pseudo-labeled data to augment supervised training. This is done by first training a teacher model with augmentation on a supervised set, then using that teacher to generate transcripts for unlabeled audio data. A heuristic filtering method based on the ratio between the number of words and audio length is used to filter the pseudo-labeled data. The pseudo-labeled data is mixed with supervised data to train the student model.

2.4 Chunk-wise Attention for Long-form ASR

In many real-world applications, ASR systems are required to transcribe minutes- or hours-long audio. This poses significant challenges to many end-to-end ASR systems, as these ASR systems

are usually trained on much shorter segments, typically less than 30 seconds. For systems that use attention-based encoders, it is impractical to use global attention to attend to the entire audio. Local self attention, which only attends to the fixed length of left and right context, is thus widely used. For example, in BEST-RQ pre-training, only 128 left and 128 right context frames are used for local self attention. However, stacking many local self attention layers creates a significant receptive field mismatch between training and inference. The left figure in Fig. 4 illustrates this issue with a network consisting of 4 local self attention layers, each using only 1 left and 1 right context frames. Since the context is leaked in every layer, the receptive field width grows linearly with respect to the number of layers; for a big encoder like that of the Conformer-2B, this means that the receptive field width for the encoder output is longer than 327 seconds. During training, the model is trained with at most 30 seconds speech segments, while at inference time, when minutes or hours long audio is fed to the model, the encoder needs to process over 300 seconds of audio to produce one encoder output—a pattern it has never trained on. Our empirical observations demonstrate that, under this train-test mismatch, these models with deep architectures and high capacity suffer from high deletion errors. We henceforth refer to this problem as the “long-form (performance) degradation” problem.

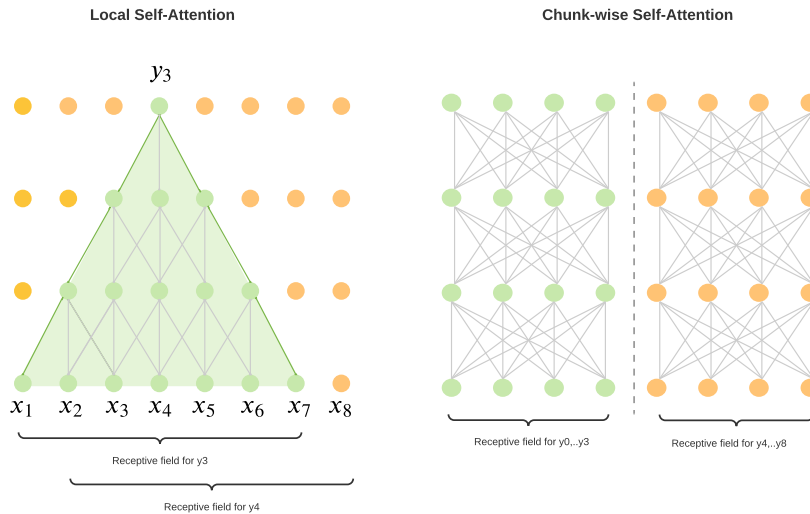


Figure 4: Comparing receptive fields of two networks with 4 layers of local self attention and chunk-wise attention.

To solve this problem, we propose a simple modification to the attention mechanism; the attention is restricted to audio chunks. This is illustrated on the right side of Fig. 4, in which 8 frames are divided into 2 chunks, and the attention is performed within each chunk. In this case, there is no context leaking in the attention layer, and thus the receptive field width is independent of the number of layers. In our experiments an 8-second chunk resulted in the best recognition quality vs. computational cost trade-off.

It is worthwhile to note there are a few other works in the literature which also modify the attention pattern to deal with the long-form audio in ASR, e.g., [61–66]. Though conceptually similar to block processing (e.g. [65, 66]), chunk-wise attention is more flexible. Block processing is performed at the input feature level, which limits the encoder layers to the context frame at the current chunk. On the other hand, chunk-wise attention allows other layers in the encoder (e.g., convolution layers) to process contextual frames beyond the current chunk. Compared with Whisper [1], which segments the audio into 30 second chunks and uses a heuristic process to carry the decoder states over, we only chunk the attention state, and allow the decoder to access the entire encoder output. We also use either a CTC or RNN-T decoder to decode on long-form audio, neither of which have been observed to hallucinate compared to attention-based sequence-to-sequence decoders. We observe our systems are robust on long-form ASR tasks with a simpler decoding process on long-form speech signals.

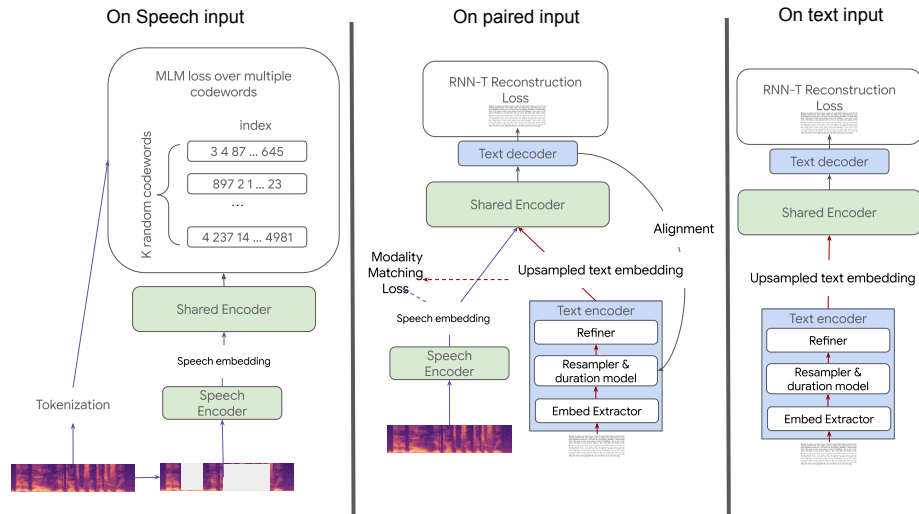


Figure 5: Overview of MOST text injection. The left-most panel depicts MOST training on unlabeled speech input; the center panel depicts training on paired speech and text input; the right-most panel depicts training on unlabeled text data.

2.5 Multi-Objective Supervised Pre-training: BEST-RQ + text-injection

In addition to pre-training with unlabeled speech, we add an additional stage of **Multi-Objective Supervised pre-Training (MOST)** as shown in Fig. 5, where we train the model jointly on unlabeled speech, unlabeled text and paired speech and text data. The training loss for this procedure is based on the text-injection loss including duration modeling and consistency regularization as in [13], to which we add a weighted BEST-RQ loss for the encoder of the model. MOST yields two benefits: (i) Training with paired speech and text data with alignment losses results in learning speech representations that are better aligned with text, improving quality on tasks like ASR and AST that require mapping the acoustics of the speech signal to text. (ii) Training simultaneously on unlabeled text in a model that learns speech and text representations jointly improves the robustness of learned representations, especially on low resource languages and domains, also generalizing to new languages with no paired data seen during training [67].

The key architectural components for constructing the text-injection loss as utilized in our approach include: (i) A speech-only encoder that utilizes a convolutional sub-sampling feature encoder and a single conformer layer. For continued pre-training the feature encoder is initialized from the BEST-RQ pre-trained checkpoint while the conformer layer is initialized randomly. (ii) A text-only encoder that consists of an embedding layer, an upsampler, and a conformer layer block. The upsampler used in this work is a learned duration based upsampling model [13], though a fixed or random repetition upsampler can also be used for text-injection [47, 53]. All components are initialized randomly. (iii) A shared conformer encoder initialized from the pre-trained BEST-RQ speech encoder. (iv) The BEST-RQ speech softmax layers initialized from the BEST-RQ checkpoint. (v) The decoder unit which is initialized randomly.

The main idea of text-injection (e.g. [13, 53, 54]) is to produce joint, co-aligned embeddings of speech and text as sequences in the same embedding space. Given this embedding space, text data with no associated audio can contribute to improving the speech task. The speech and text encoders presented above are intended to produce these embeddings, which need to be matched in the embedding space and are also required to be co-aligned in the time dimension. The embeddings enable the text data to contribute to preparing the model for downstream tasks.

To achieve these objectives, the architecture as presented above is trained using three types of data, each contributing to different types of losses:

1. The unlabeled speech passes through the shared encoder and the BEST-RQ softmax layers to contribute to the BEST-RQ loss.

2. The paired speech-text data serves multiple functions.
 - The labeled speech flows through the speech encoder, the shared encoder and the decoder unit and contributes to the standard ASR loss computed against the paired text. Here, the speech-text alignments of the paired data are extracted from the decoder unit and used to train the duration upsampler within the text encoder.
 - The text of the paired data also passes through the text encoder. The encoded text sequence is used to compute a consistency loss against the encoded speech sequence. This loss is used to train solely the text encoder—the speech encoder weights are frozen for this particular forward-propagation.
3. The unlabeled text data contributes to a reconstruction loss. This loss is constructed by passing the text through the text encoder, then masking chunks of the feature sequence produced. These masked text features live in the same embedding space as masked speech features, and thus can be passed through the shared encoder and the decoder unit to compute the ASR loss against the original text. This is the reconstruction loss used to train the model.

For training stability, MOST proceeds in two stages—we first train solely on paired data to learn stable decoder alignments for 20k steps. We then train the duration upsampler and activate the losses for unlabeled text. We refer the reader to [13] for further details.

When fine-tuning for ASR, we initialize the feature encoder of the ASR model with the speech feature encoder, initialize the conformer block with the shared conformer encoder, and add a randomly initialized task-specific transducer.

In the MOST set-up, the speech and text representations live in a shared representation space, thereby allowing us to utilize text machine translation (MT) data during the fine-tuning stage of AST tasks. We follow the same approach described in [13, 20] and report the AST results with joint fine-tuning for models prepared with MOST.

2.6 Residual Adaptation with a Frozen Encoder

Ideally, the fine-tuning process of the model should be scalable with the number of downstream tasks while in reality, fine-tuning the pre-trained USM individually for various domains and tasks becomes prohibitively expensive. In order to mitigate this issue, we explore a lightweight alternative [19] to training the full network where residual adapters with a small number of parameters are added for each individual language while the pre-trained USM is entirely frozen during fine-tuning. We experiment with adding two parallel adapters to each Conformer block, whose parameter count amounts to 2% of the original pre-trained USM, and fine-tune the adapters on downstream language tasks. When serving the model, the adapter is dynamically loaded according to the language of the input batch [68, 69]. This enables one to conduct inference on 100+ languages while keeping the total number of parameters manageable by re-using the same parameters and computation process for the majority of the time. We also find that training the adapter versus fine-tuning the entire model can reduce over-fitting especially when the training data is limited.

2.7 Training Details

Data Processing: The audio is uniformly sampled to 16 kHz quality—any audio with a different native sampling rate is either up-sampled or down-sampled. The audio is then featurized into 128-dimensional log-mel filterbank coefficients. Graphemes are used to tokenize the text for FLEURS in-domain fine-tuning, while word-piece models (WPMs) [70] are used for tokenization for all other tasks.

BEST-RQ: We follow default masking and quantization parameters of BEST-RQ as in [10]. We use a 16 codebook multi-softmax loss to stabilize training and improve performance as described in 5.1. We do not use EMA for pre-training.

MOST: We follow the text encoder and decoder architecture described in [13] but use 4k sentence-piece models (SPMs). We use a single 1536-dimensional Conformer layer as the speech encoder and Conformer-2B encoder as the shared encoder. We mix un-transcribed speech, unspoken text, and transcribed speech in each batch with fixed batch sizes of, respectively, 4096, 8192, and 1024. The model is initialized with the BEST-RQ pre-trained encoder. MOST employs a curriculum learning

schedule where training initially is conducted with un-transcribed speech and paired speech-text data, and unspoken text is utilized only after 20k steps. The joint training employing all three types of data lasts for another 100K steps.

Supervised Training: We use two separate optimizers for the encoder parameters and the decoder parameters of the network [71]. For USM-CTC and USM-LAS, we train the model for 100k steps with 2048 batch size. For in-domain experiments, the checkpoint is selected based on development set performance.

Training Large Models: We use the GShard [72] framework with the GSPMD backend [73] to train our large models on TPUs.

3 Datasets

3.1 Audio Data

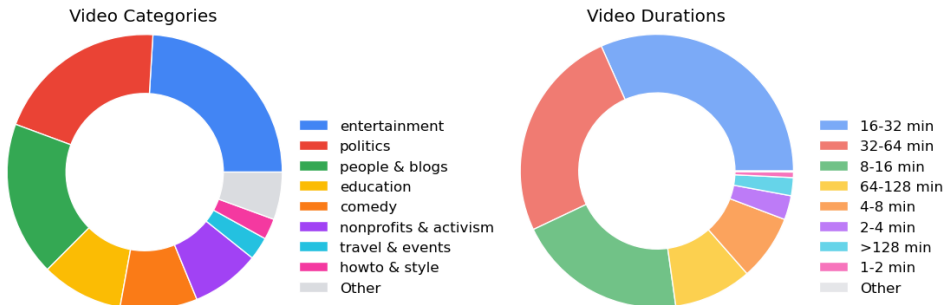


Figure 6: The video category and length distribution of YT-513-U.

The following audio datasets are used in this report to train our models:

- **YouTube SUPervised Plus (YT-SUP+):**
 - YT-SUP: 90k hours of segmented, labeled audio across 75 languages.
 - YT-Pseudo-Labeled: 100k hours of segmented, pseudo-labeled en-US audio from YT-NTL-U. The pseudo-labels are generated by a 600M CTC model trained on YT-SUP en-US data.
- **YouTube Next Thousand Languages Unsupervised (YT-NTL-U):** 12.1M hours of segmented, unlabeled audio, including:
 - YT-55-U: 12M hours of segmented, unlabeled audio on 55 rich resource languages identified by YouTube production language id models.
 - YT-513-U: 100k hours of segmented, unlabeled audio across 513 tail languages not covered by YouTube production language id models. These languages are identified by vendors.

Let us expand upon how each dataset has been constructed.

YT-SUP+: YT-SUP is a dataset with audio from videos that have user-uploaded transcripts from 75 languages. We group consecutive segments into a longer unit similar to [57]. The maximal sequence length for training is 30 seconds. The total amount of training data is 90k hours, ranging from English (en-US) (3.5k hours) to Amharic (Am-Et) (150 hours). We also introduce an additional 100k hours of en-US audio from YT-NTL-U to YT-SUP. We choose to generate pseudo-labels on this dataset using a 600M-parameter CTC YT teacher model trained on YT-SUP. Each audio is randomly segmented between 5 to 15 seconds.

YT-55-U: YT-55-U is built by first randomly collecting 3 million hours of audio from "speech-heavy" YouTube videos, filtered by language. The 3 million hours of audio is then further segmented by the YT teacher model. Instead of using a teacher model as in [3], the non-speech segments identified by a Voice Activity Detection (VAD) model are removed to yield approximately 1 million hours of

unlabeled audio data. Later, we use a YouTube production language identification model to select 55 languages from that audio.

YT-513-U: We create an additional dataset called YT-513-U to ensure coverage of lower resource languages in our pre-training dataset. We reached out to vendors and native speakers to identify YT videos containing speech in specific long tail languages, collecting a dataset of unlabeled speech in 513 languages. Vendors were tasked with ensuring a variety of domains, voices, and content in the videos that are collected in each language. These videos are segmented into speech segments using a VAD model, resulting in a total of 102k hours of speech. Our final YT-513-U dataset contains 88 languages with over 500 hours of speech each, 237 languages with between 100-500 hours, and 188 languages with less than 100 hours of data. The languages chosen for this collection are wide-ranging, with a majority of our data corresponding to languages from South Asia, Southeast Asia, West Africa, and East Africa. The distribution of video categories and lengths in our dataset are depicted in Figure 6.

In addition to YouTube data, we also include public data for **MOST** training:

- **Public Unsupervised (Pub-U):** Following [20], we use approximately 429k hours of unlabeled speech data in 51 languages. It includes: 372k hours of speech data spanning 23 languages from VoxPopuli [74], read speech data in 25 languages drawn from the v6.1 release of Common Voice [75], 50k hours of read books data in eight European languages from Multilingual LibriSpeech [76] and 1k hours of telephonic conversation data spanning 17 African and Asian languages from BABEL [77].
- **Public Supervised (Pub-S):** Similar to [20], our public supervised set includes approximately 1.3k hours of speech and transcript data spanning 14 languages from VoxPopuli, 10 hour training splits for each of the 8 MLS languages, and 1k hours of data spanning 17 languages from the Babel ASR task.

Note that the public data is only used for in-domain pre-training and is excluded for training the generic USM-LAS/CTC models. This allows us to treat the public task performance as out-of-domain benchmarks for the USM-LAS/CTC models.

3.2 Text Data

Web-NTL: For pre-training with unlabeled text, we use a web-crawled corpus of monolingual text containing over 28B sentences [78]. The dataset spans 1140 languages, 205 of which have over 1M sentences and 199 of which have between 100k and 1M sentences. We up-sample lower resource languages using temperature-based sampling [79] with $T = 3.0$. More details about the dataset and the mining approach have been described in Section 2 of [78].

3.3 Downstream Benchmarks

3.3.1 Speech Recognition (ASR)

We present our results on two public tasks, SpeechStew [2] and FLEURS [16], and an internal benchmark on YouTube.

The **SpeechStew** [2] dataset is assembled by putting together seven public speech corpora—AMI [80], Common Voice [81], English Broadcast News³, LibriSpeech [82], Switchboard/Fisher⁴, TED-LIUM v3 [83, 84] and Wall Street Journal⁵, which are all standard benchmarks [85–87] covering different domains in en-US.

The **FLEURS** [16] dataset is a publicly available, multi-way parallel dataset of 10 hours of read speech in 102 languages spanning 7 geo-groups. We restrict our use of the dataset to its ASR benchmark. Among the 102 languages present in the FLEURS benchmark, we select 62 to serve as a sub-group to compare our generic ASR system with Whisper [1], as those languages are covered by the training sets of both models. We also report full results for in-domain fine-tuning and adaptation. Unlike [16], we report both WER and CER metrics, as CER is inappropriate as an indicator of

³Linguistic data consortium (LDC) datasets LDC97S44, LDC97T22, LDC98S71 and LDC98T28.

⁴LDC datasets LDC2004T19, LDC2005T19, LDC2004S13, LDC2005S13 and LDC97S62.

⁵LDC datasets LDC93S6B and LDC94S13B.

Table 3: WERs (%) across multiple tasks for multiple settings compared against pre-existing baselines, with the exception of CoVoST 2, for which the BLEU score is presented. For the YouTube long-form set, we select the top-25 languages Whisper was trained on and exclude all languages for which Whisper produces > 40% WER to reduce the noise introduced by LAS hallucination in the Whisper model. For FLEURS, we report both the WER and the CER for our models. [†]Results omitted for the Whisper-shortform model on the YouTube long-form dataset as the model has a high deletion problem on this set. [‡]The Whisper-shortform model uses segmented decoding to reduce its hallucination problem on CORAAL. [§]Our adapter setup adds about 2.3% of the total parameters while keeping the encoder frozen from pre-training.

Task	Multilingual Long-form ASR			Multidomain en-US	Multilingual ASR		AST	
	Dataset	YouTube		CORAAL	SpeechStew	FLEURS		CoVoST 2
Langauges	en-US	18	73	en-US	en-US	62	102	21
Prior Work (single model)								
Whisper-longform	17.7	27.8	-	23.9	12.8	-	-	-
Whisper-shortform [†]	-	-	-	13.2 [‡]	11.5	36.6	-	29.1
Our Work (single model)								
USM-LAS	14.4	19.0	29.8	11.2	10.5	12.5	-	-
USM-CTC	13.7	18.7	26.7	12.1	10.8	15.5	-	-
Prior Work (in-domain fine-tuning)								
BigSSL [3]	14.8	-	-	-	7.5	-	-	-
Maestro [67]	-	-	-	-	7.2	-	-	25.2
Maestro-U [67]	-	-	-	-	-	-	26.0 (8.7)	-
Our Work (in-domain fine-tuning)								
USM	13.2	-	-	-	7.4	13.5	19.2 (6.9)	28.7
USM-M	12.5	-	-	-	7.0	11.8	17.4 (6.5)	30.7
Our Work (frozen encoder)								
USM-M-adapter [§]	-	-	-	-	7.5	12.4	17.6 (6.7)	29.6

performance for some languages. When presenting the error rate metrics, we use CER for Chinese, Japanese, Thai, Lao, and Burmese to be consistent with Whisper [1].

The test set for the **YouTube** domain consists of utterances from 73 languages with an average of 15 hours of audio per language, the audio length for each individual language ranging from 1 to 24 hours. The audio is transcribed manually from popular YouTube videos, each with a duration of up to 30 minutes.

3.3.2 Speech Translation (AST)

Following [20], we use CoVoST 2 [18] to benchmark multilingual speech translation. We evaluate the multilingual XX-to-English task that covers translation from 21 source languages into English. Depending on the language, the training data ranges in size from 1 - 264 hours.

Besides speech translation data, we also add text-to-text translation data for training the model as in [20]. This dataset includes the text translation data from CoVoST 2 combined with all data from either WMT or TED Talks, as available.

4 Key Results

4.1 Robust Speech Recognition for Massively Multilingual Tasks

In this section, we compare the performance of our models against public baselines, including Whisper large-v2⁶ [1], which has been trained on 680k hours of weakly supervised data across 100 languages.

For the massively multilingual speech recognition test dataset from YouTube, we observe that Whisper hallucinates in many languages, resulting in a WER exceeding 100%. For a reasonable comparison, we restrict the language set on which we compare the performance USM against Whisper by first selecting the top-25 languages from the training data for Whisper and further excluding languages for which Whisper produces > 40% WER. We also use segmented decoding for Whisper with 30-second segments to further reduce the effect of hallucinations. As shown in Table 3, our USM-LAS and

⁶Whisper large-v2 on Github (<https://github.com/openai/whisper.git>, revision b4308c4) is used for evaluation.

USM-CTC models outperform Whisper by a wide margin on YouTube en-US, despite training on significantly less supervised data (3.5k hours versus Whisper’s 400k hours [1]). While the USM-LAS model also requires segmented decoding to reduce long-form degradation as discussed in section 2.4, it is far more robust, out-performing Whisper by a relative 30% WER on those 18 languages. USM-CTC does not exhibit long-form performance degradation and achieves the best performance on YouTube.

On the out-of-domain long-form CORAAL set, both USM-CTC and USM-LAS outperform Whisper by more than 10% relative WER. USM-CTC and USM-LAS similarly outperform Whisper on SpeechStew, whose training data the models have not had access to.

We further compare the multilingual performance of the models on the held-out set from FLEURS. As shown in Table 3, USM-LAS and USM-CTC both outperform Whisper by 66% relative WER, despite using a smaller amount of multilingual supervised data (90k versus Whisper’s 117k, when en-US is excluded). USM-LAS consistently outperforms USM-CTC for short-form ASR tasks.

4.2 Massively Multilingual Results Beyond 100 Languages

The lower part of Table 3 shows our results for in-domain fine-tuning. Our pre-trained model improves the FLEURS benchmark significantly, even when using only 10 hours per language. Compared to the previous SoTA in [67], our model achieves a 30% relative improvement in terms of WER across 102 languages. Our results show that while generic speech models can be powerful, performance is still maximized by in-domain fine-tuning.

4.3 MOST Produces Robust Representations that Generalize to New Domains

MOST training aligns the representations of speech and text by training simultaneously on the two modalities. We investigate whether MOST representations are useful for adapting the model to new domains by freezing the entire learned encoder produced by MOST and adjusting a small amount of parameters added to the network by residual adapters. As shown in Table 3, by adding only 2% to the total number of parameters, the MOST representation model (USM-M-adapter) only performs slightly worse than the fine-tuning baselines, still showing competitive performance on downstream ASR and AST tasks. The small number of parameters being trained in this approach makes it feasible to extend our system to a large number of new domains and new tasks, even with a limited amount of training data, such as in FLEURS.

4.4 Pushing the Quality of ASR on Unseen Languages

Table 4: Noisy student training for unseen languages. WERs (%) for the teacher adapter models and the student models are presented. The relative improvement (%) of the student models can be found in the last column.

Languages	Whisper-v2	# hrs in YT-NTL	USM-LAS-Adapter	USM-M + pseudo label	Rel. Imprv.
Hausa (ha)	88.9	2175.0	24.5	22.8	7.5
Kazakh (kk)	37.7	196.0	11.8	10.9	8.3
Shona (sn)	121.0	247.0	29.1	22.2	31.1
Pashto (ps)	93.7	254.0	36.0	35.4	1.7
Yoruba (yo)	94.8	1292.0	33.4	30.6	9.2

Tail languages often do not have paired transcriptions for supervised learning—we refer to these languages as unseen languages, as the model has not seen paired data for these languages during training. To create pseudo-labels for these languages, we first build a USM-LAS-Adapter by attaching residual adapters to USM-LAS and training them using FLEURS data. By using the USM-LAS-Adapter as a teacher, we can now transcribe the unlabeled data in the unseen languages as part of the YT-NTL dataset. As shown in Table 4, we observe consistent wins for all languages on the FLEURS benchmark. For some languages, the improvement is larger than 30%. This further demonstrates the robustness of the USM-LAS model—despite using only 10 hours of out of domain data from FLEURS, the USM-LAS-Adapter is able to transcribe YouTube data to produce meaningful recognition results that lead to these improvements. We find the approach of training adapter models

on small datasets and utilizing them for pseudo-labeling to be a promising route for scaling up the number of languages that can be transcribed by USMs.

4.5 USMs are Strong AST Models

The multi-lingual speech translation performance of fine-tuned USMs are shown in Table 3. We find that we are already comparable to the CoVoST 2 SoTA BLEU score by fine-tuning the speech-only USM. We note that the previous SoTA uses 125k hours of supervised speech translation data compared to the 859 hours of data used by the USM. After MOST training, USM-M can use both speech and text as training input. By introducing text-to-text machine translation (MT) data during fine-tuning, USM-M is able to achieve an unprecedented > 30 BLEU on CoVoST (a 1 BLEU increase from SoTA).

5 Analysis and Ablations

5.1 Multi-Softmax Loss for BEST-RQ

We observe a consistent $> 5\%$ relative improvement in ASR and AST benchmarks by increasing the number of the softmax groups in the multi-softmax loss for BEST-RQ training from 1 to 16, as shown in Table 5. We also find that using multiple softmax groups significantly reduces performance variation across different pre-training runs and improves convergence speed.

Table 5: YT-55 versus YT-NTL across different domains, with and without multi-softmax groups. For simplicity, we report CER for FLEURS. For CoVoST, we report the BLEU score. YT-NTL covers 27 additional languages not covered in YT-55.

Model	pre-train Set	# Params (B)	# Softmax	FLEURS (CER)		CoVoST (BLEU)
				102 langs	27 langs	
Conformer-0.6B	YT-55	0.6	1	9.5	-	20.9
Conformer-2B	YT-55	2.0	1	7.9	9.5	26.6
Conformer-2B	YT-NTL-U	2.0	1	7.4	8.5	27.5
Conformer-2B	YT-NTL-U	2.0	16	6.9	8.1	28.7

5.2 Model and Language Scaling

We find that scaling up the model size and increasing the language coverage of the pre-training dataset greatly benefits the performance of the USMs, as demonstrated in Table 5. In particular, we find a 10% relative improvement of ASR and AST performance by using YT-NTL vs. YT-55 for pre-training, despite the fact that each newly added language in YT-NTL contains approximately 500 hours of speech—a relatively small amount. As could be expected, the relative gains on the newly covered languages are more substantial than those on other languages.

5.3 BEST-RQ is a Scalable Self-supervised Learner

BEST-RQ has been shown to outperform or be comparable to other prominent pre-training methods for speech recognition, including wav2vec 2.0 and W2v-BERT in the original work in which it was introduced [10]. Here we investigate its comparative performance and scaling properties, similar to what has been done for wav2vec 2.0 in [3] and W2v-BERT in [20]. We utilize the set-up of pre-training the model using YT-55 and fine-tuning it on CoVoST 2. As shown in Table 6, our results indicate that for the Conformer-0.6B, W2v-BERT and BEST-RQ perform similarly, but BEST-RQ obtains greater gains when scaled up. A contributing factor to this can be that W2v-BERT is more prone to codebook collapse and training instabilities at the 2B scale, while BEST-RQ by construction doesn't suffer from codebook collapse.

5.4 Chunk-wise attention for robust long-form speech recognition

Fig. 7 depicts the long-form performance degradation issue as described in section 2.4. In the figure, we see that for the shallow Conformer model with 17 layers, using a small local self attention context

Table 6: BLEU scores for the CoVoST 2 X → En task to compare BEST-RQ against W2v-BERT. Higher is better.

X → English	high	mid	low	all
Previous Work				
XLS-R (0.3B) [33]	30.6	18.9	5.1	13.2
XLS-R (1B) [33]	34.3	25.5	11.7	19.3
XLS-R (2B) [33]	36.1	27.7	15.1	22.1
Conformer-0.6B				
W2v-BERT	35.6	25.3	13.4	20.4
BEST-RQ	32.5	25.6	14.7	20.7
Conformer-2B				
W2v-BERT	36.0	27.8	15.6	22.4
BEST-RQ	35.8	31.3	21.5	26.6

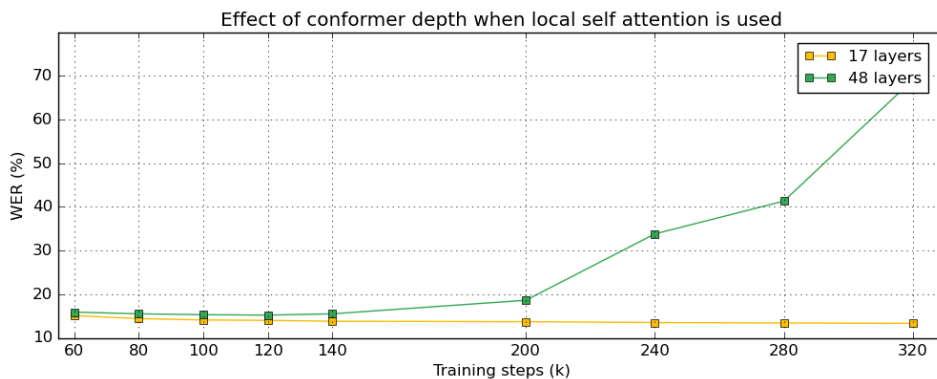


Figure 7: The word error rate measured on the YouTube en-US long-form test set for Conformer models with varying depth.

(65) length, the word error rate measured on the long-form test set gradually improves as the training progresses. With a deeper model that has 48 layers but roughly the same number of parameters, however, the larger receptive field mismatch results in higher test WERs as the training step increases.

Table 7 demonstrates that chunk-wise attention is able to address the long-form degradation issue and show robust performance across four different languages—en-US (English), ru-RU (Russian), ko-KR (Korean), and uk-UA (Ukrainian). We compare chunk-wise attention models with an 8-second chunk size (CW-8s in Table 7) against local self attention models which uses 128 context frames in each conformer layer (LSA-128). We note that further increasing the context window size of the local self attention model results in high deletion error rates on all languages of the YouTube long-form test sets. These results show that the chunk-wise attention models do not exhibit long-form performance degradation and are able to improve upon the performance of the local self attention models operating at the maximum allowed receptive field length.

Table 7: Chunk-wise attention. WER (%) is reported on the YouTube long-form set.

Model	# Params (B)	# Layers	en-US	ru-RU	ko-KR	uk-UA
LSA-128	0.6	24	16.2	16.6	26.2	15.5
CW-8s	0.6	24	12.5	14.7	19.5	15.3

5.5 TPU Serving Capacity of USM-CTC Models

In section 4, we have demonstrated that USM-CTC models are powerful generic ASR models with reliable long-form transcription performance and excellent generalization properties. Here we

Table 8: RTF for USM-2B.

Model	bf-16	Streaming	# Params (B)	TPU [88]	Batch Size	1.0/RTF
Conformer-0.1B	Y	Y	0.1	TPUv4i	64	3047
Conformer-0.6B	N	N	0.6	TPUv4i	64	1920
Conformer-2B	N	N	2.0	TPUv4i	32	827

measure the serving capacity of the USM-CTC model as represented by the real time factor (RTF) in an ideal setup where we assume that each batch sent to TPU is fully packed along the time axis. The results of these measurements are presented in Table 8. Surprisingly, we find that the 2B-parameter USM-CTC model is only $3.9\times$ slower than the 100M-parameter streaming model [89], primarily due to the fact that our models operate at batch processing mode. This result demonstrates that the USM-CTC can be used as an offline transcriber efficiently on TPUs (or GPUs).

6 Discussion

In this report, we put forward a practical and flexible approach for training speech understanding models capable of scaling speech recognition to hundreds of languages. We conclude the report with summarizing insights gained in the process:

Unlabeled versus weakly labeled data: We believe diverse unlabeled data is more practical to acquire for building usable ASR for tail languages than weakly labeled data. We have demonstrated that collaborating with native speakers to identify unsupervised data in hundreds of tail languages can be an effective route to improving recognition performance on low resource languages.

In-domain data is best: We have demonstrated that we can build a robust ASR system across many domains by utilizing a large amount of unsupervised data and a small amount of labeled data. Our results, however, also confirm that the most effective way to optimize the performance for a given domain is to use in-domain data to fine-tune the model.

CTC vs RNN-T vs LAS: The best transducer depends on the downstream task. A large pre-trained model with a frozen encoder can allow experimenters to test different transducers quickly and select the optimal transducer for their purpose.

Acknowledgments

We would like to thank Alexis Conneau, Min Ma, Shikhar Bharadwaj, Sid Dalmia, Jiahui Yu, Jian Cheng, Paul Rubenstein, Ye Jia, Justin Snyder, Vincent Tsang, Yuanzhong Xu, Tao Wang, Anusha Ramesh, Calum Barnes, Salem Haykal for useful discussions.

We appreciate valuable feedback and support from Eli Collins, Jeff Dean, Sissie Hsiao, Zoubin Ghahramani. Special thanks to Austin Tarango, Lara Tumeh, and Jason Porta for their guidance around responsible AI practices.

References

- [1] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, “Robust speech recognition via large-scale weak supervision,” *arXiv preprint arXiv:2212.04356*, 2022.
- [2] W. Chan, D. Park, C. Lee, Y. Zhang, Q. Le, and M. Norouzi, “Speechstew: Simply mix all available speech recognition data to train one large neural network,” *arXiv preprint arXiv:2104.02133*, 2021.
- [3] Y. Zhang, D. S. Park, W. Han, J. Qin, A. Gulati, J. Shor, A. Jansen, Y. Xu, Y. Huang, S. Wang, Z. Zhou, B. Li, M. Ma, W. Chan, J. Yu, Y. Wang, L. Cao, K. C. Sim, B. Ramabhadran, T. N. Sainath, F. Beaufays, Z. Chen, Q. V. Le, C.-C. Chiu, R. Pang, and Y. Wu, “Bigssl: Exploring the frontier of large-scale semi-supervised learning for automatic speech recognition,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1519–1532, 2022.
- [4] B. Li, R. Pang, T. N. Sainath, A. Gulati, Y. Zhang, J. Qin, P. Haghani, W. R. Huang, and M. Ma, “Scaling end-to-end models for large-scale multilingual asr,” *arXiv preprint arXiv:2104.14830*, 2021.

- [5] X. Li, F. Metze, D. R. Mortensen, A. W. Black, and S. Watanabe, “Asr2k: Speech recognition for around 2000 languages without audio,” *arXiv preprint arXiv:2209.02842*, 2022.
- [6] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *arXiv preprint arXiv:2006.11477*, 2020.
- [7] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, “Self-training with noisy student improves imagenet classification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 687–10 698.
- [8] 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,” *arXiv preprint arXiv:2005.09629*, 2020.
- [9] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu *et al.*, “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.
- [10] C.-C. Chiu, J. Qin, Y. Zhang, J. Yu, and Y. Wu, “Self-supervised learning with random-projection quantizer for speech recognition,” in *Proceedings of the 39th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 17–23 Jul 2022, pp. 3915–3924. [Online]. Available: <https://proceedings.mlr.press/v162/chiu22a.html>
- [11] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [12] Z. Chen, Y. Zhang, A. Rosenberg, B. Ramabhadran, G. Wang, and P. Moreno, “Injecting text in self-supervised speech pretraining,” *arXiv preprint arXiv:2108.12226*, 2021.
- [13] Z. Chen, Y. Zhang, A. Rosenberg, B. Ramabhadran, P. Moreno, A. Bapna, and H. Zen, “Maestro: Matched speech text representations through modality matching,” *arXiv preprint arXiv:2204.03409*, 2022.
- [14] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [15] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2016, pp. 4960–4964.
- [16] A. Conneau, M. Ma, S. Khanuja, Y. Zhang, V. Axelrod, S. Dalmia, J. Riesa, C. Rivera, and A. Bapna, “Fleurs: Few-shot learning evaluation of universal representations of speech,” *arXiv preprint arXiv:2205.12446*, 2022.
- [17] T. Kendall and C. Farrington, “The corpus of regional african american language. version 2021.07. eugene, or: The online resources for african american language project,” 2021.
- [18] C. Wang, A. Wu, and J. Pino, “CoVoST 2 and massively multilingual speech-to-text translation,” in *interspeech*, 2021.
- [19] J. He, C. Zhou, X. Ma, T. Berg-Kirkpatrick, and G. Neubig, “Towards a unified view of parameter-efficient transfer learning,” in *International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=0RDcd5Axok>
- [20] A. Bapna, C. Cherry, Y. Zhang, Y. Jia, M. Johnson, Y. Cheng, S. Khanuja, J. Riesa, and A. Conneau, “mslam: Massively multilingual joint pre-training for speech and text,” *arXiv preprint arXiv:2202.01374*, 2022.
- [21] Y.-A. Chung, Y. Zhang, W. Han, C.-C. Chiu, J. Qin, R. Pang, and Y. Wu, “W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training,” in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2021, pp. 244–250.
- [22] W.-N. Hsu and J. Glass, “Extracting domain invariant features by unsupervised learning for robust automatic speech recognition,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5614–5618.
- [23] Y.-A. Chung and J. Glass, “Speech2vec: A sequence-to-sequence framework for learning word embeddings from speech,” *arXiv preprint arXiv:1803.08976*, 2018.
- [24] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018.
- [25] Y.-A. Chung, W.-N. Hsu, H. Tang, and J. Glass, “An unsupervised autoregressive model for speech representation learning,” *arXiv preprint arXiv:1904.03240*, 2019.
- [26] J. Chorowski, R. J. Weiss, S. Bengio, and A. van den Oord, “Unsupervised speech representation learning using wavenet autoencoders,” *IEEE/ACM transactions on audio, speech, and language processing*, vol. 27, no. 12, pp. 2041–2053, 2019.

- [27] S. Schneider, A. Baevski, R. Collobert, and M. Auli, “wav2vec: Unsupervised pre-training for speech recognition,” *arXiv preprint arXiv:1904.05862*, 2019.
- [28] A. Baevski, S. Schneider, and M. Auli, “vq-wav2vec: Self-supervised learning of discrete speech representations,” *arXiv preprint arXiv:1910.05453*, 2019.
- [29] S. Ling, Y. Liu, J. Salazar, and K. Kirchhoff, “Deep contextualized acoustic representations for semi-supervised speech recognition,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6429–6433.
- [30] A. Baevski, M. Auli, and A. Mohamed, “Effectiveness of self-supervised pre-training for speech recognition,” *arXiv preprint arXiv:1911.03912*, 2019.
- [31] M. Riviere, A. Joulin, P.-E. Mazaré, and E. Dupoux, “Unsupervised pretraining transfers well across languages,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7414–7418.
- [32] K. Kawakami, L. Wang, C. Dyer, P. Blunsom, and A. v. d. Oord, “Learning robust and multilingual speech representations,” *arXiv preprint arXiv:2001.11128*, 2020.
- [33] A. Babu, C. Wang, A. Tjandra, K. Lakhotia, Q. Xu, N. Goyal, K. Singh, P. von Platen, Y. Saraf, J. Pino *et al.*, “Xls-r: Self-supervised cross-lingual speech representation learning at scale,” *arXiv preprint arXiv:2111.09296*, 2021.
- [34] G. Zavaliagkos and T. Colthurst, “Utilizing untranscribed training data to improve performance,” in *DARPA Broadcast News Transcription and Understanding Workshop, Landsdowne*, 1998, pp. 301–305.
- [35] L. Lamel, J. luc Gauvain, and G. Adda, “Lightly supervised acoustic model training,” in *Proc. ISCA ITRW ASR2000*, 2000, pp. 150–154.
- [36] S. Novotney and R. Schwartz, “Analysis of low-resource acoustic model self-training,” in *Tenth Annual Conference of the International Speech Communication Association*, 2009.
- [37] S. Thomas, M. L. Seltzer, K. Church, and H. Hermansky, “Deep neural network features and semi-supervised training for low resource speech recognition,” in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 6704–6708.
- [38] B. Li, T. N. Sainath, R. Pang, and Z. Wu, “Semi-supervised training for end-to-end models via weak distillation,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 2837–2841.
- [39] J. Kahn, A. Lee, and A. Hannun, “Self-training for end-to-end speech recognition,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7084–7088.
- [40] G. Synnaeve, Q. Xu, J. Kahn, T. Likhomanenko, E. Grave, V. Pratap, A. Sriram, V. Liptchinsky, and R. Collobert, “End-to-end asr: from supervised to semi-supervised learning with modern architectures,” in *arXiv*, 2019.
- [41] S. H. K. Parthasarathi and N. Strom, “Lessons from building acoustic models with a million hours of speech,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6670–6674.
- [42] W.-N. Hsu, A. Lee, G. Synnaeve, and A. Hannun, “Semi-supervised speech recognition via local prior matching,” *arXiv preprint arXiv:2002.10336*, 2020.
- [43] Q. Xu, T. Likhomanenko, J. Kahn, A. Hannun, G. Synnaeve, and R. Collobert, “Iterative pseudo-labeling for speech recognition,” *arXiv preprint arXiv:2005.09267*, 2020.
- [44] Z. Chen, A. Rosenberg, Y. Zhang, H. Zen, M. Ghods, Y. Huang, J. Emond, G. Wang, B. Ramabhadran, and P. J. Moreno, “Semi-Supervision in ASR: Sequential MixMatch and Factorized TTS-Based Augmentation,” in *Proc. Interspeech 2021*, 2021, pp. 736–740.
- [45] A. Renduchintala, S. Ding, M. Wiesner, and S. Watanabe, “Multi-modal data augmentation for end-to-end asr,” *arXiv preprint arXiv:1803.10299*, 2018.
- [46] A. Bapna, Y.-a. Chung, N. Wu, A. Gulati, Y. Jia, J. H. Clark, M. Johnson, J. Riesa, A. Conneau, and Y. Zhang, “Slam: A unified encoder for speech and language modeling via speech-text joint pre-training,” *arXiv preprint arXiv:2110.10329*, 2021.
- [47] S. Thomas, B. Kingsbury, G. Saon, and H.-K. J. Kuo, “Integrating text inputs for training and adapting rnn transducer asr models,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 8127–8131.
- [48] Y. Cheng, Y. Zhang, M. Johnson, W. Macherey, and A. Bapna, “Mu²slam: Multitask, multilingual speech and language models,” *arXiv preprint arXiv:2212.09553*, 2022.

- [49] Z.-H. Zhang, L. Zhou, J. Ao, S. Liu, L. Dai, J. Li, and F. Wei, "Speechut: Bridging speech and text with hidden-unit for encoder-decoder based speech-text pre-training," in *Conference on Empirical Methods in Natural Language Processing*, 2022.
- [50] Z.-H. Zhang, S. Chen, L. Zhou, Y. Wu, S. Ren, S. Liu, Z. Yao, X. Gong, L. Dai, J. Li, and F. Wei, "Speechlm: Enhanced speech pre-training with unpaired textual data," *ArXiv*, vol. abs/2209.15329, 2022.
- [51] S. Khurana, A. Laurent, and J. R. Glass, "Samu-xlsr: Semantically-aligned multimodal utterance-level cross-lingual speech representation," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, pp. 1493–1504, 2022.
- [52] X. Zhou, J. Wang, Z. Cui, S. Zhang, Z. Yan, J. Zhou, and C. Zhou, "Mmspeech: Multi-modal multi-task encoder-decoder pre-training for speech recognition," *ArXiv*, vol. abs/2212.00500, 2022.
- [53] T. N. Sainath, R. Prabhavalkar, A. Bapna, Y. Zhang, Z. Huo, Z. Chen, B. Li, W. Wang, and T. Strohmaier, "Joist: A joint speech and text streaming model for asr," in *2022 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2023, pp. 52–59.
- [54] Z. Meng, W. Wang, R. Prabhavalkar, T. N. Sainath, T. Chen, E. Variiani, Y. Zhang, B. Li, A. Rosenberg, and B. Ramabhadran, "Jeit: Joint end-to-end model and internal language model training for speech recognition," in *ICASSP, 2023*, 2023.
- [55] Z. Meng, T. Chen, R. Prabhavalkar, Y. Zhang, G. Wang, K. Audhkhasi, J. Emond, T. Strohmaier, B. Ramabhadran, W. R. Huang *et al.*, "Modular hybrid autoregressive transducer," in *2022 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2023, pp. 197–204.
- [56] C.-C. Chiu, W. Han, Y. Zhang, R. Pang, S. Kishchenko, P. Nguyen, A. Narayanan, H. Liao, S. Zhang, A. Kannan *et al.*, "A comparison of end-to-end models for long-form speech recognition," in *2019 IEEE automatic speech recognition and understanding workshop (ASRU)*. IEEE, 2019, pp. 889–896.
- [57] Z. Lu, Y. Pan, T. Dautre, P. Haghani, L. Cao, R. Prabhavalkar, C. Zhang, and T. Strohmaier, "Input length matters: Improving rnn-t and mwer training for long-form telephony speech recognition," *arXiv preprint arXiv:2110.03841*, 2021.
- [58] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov, "Transformer-xl: Attentive language models beyond a fixed-length context," *arXiv preprint arXiv:1901.02860*, 2019.
- [59] A. Graves, "Sequence transduction with recurrent neural networks," *arXiv preprint arXiv:1211.3711*, 2012.
- [60] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 5485–5551, 2020.
- [61] B. Ramabhadran, K. Audhkhasi, P. J. M. Mengibar, and T. Chen, "Mixture model attention: Flexible streaming and non-streaming automatic speech recognition," in *Proceedings of Interspeech, 2021*, 2021.
- [62] L. Lu, C. Liu, J. Li, and Y. Gong, "Exploring transformers for large-scale speech recognition," *arXiv preprint arXiv:2005.09684*, 2020.
- [63] X. Chen, Y. Wu, Z. Wang, S. Liu, and J. Li, "Developing real-time streaming transformer transducer for speech recognition on large-scale dataset," in *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 5904–5908.
- [64] C. Wu, Y. Wang, Y. Shi, C.-F. Yeh, and F. Zhang, "Streaming transformer-based acoustic models using self-attention with augmented memory," *arXiv preprint arXiv:2005.08042*, 2020.
- [65] Y. Shi, Y. Wang, C. Wu, C.-F. Yeh, J. Chan, F. Zhang, D. Le, and M. Seltzer, "Emformer: Efficient memory transformer based acoustic model for low latency streaming speech recognition," in *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6783–6787.
- [66] E. Tsunoo, Y. Kashiwagi, T. Kumakura, and S. Watanabe, "Transformer asr with contextual block processing," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 427–433.
- [67] Z. Chen, A. Bapna, A. Rosenberg, Y. Zhang, B. Ramabhadran, P. Moreno, and N. Chen, "Maestro: Leveraging joint speech-text representation learning for zero supervised speech asr," *arXiv preprint arXiv:2210.10027*, 2022.
- [68] F. Biadsy, Y. Chen, X. Zhang, O. Rybakov, A. Rosenberg, and P. J. Moreno, "A scalable model specialization framework for training and inference using submodels and its application to speech model personalization," in *Proc. Interspeech 2022*. ISCA, 2022, pp. 5125–5129.
- [69] K. Tomanek, V. Zayats, D. Padfield, K. Vaillancourt, and F. Biadsy, "Residual adapters for parameter-efficient asr adaptation to atypical and accented speech," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021, pp. 6751–6760.
- [70] M. Schuster and K. Nakajima, "Japanese and korean voice search," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012, pp. 5149–5152.

- [71] Y. Zhang, J. Qin, D. S. Park, W. Han, C.-C. Chiu, R. Pang, Q. V. Le, and Y. Wu, "Pushing the limits of semi-supervised learning for automatic speech recognition," *arXiv preprint arXiv:2010.10504*, 2020.
- [72] D. Lepikhin, H. Lee, Y. Xu, D. Chen, O. Firat, Y. Huang, M. Krikun, N. Shazeer, and Z. Chen, "GShard: Scaling giant models with conditional computation and automatic sharding," *CoRR*, vol. abs/2006.16668, 2020. [Online]. Available: <https://arxiv.org/abs/2006.16668>
- [73] Y. Xu, H. Lee, D. Chen, B. A. Hechtman, Y. Huang, R. Joshi, M. Krikun, D. Lepikhin, A. Ly, M. Maggioni, R. Pang, N. Shazeer, S. Wang, T. Wang, Y. Wu, and Z. Chen, "GSPMD: general and scalable parallelization for ML computation graphs," *CoRR*, vol. abs/2105.04663, 2021. [Online]. Available: <https://arxiv.org/abs/2105.04663>
- [74] C. Wang, M. Rivière, A. Lee, A. Wu, C. Talnikar, D. Haziza, M. Williamson, J. Pino, and E. Dupoux, "Voxpopuli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation," *arXiv preprint arXiv:2101.00390*, 2021.
- [75] R. Ardila, M. Branson, K. Davis, M. Henretty, M. Kohler, J. Meyer, R. Morais, L. Saunders, F. M. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," *arXiv preprint arXiv:1912.06670*, 2019.
- [76] V. Pratap, Q. Xu, A. Sriram, G. Synnaeve, and R. Collobert, "Mls: A large-scale multilingual dataset for speech research," *arXiv preprint arXiv:2012.03411*, 2020.
- [77] M. J. F. Gales, K. Knill, A. Ragni, and S. P. Rath, "Speech recognition and keyword spotting for low-resource languages: Babel project research at cued," in *SLTU*, 2014.
- [78] A. Bapna, I. Caswell, J. Kreutzer, O. Firat, D. van Esch, A. Siddhant, M. Niu, P. Baljekar, X. Garcia, W. Macherey *et al.*, "Building machine translation systems for the next thousand languages," *arXiv preprint arXiv:2205.03983*, 2022.
- [79] N. Arivazhagan, A. Bapna, O. Firat, D. Lepikhin, M. Johnson, M. Krikun, M. X. Chen, Y. Cao, G. Foster, C. Cherry *et al.*, "Massively multilingual neural machine translation in the wild: Findings and challenges," *arXiv preprint arXiv:1907.05019*, 2019.
- [80] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, M. Kronenthal *et al.*, "The ami meeting corpus: A pre-announcement," in *International workshop on machine learning for multimodal interaction*. Springer, 2005, pp. 28–39.
- [81] R. Ardila, M. Branson, K. Davis, M. Henretty, M. Kohler, J. Meyer, R. Morais, L. Saunders, F. M. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," *arXiv preprint arXiv:1912.06670*, 2019.
- [82] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [83] A. Rousseau, P. Deléglise, and Y. Esteve, "Ted-lium: an automatic speech recognition dedicated corpus." in *LREC*, 2012, pp. 125–129.
- [84] F. Hernandez, V. Nguyen, S. Ghannay, N. Tomashenko, and Y. Esteve, "Ted-lium 3: twice as much data and corpus repartition for experiments on speaker adaptation," in *International conference on speech and computer*. Springer, 2018, pp. 198–208.
- [85] J.-L. Gauvain, L. F. Lamel, G. Adda, and M. Adda-Decker, "The limsi continuous speech dictation system: evaluation on the arpa wall street journal task," in *Proceedings of ICASSP'94. IEEE International Conference on Acoustics, Speech and Signal Processing*, vol. 1. IEEE, 1994, pp. 1–557.
- [86] F. Kubala, J. Davenport, H. Jin, D. Liu, T. Leek, S. Matsoukas, D. Miller, L. Nguyen, F. Richardson, R. Schwartz *et al.*, "The 1997 bbn byblos system applied to broadcast news transcription," in *Proc. DARPA Broadcast News Transcription and Understanding Workshop*. Morgan Kaufmann, 1998, pp. 35–40.
- [87] S. Chen, M. Gales, P. Gopalakrishnan, R. Gopinath, H. Printz, D. Kanevsky, P. Olsen, and L. Polymenakos, "Ibm's lvsr system for transcription of broadcast news used in the 1997 hub4 english evaluation," in *Proceedings of the Speech Recognition Workshop*. Citeseer, 1998.
- [88] N. P. Jouppi, D. H. Yoon, M. Ashcraft, M. Gottscho, T. B. Jablin, G. Kurian, J. Laudon, S. Li, P. Ma, X. Ma *et al.*, "Ten lessons from three generations shaped google's tpuv4i: Industrial product," in *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2021, pp. 1–14.
- [89] B. Li, A. Gulati, J. Yu, T. N. Sainath, C.-C. Chiu, A. Narayanan, S.-Y. Chang, R. Pang, Y. He, J. Qin *et al.*, "A better and faster end-to-end model for streaming asr," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 5634–5638.