# REPRODUCING WHISPER-STYLE TRAINING USING AN OPEN-SOURCE TOOLKIT AND PUBLICLY AVAILABLE DATA

*Yifan Peng[1], Jinchuan Tian[1], Brian Yan[1], Dan Berrebbi[1], Xuankai Chang[1], Xinjian Li[1], Jiatong Shi[1], Siddhant Arora[1], William Chen[1], Roshan Sharma[1], Wangyou Zhang[1,2], Yui Sudo[3], Muhammad Shakeel[3], Jee-weon Jung[1], Soumi Maiti[1], Shinji Watanabe[1]*

[1]Carnegie Mellon University, USA
[2]Shanghai Jiao Tong University, China
[3]Honda Research Institute Japan, Japan

## ABSTRACT

Pre-training speech models on large volumes of data has achieved remarkable success. OpenAI Whisper is a multilingual multitask model trained on 680k hours of supervised speech data. It generalizes well to various speech recognition and translation benchmarks even in a zero-shot setup. However, the full pipeline for developing such models (from data collection to training) is not publicly accessible, which makes it difficult for researchers to further improve its performance and address training-related issues such as efficiency, robustness, fairness, and bias. This work presents an Open Whisper-style Speech Model (OWSM), which reproduces Whisper-style training using an open-source toolkit and publicly available data. OWSM even supports more translation directions and can be more efficient to train. We will publicly release all scripts used for data preparation, training, inference, and scoring as well as pre-trained models and training logs to promote open science.[1]

***Index Terms***— Pre-training, whisper, speech recognition, speech translation

## 1. INTRODUCTION

Large-scale Transformers [1] have garnered significant attention in natural language processing (NLP) [2–7]. These models, trained on extensive datasets, have showcased remarkable emergent capabilities in diverse downstream tasks. Notably, the application of similar pre-training techniques has also found success in the domain of speech processing. Self-supervised learning (SSL) techniques have demonstrated impressive achievements [8–14]. Furthermore, large-scale supervised learning has emerged as a promising avenue for the development of universal speech models capable of performing multiple speech tasks within a single model [15–18]. OpenAI Whisper [15] is a series of multilingual multitask models trained on 680k hours of labeled speech data which is carefully curated from diverse sources on the Internet.

Despite the release of pre-trained Whisper models and inference code, the comprehensive pipeline for model development (from data preparation to training) remains inaccessible to the public, which has been a common situation for large language models (LLMs). This limitation engenders several concerns. Firstly, the utilization of pre-trained models on novel benchmarks has the potential risk of data leakage, as users are deprived of knowledge regarding the actual training data. Secondly, researchers face significant difficulties in comprehending the underlying mechanisms and elucidating methods for enhancing the model's performance, given their lack of access to the training dynamics. Thirdly, the absence of access to the complete model development pipeline poses notable challenges in effectively tackling issues related to robustness, fairness, bias, and toxicity, all of which frequently arise as a result of the data and training procedure [19–21].

Recently, there has been a concerted effort to foster open science in the realm of LLM research by advocating for the release of complete training pipelines [5]. Inspired by this, we present the Open Whisper-style Speech Model (OWSM)[2], which reproduces Whisper-style training using an open-source toolkit and publicly available data. OWSM follows the design of Whisper [15] to support essential tasks such as language identification (LID), multilingual automatic speech recognition (ASR), and utterance-level segmentation. Notably, OWSM also exhibits several technical novelties. It is designed to support any-to-any speech translation as opposed to solely any-to-English translation (see Section 3.4 for results). OWSM also adopts multiple strategies to enhance the efficiency (see Section 2.5 for discussions).

We will provide reproducible recipes encompassing the entire pipeline, including data preparation, training, inference, and scoring. Furthermore, we will release pre-trained models and training logs, enabling researchers to delve into the specifics of the training process and gain valuable insights for their own investigations. While OWSM shows competitive or even superior performance compared to Whisper in certain benchmarks, it is essential to clarify that our objective is not to engage in a comprehensive competition with Whisper. The scope of our endeavor is constrained by the fact that our largest dataset comprises only a quarter of the training set used by Whisper, and our resource limitations restrict us from conducting multiple trial runs. Instead, by sharing these resources, we aim to promote transparency and facilitate progress and advancements in the field of large-scale pre-training for speech processing.

## 2. WHISPER-STYLE TRAINING

### 2.1. Multitask data format

OpenAI Whisper [15] employs a single sequence-to-sequence model to perform multiple speech processing tasks, including LID, multilingual ASR, any-to-English ST, and utterance-level segmentation.

---

[1]https://github.com/espnet/espnet
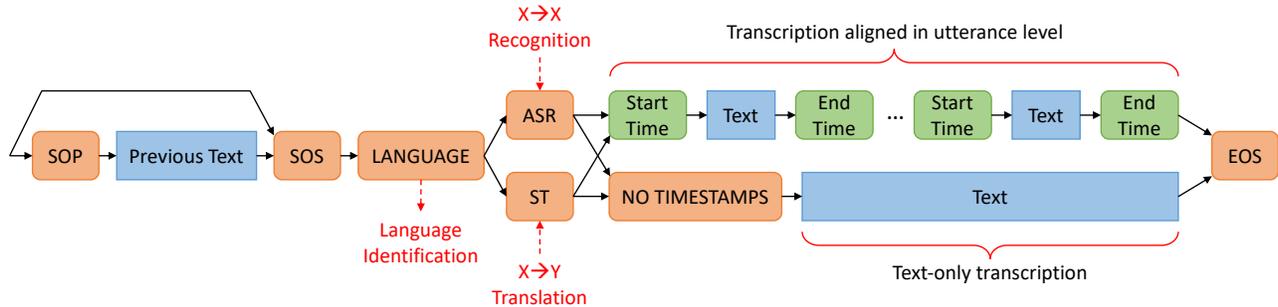
[2]OWSM is pronounced as "awesome".

**Fig. 1**: Multitask data format used by our OWSM, which mostly follows OpenAI Whisper [15]. Different speech processing tasks are represented in a unified format, which can be predicted by an autoregressive decoder. Note that OWSM is designed to support any-to-any speech-to-text translation, whereas Whisper can only perform any-to-English translation. Blue boxes denote standard text tokens, while orange and green boxes are special tokens. SOP, SOS, and EOS represent start-of-prompt, start-of-sentence, and end-of-sentence, respectively.

Our OWSM mostly follows this design, but extends it to potentially support any-to-any ST. Figure 1 illustrates the multitask data format. Data samples from different tasks are represented in a unified format, which can be predicted by the decoder in an autoregressive manner. Specifically, each sample is converted to a sequence of tokens with two segments separated by special tokens. The first segment (before "SOS") is an optional text prompt used as a condition, while the second segment is the actual target. The target starts with a special token denoting the language of the input speech. Then, it uses a task token to distinguish between ASR and ST. There is a separate ST token for each target language, which enables translation to any language. Finally, it appends the text transcription either with or without utterance-level timestamps. All timestamps are quantized and represented as special tokens.

### 2.2. Data preparation

Whisper is pre-trained on 680k hours of labeled audio data sourced from the Internet, which is not publicly accessible. To construct a speech dataset for large-scale supervised learning, we combine training sets from various publicly available ASR and ST corpora. These diverse corpora encompass a wide range of speaking styles, recording environments, and languages. Our datasets are prepared using an open-source toolkit, ESPnet [22]. However, OWSM is trained on long-form audio data, which deviates from previous recipes in ESPnet. Consequently, we have developed new data preparation scripts tailored specifically for Whisper-style training. We concatenate consecutive utterances within the same long talk based on their original timestamps. Each long-form utterance is limited to a maximum duration of 30 seconds. During training, all utterances are padded to precisely 30 seconds, optimizing the utilization of computational resources.

To date, we have developed three versions at different scales, denoted as OWSM v1, v2, and v3 in Table 1. Our largest dataset, v3, comprises 180k hours of labeled audio data. This constitutes approximately one quarter of the total data employed by OpenAI Whisper in its training process [15]. The individual datasets utilized by our models are listed below:

- OWSM v1: AISHELL-1 [23], CoVoST2 [24], GigaSpeech [25], LibriSpeech [26], MuST-C [27], SPGISpeech [28], and TEDLIUM3 [29].
- OWSM v2: all data in v1, GigaST [30], Multilingual LibriSpeech [31], and WenetSpeech [32].

- OWSM v3: all data in v2, AIDATATANG [33], AMI [34], Babel [35], Common Voice [36], Fisher (Switchboard) [37], Fisher Callhome Spanish [38], FLEURS [39], Googlei18n[3], KsponSpeech [40], MagicData [41], ReazonSpeech [42], Russian Open STT [43], VCTK [44], VoxForge [45], VoxPopuli [46], and WSJ [47].

### 2.3. Model architectures

OWSM follows Whisper to utilize a Transformer encoder-decoder architecture [1], where the encoder and decoder have the same number of layers. However, OWSM additionally employs a joint CTC loss for ASR targets [48], which was empirically shown to stabilize our training process. The input waveforms are converted to 80-dimensional log Mel filterbanks with a window length of 25ms and a hop length of 10ms. The extracted features are augmented using SpecAugment [49] and normalized by their global mean and variance. The features are then processed by a two-dimensional convolution module to reduce the sequence length. OpenAI Whisper [15] always downsamples the sequence by 2, resulting in a time resolution of 20ms. Our OWSM v2 and v3 perform 4 times downsampling, which further improves efficiency. The detailed configurations of Transformer encoder and decoder layers are summarized in Table 1. OWSM v1 and v3 use the same configurations as Whisper small and medium, respectively, while OWSM v2 is slightly smaller than v3. [4]

For inference, OpenAI Whisper implements both greedy decoding and beam search with temperature fallback. The latter is a complicated procedure relying on many heuristics and hyperparameters such as beam sizes, temperatures, log probability threshold, and gzip compression rate threshold. Our OWSM utilizes the ESPnet framework [22], thereby ensuring compatibility with various decoding algorithms originally supported by ESPnet, including greedy search, beam search, and joint CTC/attention decoding (for ASR only) [50].

---

[3]Resources 32, 35, 36, 37, 41, 42, 43, 44, 52, 53, 54, 61, 63, 64, 65, 66, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, and 86 from openslr.org.

[4]OWSM has slightly more parameters than Whisper under the same configuration, because the ESPnet model has a larger convolution downsampling module and does not share the input embedding and output projection in its decoder.

**Table 1**: Details of data, model architectures, and training configurations. We gradually increase data and model sizes from v1 to v3. The model configurations of OWSM v1 and v3 match those of Whisper small and medium, respectively. Although OWSM v3 covers more languages than Whisper, our data size remains significantly smaller, making our task much more challenging. *Our v3 model is initialized with the pre-trained v2 to reduce training time (see Section 2.5).

| | OpenAI Whisper | | | OWSM (ours) | | |
|---|---|---|---|---|---|---|
| | small | medium | large | v1 | v2 | v3* |
| *Data* | | | | | | |
| Total hours (k) | | 680 | | 38 | 129 | 180 |
| - English ASR | | 438 | | 22 | 67 | 73 |
| - Multilingual ASR | | 117 | | 1 | 22 | 67 |
| - Translation | | 125 | | 15 | 40 | 40 |
| Languages | | 99 | | 22 | 23 | 151 |
| BPE vocabulary size | | 51,865 | | 20k | 50k | 50k |
| *Model architectures* | | | | | | |
| Parameters (M) | 244 | 769 | 1550 | 272 | 712 | 889 |
| Hidden size | 768 | 1024 | 1280 | 768 | 1024 | 1024 |
| Layers | 12 | 24 | 32 | 12 | 18 | 24 |
| Attention heads | 12 | 16 | 20 | 12 | 16 | 16 |
| Time resolution (ms) | 20 | 20 | 20 | 20 | 40 | 40 |
| *Training configurations* | | | | | | |
| Batch size | | 256 | | | 256 | |
| Total updates | | 1,048,576 | | 300k | 500k | 470k |
| Warmup updates | | 2048 | | 10k | 20k | 10k |
| Learning rate | 5e-4 | 2.5e-4 | 1.75e-4 | 1e-3 | 5e-4 | 2.5e-4 |
| Optimizer | | AdamW | | | AdamW | |
| Joint CTC weight | | NA | | | 0.3 | |

**Table 2**: WER % (↓) of English ASR using greedy search. OpenAI Whisper uses 438k hours of English ASR data, while OWSM uses at most 73k hours. As shown in Table 1, the configurations of OWSM v1 and v3 match those of Whisper small and medium, respectively. Whisper large is significantly larger, so it is not included in the comparison. †The larger degradation of OWSM v3 on WSJ is likely caused by inconsistent case and punctuation of the training data (see the last paragraph in Section 2.5). The gray color means OWSM is better than Whisper small and medium.

| Dataset | OpenAI Whisper | | OWSM (ours) | | |
|---|---|---|---|---|---|
| | small | medium | v1 | v2 | v3 |
| Common Voice en | 15.7 | **11.9** | 20.1 | 14.4 | 14.5 |
| FLEURS en | 9.6 | **6.4** | 13.2 | 10.9 | 10.9 |
| LibriSpeech test-clean | 3.3 | 2.8 | 5.4 | **2.2** | 2.7 |
| LibriSpeech test-other | 7.7 | 6.5 | 10.9 | **5.1** | 6.0 |
| Switchboard eval2000 | 22.2 | 19.4 | 28.7 | 20.4 | **17.2** |
| TEDLIUM test | **4.6** | 5.1 | 6.6 | **4.6** | 4.8 |
| VoxPopuli en | 8.5 | **7.6** | 14.2 | 10.3 | 9.2 |
| WSJ eval92 | 4.3 | **2.9** | 4.3 | 3.7 | 13.4† |

## 2.4. Training details

OWSM is implemented in ESPnet [22] based on PyTorch [51]. Table 1 compares the training hyperparameters of different models. OWSM uses the same batch size as Whisper, but the number of total updates is smaller. OWSM is trained on NVIDIA A100 (40GB) GPUs. Each GPU takes two samples, and gradient accumulation is applied whenever necessary to ensure the total batch size is 256. Specifically, OWSM v1 is trained for around 7 days on 32 A100 GPUs and OWSM v2 and v3 are trained for around 10 days on 64 A100 GPUs. After training, five checkpoints with the highest validation accuracies are averaged to generate the final checkpoint.

## 2.5. Challenges and training tips

Large-scale distributed training presents significant challenges, particularly when the computation budget is limited. As we scale up from a few thousand hours of data to nearly 200 thousand hours, we have encountered a range of issues. Here, we discuss some of the challenges and provide valuable training tips to help overcome these obstacles effectively. We will release our scripts that support these techniques.

**Time resolution:** Whisper employs a time resolution of 20ms within its encoder module, resulting in a sequence length of 1500 for 30-second inputs. This significantly increases GPU memory consumption and makes training slower and more difficult. In contrast, contemporary state-of-the-art ASR and ST models [52–55] adopt larger downsampling rates. Starting from OWSM v2, we have adopted a time resolution of 40ms, effectively reducing the sequence

length and mitigating the associated computational demands. We have also found that a shorter sequence length facilitates easier convergence of the model.

**Joint ASR CTC loss:** In our preliminary experiments, we observed suboptimal convergence of the attention-based encoder-decoder model trained on multiple tasks and diverse data. Incorporating a joint ASR CTC loss [48] to the encoder output can stabilize training and expedite convergence.

**Warm initialization:** When training our largest model, OWSM v3, we employ a warm initialization technique by leveraging the pre-trained OWSM v2. Specifically, the first 18 layers of OWSM v3 are initialized with v2 (which has precisely 18 layers), whereas the remaining 6 layers are initialized randomly. This v3 model converges much faster than training from scratch. However, it remains to be investigated whether a warm initialization adversely affects the final performance of the model.

**Memory and efficiency issues:** We have developed several strategies to address memory and efficiency issues caused by large data. To train the BPE tokenization models using SentencePiece [56], we randomly select 10 million text transcriptions instead of using the whole set to reduce memory usage. For training, the entire text file is too large to be distributed across different workers. We partition the training set into 5 to 12 non-overlapping subsets and use multiple data iterators to construct mini-batches. We further filter out samples with extremely long transcriptions (e.g., greater than 600 tokens including both prompt and target) which are caused by incorrect alignments in the original corpus (e.g., Common Voice). Without such filtering, the training will occasionally encounter out-of-memory errors. Additionally, we validate intermediate checkpoints using only 10% of the full validation set. This might generate slightly inaccurate estimates of the actual performance, but it significantly reduces the validation time and thus allows for more frequent validation and checkpoint saving, which is crucial for large-scale distributed training. In fact, we encountered various failures mainly due to file system or communication errors, and we had to manually resume from previous checkpoints.

**Inconsistent case and punctuation.** Our training data is gathered from many public corpora. Some of them provide raw transcripts in true case with punctuation, but the others only provide

**Table 3**: WER/CER % (↓) of multilingual ASR using greedy search. Training data sizes (in hours) are also provided. The gray color means OWSM is better than Whisper small and medium.

| Dataset | Language | Metric | OpenAI Whisper | | | OWSM v1 | | OWSM v2 | | OWSM v3 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | hours | small | medium | hours | result | hours | result | hours | result |
| Multilingual LibriSpeech | English | WER | 438k | 9.1 | 10.2 | 22k | 13.7 | 67k | **6.7** | 73k | 7.4 |
| | Spanish | | 11k | 9.1 | **6.1** | 0.1k | 37.2 | 1.0k | 11.7 | 2.0k | 11.7 |
| | French | | 10k | 13.6 | **9.7** | 0.3k | 41.8 | 1.3k | 13.0 | 2.5k | 14.1 |
| | German | | 13k | 11.5 | **8.1** | 0.2k | 43.3 | 2.2k | 11.8 | 3.7k | 11.9 |
| | Dutch | | 2.1k | 18.2 | **12.2** | 0.007k | 78.7 | 1.6k | 16.9 | 1.7k | 17.7 |
| | Italian | | 2.6k | 21.3 | **15.6** | 0.04k | 54.9 | 0.3k | 23.1 | 0.7k | 24.5 |
| | Portuguese | | 8.6k | 13.8 | **8.9** | 0.009k | 90.9 | 0.2k | 31.8 | 0.3k | 28.2 |
| | Polish | | 4.3k | 12.5 | **6.8** | 0 | NA | 0.1k | 89.7 | 0.3k | 37.0 |
| AISHELL-1 | Chinese | CER | 23k | 25.1 | 15.7 | 0.2k | 22.6 | 15k | **5.9** | 16k | 7.1 |
| KsponSpeech eval-clean | Korean | | 8k | 24.0 | **17.6** | 0 | NA | 0 | NA | 1.0k | 20.5 |
| KsponSpeech eval-other | | | | 15.4 | **12.8** | | | | | | 22.6 |
| ReazonSpeech | Japanese | | 7k | 32.5 | 25.3 | ≈0 | NA | ≈0 | NA | 19k | **11.3** |

**Table 4**: WER % (↓) of long-form ASR on the TEDLIUM2 test set. Unsegmented long talks are transcribed in chunks of 30 seconds. It is shifted based on predicted timestamps.

| Beam size | OpenAI Whisper | | OWSM (ours) | |
|---|---|---|---|---|
| | small | medium | v2 | v3 |
| 1 | 4.4 | **3.8** | 7.2 | 9.2 |
| 5 | 4.2 | **3.8** | 6.6 | 7.6 |

normalized transcripts in lower or upper case without any punctuation. During inference, we find that OWSM models are so powerful that they are able to recognize the corpus and generate outputs that are consistent with the training data format. For example, the training data of WSJ is in upper case. When tested on WSJ test sets, OWSM also mostly generates text in upper case. Since only a very small portion of training data is in upper case, OWSM v3 performs poorly on WSJ (see Table 2). In the future, we will normalize the text to address this issue. Note that this analysis demonstrates the benefit of using public data and open-source code, without which we cannot discover such issues.

## 3. EXPERIMENTS

### 3.1. English speech recognition

Table 2 presents word error rates (WER) on standard English ASR benchmarks. Greedy search is employed without any external language models. To ensure fair comparison, we prepare all test data in ESPnet and evaluate Whisper in the same setup instead of reporting results from their paper [15]. The text is normalized using the English or basic normalizer provided by Whisper. Whisper large is not included since it is significantly larger than the other models. Although many public ASR corpora are combined, our English training data is still significantly smaller than that of Whisper (73k vs 438k hours). However, our OWSM models achieve competitive results in most benchmarks. OWSM models even outperform Whisper on LibriSpeech and Switchboard.

By comparing different versions of OWSM, we observe that its English ASR capability is largely improved from v1 to v2, demonstrating the effectiveness of scaling up in terms of the number of model parameters and the amount of training data. However,

OWSM v3 does not show a consistent improvement over v2 in all benchmarks. OWSM v3 achieves lower WERs on Switchboard and VoxPopuli test sets, likely because their training sets are newly added (see Section 2.2). OWSM v3 has slight degradations on LibriSpeech and a large degradation on WSJ. This is probably due to the shift of data distributions from v2 to v3. As shown in Table 1, our v3 dataset contains significantly more languages compared to v2 (151 vs 23), but the model size is only slightly increased (889M vs 712M). Hence, the model has to adjust its capacity from English to other languages or from one type of speech to another type. This issue might be mitigated with larger models and more diverse data. We will explore it in the future. Please refer to the last paragraph in Section 2.5 for more discussions.

We have also investigated the inference speed. Specifically, we select 50 utterances of 30 seconds from our prepared TEDLIUM dev set, and decode OWSM v3 with greedy search using a single NVIDIA A40 GPU. The average decoding time for each 30-second utterance is 2.3 seconds.

### 3.2. Multilingual speech recognition

Table 3 shows the ASR results on multilingual benchmarks. In general, OpenAI Whisper achieves better performance than our OWSM, because Whisper employs significantly more training data in all languages except Japanese. For Japanese, OWSM v3 outperforms Whisper by a large margin (CER: 11.3 vs 25.3) thanks to the larger amount of training data (19k vs 7k hours) from ReazonSpeech [42]. Notably, OWSM v2 achieves the best results on the English and Chinese test sets from Multilingual LibriSpeech and AISHELL, respectively, despite being trained on less data.

The trend across different versions of OWSM is consistent with that in Section 3.1. OWSM v2 is drastically improved compared to v1 in all languages, which verifies the benefits of scaling up. OWSM v3 outperforms v2 in a few languages but achieves comparable or slightly worse results in the others. Again, this is likely because the model needs to adjust its capacity to support much more languages in v3.

### 3.3. Long-form speech recognition

Similar to Whisper, OWSM performs long-form ASR by consecutively transcribing 30-second audio segments and shifting the win-

**Table 5**: Examples of ASR on 30-second audio segments, generated by OWSM v2 using greedy search. Utterances can be segmented in different ways, but the predicted timestamps are usually accurate. Differences between the reference and prediction are marked in red.

| # | Groundtruth from the dev set of MuST-C v2 | Prediction by OWSM v2 |
|---|---|---|
| 1 | \<en>\<asr>\<0.00> I'm going to talk today about energy and climate.\<3.50>\<4.28> And that might seem a bit surprising, because my full-time work at the foundation is mostly about vaccines and seeds, about the things that we need to invent and deliver to help the poorest two billion live better lives.\<18.38>\<19.64> But energy and climate are extremely important to these people; in fact, more important than to anyone else on the planet.\<28.52> | \<en>\<asr>\<0.00> I'm going to talk today about energy and climate.\<3.52>\<4.26> And that might seem a bit surprising, because my full-time work at the foundation is mostly about vaccines and seeds, about the things that we need to invent and deliver to help the poorest two billion live better lives.\<18.40>\<19.62> But energy and climate are extremely important to these people, in fact more important than to anyone else on the planet.\<28.52> |
| 2 | \<en>\<asr>\<0.00> Several years ago here at TED, Peter Skillman introduced a design challenge called the marshmallow challenge.\<5.60>\<5.80> And the idea's pretty simple: Teams of four have to build the tallest free-standing structure out of 20 sticks of spaghetti, one yard of tape, one yard of string and a marshmallow.\<16.52>\<16.52> The marshmallow has to be on top.\<18.18>\<18.54> And, though it seems really simple, it's actually pretty hard because it forces people to collaborate very quickly.\<25.04>\<25.42> And so, I thought this was an interesting idea, and I incorporated it into a design workshop.\<29.72> | \<en>\<asr>\<0.00> Several years ago here at TED, Peter Skillman introduced a design challenge called the Marshmellow Challenge, and the idea is pretty simple.\<7.32>\<7.50> Teams of four have to build the tallest freestanding structure out of 26 of spaghetti, one yard of tape, one yard of string and a marshmallow.\<16.50>\<16.54> The marshmallow has to be on top.\<18.20>\<18.54> And though it seems really simple, it's actually pretty hard because it forces people to collaborate very quickly.\<25.04>\<25.44> And so I thought this was an interesting idea, and I incorporated it into a design workshop.\<30.00> |

**Table 6**: BLEU % (↑) of speech translation. OpenAI Whisper supports any-to-English translation. OWSM can support more directions. The sizes of training sets (in hours) are also provided.

| Dataset | Source | Target | OpenAI Whisper | | OWSM v2 | | OWSM v3 | |
|---|---|---|---|---|---|---|---|---|
| | | | small | medium | hours | result | hours | result |
| MuST-C | English | German | | | 14k | **28.5** | 14k | 27.9 |
| | | Chinese | | NA | 14k | 20.5 | 14k | **20.7** |
| | | Japanese | | | 1.0k | **10.5** | 1.0k | 9.4 |
| | | Spanish | | | 0.5k | **23.4** | 0.5k | 22.5 |
| | | French | | | 0.5k | **28.5** | 0.5k | 26.2 |
| CoVoST | German | English | 4.3k | 26.2 | **34.8** | 0.2k | 18.6 | 0.2k | 18.0 |
| | Chinese | | 12k | 6.3 | **13.6** | 0.01k | 3.0 | 0.01k | 3.3 |
| | Japanese | | 8.9k | 15.9 | **22.9** | 0.001k | 0.1 | 0.001k | 0.1 |
| | Spanish | | 6.7k | 34.2 | **40.2** | 0.1k | 24.9 | 0.1k | 22.7 |
| | French | | 4.5k | 27.8 | **34.8** | 0.3k | 26.0 | 0.3k | 23.7 |

**Table 7**: Accuracy % (↑) of language identification. OWSM v3 supports 151 languages, whereas Whisper supports 99 languages.

| Dataset | OpenAI Whisper | | OWSM (ours) |
|---|---|---|---|
| | small | medium | v3 |
| FLEURS | 53.1 | 54.8 | **81.4** |

**Table 8**: WER/CER % (↓) of OWSM v3 using different decoding algorithms in ESPnet.

| Dataset | Metric | CTC | Attention | Joint CTC/attention |
|---|---|---|---|---|
| Common Voice en | | 18.6 | 14.5 | **12.9** |
| FLEURS en | | 17.3 | 10.9 | **9.7** |
| LibriSpeech test-clean | | 4.5 | 2.7 | **2.6** |
| LibriSpeech test-other | WER | 8.1 | 6.0 | **5.4** |
| Switchboard eval2000 | | 19.4 | 17.2 | **16.6** |
| TEDLIUM test | | 6.7 | 4.8 | **4.7** |
| VoxPopuli en | | 12.3 | 9.2 | **8.7** |
| WSJ eval92 | | 32.0 | 13.4 | **11.4** |
| AISHELL-1 test | CER | 9.2 | 7.1 | **6.5** |
| ReazonSpeech test | | 17.0 | 11.3 | **10.3** |

dow based on predicted timestamps. Table 4 presents the long-form ASR results on the TEDLIUM test set, where each input audio is an unsegmented long talk. OWSM v2 achieves 7.2% WER with greedy decoding and 6.6% WER with beam search. Whisper models achieve lower WERs in both cases, likely because: (1) their training set is larger; (2) their data, collected from the Internet, is originally in a long form, which can be more realistic than ours; (3) they apply various heuristics to improve the timestamp prediction and also the quality of text (see Section 4.5 in their official report [15]). In our future work, we will explore more strategies to enhance the long-form performance.

Table 5 shows two examples from TED talks, where timestamps are generated along with text tokens. Although utterances can be segmented in different ways, the boundaries predicted by OWSM are usually very close to the reference.

### 3.4. Speech translation

Table 6 compares different models on two ST benchmarks: MuST-C (English-to-X) and CoVoST (X-to-English). Whisper only supports the latter, while OWSM supports both directions.

OWSM models achieve notable results on MuST-C thanks to the sufficient amount of training data (more than 500 hours for each language). The BLEU scores of Chinese and Japanese are lower than those of European languages even with enough training data. This indicates that the model has difficulty in translating between very different languages.

On CoVoST, the performance of OWSM ranges across language pairs as the amount of training data varies from 1 to 300 hours. On Chinese and Japanese, OWSM outputs have low intelligibility while on the European languages OWSM outputs are moderately intelligible. On the other hand, Whisper is trained on 4k to 12k hours and thus achieves greater BLEU scores on X-to-English in general.

Similar to the findings in Section 3.1 and Section 3.2, OWSM v3 shows comparable or slightly worse performance than OWSM v2. This is because OWSM v3 employs almost the same amount of ST

data but it has to recognize drastically more languages (see Table 1). Some of its capacity needs to be assigned to these additional languages.

### 3.5. Language identification

As described in Section 2.1 and Figure 1, OWSM predicts a language token at the beginning of decoding, which effectively performs the LID task. Table 7 compares Whisper and OWSM on the FLERUS test set prepared in ESPnet. OWSM v3 achieves a top-1 accuracy of 81.4%, which outperforms Whisper small and medium by a large margin. This is because OWSM v3 utilizes the training data from Common Voice and FLEURS, containing 151 languages in total, whereas Whisper supports 99 languages that only cover a subset of the languages in FLEURS. Nevertheless, this result demonstrates that OWSM has a strong capability in speech classification although it is designed as a sequence-to-sequence model.

### 3.6. Comparison of decoding algorithms

OWSM is compatible with various decoding algorithms in ESPnet. Table 8 compares three commonly used algorithms: CTC only (greedy), attention only (greedy), and joint CTC/attention (with beam size 10 and CTC weight 0.3). Beam search with joint CTC/attention achieves the best results in all test sets. Attention-only decoding outperforms CTC-only, indicating that the decoder has strong capacity.

### 4. DISCUSSIONS AND FUTURE DIRECTIONS

This work serves as an exploratory endeavor in reproducing Whisper-style training using open-source resources. Moving forward, we will delve into the following directions.

Firstly, the current OWSM still falls behind Whisper in many benchmarks, likely because: (1) OWSM supports more languages and more translation directions, which increases the difficulty of multitask learning; (2) our training set is significantly smaller than that of Whisper in nearly all languages and tasks; (3) we directly leverage public ASR and ST corpora which may be less diverse than Whisper's data collected from the Internet. These issues can probably be addressed by utilizing more advanced encoder [52–54, 57] or decoder [58] architectures, collecting more diverse ASR and ST data from public sources, and incorporating self-supervised speech representations [8, 9] as in Google USM [18].

Secondly, we plan to incorporate additional speech processing tasks into the multitask framework, including spoken language understanding and speech generation based on discrete representations, thereby working towards the development of "universal speech models".

Thirdly, these large pre-trained models are unsuitable for deployment in real-world applications. Various compression techniques [59–64] can be applied to reduce the model size and computation.

Fourthly, OWSM provides a valuable testbed for investigating and exploring various machine learning problems such as data imbalance, continual learning [65], adversarial robustness [66], and machine unlearning [67].

### 5. CONCLUSION

This work presents OWSM, which reproduces Whisper-style training using an open-source toolkit and publicly available data. OWSM follows the multitask framework of OpenAI Whisper, but extends it to support more translation directions. Several strategies are developed to improve efficiency. We will open-source all scripts for data preparation, training, inference, and scoring as well as pre-trained models and training logs. We believe this can promote transparency and facilitate advancements in the large-scale pre-training of speech models.

### 7. REFERENCES

[1] A. Vaswani et al., "Attention is all you need," in *Proc. NeurIPS*, 2017.

[2] Tom Brown et al., "Language models are few-shot learners," 2020.

[3] Jack W Rae et al., "Scaling language models: Methods, analysis & insights from training gopher," *arXiv:2112.11446*, 2021.

[4] Aakanksha Chowdhery et al., "Palm: Scaling language modeling with pathways," *arXiv:2204.02311*, 2022.

[5] Susan Zhang et al., "Opt: Open pre-trained transformer language models," *arXiv:2205.01068*, 2022.

[6] Hugo Touvron et al., "Llama: Open and efficient foundation language models," *arXiv:2302.13971*, 2023.

[7] OpenAI, "GPT-4 Technical Report," *arXiv:2303.08774*, 2023.

[8] Alexei Baevski, Yuhao Zhou, et al., "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. NeurIPS*, 2020.

[9] Wei-Ning Hsu et al., "HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 3451–3460, 2021.

[10] Arun Babu, Changhan Wang, Andros Tjandra, et al., "XLS-R: Self-supervised Cross-lingual Speech Representation Learning at Scale," in *Proc. Interspeech*, 2022.

[11] Shu wen Yang et al., "SUPERB: Speech Processing Universal PERformance Benchmark," in *Proc. Interspeech*, 2021.

[12] Abdelrahman Mohamed et al., "Self-supervised speech representation learning: A review," *IEEE J. Sel. Topics Signal Process.*, vol. 16, no. 6, pp. 1179–1210, 2022.

[13] Xuankai Chang, Takashi Maekaku, et al., "An exploration of self-supervised pretrained representations for end-to-end speech recognition," in *Proc. ASRU*, 2021.

[14] Yifan Peng et al., "A Study on the Integration of Pre-trained SSL, ASR, LM and SLU Models for Spoken Language Understanding," in *Proc. SLT*, 2022.

[15] Alec Radford et al., "Robust speech recognition via large-scale weak supervision," *arXiv:2212.04356*, 2022.

[16] William Chan et al., "Speechstew: Simply mix all available speech recognition data to train one large neural network," *arXiv:2104.02133*, 2021.

[17] Bo Li et al., "Scaling end-to-end models for large-scale multilingual asr," in *Proc. ASRU*, 2021.

[18] Yu Zhang et al., "Google USM: Scaling automatic speech recognition beyond 100 languages," *arXiv:2303.01037*, 2023.

[19] Paul Pu Liang et al., "Towards understanding and mitigating social biases in language models," in *Proc. ICML*, 2021.

[20] Jindong Wang et al., "On the robustness of chatgpt: An adversarial and out-of-distribution perspective," *arXiv:2302.12095*, 2023.

[21] Sébastien Bubeck et al., "Sparks of artificial general intelligence: Early experiments with gpt-4," *arXiv:2303.12712*, 2023.

[22] Shinji Watanabe et al., "ESPnet: End-to-End Speech Processing Toolkit," in *Proc. Interspeech*, 2018.

[23] Hui Bu et al., "AISHELL-1: An open-source Mandarin speech corpus and a speech recognition baseline," in *Proc. O-COCOSDA*, 2017.

[24] Changhan Wang et al., "CoVoST 2 and Massively Multilingual Speech Translation," in *Interspeech*, 2021.

[25] Guoguo Chen et al., "GigaSpeech: An Evolving, Multi-Domain ASR Corpus with 10,000 Hours of Transcribed Audio," in *Proc. Interspeech*, 2021.

[26] Vassil Panayotov et al., "Librispeech: An ASR corpus based on public domain audio books," in *ICASSP*, 2015.

[27] Roldano Cattoni et al., "Must-c: A multilingual corpus for end-to-end speech translation," *Computer speech & language*, vol. 66, pp. 101155, 2021.

[28] Patrick K O'Neill et al., "Spgispeech: 5,000 hours of transcribed financial audio for fully formatted end-to-end speech recognition," *arXiv:2104.02014*, 2021.

[29] François Hernandez et al., "Ted-lium 3: Twice as much data and corpus repartition for experiments on speaker adaptation," in *Speech & Computer*, 2018, pp. 198–208.

[30] Rong Ye et al., "Gigast: A 10,000-hour pseudo speech translation corpus," *arXiv:2204.03939*, 2022.

[31] Vineel Pratap et al., "Mls: A large-scale multilingual dataset for speech research," *arXiv:2012.03411*, 2020.

[32] Binbin Zhang et al., "Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition," in *Proc. ICASSP*, 2022.

[33] "aidatatang_200zh, a free Chinese Mandarin speech corpus by Beijing DataTang Technology Co., Ltd," .

[34] Jean Carletta, "Unleashing the killer corpus: experiences in creating the multi-everything AMI Meeting Corpus," *Lang. Res. Eval.*, vol. 41, pp. 181–190, 2007.

[35] "The babel program: https://www.iarpa.gov/index.php/research-programs/babel," .

[36] Rosana Ardila et al., "Common voice: A massively-multilingual speech corpus," *arXiv:1912.06670*, 2019.

[37] J.J. Godfrey et al., "SWITCHBOARD: telephone speech corpus for research and development," in *Proc. ICASSP*, 1992.

[38] Matt Post et al., "Improved speech-to-text translation with the fisher and callhome Spanish-English speech translation corpus," 2013.

[39] Alexis Conneau et al., "FLEURS: Few-Shot Learning Evaluation of Universal Representations of Speech," in *Proc. SLT*, 2022.

[40] Jeong-Uk Bang et al., "Ksponspeech: Korean spontaneous speech corpus for automatic speech recognition," *Applied Sciences*, vol. 10, no. 19, pp. 6936, 2020.

[41] Zehui Yang et al., "Open source magicdata-ramc: A rich annotated mandarin conversational (ramc) speech dataset," *arXiv:2203.16844*, 2022.

[42] Yue Yin, Daijiro Mori, et al., "ReazonSpeech: A Free and Massive Corpus for Japanese ASR," 2023.

[43] Anna Slizhikova et al., "Russian Open Speech To Text (STT/ASR) Dataset," 2020.

[44] Junichi Yamagishi et al., "CSTR VCTK Corpus: English Multi-speaker Corpus for CSTR Voice Cloning Toolkit," 2019.

[45] "VoxForge: http://www.voxforge.org/," .

[46] Changhan Wang et al., "VoxPopuli: A Large-Scale Multilingual Speech Corpus for Representation Learning, Semi-Supervised Learning and Interpretation," in *Proc. ACL*, 2021.

[47] Douglas B Paul and Janet Baker, "The design for the Wall Street Journal-based CSR corpus," in *Proc. Workshop on Speech and Natural Language*, 1992.

[48] Suyoun Kim, Takaaki Hori, and Shinji Watanabe, "Joint ctc-attention based end-to-end speech recognition using multi-task learning," in *Proc. ICASSP*, 2017.

[49] Daniel S. Park, William Chan, et al., "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," in *Proc. Interspeech*, 2019.

[50] Takaaki Hori, Shinji Watanabe, and John R Hershey, "Joint CTC/attention decoding for end-to-end speech recognition," in *Proc. ACL*, 2017.

[51] A. Paszke et al., "Pytorch: An imperative style, high-performance deep learning library," in *NeurIPS*, 2019.

[52] Anmol Gulati et al., "Conformer: Convolution-augmented Transformer for Speech Recognition," in *Proc. Interspeech*, 2020.

[53] Yifan Peng, Siddharth Dalmia, Ian Lane, and Shinji Watanabe, "Branchformer: Parallel MLP-attention architectures to capture local and global context for speech recognition and understanding," in *Proc. ICML*, 2022.

[54] Kwangyoun Kim et al., "E-Branchformer: Branchformer with Enhanced Merging for Speech Recognition," in *Proc. SLT*, 2022.

[55] Sehoon Kim et al., "Squeezeformer: An efficient transformer for automatic speech recognition," in *Proc. NeurIPS*, 2022.

[56] Taku Kudo and John Richardson, "SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing," 2018.

[57] Yifan Peng et al., "A Comparative Study on E-Branchformer vs Conformer in Speech Recognition, Translation, and Understanding Tasks," in *Proc. Interspeech*, 2023.

[58] Koichi Miyazaki, Masato Murata, and Tomoki Koriyama, "Structured State Space Decoder for Speech Recognition and Synthesis," in *Proc. ICASSP*, 2023.

[59] Heng-Jui Chang et al., "DistilHuBERT: Speech representation learning by layer-wise distillation of hidden-unit BERT," in *Proc. ICASSP*, 2022.

[60] Cheng-I Jeff Lai et al., "PARP: Prune, Adjust and Re-Prune for Self-Supervised Speech Recognition," in *Proc. NeurIPS*, 2021.

[61] Yifan Peng et al., "Structured Pruning of Self-Supervised Pre-trained Models for Speech Recognition and Understanding," in *Proc. ICASSP*, 2023.

[62] Yifan Peng, Yui Sudo, et al., "DPHuBERT: Joint Distillation and Pruning of Self-Supervised Speech Models," in *Proc. Interspeech*, 2023.

[63] Yizeng Han, Gao Huang, et al., "Dynamic neural networks: A survey," vol. 44, no. 11, pp. 7436–7456, 2021.

[64] Yifan Peng, Jaesong Lee, et al., "I3D: Transformer architectures with input-dependent dynamic depth for speech recognition," in *Proc. ICASSP*, 2023.

[65] German I Parisi, Ronald Kemker, et al., "Continual lifelong learning with neural networks: A review," *Neural networks*, vol. 113, pp. 54–71, 2019.

[66] Raphael Olivier and Bhiksha Raj, "There is more than one kind of robustness: Fooling whisper with adversarial examples," *arXiv:2210.17316*, 2022.

[67] Thanh Tam Nguyen et al., "A survey of machine unlearning," *arXiv:2209.02299*, 2022.