

# On the Robust Approximation of ASR Metrics

Anonymous ACL submission

## Abstract

Recent advances in speech foundation models are largely driven by scaling both model size and data, enabling them to perform a wide range of tasks, including speech recognition. Traditionally, ASR models are evaluated using metrics like Word Error Rate (WER) and Character Error Rate (CER), which depend on ground truth labels. As a result of limited labeled data from diverse domains and testing conditions, the true generalization capabilities of these models beyond standard benchmarks remain unclear. Moreover, labeling data is both costly and time-consuming. To address this, we propose a novel label-free approach for approximating ASR performance metrics, eliminating the need for ground truth labels. Our method utilizes multimodal embeddings in a unified space for speech and transcription representations, combined with a high-quality proxy model to compute proxy metrics. These features are used to train a regression model to predict key ASR metrics like Word Error Rate (WER) and Character Error Rate (CER). We experiment with over 40 models across 14 datasets representing both standard and in-the-wild testing conditions. Our results show that we approximate the metrics within a single-digit absolute difference across all experimental configurations, outperforming the most recent baseline by more than 50%.

## 1 Introduction

Automatic Speech Recognition (ASR) models have made significant advancements in recent years, achieving near-human performance on several standard evaluation benchmarks (Radford et al., 2022; Seamless Communication, 2023; Communication et al., 2023; Harper et al., *inter alia*). These models are typically evaluated using metrics like Word Error Rate (WER) and Character Error Rate (CER) (Likhomanenko et al., 2020), which are essential for assessing model performance.

However, these metrics are dependent on ground truths, which are often scarce in resource-constrained environments, and human labeling is both costly and time-consuming. To mitigate this challenge, several reference-free evaluation methods are proposed (Yuksel et al., 2023b; Kalgaonkar et al., 2015; Swarup et al., 2019; Qiu et al., 2021; Del-Agua et al., 2018; Raj et al., 2011). While these approaches eliminate the reliance on labeled data, they primarily offer relative assessments of transcription quality, rather than providing precise error counts or rates. As a result, their applicability in real-world scenarios, where actionable performance metrics are crucial for further model refinement and deployment, is limited.

Given the limitations of both methods, approximating ASR metrics has emerged as a promising alternative for label-free evaluation (Chowdhury and Ali, 2023; Sheshadri et al., 2021b; Ali and Renals, 2018). This approach typically involves training regression (Jalalvand et al., 2016) and/or classification models (Sheshadri et al., 2021a) on top of speech and text encoders. While this method offers a close approximation of error metrics, several important questions remain unresolved. Specifically, an approximation model trained on dataset sampled from  $D$  to predict ASR metrics for a source model  $M$  must be evaluated under diverse conditions: 1) on test data that is IID (independent and identically distributed) sampled from  $D$ ; 2) on out-of-distribution (OOD) data representing diverse domains and recording conditions; 3) on IID data, but transcription from a target model  $T$ ; and 4) on OOD data with transcriptions from a target model  $T$ . Most prior works (Chowdhury and Ali, 2023; Sheshadri et al., 2021b) focus primarily on the first condition. Moreover, recent advancements in multimodal foundation models offer new opportunities to directly train regression models on unified speech and text embeddings.

To address these critical research gaps, we pro-

pose a novel framework for approximating the performance of a wide range of ASR models, both on standard benchmarks and in-the-wild scenarios. Specifically, we compute the similarity between speech and text embeddings in a unified space, capturing the semantic alignment between the two modalities. Additionally, we incorporate a high-quality reference model as a proxy, based on the intuition that agreement with a reliable proxy correlates with transcription quality, as shown in prior works (Waheed et al., 2025). These features are then used to train a regression model to predict key ASR metrics, such as WER, CER, and absolute word and character error counts.

In summary, our work represents one of the most comprehensive studies to date on approximating ASR metrics at scale, in terms of both data and model coverage. Our proposed approach serves as a reference-free evaluation particularly suited for label-scarce scenarios. Beyond evaluation, our method is especially valuable for tasks such as pseudo-labeling, where high-quality transcriptions are essential for downstream applications like knowledge distillation (Waheed et al., 2024; Gandhi et al., 2023).

Our contributions are as follows:

- We evaluate over 40 ASR models across 14 diverse evaluation setups, including both standard benchmarks and domain-specific, unseen conditions followed by training regression models to approximate ASR metrics.
- We compare our approach with the most recent work on approximating ASR metrics and demonstrate over a 100% improvement against the strong baseline.
- We conduct a rigorous ablation study to analyze the impact of different experimental configurations, providing deeper insights into the robustness of our approach. Our findings show that our method is resilient to diverse evaluation setups and requires only a small amount of training data.

**Outline.** The remainder of this paper is organized as follows: Section 2 reviews related work. Section 3 presents our proposed methodology. Sections 4 and 5 detail our experimental setup, results, and ablation study, respectively. Section 6 concludes the paper and outlines future directions.

**Reproducibility.** We are committed to making all code, data, configurations, and logs available

upon acceptance. Additionally, we will provide a lightweight Python package to seamlessly use our trained approximators.

## 2 Related Work

Automatic speech recognition (ASR) has seen remarkable progress in recent years, driven by advances in deep learning and the availability of extensive training datasets (Radford et al., 2022; Communication et al., 2023). Transformer (Vaswani et al., 2023) based models, in particular, have significantly contributed to these developments by effectively capturing long-range dependencies and contextual nuances in speech, achieving state-of-the-art (SOTA) performance across diverse benchmarks (Kheddar et al., 2024; Dhanjal and Singh, 2024; Zimerman and Wolf, 2023). While traditional evaluation metrics like Word Error Rate (WER) and Character Error Rate (CER) are de-facto evaluation metrics in benchmarking ASR systems (Lin et al., 2021; Park et al., 2024), scenarios where ground truth transcriptions are unavailable have caught interest in reference-free ASR evaluation methods (Karbasi and Kolossa, 2022; Wang et al., 2024; Kuhn et al., 2024).

Reference-free ASR evaluation methods aim to estimate ASR performance without requiring ground truth transcriptions (Ospanov et al., 2024). Earlier approaches rely on heuristic features or metadata such as speaker demographics, background noise, and linguistic characteristics (Litman et al., 2000; Yoon et al., 2010), limiting their applicability across varied contexts. However, recent advancements focus on deep learning-based frameworks, such as convolutional neural networks (CNNs) (Elloumi et al., 2018) and contrastive learning methods (Yuksel et al., 2023a), to predict ASR quality directly from encoded speech and text. For instance, methods like NoRefER (Yuksel et al., 2023b) employ Siamese architectures fine-tuned on ASR hypotheses, achieving high correlation with traditional metrics and improving WER by optimizing hypothesis ensembling (Park et al., 2024).

Efforts to approximate ASR metrics have explored hybrid approaches that combine traditional and reference-free methods, such as leveraging word confidence scores, linguistic embeddings, or post-processing adaptations to estimate WER and CER without explicit references (Ali and Renals, 2020, 2018; Kuhn et al., 2024; Negri et al., 2014). However, these approaches often suffer from re-

liance on specific ASR models or domain characteristics, limiting their generalizability. Unlike existing methods, our work addresses these limitations by introducing a robust, model and data-agnostic framework that evaluates ASR outputs across diverse datasets and configurations, emphasizing adaptability to unseen domains and variations.

### 3 Methodology

We present a scalable and robust method to approximate ASR performance metrics using multi-modal unified embeddings, proxy references, and regression models. The primary goal is to eliminate reliance on ground-truth labels, enabling performance evaluation in label-scarce scenarios. The pipeline consists of three components: representation similarity in a unified speech-text embedding space, agreement with a high-quality proxy reference, and a regression model trained on these features to predict ASR metrics. Our pipeline diagram is shown in Appendix 3 Figure 3.

#### 3.1 Similarity in Unified Representation Space

The foundation of our approach is the SONAR model (Duquenne et al., 2023), a state-of-the-art multimodal (speech-text) model trained to produce unified embeddings for both speech and text inputs. Let  $x_{\text{speech}}$  represent the input speech signal and  $x_{\text{text}}$  denote the corresponding ASR-generated transcription. SONAR maps these inputs to a shared embedding space, generating  $e_{\text{speech}}$  and  $e_{\text{text}}$ :

$$e_{\text{speech}} = f_{\text{SONAR}}(x_{\text{speech}}), \quad e_{\text{text}} = f_{\text{SONAR}}(x_{\text{text}}) \quad (1)$$

where  $f_{\text{SONAR}}$  represents the embedding model. The alignment between these embeddings is quantified using cosine similarity:

$$\text{Similarity}(x_{\text{speech}}, x_{\text{text}}) = \frac{e_{\text{speech}} \cdot e_{\text{text}}}{\|e_{\text{speech}}\| \|e_{\text{text}}\|} \quad (2)$$

This similarity metric serves as an initial indicator of transcription quality, with higher values suggesting better alignment between the speech and text representations.

#### 3.2 Agreement with a Proxy Reference

To complement the similarity score, we utilize proxy references generated by a high-quality ASR model, denoted as  $x_{\text{proxy}}$ . The comparison between the ASR-generated transcription  $x_{\text{text}}$  and the proxy

reference  $x_{\text{proxy}}$  is quantified using Word Error Rate (pWER) and Character Error Rate (pCER) as defined in Appendix A.1.

These metrics assess transcription quality by comparing it with a reliable proxy reference, without using ground-truth labels at any stage. Proxy references are dynamically selected by profiling 41 models across datasets and ranking them by average performance. For each target ASR model, the reference is the highest-ranking model other than the target itself. For example, if `whisper-large-v3` ranks highest, the reference for `whisper` will be the second-best model. This ensures the proxy reference is both relevant and reliable for evaluating the target model.

#### 3.3 Regression Model for Metric Prediction

The extracted features, including similarity scores and proxy metrics, are concatenated to form the input to a regression model. Let  $z = [\text{Similarity}, \text{pWER}/\text{pCER}]$  represent the feature vector. The regression model  $g$  estimates the ASR metrics  $\hat{y}$ , denoted as aWER and/or aCER:

$$\hat{y} = g(z) \quad (3)$$

The regression model is an ensemble of Random Forest, Gradient Boosting, and Histogram-based Gradient Boosting regressors. Each base model is fine-tuned via grid search for hyperparameter optimization. The ensemble is trained to minimize the mean absolute error between predicted and ground-truth metrics. Additionally, a ridge regression model with non-negativity constraints is included in the ensemble to ensure predictions remain within valid ranges. Additional details of our regression pipeline are provided in Section 4, with hyperparameter details in Appendix A.4.

#### 3.4 Evaluation

We evaluate the regression model’s performance across four setups, including IID and OOD data and different model configurations. Specifically, we train our regression model on one ASR system (source) on one dataset and evaluate it on both IID and OOD data for the source model and for a target model. These scenarios assess the model’s robustness and generalization under diverse real-world conditions.

Let  $\mathcal{D}_{M,B}$  denote the 10 benchmark datasets, and  $\mathcal{D}_{M,W}$  represent the four in-the-wild datasets, as described in Section 4.1, where  $M \in \{S, T\}$  refers to either the source model  $S$  or the target model  $T$ .

The regression model is trained on data  $\mathcal{D}_{S,B}^{\text{train}} \sim \mathcal{D}_{S,B}$  and evaluated on the IID test set  $\mathcal{D}_{S,B}^{\text{test-IID}} \sim \mathcal{D}_{S,B}$ , consisting of 80% and 20% of the data, respectively. Additionally, the model is evaluated on  $\mathcal{D}_{T,B}^{\text{test-IID}}$ ,  $\mathcal{D}_{S,W}$ , and  $\mathcal{D}_{T,W}$ . Below, we detail the formulation of each evaluation setup.

**Case 1: IID Evaluation (Source  $S$ )** The regression model is trained on  $\mathcal{D}_{S,B}^{\text{train}}$  and evaluated on  $\mathcal{D}_{S,B}^{\text{test-IID}}$ . Let  $x_1^S = f(s, o^S)$  represent the similarity between input speech  $s$  and the ASR output  $o^S$ , and  $x_2^S = g(o^S, r)$  represent the agreement with the proxy reference  $r$ , where  $o^S$  is the ASR output produced by the source model  $S$ . The evaluation is formulated as:

$$\mathcal{L}_{\text{IID}}^S = E_{(x_1^S, x_2^S, y) \sim \mathcal{D}_{S,B}^{\text{test-IID}}} [\mathcal{L}(h(x_1^S, x_2^S), y)] \quad (4)$$

**Case 2: IID Evaluation (Target  $T$ )** The regression model trained on  $\mathcal{D}_{S,B}^{\text{train}}$  is evaluated on the IID test set  $\mathcal{D}_{T,B}^{\text{test-IID}}$ . Let  $x_1^T = f(s, o^T)$  represent the similarity between input speech  $s$  and the ASR output  $o^T$ , and  $x_2^T = g(o^T, r)$  represent the agreement with the proxy reference  $r$ , where  $o^T$  is the ASR output produced by the target model  $T$ . The evaluation is expressed as:

$$\mathcal{L}_{\text{IID}}^T = E_{(x_1^T, x_2^T, y) \sim \mathcal{D}_{T,B}^{\text{test-IID}}} [\mathcal{L}(h(x_1^T, x_2^T), y)] \quad (5)$$

**Case 3: OOD Evaluation (Source  $S$ )** The regression model trained on  $\mathcal{D}_{S,B}^{\text{train}}$  is evaluated on the out-of-distribution set  $\mathcal{D}_{S,W}$ . Let  $x_1^S = f(s, o^S)$  represent the similarity between the input speech  $s$  and the ASR output  $o^S$ , and  $x_2^S = g(o^S, r)$  represent the agreement with the proxy reference  $r$ , where  $o^S$  is the ASR output produced by the source model  $S$ . The evaluation is defined as:

$$\mathcal{L}_{\text{OOD}}^S = E_{(x_1^S, x_2^S, y) \sim \mathcal{D}_{S,W}} [\mathcal{L}(h(x_1^S, x_2^S), y)] \quad (6)$$

**Case 4: OOD Evaluation (Target  $T$ )** The regression model trained on  $\mathcal{D}_{S,B}^{\text{train}}$  is evaluated on the out-of-distribution set  $\mathcal{D}_{T,W}$ , using the ASR output produced by the target model  $T$ . Let  $x_1^T = f(s, o^T)$  represent the similarity between the input speech  $s$  and the ASR output  $o^T$ , and  $x_2^T = g(o^T, r)$  represent the agreement with the proxy reference  $r$ , where  $o^T$  is the ASR output produced by the target model  $T$ . The evaluation is expressed as:

$$\mathcal{L}_{\text{OOD}}^T = E_{(x_1^T, x_2^T, y) \sim \mathcal{D}_{T,W}} [\mathcal{L}(h(x_1^T, x_2^T), y)] \quad (7)$$

**Note.** For computational feasibility, the primary experiments train the regression model on 9 out

of the 10 datasets in  $\mathcal{D}_{S,B}^{\text{train}}$  and evaluate it on the remaining dataset, as well as on all four datasets in  $\mathcal{D}_{S,B}^{\text{OOD}}$ . This process is repeated for each dataset in  $\mathcal{D}_{S,B}^{\text{train}}$ , ensuring robust evaluation across various testing conditions. No examples from  $\mathcal{D}_{M,\text{OOD}}$  are used at any stage for training the regression model.

## 4 Experiments

In this section, we present the experimental setup used to evaluate our ASR metrics approximation tool. We describe the datasets, models, and regression pipeline used in our experiments, highlighting the diversity of ASR systems and testing conditions.

### 4.1 Datasets

To evaluate the robustness and generalizability of our ASR metrics approximation tool, we use datasets sourced from multiple distributions, divided into two types: **Standard Benchmark** and **Wild Challenge** datasets. Below we describe the datasets and provide additional details in Appendix A.2 Table 3.

**Standard Benchmark Datasets.** We include widely used datasets representing diverse domains and acoustic conditions. *LibriSpeech* (Panayotov et al., 2015) provides 1,000 hours of English read audiobooks, covering both clean and noisy conditions. *TED-LIUM* (Rousseau et al., 2014) consists of TED talks from 2,000 speakers. *GigaSpeech* (Chen et al., 2021) spans audiobooks, podcasts, and YouTube, incorporating both read and spontaneous speech. *SPGISpeech* (Technologies, 2021) features 5,000 hours of earnings calls with a focus on orthographic accuracy. *Common Voice* (Ardila et al., 2020) is a multilingual, crowdsourced corpus with diverse accents. *Earnings22* (Rio et al., 2022) provides 119 hours of accented, real-world earnings calls. Additional datasets include *AMI (IHM)* (Carletta et al., 2005), with 100 hours of English meeting recordings from non-native speakers, and *People’s Speech* (Galvez et al., 2021), emphasizing inclusivity and linguistic diversity. *SLUE-VoXCeleb* (Shon et al., 2022) contains conversational voice snippets, capturing diverse speaking styles and emotions.

**Wild Datasets.** The wild set focuses on real-world variability and challenging scenarios. *Pri-mock57* (Papadopoulos Korfiatis et al., 2022) includes telemedicine consultations with diverse accents, ages, and scenarios, recorded by clinicians

371 and actors. *VoxPopuli Accented* (Wang et al., 2021) 422  
372 contains multilingual speeches from European Parli- 423  
373 ament recordings, rich in named entities. *AT-* 424  
374 *COsim* (Jan van Doorn, 2023) features 10 hours of 425  
375 non-native English speech from air traffic control 426  
376 simulations with clean utterance-level transcrip- 427  
377 tions. Additionally, we include a noisy subset of  
378 *LibriSpeech* (Panayotov et al., 2015), which reflects  
379 challenging real-world conditions.

## 380 4.2 Models

381 We evaluate our ASR metrics approximation for 429  
382 a range of state-of-the-art ASR models, put into 430  
383 three categories based on their architecture and 431  
384 functionality. Below we describe the datasets and 432  
385 provide additional details in Appendix A.2 and in 433  
386 Table 4. 434

387 **Encoder-Decoder Models.** We include multiple 435  
388 encoder-decoder families of models capable of per- 436  
389 forming ASR tasks in a zero-shot setting. More 437  
390 specifically, we include whisper (Radford et al., 438  
391 2023) and distil-whisper (Gandhi et al., 2023) 439  
392 models which perform really well across diverse 440  
393 testing settings. We also include seamless (Com- 441  
394 munication et al., 2023; Seamless Communica- 442  
395 tion, 2023; Barrault et al., 2025), SpeechT5 (Ao 443  
396 et al., 2022) which are unified encoder-decoder 444  
397 framework for tasks such as ASR, speech synthe- 445  
398 sis, translation, and voice conversion. MMS (Pratap 446  
399 et al., 2023) supports hundreds of languages and 447  
400 excels in resource-constrained scenarios. *Moon-* 448  
401 *shine* (2) (Jeffries et al., 2024), a lightweight and 449  
402 efficient model, is designed for edge deployments 450  
403 with strong performance. 451

404 **NeMo-ASR Models.** We use multiple models 452  
405 from the NeMo-ASR (Gulati et al., 2020; Variani 453  
406 et al., 2020; Noroozi et al., 2024; Tang et al., 2023; 454  
407 Harper et al.) toolkit by NVIDIA. These models 455  
408 include architectures such as Canary and Parakeet, 456  
409 which use highly efficient speech encoders like 457  
410 Fast-Conformer (Rekesh et al., 2023) in combina- 458  
411 tion with various decoders (CTC, RNN-T, TDT) 459  
412 and Conformer-CTC (Guo et al., 2021), making 460  
413 them suitable for a wide range of ASR tasks. In our 461  
414 work, we evaluate 11 models from the NeMo-ASR 462  
415 toolkit. 463

416 **Encoder-Only and Decoder-Only Models.** We 464  
417 include self-supervised encoder-only models 465  
418 and their derivatives, as well as decoder-only 466  
419 models like SpeechLLM. Specifically, we use 467  
420 Wav2Vec2 (Schneider et al., 2019; Baevski 468  
421 et al., 2020), *HuBERT* (Hsu et al., 2021), and

*Data2Vec* (Baevski et al., 2022). Additionally, 422  
we include speech language models like *Speech-* 423  
*LLM* (Rajaa and Tushar), which combines speech 424  
embeddings with language models to predict meta- 425  
data such as speaker attributes, emotions, and ac- 426  
cents, offering robust multimodal capabilities. 427

## 428 4.3 Experimental Setup

429 We evaluate all models listed in Section 4.2 on 430  
1000 examples sampled randomly from the *test* 431  
split of each dataset, as described in Section 4.1. 432  
Since all models are trained at a 16 kHz sampling 433  
rate, we (re)sample the audio inputs accordingly. 434  
For ASR, we employ greedy decoding without us- 435  
ing specialized decoding strategies, and all other 436  
parameters are default unless otherwise specified. 437  
Orthographic transcriptions undergo basic text post- 438  
processing before computing ASR metrics, using 439  
the implementation from whisper<sup>1</sup>. We obtain all 440  
models from Huggingface Hub<sup>2</sup> and implement 441  
the ASR pipeline using the Transformers (Wolf 442  
et al., 2020) library. 443

444 For multimodal embeddings, we use 445  
SONAR (Duquenne et al., 2023), a 1024- 446  
dimensional sentence-level multilingual 447  
embedding model. Specifically, we utilize 448  
text\_sonar\_basic\_encoder for text encoding 449  
and speech\_sonar\_basic\_encoder for speech 450  
encoding. These encoders provide unified rep- 451  
resentations, enabling text reconstruction from 452  
speech. 453

454 The regression framework uses a stacking en- 455  
semble with base regressors and a final estima- 456  
tor. Hyperparameter tuning is performed with 457  
RandomizedSearchCV to minimize MAE. The 458  
model is trained on 9 benchmark datasets and eval- 459  
uated on the remaining benchmark dataset and 460  
four in-the-wild datasets. This process is repeated 461  
for all 10 benchmark datasets. Additional details 462  
of the regression pipeline are provided in Section 3 463  
and low-level details in Appendix A.4.1. 464

465 We conduct ASR experiments on a single 466  
A100/H100 GPU, while the regression model train- 467  
ing runs on CPUs. Although ASR time and mem- 468  
ory consumption depend on the model size, em- 469  
bedding extraction for 1000 audio-text pairs takes 470  
approximately one minute on a single consumer-  
grade GPU without parallelization or additional  
efficiency measures. Appendix A.4 provides fur-  
ther experimental setup details.

<sup>1</sup><https://bit.ly/enormwhisper>

<sup>2</sup><https://huggingface.co/models>

**Baselines.** The recent literature directly aligned with our approach is limited. For instance, eWER (Ali and Renals, 2018) and eWER2 (Ali and Renals, 2020) estimate error rates based on the input signal, which differs from our approach. In contrast, we incorporate the model’s output transcript into the error rate approximation function. The most closely related recent works are WERBERT (Sheshadri et al., 2021a) and eWER3 (Chowdhury and Ali, 2023), which share a similar overall architecture. Both use encoders for text, speech, and other modalities, followed by a regression model trained in an end-to-end setting. Since eWER3 is the more recent of the two, we use it as our baseline. In eWER3, the speech encoder is *wav2vec2* (Baevski et al., 2020), and the text encoder is *roberta-base* (Liu et al., 2019), with a regression model trained on top while both encoders remain frozen. Given the unavailability of public code or pretrained models for evaluation, we implement eWER3 with some modifications to ensure a fair comparison. Specifically, we extract features from both encoders and apply PCA for dimensionality reduction on each modality before training our regression pipeline. For both speech and text, we experiment with 32 and 64 PCA components (referred to as *nc* in Table 2).

## 5 Results

We conduct experiments using two dataset categories: standard benchmarks and in-the-wild, as described in Section 4.1. For each ASR model, a leave-one-out strategy is used, training the regression model on 9 benchmark datasets and testing it on the remaining benchmark dataset and all four in-the-wild datasets. This ensures comprehensive evaluation exclusively on out-of-domain data. Additionally, in-domain testing is included in ablation studies, as detailed in Section 5.3. The regression model is trained to predict absolute error counts (word and character levels), which are normalized by the reference length to compute approximate error rates (*aWER* and *aCER*). We also train regression models to directly predict WER and CER.

### 5.1 Evaluation on In-the-Wild Datasets

The wild datasets provide a realistic testbed for evaluating the regression model’s ability to approximate error rates under real-world conditions. The results are presented in Table 1. High-performing

Model	LS_Noise	Primock57	ATCosim	VP_Acc
w2v2-ls	8.8/10.2	32.8/35.6	43.0/49.5	20.4/26.4
can-1b	4.1/6.4	16.2/13.4	30.4/35.5	23.2/12.1
d2v-base	14.9/16.4	39.6/41.7	66.0/71.2	28.4/33.8
d2v-large	7.2/8.6	28.3/30.7	44.0/51.1	21.4/26.5
distil-l-v2	7.3/9.2	18.3/13.0	69.5/66.7	14.9/14.5
distil-l-v3	6.1/8.3	18.4/12.9	69.0/63.6	14.8/14.0
distil-s.en	9.1/10.6	19.3/14.7	74.9/69.1	14.6/14.7
sm4t-l	11.2/12.3	41.7/37.8	75.0/82.5	29.3/19.9
sm4t-m	14.9/15.6	44.1/39.7	54.6/60.4	30.5/22.5
hub-l-ls-ft	7.3/8.8	29.5/32.0	50.4/56.9	21.4/26.6
hub-xl-ls-ft	6.8/8.3	31.1/32.9	46.7/53.0	21.8/27.7
mms-1b-a	9.5/11.1	36.2/34.4	63.4/71.8	29.9/23.8
mms-1b-f102	24.0/24.9	70.2/67.8	93.4/99.0	39.4/38.2
moon-b	11.3/12.4	19.9/18.5	65.5/66.2	17.1/20.8
moon-t	15.5/17.4	29.2/29.5	62.9/68.5	22.1/26.2
par-ctc-0.6b	4.6/7.4	16.3/13.8	32.9/42.9	16.3/13.8
par-ctc-1.1b	4.5/6.9	16.6/14.1	30.9/39.9	16.4/12.4
par-rnnt-0.6b	3.8/6.9	16.3/13.2	31.6/41.8	17.3/12.6
par-rnnt-1.1b	3.5/6.1	14.6/13.3	27.3/37.6	18.1/10.4
par-tdt-1.1b	3.4/6.0	13.5/13.2	28.3/35.7	17.9/10.2
pkt-ctc-110m	6.1/8.6	16.7/13.0	39.9/42.4	19.2/12.5
sm4t-v2-l	7.2/8.4	34.6/31.7	52.4/57.6	33.8/24.5
spchllm-1.5B	15.3/16.6	42.0/41.8	121.1/125.4	57.0/59.3
spchllm-2B	13.9/15.6	39.4/40.3	60.6/64.1	39.2/44.1
stt-cfc-l	5.8/6.8	16.1/17.6	35.9/38.0	18.6/11.5
stt-cfc-s	9.7/11.2	22.2/24.6	43.7/47.7	16.4/15.6
stt-fc-cfc-l	6.8/10.0	17.6/23.9	34.9/47.6	18.9/13.3
stt-fc-td-l	6.0/8.8	17.0/20.6	34.5/46.5	21.1/15.1
w2v2-960h	17.4/18.5	44.7/47.1	68.4/74.0	29.9/36.5
w2v2-crelpos	5.9/7.4	28.5/30.3	47.2/54.0	22.4/26.7
w2v2-crope	6.6/8.1	31.7/33.4	49.8/56.9	21.9/26.3
w2v2-l-960h	11.6/12.6	37.8/40.2	66.4/72.7	26.3/33.3
w2v2-l-lv60-s	7.8/9.4	33.1/35.5	40.5/48.8	19.3/24.9
w2v2-l-rft-ls	10.0/11.5	32.2/34.6	48.9/55.7	22.0/28.6
whisper-l	6.2/8.1	18.8/13.9	65.3/66.9	18.7/15.9
whisper-l-v2	5.4/6.6	19.0/13.1	64.8/74.8	20.0/18.1
whisper-l-v3	4.6/5.9	18.7/12.0	64.7/73.9	19.2/18.1
whisper-l-v3-t	4.9/6.0	18.5/12.3	66.0/72.5	24.3/23.2
whisper-m.en	6.5/7.9	19.5/14.0	66.2/73.8	27.6/26.4
whisper-s.en	8.2/9.7	20.0/15.1	67.1/73.8	17.3/17.5
whisper-tiny	18.5/20.7	30.0/26.6	97.6/102.5	29.8/33.2

Table 1: Actual and approximated WER ( $\downarrow$ ), separated by a slash, on out-of-distribution wild datasets. The regression model is trained independently for each ASR model on standard benchmarks, making the wild datasets out-of-distribution. Model names are shortened due to space. See Table 7 for full names.

models, like *canary-1b*, demonstrate strong agreement between predicted and actual error rates. For example, on VP\_Accented, *canary-1b* achieves a WER of 23.2% and an *aWER* of 12.1%, with a minimal difference of 1.1%. On Primock57, a clinical consultation dataset, the model shows robustness with a WER of 16.2% and an *aWER* of 13.4%, highlighting its effective generalization across diverse and domain-specific contexts.

Models like *data2vec-audio-large-960h* also maintain strong performance, with deviations consistently under 2% on various datasets. For example, on LibriSpeech-test-noise, the model’s actual WER is 7.2% while the approximated *aWER* is 8.6%, showcasing its reliability in noisy con-

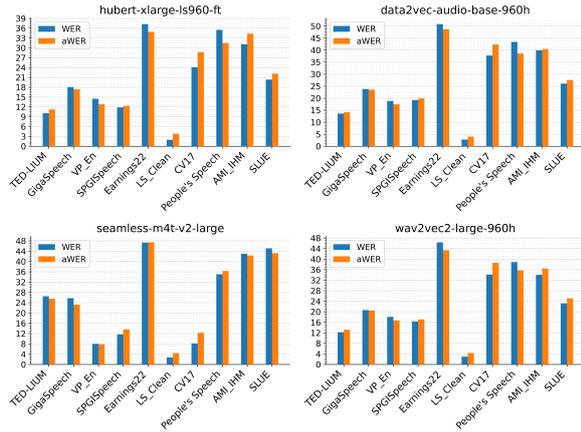


Figure 1: Actual and approximated WER for four models across standard benchmark.

ditions. Even on acoustically complex datasets like ATCOsim, where the WER is 44.0% and the  $aWER$  is 51.1%, the model exhibits a reasonable alignment between approximated and actual error rates.

In contrast, models with high actual error rates, such as *mms-1b-fl102*, show slightly larger deviations, particularly on datasets with challenging conditions. For instance, on ATCOsim, the WER is 93.4% and the  $aWER$  is 99.0%, resulting in a significant deviation of 5.6%, the highest observed across all in-the-wild datasets. Similarly, on Pri-mock57, where the WER is 70.2% and the  $aWER$  is 67.8%, the approximation also struggles to align due to the inherently high error rates. This highlights that extreme error cases often correspond to semantically nonsensical outputs, where the distinction between high and extremely high error rates becomes less relevant.

## 5.2 Evaluation on Benchmark Datasets

We summarize results on 10 standard benchmark datasets in Appendix A.5 Tables 8 and 9. Each table reports actual WER/CER alongside the approximated WER/CER (denoted by  $aWER/aCER$ ).

Overall, models such as *parakeet-tdt-1.1b* and *whisper-large-v3* show relatively small differences between WER and  $aWER$ , indicating reliable approximations. For instance, the actual WER for *whisper-large-v3* on **AMI\_IHM** is 19.0% compared to  $aWER$  of 17.1%, that’s only a 1.9% gap. Conversely, some challenging datasets (e.g., **CV11** and **Earnings22**) reveal larger discrepancies for specific models, particularly those with higher overall error rates. For example, *mms-1b-fl102* exhibits a wide WER/ $aWER$  gap on **Earnings22**, suggest-

ing difficulty handling accented or domain-specific speech.

In general, high-performing ASR models demonstrate small WER– $aWER$  gaps, indicating that it’s easy to approximate when error rates are low. However, models with higher WERs or faced with more acoustically or linguistically challenging test sets tend to show wider divergences between WER and  $aWER$ . Despite these variations, most results remain within a reasonable margin, highlighting the robustness of our approximation model in diverse out-of-distribution scenarios.

These results underscore the critical role of model quality in achieving reliable approximations. The approximation framework remains effective for high-performing models, while deviations tend to increase in cases of semantically divergent or poorly structured outputs, reflecting the inherent challenges in approximating errors for low-performing systems.

## 5.3 Ablation

We conduct ablation experiments to evaluate the robustness of the approximation model and the contributions of its individual components. Using the evaluation setup outlined in Section 3.4, we select *data2vec-audio-base-960h* as the source model ( $S$ ) and *wav2vec2-base-960h* as the target model ( $T$ ). The results are summarized in Table 2, where IID results correspond to Case-I 3.4, and  $D$ ,  $M$ , and  $D + M$  under OOD represent Case II 3.4, Case-III 3.4, and Case-IV 3.4, respectively. The reference model’s  $r$  value represents the average WER across all datasets. We include reference models with varying  $r$  values, such as *whisper-large-v3* ( $r = 17.8$ ), *whisper-medium.en* ( $r = 20.1$ ), *whisper-tiny* ( $r = 33.4$ ), and *mms-1b-fl102* ( $r = 51.0$ ).

The results in Table 2 demonstrate the importance of proxy references in improving the regression model’s performance. Training without proxy references (*w/o PR*) significantly increases the mean absolute error (MAE) across all conditions. For instance, the IID MAE increases from 1.03 (Base) to 3.13, and the OOD  $D + M$  MAE rises from 1.07 (Base) to 3.33, highlighting the essential role of proxy references in approximation.

Increasing the number of high-quality proxy references ( $MPR$ ) further reduces errors. Under IID conditions, the MAE decreases from 1.00 with  $n = 2$  to 0.93 with  $n = 5$ . Similarly, in OOD  $D + M$ , the error drops from 1.06 ( $MPR, n = 2$ )

Method	IID	D	OOD M	D + M
eWER3(nc=32)	2.03 <sup>0.07</sup>	2.09 <sup>0.04</sup>	2.06 ± 0.03	2.12 <sup>0.04</sup>
eWER3(nc=64)	1.98 <sup>0.06</sup>	2.07 <sup>0.05</sup>	2.00 <sup>0.04</sup>	2.09 <sup>0.05</sup>
Base	1.03 <sup>0.03</sup>	1.05 <sup>0.01</sup>	1.03 <sup>0.02</sup>	1.07 <sup>0.01</sup>
w/o S	1.04 <sup>0.03</sup>	1.05 <sup>0.01</sup>	1.04 <sup>0.03</sup>	1.05 <sup>0.01</sup>
w/o PR	3.13 <sup>0.07</sup>	3.22 <sup>0.02</sup>	3.23 <sup>0.05</sup>	3.33 <sup>0.02</sup>
w/ MPR (n=2)	1.00 <sup>0.02</sup>	1.04 <sup>0.02</sup>	0.99 <sup>0.02</sup>	1.06 <sup>0.02</sup>
w/ MPR (n=3)	0.96 <sup>0.02</sup>	0.97 <sup>0.01</sup>	0.95 <sup>0.02</sup>	0.99 <sup>0.01</sup>
w/ MPR (n=4)	0.95 <sup>0.02</sup>	0.96 <sup>0.02</sup>	0.94 <sup>0.02</sup>	0.98 <sup>0.02</sup>
w/ MPR (n=5)	0.93 <sup>0.02</sup>	0.93 <sup>0.01</sup>	0.92 <sup>0.02</sup>	0.95 <sup>0.01</sup>
w/MPR (n=10)	0.90 <sup>0.02</sup>	0.93 <sup>0.01</sup>	0.88 <sup>0.02</sup>	0.95 <sup>0.01</sup>
w/MPR (n=20)	0.89 <sup>0.02</sup>	0.96 <sup>0.02</sup>	0.87 <sup>0.02</sup>	0.96 <sup>0.02</sup>
w/ mMPR (n=3)	0.98 <sup>0.02</sup>	0.96 <sup>0.02</sup>	0.97 <sup>0.02</sup>	0.98 <sup>0.02</sup>
w/ mMPR (n=5)	0.94 <sup>0.02</sup>	0.94 <sup>0.02</sup>	0.93 <sup>0.01</sup>	0.96 <sup>0.02</sup>
w/mMPR (n=10)	0.92 <sup>0.02</sup>	0.94 <sup>0.02</sup>	0.91 <sup>0.02</sup>	0.96 <sup>0.02</sup>
w/mMPR (n=20)	1.04 <sup>0.02</sup>	1.05 <sup>0.01</sup>	1.02 <sup>0.02</sup>	1.04 <sup>0.01</sup>
Base (r=17.8)	1.31 <sup>0.04</sup>	1.44 <sup>0.02</sup>	1.31 <sup>0.04</sup>	1.40 <sup>0.01</sup>
Base (r=20.1)	1.36 <sup>0.04</sup>	1.36 <sup>0.01</sup>	1.34 <sup>0.03</sup>	1.34 <sup>0.01</sup>
Base (r=33.4)	1.55 <sup>0.04</sup>	1.69 <sup>0.02</sup>	1.55 <sup>0.04</sup>	1.63 <sup>0.02</sup>
Base (r=51.0)	2.03 <sup>0.02</sup>	2.10 <sup>0.01</sup>	2.08 <sup>0.05</sup>	2.09 <sup>0.01</sup>
w/o S (r=17.8)	1.47 <sup>0.04</sup>	1.56 <sup>0.01</sup>	1.48 <sup>0.04</sup>	1.54 <sup>0.01</sup>
w/o S (r=20.1)	1.55 <sup>0.02</sup>	1.50 <sup>0.01</sup>	1.55 <sup>0.03</sup>	1.50 <sup>0.01</sup>
w/o S (r=33.4)	1.79 <sup>0.07</sup>	1.89 <sup>0.02</sup>	1.78 <sup>0.06</sup>	1.82 <sup>0.02</sup>
w/o S (r=51.0)	2.23 <sup>0.02</sup>	2.24 <sup>0.01</sup>	2.28 <sup>0.04</sup>	2.21 <sup>0.01</sup>

Table 2: Mean absolute error ( $\downarrow$ ) between predicted word error count and actual error count (in absolute terms) across different configurations. PR - Proxy Reference, S - Similarity, MPR - Multiple PR, D - Data, M - Model. The OOD results are averaged across four wild datasets.  $n$  is the number of proxy references.  $r\downarrow$  is the average WER for proxy reference across 14 datasets. Superscript represents the standard deviation across five runs.

to 0.95 (*MPR*,  $n = 5$ ), demonstrating that multiple high-quality references enhance model robustness.

The quality of references, quantified by the  $r$ -value, also plays a critical role. For example, in IID conditions, the MAE increases from 1.31 for  $r = 17.8$  to 2.03 for  $r = 51.0$ . A similar trend is observed in OOD  $D + M$ , where the MAE rises from 1.40 ( $r = 17.8$ ) to 2.09 ( $r = 51.0$ ). The absence of similarity (*w/o S*) combined with low-quality proxies further degrades performance, underscoring the importance of both high-quality references and similarity measures. These trends are similarly observed for character-level error count approximation, as detailed in Appendix Table 6.

**Scaling Training Data for Regression.** To evaluate the impact of training data size on the regression model, we scale the data from 1K to 10K examples in increments of 1K. As shown in Figure 2, the

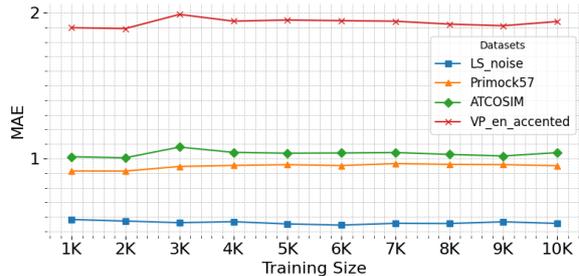


Figure 2: Mean absolute error ( $\downarrow$ ) between predicted and actual word error counts across varying training data sizes for the regression model. The model was trained on 10 standard benchmarks and evaluated on four in-the-wild test sets.

model’s performance does not exhibit a clear trend with increasing training data size. Some datasets show slight improvements with more data; others show minimal improvement. This suggests that the regression model is largely agnostic to the size of the training data. In fact, it appears that a relatively small dataset of just 1,000 examples is sufficient to train a robust approximation model. This underscores the model’s ability to generalize effectively with limited data, making it an efficient choice for scenarios with constrained datasets.

## 6 Conclusion

We present a framework for approximating ASR metrics, demonstrating its effectiveness in generalizing to unseen, in-the-wild, and challenging conditions. Our results show that the model performs well with absolute error counts, consistently outperforming strong baseline, with error rates remaining relatively low. We show that our proposed method achieves consistent performance across 40 ASR models and 14 evaluation setups, including both standard benchmarks and domain-specific conditions. The trained regression model can be efficiently used to approximate ASR metrics, particularly in data-constrained environments, such as critical domains with limited labeled data. In summary, our work bridges the gap between theoretical advancements and real-world applications, paving the way for more reliable and scalable ASR systems. While in this work, we explore the impact of training data size within a single language, future work will focus on extending this framework to support multiple languages and exploring language-agnostic ASR metric approximation.

## 7 Limitations

In this work, we introduced a framework for approximating ASR metrics, evaluated across various ASR models and datasets. Despite the promising results, there are several limitations to consider.

**Evaluation.** While our evaluation setup is comprehensive, consisting of over 40 models and 14 datasets representing various acoustic and linguistic conditions such as natural noise, dialects, and accents—far surpassing previous works—we have not explored more nuanced conditions such as gender, non-native speech, and approximation across various age groups. Additionally, while the framework has demonstrated strong performance in approximating ASR metrics across multiple datasets, its generalization to highly diverse or extreme real-world conditions might still require further investigation.

**Language.** Furthermore, the evaluation is currently limited to a single language; expanding this framework to multiple languages or achieving language-agnostic ASR metric approximation remains an important direction for future work.

**Compute.** While, unlike previous works, our final approximator is a simple machine learning model that does not require GPUs to run, we do utilize a single GPU for multimodal embedding extraction, which could be performed on any consumer-grade GPU.

## 8 Ethics Statement

**Data Collection and Release.** The datasets used in our experiments consist of publicly available ASR data from both benchmark and in-the-wild sources, as detailed in Section 4.1. We ensure that the use of these datasets aligns with the principles of fair use, specifically in a non-commercial academic context or as specified in their original license. All datasets are openly accessible, and no private or confidential data is included in this work to the best of our knowledge.

**Intended Use.** By enabling the approximation of ASR performance metrics with minimal data, our work has the potential to impact applications in domains with limited data availability, such as healthcare, emergency response, and low-resource language research. We believe our approach will foster further research in scalable, low-cost ASR systems with comprehensive evaluation, benefiting industries and research areas that serve underrepre-

sented or resource-limited populations.

**Potential Misuse and Bias.** While our regression model has demonstrated effectiveness in approximating ASR metrics, it is important to consider potential misuse and bias. Given that our model is trained on diverse datasets, including those with various linguistic, acoustic, and demographic variations, there is a risk that the model may inherit biases present in the data, particularly with respect to accents, dialects, and socio-linguistic factors. Additionally, as our model approximates error rates, it could be misused in applications where the approximation may not be sufficient for real-world critical tasks. We recommend cautious deployment and further evaluation in sensitive applications, especially those where fairness and accuracy are critical.

## References

- Ahmed Ali and Steve Renals. 2018. [Word error rate estimation for speech recognition: e-WER](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 20–24, Melbourne, Australia. Association for Computational Linguistics.
- Ahmed M. Ali and Steve Renals. 2020. [Word error rate estimation without asr output: e-wer2](#). In *Inter-speech*.
- Junyi Ao, Rui Wang, Long Zhou, Chengyi Wang, Shuo Ren, Yu Wu, Shujie Liu, Tom Ko, Qing Li, Yu Zhang, Zhihua Wei, Yao Qian, Jinyu Li, and Furu Wei. 2022. [Specht5: Unified-modal encoder-decoder pre-training for spoken language processing](#). *Preprint*, arXiv:2110.07205.
- Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus. <https://commonvoice.mozilla.org/en/datasets>. Accessed: [Insert Date].
- Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. 2022. [data2vec: A general framework for self-supervised learning in speech, vision and language](#). *Preprint*, arXiv:2202.03555.
- Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. [wav2vec 2.0: A framework for self-supervised learning of speech representations](#). *Preprint*, arXiv:2006.11477.
- Loïc Barrault, Yu-An Chung, Mariano Coria Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, Christopher Klaiber, Pengwei Li,

774	Daniel Licht, Jean Maillard, Alice Rakotoarison,	Miguel Ángel Del-Agua, Adrià Giménez, Albert San-	834
775	Kaushik Ram Sadagopan, Guillaume Wenzek, Ethan	chis, Jorge Civera, and Alfons Juan. 2018. <a href="#">Speaker-</a>	835
776	Ye, Bapi Akula, Peng-Jen Chen, Naji El Hachem,	<a href="#">adapted confidence measures for asr using deep bidi-</a>	836
777	Brian Ellis, Gabriel Mejia Gonzalez, Justin Haaheim,	<a href="#">rectional recurrent neural networks.</a>	837
778	Prangthip Hansanti, Russ Howes, Bernie Huang,	<i>IEEE/ACM</i>	838
779	Min-Jae Hwang, Hirofumi Inaguma, Somya Jain,	<i>Transactions on Audio, Speech, and Language Pro-</i>	839
780	Elahe Kalbassi, Amanda Kallet, Iliia Kulikov, Jan-	<i>cessing</i> , 26(7):1198–1206.	
781	ice Lam, Daniel Li, Xutai Ma, Ruslan Mavlyut-	Amandeep Singh Dhanjal and Williamjeet Singh. 2024.	840
782	ov, Benjamin Peloquin, Mohamed Ramadan, Abi-	A comprehensive survey on automatic speech recog-	841
783	nesh Ramakrishnan, Anna Sun, Kevin Tran, Tuan	nition using neural networks. <i>Multimedia Tools and</i>	842
784	Tran, Igor Tufanov, Vish Vogeti, Carleigh Wood,	<i>Applications</i> , 83(8):23367–23412.	843
785	Yilin Yang, Bokai Yu, Pierre Andrews, Can Bali-	Paul-Ambroise Duquenne, Holger Schwenk, and Benoît	844
786	oglu, Marta R. Costa-jussà, Onur Çelebi, Maha El-	Sagot. 2023. Sonar: sentence-level multimodal and	845
787	bayad, Cynthia Gao, Francisco Guzmán, Justine Kao,	language-agnostic representations. <i>arXiv e-prints</i> ,	846
788	Ann Lee, Alexandre Mourachko, Juan Pino, Sravya	pages arXiv–2308.	847
789	Popuri, Christophe Ropers, Safiyyah Saleem, Hol-	Zied Elloumi, Laurent Besacier, Olivier Galibert, Juli-	848
790	ger Schwenk, Paden Tomasello, Changhan Wang,	ette Kahn, and Benjamin Lecouteux. 2018. <a href="#">Asr</a>	849
791	Jeff Wang, Skyler Wang, and SEAMLESS Com-	<a href="#">performance prediction on unseen broadcast pro-</a>	850
792	munication Team. 2025. <a href="#">Joint speech and text machine</a>	<a href="#">grams using convolutional neural networks.</a>	851
793	<a href="#">translation for up to 100 languages.</a>	<i>Preprint</i> ,	852
794	<i>Nature</i> ,	arXiv:1804.08477.	
	637(8046):587–593.		
795	Jean Carletta, Simone Ashby, Sebastien Bourban, Mike	Daniel Galvez, Greg Diamos, Juan Ciro, Juan Felipe	853
796	Flynn, Mael Guillemot, Thomas Hain, Jaroslav	Cerón, Keith Achorn, Anjali Gopi, David Kanter,	854
797	Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa	Maximilian Lam, Mark Mazumder, and Vijay Janapa	855
798	Kronenthal, et al. 2005. The ami meeting corpus: A	Reddi. 2021. <a href="#">The people’s speech: A large-scale</a>	856
799	pre-announcement. <a href="https://groups.inf.ed.ac.uk/ami/corpus/">https://groups.inf.ed.ac.</a>	<a href="#">diverse english speech recognition dataset for com-</a>	857
800	<a href="https://groups.inf.ed.ac.uk/ami/corpus/">uk/ami/corpus/</a> . Accessed: [Insert Date].	<a href="#">mercial usage.</a>	858
		<i>Preprint</i> , arXiv:2111.09344.	
801	Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu	Sanchit Gandhi, Patrick von Platen, and Alexander M.	859
802	Du, Wei-Qiang Zhang, Chao Weng, Dan Su, Daniel	Rush. 2023. <a href="#">Distil-whisper: Robust knowledge dis-</a>	860
803	Povey, Jan Trmal, Junbo Zhang, et al. 2021. Gi-	<a href="#">tillation via large-scale pseudo labelling.</a>	861
804	gaspeech: An evolving, multi-domain asr corpus with	<i>Preprint</i> ,	862
805	10,000 hours of transcribed audio. <i>arXiv preprint</i>	arXiv:2311.00430.	
806	<i>arXiv:2106.06909</i> .	Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki	863
		Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo	864
807	Shammur Absar Chowdhury and Ahmed Ali. 2023.	Wang, Zhengdong Zhang, Yonghui Wu, and Ruom-	865
808	<a href="#">Multilingual word error rate estimation: e-wer3.</a>	ing Pang. 2020. <a href="#">Conformer: Convolution-augmented</a>	866
809	<i>Preprint</i> , arXiv:2304.00649.	<a href="#">transformer for speech recognition.</a>	867
		<i>Preprint</i> ,	868
		arXiv:2005.08100.	
810	Seamless Communication, Loïc Barrault, Yu-An Chung,	Pengcheng Guo, Xuankai Chang, Shinji Watanabe, and	869
811	Mariano Cora Meglioli, David Dale, Ning Dong,	Lei Xie. 2021. <a href="#">Multi-speaker asr combining non-</a>	870
812	Paul-Ambroise Duquenne, Hady Elsahar, Hongyu	<a href="#">autoregressive conformer ctc and conditional speaker</a>	871
813	Gong, Kevin Heffernan, John Hoffman, Christopher	<a href="#">chain.</a>	872
814	Klaiber, Pengwei Li, Daniel Licht, Jean Maillard,	<i>Preprint</i> , arXiv:2106.08595.	
815	Alice Rakotoarison, Kaushik Ram Sadagopan, Guil-	Eric Harper, Somshubra Majumdar, Oleksii Kuchaiev,	873
816	laume Wenzek, Ethan Ye, Bapi Akula, Peng-Jen	Li Jason, Yang Zhang, Evelina Bakhturina, Vahid	874
817	Chen, Naji El Hachem, Brian Ellis, Gabriel Mejia	Noroozi, Sandeep Subramanian, Koluguri Nithin,	875
818	Gonzalez, Justin Haaheim, Prangthip Hansanti, Russ	Huang Jocelyn, Fei Jia, Jagadeesh Balam, Xuesong	876
819	Howes, Bernie Huang, Min-Jae Hwang, Hirofumi In-	Yang, Micha Livne, Yi Dong, Sean Naren, and Boris	877
820	aguma, Somya Jain, Elahe Kalbassi, Amanda Kallet,	Ginsburg. <a href="#">NeMo: a toolkit for Conversational AI</a>	878
821	Iliia Kulikov, Janice Lam, Daniel Li, Xutai Ma, Rus-	<a href="#">and Large Language Models.</a>	879
822	lan Mavlyutov, Benjamin Peloquin, Mohamed Ram-	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai,	880
823	adan, Abinesh Ramakrishnan, Anna Sun, Kevin	Kushal Lakhota, Ruslan Salakhutdinov, and Abdel-	881
824	Tran, Tuan Tran, Igor Tufanov, Vish Vogeti, Carleigh	rahman Mohamed. 2021. <a href="#">Hubert: Self-supervised</a>	882
825	Wood, Yilin Yang, Bokai Yu, Pierre Andrews, Can	<a href="#">speech representation learning by masked prediction</a>	883
826	Balioglu, Marta R. Costa-jussà, Onur Celebi, Maha	<a href="#">of hidden units.</a>	884
827	Elbayad, Cynthia Gao, Francisco Guzmán, Justine	<i>Preprint</i> , arXiv:2106.07447.	
828	Kao, Ann Lee, Alexandre Mourachko, Juan Pino,	Shahab Jalalvand, Matteo Negri, Marco Turchi, José G.	885
829	Sravya Popuri, Christophe Ropers, Safiyyah Saleem,	C. de Souza, Daniele Falavigna, and Mohammed	886
830	Holger Schwenk, Paden Tomasello, Changhan Wang,	R. H. Qwaider. 2016. <a href="#">TranscRater: a tool for au-</a>	887
831	Jeff Wang, and Skyler Wang. 2023. <a href="#">Seamlessm4t:</a>	<a href="#">tomatic speech recognition quality estimation.</a>	888
832	<a href="#">Massively multilingual multimodal machine transla-</a>	<i>In</i>	889
833	<a href="#">tion.</a>	<i>Proceedings of ACL-2016 System Demonstrations</i> ,	
	<i>Preprint</i> , arXiv:2308.11596.		

890	pages 43–48, Berlin, Germany. Association for Computational Linguistics.	Vahid Noroozi, Somshubra Majumdar, Ankur Kumar, Jagadeesh Balam, and Boris Ginsburg. 2024. <a href="#">Stateful conformer with cache-based inference for streaming automatic speech recognition</a> . <i>Preprint</i> , arXiv:2312.17279.	943
891			944
892	Jan van Doorn. 2023. <a href="#">atcosim (revision b5839d9)</a> .		945
893	Nat Jeffries, Evan King, Manjunath Kudlur, Guy Nicholson, James Wang, and Pete Warden. 2024. <a href="#">Moonshine: Speech recognition for live transcription and voice commands</a> . <i>Preprint</i> , arXiv:2410.15608.		946
894		Azim Ospanov, Jingwei Zhang, Mohammad Jalali, Xue-nan Cao, Andrej Bogdanov, and Farzan Farnia. 2024. Towards a scalable reference-free evaluation of generative models. <i>arXiv preprint arXiv:2407.02961</i> .	947
895			948
896			949
897	Kaustubh Kalgaonkar, Chaojun Liu, Yifan Gong, and Kaisheng Yao. 2015. <a href="#">Estimating confidence scores on asr results using recurrent neural networks</a> . In <i>2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 4999–5003.		950
898		Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. Librispeech: An asr corpus based on public domain audio books. <a href="http://www.openslr.org/12/">http://www.openslr.org/12/</a> . Accessed: [Insert Date].	951
899			952
900			953
901			954
902			955
903	Mahdie Karbasi and Dorothea Kolossa. 2022. Asr-based speech intelligibility prediction: A review. <i>Hearing Research</i> , 426:108606.	Alex Papadopoulos Korfiatis, Francesco Moramarco, Radmila Sarac, and Aleksandar Savkov. 2022. <a href="#">PriMock57: A dataset of primary care mock consultations</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 588–598, Dublin, Ireland. Association for Computational Linguistics.	956
904			957
905			958
906	Hamza Kheddar, Mustapha Hemis, and Yassine Himeur. 2024. Automatic speech recognition using advanced deep learning approaches: A survey. <i>Information Fusion</i> , page 102422.		959
907			960
908			961
909			962
910	Korbinian Kuhn, Verena Kersken, Benedikt Reuter, Niklas Egger, and Gottfried Zimmermann. 2024. Measuring the accuracy of automatic speech recognition solutions. <i>ACM Transactions on Accessible Computing</i> , 16(4):1–23.	Chanho Park, Hyunsik Kang, and Thomas Hain. 2024. Character error rate estimation for automatic speech recognition of short utterances. In <i>2024 32nd European Signal Processing Conference (EUSIPCO)</i> , pages 131–135. IEEE.	963
911			964
912			965
913			966
914			967
915	Tatiana Likhomanenko, Qiantong Xu, Vineel Pratap, Paden Tomasello, Jacob Kahn, Gilad Avidov, Ronan Collobert, and Gabriel Synnaeve. 2020. Rethinking evaluation in asr: Are our models robust enough? <i>arXiv preprint arXiv:2010.11745</i> .	Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, Sayani Kundu, Ali Elkahky, Zhaoheng Ni, Apoorv Vyas, Maryam Fazel-Zarandi, Alexei Baevski, Yossi Adi, Xiaohui Zhang, Wei-Ning Hsu, Alexis Conneau, and Michael Auli. 2023. <a href="#">Scaling speech technology to 1,000+ languages</a> . <i>Preprint</i> , arXiv:2305.13516.	968
916			969
917			970
918			971
919			972
920	Yu-Yi Lin, Wei-Zhong Zheng, Wei Chung Chu, Ji-Yan Han, Ying-Hsiu Hung, Guan-Min Ho, Chia-Yuan Chang, and Ying-Hui Lai. 2021. A speech command control-based recognition system for dysarthric patients based on deep learning technology. <i>Applied Sciences</i> , 11(6):2477.		973
921			974
922			975
923			976
924			977
925			978
926	Diane J. Litman, Julia B. Hirschberg, and Marc Swerts. 2000. <a href="#">Predicting automatic speech recognition performance using prosodic cues</a> . In <i>1st Meeting of the North American Chapter of the Association for Computational Linguistics</i> .	David Qiu, Qiuqia Li, Yanzhang He, Yu Zhang, Bo Li, Liangliang Cao, Rohit Prabhavalkar, Deepti Bhatia, Wei Li, Ke Hu, Tara N. Sainath, and Ian McGraw. 2021. <a href="#">Learning word-level confidence for subword end-to-end asr</a> . In <i>ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 6393–6397.	979
927			980
928			981
929			982
930			983
931	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. <a href="#">Roberta: A robustly optimized bert pretraining approach</a> . <i>Preprint</i> , arXiv:1907.11692.	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. <a href="#">Robust speech recognition via large-scale weak supervision</a> . <i>Preprint</i> , arXiv:2212.04356.	984
932			985
933			986
934			987
935			988
936	Matteo Negri, Marco Turchi, José G. C. de Souza, and Daniele Falavigna. 2014. <a href="#">Quality estimation for automatic speech recognition</a> . In <i>Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers</i> , pages 1813–1823, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In <i>International conference on machine learning</i> , pages 28492–28518. PMLR.	989
937			990
938			991
939			992
940			993
941			994
942			995
		Bhiksha Raj, Rita Singh, and James Baker. 2011. <a href="#">A paired test for recognizer selection with untranscribed data</a> . In <i>2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5676–5679.	996
			997
		Shangeth Rajaa and Abhinav Tushar. <a href="#">SpeechLLM: Multi-Modal LLM for Speech Understanding</a> .	997

998	Dima Rekesh, Nithin Rao Koluguri, Samuel Kriman, Somshubra Majumdar, Vahid Noroozi, He Huang, Oleksii Hrinchuk, Krishna Puvvada, Ankur Kumar, Jagadeesh Balam, and Boris Ginsburg. 2023. <a href="#">Fast conformer with linearly scalable attention for efficient speech recognition</a> . <i>Preprint</i> , arXiv:2305.05084.	2022-2022 <i>IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 7927–7931. IEEE.	1056 1057 1058
1000	MD Rio, Peter Ha, Quinten McNamara, Corey Miller, and Shipra Chandra. 2022. Earnings-22: A practical benchmark for accents in the wild.	Prakhar Swarup, Roland Maas, Sri Garimella, Sri Harish Mallidi, and Björn Hoffmeister. 2019. <a href="#">Improving asr confidence scores for alexa using acoustic and hypothesis embeddings</a> .	1059 1060 1061 1062
1001	Anthony Rousseau, Paul Deléglise, and Yann Es-tève. 2014. Ted-lium 3: Twice as much data and corpus repartition for experiments on speaker adaptation. <a href="https://lium.univ-lemans.fr/ted-lium3/">https://lium.univ-lemans.fr/ted-lium3/</a> . Accessed: [Insert Date].	Yun Tang, Anna Y. Sun, Hirofumi Inaguma, Xinyue Chen, Ning Dong, Xutai Ma, Paden D. Tomasello, and Juan Pino. 2023. <a href="#">Hybrid transducer and attention based encoder-decoder modeling for speech-to-text tasks</a> . <i>Preprint</i> , arXiv:2305.03101.	1063 1064 1065 1066 1067
1002	Steffen Schneider, Alexei Baeovski, Ronan Collobert, and Michael Auli. 2019. <a href="#">wav2vec: Unsupervised pre-training for speech recognition</a> . <i>Preprint</i> , arXiv:1904.05862.	Kensho Technologies. 2021. Spgispeech: A large-scale, high-quality dataset for speech recognition in financial earnings calls. <a href="https://datasets.kensho.com/datasets/spgispeech">https://datasets.kensho.com/datasets/spgispeech</a> . Accessed: [Insert Date].	1068 1069 1070 1071 1072
1003	Yu-An Chung Mariano Coria Meglioli David Dale Ning Dong Mark Duppenthaler Paul-Ambroise Duquenne Brian Ellis Hady Elsahar Justin Haaheim John Hoffman Min-Jae Hwang Hirofumi Inaguma Christopher Klaiber Ilia Kulikov Pengwei Li Daniel Licht Jean Maillard Ruslan Mavlyutov Alice Rakotoarison Kaushik Ram Sadagopan Abinesh Ramakrishnan Tuan Tran Guillaume Wenzek Yilin Yang Ethan Ye Ivan Evtimov Pierre Fernandez Cynthia Gao Prangthip Hansanti Elahe Kalbassi Amanda Kallet Artyom Kozhevnikov Gabriel Mejia Robin San Roman Christophe Touret Corinne Wong Carleigh Wood Bokai Yu Pierre Andrews Can Bahiloglu Peng-Jen Chen Marta R. Costa-jussà Maha Elbayad Hongyu Gong Francisco Guzmán Kevin Heffernan Somya Jain Justine Kao Ann Lee Xutai Ma Alex Mourachko Benjamin Peloquin Juan Pino Sravya Popuri Christophe Ropers Safiyah Saleem Holger Schwenk Anna Sun Paden Tomasello Changhan Wang Jeff Wang Skyler Wang Mary Williamson Seamless Communication, Loïc Barrault. 2023. Seamless: Multilingual expressive and streaming speech translation.	Ehsan Variani, David Rybach, Cyril Allauzen, and Michael Riley. 2020. <a href="#">Hybrid autoregressive transducer (hat)</a> . <i>Preprint</i> , arXiv:2003.07705.	1073 1074 1075
1004	Akshay Krishna Sheshadri, Anvesh Rao Vijjini, and Sukhdeep Kharbanda. 2021a. <a href="#">WER-BERT: Automatic WER estimation with BERT in a balanced ordinal classification paradigm</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 3661–3672, Online. Association for Computational Linguistics.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. <a href="#">Attention is all you need</a> . <i>Preprint</i> , arXiv:1706.03762.	1076 1077 1078 1079
1005	Akshay Krishna Sheshadri, Anvesh Rao Vijjini, and Sukhdeep Kharbanda. 2021b. <a href="#">Wer-bert: Automatic wer estimation with bert in a balanced ordinal classification paradigm</a> . <i>Preprint</i> , arXiv:2101.05478.	Abdul Waheed, Karima Kadaoui, and Muhammad Abdul-Mageed. 2024. <a href="#">To distill or not to distill? on the robustness of robust knowledge distillation</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12603–12621, Bangkok, Thailand. Association for Computational Linguistics.	1080 1081 1082 1083 1084 1085 1086
1006	Suwon Shon, Ankita Pasad, Felix Wu, Pablo Brusco, Yoav Artzi, Karen Livescu, and Kyu J Han. 2022. Slue: New benchmark tasks for spoken language understanding evaluation on natural speech. In <i>ICASSP</i>	Abdul Waheed, Karima Kadaoui, Bhiksha Raj, and Muhammad Abdul-Mageed. 2025. <a href="#">udistil-whisper: Label-free data filtering for knowledge distillation in low-data regimes</a> . <i>Preprint</i> , arXiv:2407.01257.	1087 1088 1089 1090
1007	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pi-er-ic Cistac, Tim Rault, Rémi Louf, Morgan Funtow-icz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	Changhan Wang, Morgane Riviere, Ann Lee, Anne Wu, Chaitanya Talnikar, Daniel Haziza, Mary Williamson, Juan Pino, and Emmanuel Dupoux. 2021. <a href="#">VoxPop-uli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 993–1003, Online. Association for Computational Linguistics.	1091 1092 1093 1094 1095 1096 1097 1098 1099 1100 1101
1008		Haolan Wang, Amin Edraki, Wai-Yip Chan, Iván López-Espejo, and Jesper Jensen. 2024. No-reference speech intelligibility prediction leveraging a noisyy-speech asr pre-trained model. In <i>Proc. Interspeech 2024</i> , pages 3849–3853.	1102 1103 1104 1105 1106
1009			
1010			
1011			
1012			
1013			
1014			
1015			
1016			
1017			
1018			
1019			
1020			
1021			
1022			
1023			
1024			
1025			
1026			
1027			
1028			
1029			
1030			
1031			
1032			
1033			
1034			
1035			
1036			
1037			
1038			
1039			
1040			
1041			
1042			
1043			
1044			
1045			
1046			
1047			
1048			
1049			
1050			
1051			
1052			
1053			
1054			
1055			

1112 Teven Le Scao, Sylvain Gugger, Mariama Drame,  
1113 Quentin Lhoest, and Alexander M. Rush. 2020. [Huggingface’s transformers: State-of-the-art natural language processing](#). *Preprint*, arXiv:1910.03771.  
1114  
1115

1116 Su-Youn Yoon, Lei Chen, and Klaus Zechner. 2010.  
1117 [Predicting word accuracy for the automatic speech recognition of non-native speech](#). In *Interspeech 2010*, pages 773–776.  
1118  
1119

1120 Kamer Ali Yuksel, Thiago Ferreira, Ahmet Gunduz, Mohamed Al-Badrashiny, and Golara Javadi.  
1121 2023a. [A reference-less quality metric for automatic speech recognition via contrastive-learning of a multi-language model with self-supervision](#). *Preprint*,  
1122 arXiv:2306.13114.  
1123  
1124

1126 Kamer Ali Yuksel, Thiago Ferreira, Golara Javadi, Mohamed El-Badrashiny, and Ahmet Gunduz. 2023b.  
1127 [Norefer: a referenceless quality metric for automatic speech recognition via semi-supervised language model fine-tuning with contrastive learning](#).  
1128 *Preprint*, arXiv:2306.12577.  
1129  
1130  
1131

1132 Itamar Zimmerman and Lior Wolf. 2023. [On the long range abilities of transformers](#). *arXiv preprint arXiv:2311.16620*.  
1133  
1134

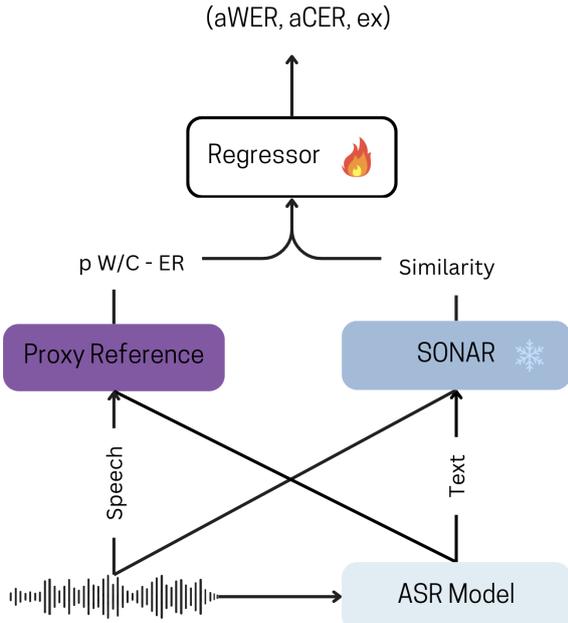


Figure 3: Pipeline diagram for our framework. The proxy reference is an ASR model that takes input speech and generates a transcription. We use the output from the source model as a hypothesis and the output from the proxy reference as a reference ground truth to calculate the WER and CER, which we denote as pWER and pCER. We then use this along with similarity in SONAR embeddings for input speech and hypothesis to train the regression model.

## A Appendix

### A.1 Methodology

$$\text{pWER}(x_{\text{text}}, x_{\text{proxy}}) = \frac{\text{EditDistance}(x_{\text{text}}, x_{\text{proxy}})}{\text{WordCount}(x_{\text{proxy}})} \quad (8)$$

$$\text{pCER}(x_{\text{text}}, x_{\text{proxy}}) = \frac{\text{EditDistance}(x_{\text{text}}, x_{\text{proxy}})}{\text{CharCount}(x_{\text{proxy}})} \quad (9)$$

### A.2 Datasets

To evaluate the robustness and generalizability of our ASR metrics approximation tool, data were sourced from multiple repositories, which we divided into two distinct groups: Standard Benchmark and Wild Challenge dataset.

#### A.2.1 Standard Benchmark Datasets

There are six datasets in total that fall under the benchmark group. These datasets are categorized based on their frequent use in ASR model training and their representation of commonly encountered domains in real-world applications.

**LibriSpeech (Panayotov et al., 2015).** prioritized speaker and content balance over explicit consideration of speech characteristics. It comprises approximately 1000 hours of English read audiobooks, with subsets featuring both clean and noisy speech conditions to simulate different acoustic environments. While the dataset covers diverse subject matter, its focus on formal, clear speech from public domain books means it lacks the natural variability of spontaneous speech, limiting its representation of conversational or informal dialogue.

**TED-LIUM (Rousseau et al., 2014).** contains TED Talks totaling 452 hours of English speech data from approximately 2,000 speakers, recorded in close-talk microphone conditions. The corpus features narrated speaking styles, capturing clear and articulate speech. While it provides non-orthographic transcriptions, lacking formatting such as capitalization and punctuation, it remains a valuable resource for training and benchmarking automatic speech recognition (ASR) models.

**GigaSpeech (Chen et al., 2021).** is a multi-domain, multi-style speech recognition corpus incorporating diverse acoustic and linguistic conditions. It sources audio from three primary domains: audiobooks, podcasts, and YouTube, covering a wide range of speaking styles, including both read and spontaneous speech. The dataset covers a broad spectrum of topics, such as arts, science, sports, and more, making it highly versatile for training robust speech recognition models.

**SPGISpeech (Technologies, 2021).** contains 5,000 hours of professionally transcribed audio from corporate earnings calls, featuring both spontaneous and narrated speaking styles. It emphasizes orthographic accuracy, providing fully formatted text with capitalization, punctuation, and denormalization of non-standard words.

**Common Voice (Ardila et al., 2020).** (a multi-lingual corpus of narrated prompts built through crowdsourcing. Recorded in teleconference conditions, the corpus features narrated speaking styles and emphasizes inclusivity by covering a wide range of accents and languages, including low-resource ones.

**Earnings22 (Rio et al., 2022).** is a 119-hour corpus of English-language earnings calls from global companies, designed to address the lack of real-world, accented speech data in ASR benchmarking

1203  
1204 **AMI (IHM) (Carletta et al., 2005).** The AMI  
1205 Meeting Corpus is a 100-hour dataset of English  
1206 meeting recordings, featuring multimodal data  
1207 synchronized across close-talking and far-field  
1208 microphones, room-view and individual cameras,  
1209 slide projectors, and whiteboards. It includes  
1210 mostly non-native speakers recorded in three  
1211 rooms with varying acoustics. Digital pens  
1212 capture unsynchronized handwritten notes,  
1213 supporting research in speech recognition,  
1214 diarization, and multimodal interaction. Avail-  
1215 able under edinburghcstr/ami, it is widely used  
1216 for meeting analysis and speech processing studies.

1217  
1218 **People’s Speech (Galvez et al., 2021).** Thousands  
1219 of hours of labeled speech data collected from  
1220 diverse speakers, covering a wide range of  
1221 topics, accents, and speaking styles. The dataset  
1222 emphasizes inclusivity and linguistic diversity,  
1223 making it suitable for developing robust and  
1224 generalized speech models. It is widely used  
1225 in academic and industrial research to advance  
1226 the state-of-the-art in automatic speech recog-  
1227 nition (ASR) and other speech-related applications.

1228  
1229 **SLUE - VolxCeleb (Shon et al., 2022).**consists  
1230 of single-sided conversational voice snippets ex-  
1231 tracted from YouTube videos, originally designed  
1232 for speaker recognition. The dataset represents  
1233 natural, unscripted speech in diverse real-world  
1234 settings, capturing a wide range of speaking styles,  
1235 emotions, and acoustic conditions. Utterances  
1236 containing slurs were excluded, and partial words  
1237 were trimmed using a forced aligner to ensure  
1238 clean, usable segments.

### 1240 A.2.2 Wild Challenge Set

1241 **Primock57 (Papadopoulos Korfiatis et al.,**  
1242 **2022).** contains mock consultations conducted by  
1243 seven clinicians and 57 actors posing as patients,  
1244 representing a diverse range of ethnicities, accents,  
1245 and ages. Each actor was provided with a detailed  
1246 case card outlining a primary care scenario, such  
1247 as urinary tract infections, cardiovascular issues, or  
1248 mental health concerns, ensuring the conversations  
1249 were realistic and clinically relevant. The consulta-  
1250 tions were recorded using telemedicine software,  
1251 capturing separate audio channels for clinicians  
1252 and patients, and transcribed by experienced  
1253 professionals to ensure accuracy.

1254 **VoxPopuli Accented (Wang et al., 2021).** is a  
1255 comprehensive multilingual speech corpus derived  
1256 from European Parliament event recordings. It  
1257 includes audio, transcripts, and timestamps sourced  
1258 directly from the official Parliament website. Due  
1259 to its origin, the dataset features a rich collection  
1260 of named entities, making it particularly suitable  
1261 for tasks like Named Entity Recognition (NER).  
1262 **ATCOsim (Jan van Doorn, 2023).**is a specialized  
1263 database containing ten hours of English speech  
1264 from ten non-native speakers, recorded during  
1265 real-time ATC simulations using close-talk  
1266 headset microphones. It features orthographic  
1267 transcriptions, speaker metadata, and session  
1268 details. With a 32 kHz sampling frequency and  
1269 10,078 clean, utterance-level recordings.

### 1270 A.3 Models

1271  
1272 **Whisper Models (Radford et al., 2023).** is  
1273 a transformer-based model that processes 80-  
1274 dimensional log-mel filter bank features from 16  
1275 kHz audio, utilizing a 2D CNN stack followed by a  
1276 transformer encoder-decoder architecture. Trained  
1277 on a vast multilingual dataset of 680,000 hours,  
1278 it incorporates timestamp tokens into its vocabu-  
1279 lary and operates on 30-second audio windows  
1280 during inference, auto-regressively generating text  
1281 sequences while leveraging encoder outputs as con-  
1282 text. Variants of Whisper, such as Distilled, Large,  
1283 Base, and Medium, offer different trade-offs in  
1284 model size and performance, catering to diverse  
1285 computational and accuracy requirements.

1286 **Seamless Models (Communication et al., 2023;**  
1287 **Seamless Communication, 2023; Barrault et al.,**  
1288 **2025).** is a cutting-edge multilingual and multitask  
1289 model for speech and text translation. Built on  
1290 the UnitY architecture, it uses w2v-BERT 2.0 for  
1291 speech encoding and NLLB for text encoding,  
1292 supporting nearly 100 languages. A text decoder  
1293 handles ASR and translation, while a text-to-unit  
1294 (T2U) model and multilingual HiFi-GAN vocoder  
1295 generate speech. Leveraging SONAR embeddings  
1296 and SeamlessAlign (443,000 hours of aligned  
1297 speech/text data), it achieves SOTA results in ASR,  
1298 speech-to-text, speech-to-speech, and text-to-text  
1299 translation, excelling in low-resource languages. It  
1300 introduces BLASER 2.0 for robust evaluation and  
1301 outperforms competitors in noisy environments.

1302  
1303 **Nemo-ASR-Models (Gulati et al., 2020; Vari-**  
1304 **ani et al., 2020; Rekesh et al., 2023; Noroozi**

Name	Type	Description
<b>LibriSpeech</b>	Bench	A corpus of approximately 1,000 hours of 16kHz read English speech, derived from LibriVox audiobooks, segmented and aligned for ASR tasks.
<b>TED-LIUM</b>	Bench	Contains TED Talks totaling 452 hours of English speech data from approximately 2,000 speakers, recorded in close-talk microphone conditions.
<b>GigaSpeech</b>	Bench	A multi-domain, multi-style speech recognition corpus incorporating diverse acoustic and linguistic conditions, sourced from audiobooks, podcasts, and YouTube.
<b>SPGISpeech</b>	Bench	Contains 5,000 hours of professionally transcribed audio from corporate earnings calls, featuring both spontaneous and narrated speaking styles.
<b>Common Voice</b>	Bench	A multilingual corpus of narrated prompts built through crowdsourcing, recorded in teleconference conditions, covering a wide range of accents and languages.
<b>Earnings22</b>	Bench	A 119-hour corpus of English-language earnings calls from global companies, designed to address the lack of real-world, accented speech data in ASR benchmarking.
<b>AMI (IHM)</b>	Bench	The AMI Meeting Corpus is a 100-hour dataset of English meeting recordings, featuring multimodal data synchronized across various devices.
<b>People’s Speech</b>	Bench	Contains thousands of hours of labeled speech data collected from diverse speakers, covering a wide range of topics, accents, and speaking styles.
<b>SLUE - VoxCeleb</b>	Wild	Consists of single-sided conversational voice snippets extracted from YouTube videos, originally designed for speaker recognition.
<b>Primock57</b>	Wild	Contains mock consultations conducted by seven clinicians and 57 actors posing as patients, representing a diverse range of ethnicities, accents, and ages.
<b>VoxPopuli Accented</b>	Wild	A comprehensive multilingual speech corpus derived from European Parliament event recordings, featuring a rich collection of named entities.
<b>ATCOsim</b>	Wild	A specialized database containing ten hours of English speech from ten non-native speakers, recorded during real-time air traffic control simulations.

Table 3: Overview of various ASR along with brief description.

1305 **et al., 2024; Tang et al., 2023; Harper et al.)**  
1306 We included several NVIDIA’s NeMo advanced  
1307 automatic speech recognition (ASR) models, in-  
1308 cluding Canary, Parakeet (110M, 0.6B, and 1.1b),  
1309 Conformer-CTC, and Fast-Conformer, as each is  
1310 designed for specific use cases and optimized for

performance. Canary-1B is a state-of-the-art multi-  
lingual, multitask model featuring a FastConformer  
encoder and Transformer decoder. The Parakeet  
family includes models with a FastConformer en-  
coder paired with different decoders: CTC, RNN-T,  
or TDT. Conformer-CTC is a non-autoregressive

1311  
1312  
1313  
1314  
1315  
1316

1317	model based on the Conformer architecture, combining self-attention and convolution for global and local feature extraction. It uses CTC loss and a linear decoder, supporting both sub-word (BPE) and character-level encodings. While Fast-Conformer is an optimized version of the Conformer architecture, offering significant speed improvements (2.4x faster) with minimal quality degradation. It uses 8x depthwise convolutional subsampling and reduced kernel sizes for efficiency.	1369
1318		1370
1319		1371
1320		1372
1321		
1322		
1323		
1324		
1325		
1326		
1327	<b>Wav2Vec2 Models (Schneider et al., 2019; Baeovski et al., 2020).</b> is a self-supervised pre-trained model designed to process raw audio inputs and generate speech representations. The model architecture consists of three key components: a convolutional feature encoder, a context network, and a quantization module. The convolutional feature encoder converts raw waveforms into latent representations, which are then processed by the context network a transformer based stack with 24 blocks, a hidden size of 1024, 16 attention heads, and a feed-forward dimension of 4096 to capture contextual information. The quantization module maps these latent representations to quantized forms.	1373
1328		1374
1329		1375
1330		1376
1331		1377
1332		1378
1333		1379
1334		1380
1335		1381
1336		1382
1337		1383
1338		1384
1339		1385
1340		1386
1341	<b>HuBERT Models (Hsu et al., 2021).</b> is a self-supervised learning framework designed for speech representation learning where CNN-encoded audio features are randomly masked. During training, the model predicts cluster assignments for masked regions of the input speech, forcing it to learn both acoustic and language models from continuous inputs.	1387
1342		1388
1343		1389
1344		1390
1345		1391
1346		1392
1347		1393
1348		
1349	<b>Audio/Speech Language Models 1.5B and 2B (Rajaa and Tushar)</b> is a multi-modal Language Model designed to analyze and predict metadata from a speaker’s turn in a conversation. It integrates a speech encoder to convert speech signals into meaningful embeddings, which are then processed alongside text instructions by TinyLlama-1.1B-Chat-v1.0 to generate predictions. The model accepts 16 KHz audio inputs and predicts metadata such as SpeechActivity, Transcript, Gender, Age, Accent, and Emotion.	1394
1350		
1351		
1352		
1353		
1354		
1355		
1356		
1357		
1358		
1359		
1360	<b>SpeechT5 (Ao et al., 2022).</b> unified modal framework capable of handling a wide range of tasks, including automatic speech recognition (ASR), speech synthesis, speech translation, voice conversion, speech enhancement, and speaker identification. Its audio post-net, which can incorporate speaker embeddings to enable prosody transfer, making it effective for tasks like voice conversion and speech synthesis. By leveraging its encoder-	1395
1361		
1362		
1363		
1364		
1365		
1366		
1367		
1368		
	decoder architecture, SpeechT5 can generate high-quality mel-spectrograms from text input while preserving speaker-specific characteristics like emotion and gender.	1369
		1370
		1371
		1372
	<b>A.4 Experiments</b>	1373
	<b>A.4.1 Regression Pipeline.</b>	1374
	The regression framework is a stacking ensemble comprising multiple base regressors and a final estimator. We perform basic hyperparameter tuning using RandomizedSearchCV with 5-fold cross-validation, with the objective to minimize <i>mean absolute error (MAE)</i> . The search explores key hyperparameters such as <code>n_estimators</code> , <code>max_depth</code> , <code>learning_rate</code> , and <code>min_samples_split</code> , balancing model complexity and generalization. We provide hyperparameter and other details in 5. The model is trained on 14 datasets divided into two groups: <i>bench</i> (10 standard benchmark datasets) and <i>in-the-wild</i> (4 diverse, real-world datasets). A leave-one-out strategy is applied to the <i>bench</i> set, where the model is trained on 9 datasets and evaluated on the remaining one. All trained models are also evaluated on the <i>in-the-wild</i> set, which remains isolated during training to assess out-of-domain generalization.	1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
	<b>A.5 Results</b>	1394
	This is an appendix.	1395

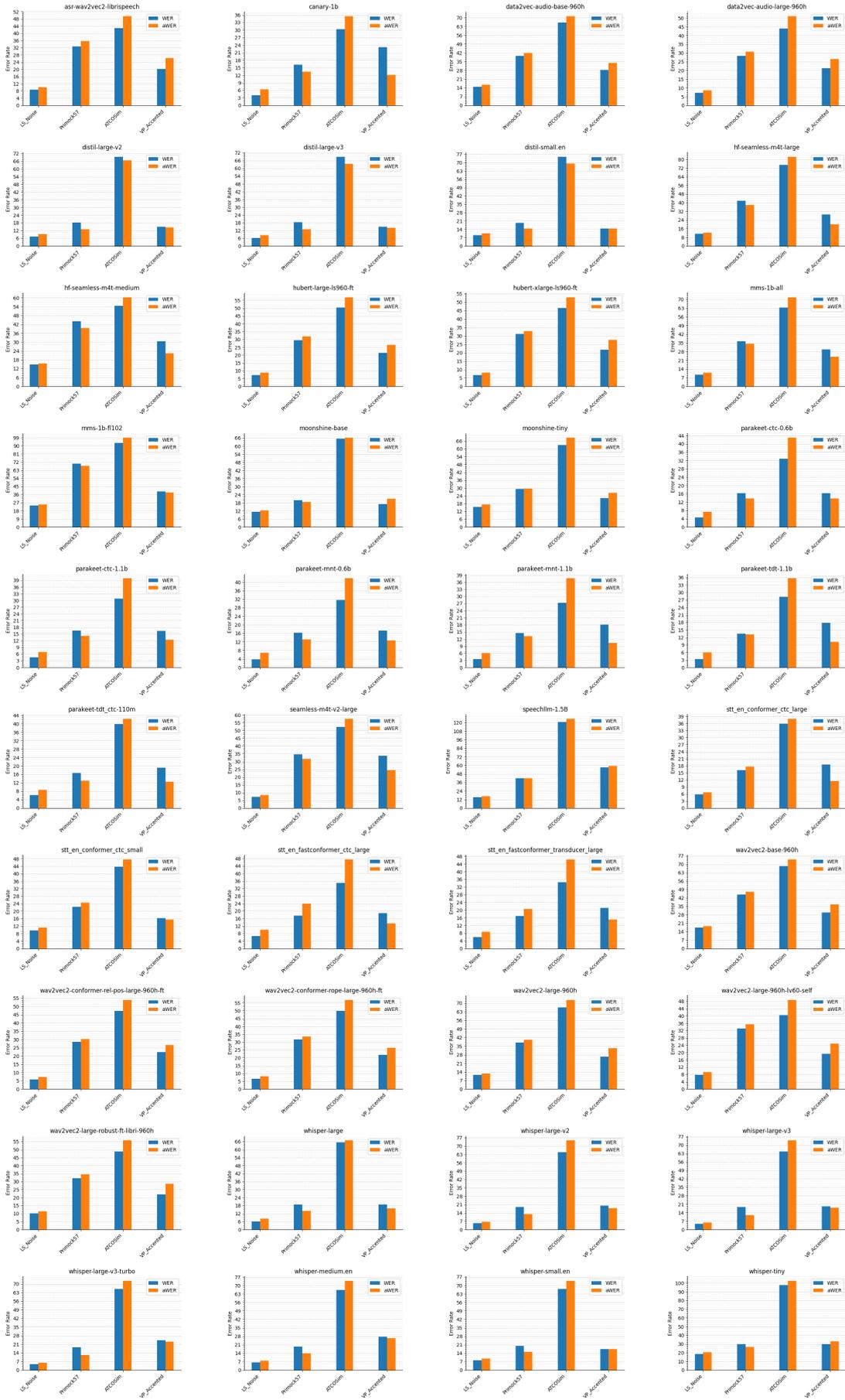


Figure 4: Comparison of Actual vs Approximated WER across models.

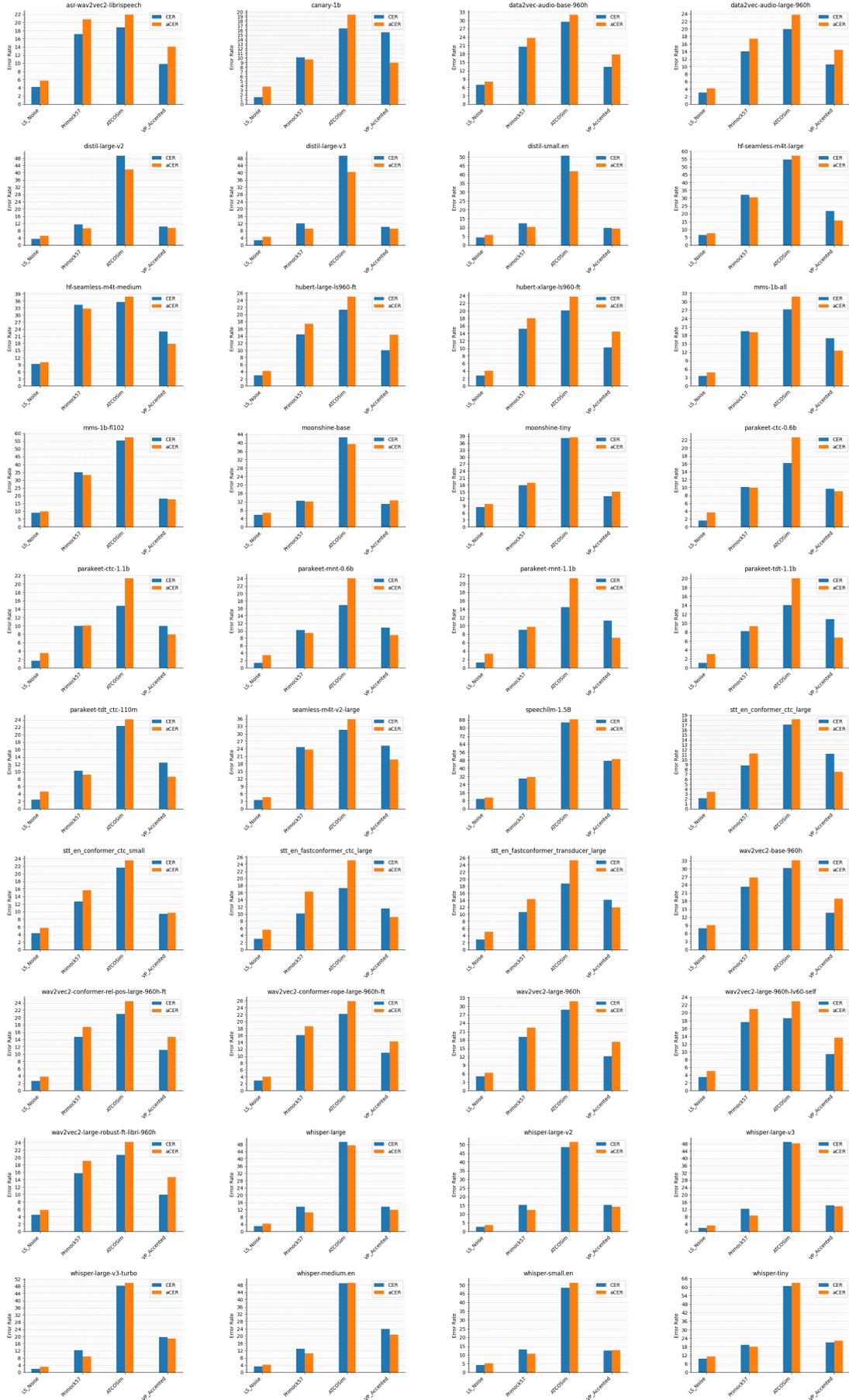


Figure 5: Comparison of Actual vs Approximated CER across models.

Model Type and Models	Description
<b>nemo_asr</b> – parakeet-ctc-1.1b – parakeet-ctc-0.6b – stt_en_conformer_ctc_large – stt_en_fastconformer_ctc_large – stt_en_conformer_ctc_small – parakeet-tdt-1.1b – parakeet-rnnt-1.1b – parakeet-rnnt-0.6b – stt_en_fastconformer_transducer_large – parakeet-tdt_ctc-110m – canary-1b	NVIDIA’s NeMo ASR models offer diverse architectures for speech-to-text applications. The Conformer-CTC model combines self-attention and convolutional operations, using Connectionist Temporal Classification (CTC) loss for efficient transcription. The Conformer-Transducer extends this by incorporating a Recurrent Neural Network Transducer (RNNT) decoder for autoregressive modeling. The Conformer-HAT variant separates label and blank score predictions, enhancing integration with external language models. For improved performance, the Fast-Conformer introduces depthwise convolutional subsampling, achieving approximately 2.4x faster encoding with minimal accuracy loss.
<b>speechbrain</b> – asr-wav2vec2-librispeech	SpeechBrain provides robust models for ASR and speaker recognition.
<b>data2vec</b> – data2vec-audio-large-960h – data2vec-audio-base-960h	Data2Vec models by Facebook are designed for speech representation learning and ASR. These models use a unified learning framework for multiple modalities.
<b>wav2vec2</b> – wav2vec2-large-960h-lv60-self – wav2vec2-large-robust-ft-libri-960h – wav2vec2-large-960h – wav2vec2-base-960h – wav2vec2-conformer-rope-large-960h-ft – wav2vec2-conformer-rel-pos-large-960h-ft	Wav2Vec2 models leverage self-supervised learning on raw audio for ASR. With advanced configurations, these models provide high accuracy for diverse speech-to-text tasks.
<b>mms</b> – mms-1b-all – mms-1b-fl102	The Multilingual Speech (MMS) models by Facebook excel at speech recognition for multiple languages and accents.
<b>hubert</b> – hubert-xlarge-ls960-ft – hubert-large-ls960-ft	HuBERT models provide high-quality speech representations for ASR and other downstream speech tasks.
<b>seamless</b> – hf-seamless-m4t-large – hf-seamless-m4t-medium – seamless-m4t-v2-large	Seamless models focus on multilingual transcription and translation, offering robust real-time speech processing solutions.
<b>speechllm</b> – speechllm-1.5B – speechllm-2B	SpeechLLM models are fine-tuned for ASR and text generation, leveraging billions of parameters for high performance.
<b>whisper</b> – whisper-large-v3 – distil-large-v3 – whisper-large-v2 – whisper-large-v3-turbo – distil-large-v2 – whisper-large – whisper-tiny – whisper-medium.en – distil-small.en – whisper-small.en	Whisper models by OpenAI provide state-of-the-art transcription and translation capabilities for multilingual ASR. These models range from tiny to large configurations.
<b>moonshine</b> – moonshine-base – moonshine-tiny	Moonshine models are lightweight and optimized for efficient ASR on edge devices with minimal computational resources.

Table 4: Overview of various ASR along with brief description.

<b>Model</b>	<b>Hyperparameter</b>	<b>Values</b>
Random Forest (RF)	n_estimators	{100, 200, 300, 500, 700, 1000}
	max_depth	{5, 10, 15, 20, 25, 30}
	min_samples_split	{2, 5, 10, 15, 20}
	min_samples_leaf	{1, 2, 4, 8}
Gradient Boosting (GBR)	n_estimators	{100, 200, 400, 600, 800}
	learning_rate	{0.001, 0.01, 0.05, 0.1, 0.2}
	max_depth	{3, 5, 7, 10}
	min_impurity_decrease	{0.0, 0.001, 0.01, 0.1, 0.2}
HistGradientBoosting (HGB)	max_iter	{100, 200, 300, 400, 500}
	learning_rate	{0.001, 0.01, 0.05, 0.1, 0.2}
	max_depth	{3, 5, 7, 10, 15}
	loss	{Poisson}
Ridge Regression (Final Estimator)	alpha	{1e-3, 1e-2, 0.1, 1, 10, 100, 1000}
	positive	{True}
Pipeline	passthrough	{True}

Table 5: Hyperparameter details for regression model.

Method	IID	D	OOD M	D + M
Base	3.79 <sup>0.16</sup>	3.56 <sup>0.06</sup>	3.76 <sup>0.18</sup>	3.69 <sup>0.06</sup>
w/o S	3.83 <sup>0.14</sup>	3.65 <sup>0.06</sup>	3.82 <sup>0.16</sup>	3.73 <sup>0.07</sup>
w/o PR	8.43 <sup>0.28</sup>	8.36 <sup>0.08</sup>	8.67 <sup>0.24</sup>	8.66 <sup>0.08</sup>
w/ MPR (n=2)	3.69 <sup>0.14</sup>	3.57 <sup>0.06</sup>	3.66 <sup>0.17</sup>	3.69 <sup>0.06</sup>
w/ MPR (n=3)	3.62 <sup>0.13</sup>	3.44 <sup>0.07</sup>	3.58 <sup>0.15</sup>	3.56 <sup>0.07</sup>
w/ MPR (n=4)	3.57 <sup>0.13</sup>	3.40 <sup>0.06</sup>	3.53 <sup>0.13</sup>	3.52 <sup>0.06</sup>
w/ MPR (n=5)	3.49 <sup>0.13</sup>	3.37 <sup>0.06</sup>	3.47 <sup>0.12</sup>	3.49 <sup>0.07</sup>
w/ mMPR (n=3)	3.61 <sup>0.15</sup>	3.40 <sup>0.09</sup>	3.57 <sup>0.13</sup>	3.51 <sup>0.09</sup>
w/ mMPR (n=5)	3.80 <sup>0.15</sup>	3.47 <sup>0.03</sup>	3.77 <sup>0.13</sup>	3.56 <sup>0.04</sup>
Base (r=11.9)	4.68 <sup>0.17</sup>	5.16 <sup>0.06</sup>	4.64 <sup>0.16</sup>	5.06 <sup>0.05</sup>
Base (r=14.0)	4.84 <sup>0.18</sup>	4.88 <sup>0.07</sup>	4.75 <sup>0.17</sup>	4.77 <sup>0.07</sup>
Base (r=20.2)	5.13 <sup>0.12</sup>	5.38 <sup>0.07</sup>	5.12 <sup>0.10</sup>	5.30 <sup>0.07</sup>
Base (r=23.5)	5.60 <sup>0.13</sup>	6.12 <sup>0.07</sup>	5.69 <sup>0.21</sup>	6.03 <sup>0.05</sup>
w/o S (r=11.9)	5.50 <sup>0.21</sup>	5.84 <sup>0.06</sup>	5.55 <sup>0.21</sup>	5.65 <sup>0.05</sup>
w/o S (r=14.0)	5.73 <sup>0.12</sup>	5.50 <sup>0.05</sup>	5.71 <sup>0.13</sup>	5.37 <sup>0.06</sup>
w/o S (r=20.2)	6.16 <sup>0.18</sup>	6.24 <sup>0.08</sup>	6.13 <sup>0.10</sup>	5.97 <sup>0.09</sup>
w/o S (r=23.5)	6.38 <sup>0.09</sup>	6.77 <sup>0.08</sup>	6.43 <sup>0.16</sup>	6.58 <sup>0.08</sup>

Table 6: Mean absolute error between predicted character count and actual character error count (in absolute terms) across different configurations. R - Regression, C - Classification, PR - Proxy Reference, S - Silarity, MPR - Multiple PR. The OOD results are averaged across five wild datasets. Superscript represents the standard deviation across five runs.

Model	LS_Noise	Primock57	Atcosim	VP_accented
asr-wav2vec2-librispeech	4.2/5.8	17.2/20.8	18.8/21.9	9.8/14.0
canary-1b	1.5/3.8	10.1/9.7	16.4/19.4	15.6/9.0
data2vec-audio-base-960h	7.0/8.1	20.5/23.7	29.5/32.0	13.3/17.8
data2vec-audio-large-960h	3.1/4.2	14.1/17.4	20.0/23.8	10.6/14.4
distil-large-v2	3.5/5.2	11.5/9.2	49.5/41.8	10.2/9.4
distil-large-v3	2.7/4.6	11.9/9.1	49.4/40.5	10.1/9.0
distil-small.en	4.2/5.8	12.2/10.4	50.7/41.8	9.7/9.2
hf-seamless-m4t-large	6.5/7.5	32.1/30.6	54.7/57.2	21.8/15.8
hf-seamless-m4t-medium	9.4/10.1	34.4/32.7	35.5/37.9	23.1/17.9
hubert-large-ls960-ft	3.0/4.2	14.4/17.4	21.3/25.0	10.0/14.3
hubert-xlarge-ls960-ft	2.7/4.1	15.3/18.1	20.1/23.8	10.2/14.5
mms-1b-all	3.6/4.8	19.5/19.1	27.2/31.8	17.0/12.6
mms-1b-fl102	9.0/10.0	35.0/33.2	55.4/57.3	18.2/17.6
moonshine-base	5.7/6.8	12.4/12.1	42.6/39.5	10.9/12.6
moonshine-tiny	8.5/9.9	17.9/19.0	38.2/38.4	13.2/15.1
parakeet-ctc-0.6b	1.7/3.7	10.1/9.9	16.2/22.7	9.7/9.0
parakeet-ctc-1.1b	1.7/3.6	10.0/10.1	14.8/21.4	10.0/8.0
parakeet-rnnt-0.6b	1.3/3.4	10.1/9.4	16.9/24.1	10.9/8.8
parakeet-rnnt-1.1b	1.3/3.3	9.1/9.7	14.5/21.3	11.2/7.2
parakeet-tdt-1.1b	1.1/3.1	8.2/9.4	14.0/20.0	10.9/6.8
parakeet-tdt_ctc-110m	2.5/4.7	10.3/9.2	22.3/24.2	12.4/8.6
seamless-m4t-v2-large	3.5/4.6	24.6/23.7	31.6/35.8	25.2/19.8
speechllm-1.5B	9.9/11.2	30.1/31.4	85.4/88.7	47.3/49.0
speechllm-2B	8.4/9.3	25.3/27.7	33.5/36.1	24.0/28.3
stt_en_conformer_ctc_large	2.1/3.4	8.8/11.2	17.1/18.2	11.1/7.5
stt_en_conformer_ctc_small	4.3/5.7	12.7/15.6	21.6/23.6	9.5/9.7
stt_en_fastconformer_ctc_large	3.0/5.6	10.1/16.3	17.3/25.1	11.5/9.2
stt_en_fastconformer_transducer_large	2.8/5.0	10.6/14.3	18.7/25.3	14.2/11.9
wav2vec2-base-960h	7.9/9.1	23.3/26.7	30.3/33.2	13.7/18.9
wav2vec2-conformer-rel-pos-large-960h-ft	2.6/3.8	14.7/17.4	21.0/24.5	11.2/14.7
wav2vec2-conformer-rope-large-960h-ft	2.9/4.0	16.1/18.7	22.2/25.9	11.0/14.2
wav2vec2-large-960h	5.1/6.3	19.1/22.4	28.8/31.8	12.2/17.4
wav2vec2-large-960h-lv60-self	3.5/5.0	17.6/21.0	18.6/23.0	9.3/13.6
wav2vec2-large-robust-ft-libri-960h	4.5/5.8	15.7/19.0	20.7/24.2	10.0/14.7
whisper-large	2.9/4.2	13.7/10.6	49.3/47.5	13.7/11.9
whisper-large-v2	2.6/3.8	15.3/12.5	48.6/51.5	15.3/14.2
whisper-large-v3	2.0/3.3	12.3/8.7	48.9/48.3	14.3/13.6
whisper-large-v3-turbo	2.0/3.2	12.4/8.8	48.3/49.9	19.7/19.0
whisper-medium.en	3.3/4.3	13.1/10.5	49.1/49.2	23.8/20.6
whisper-small.en	4.2/5.3	13.1/10.8	48.4/51.2	12.5/12.7
whisper-tiny	9.8/11.3	19.3/18.2	60.8/63.0	21.0/22.3

Table 7: Actual and approximated CER ( $\downarrow$ ), separated by a slash, on out-of-distribution wild datasets. The regression model is trained independently for each ASR model on standard benchmarks, making the wild datasets out-of-distribution.

Model	AMI_IHM		CV11		Earnings22		Gigaspeech		LibriSpeech <sub>clean</sub>	
	WER/aWER	CER/aCER	WER/aWER	CER/aCER	WER/aWER	CER/aCER	WER/aWER	CER/aCER	WER/aWER	CER/aCER
asr-wav2vec2-librispeech	28.4/30.5	13.8/17.6	25.0/29.7	11.7/15.0	37.3/33.2	21.3/16.1	16.6/16.5	6.9/7.4	1.8/3.8	0.5/2.2
canary-1b	15.4/17.6	9.2/12.7	8.7/14.2	4.1/8.5	21.8/16.0	15.8/9.1	11.1/6.9	5.5/4.3	1.5/5.7	0.5/3.5
data2vec-audio-base-960h	39.9/40.4	19.9/23.5	37.8/42.3	18.3/21.7	50.8/48.6	28.0/25.0	23.8/23.5	10.1/10.8	2.8/4.0	0.9/1.6
data2vec-audio-large-960h	34.1/36.1	16.9/21.2	23.3/27.9	10.9/14.1	37.7/34.5	21.2/16.7	17.0/16.6	7.2/7.4	1.8/3.9	0.5/1.7
distil-large-v2	17.8/16.8	11.2/11.5	14.2/19.7	7.1/10.6	19.3/20.0	12.5/13.7	12.8/8.2	7.1/5.4	3.4/6.7	1.5/4.2
distil-large-v3	18.5/17.3	11.6/11.7	13.7/19.4	6.6/10.3	18.4/19.8	12.1/13.0	12.2/7.9	6.9/5.3	2.8/6.6	1.2/4.1
distil-small.en	18.5/18.4	11.1/12.6	18.5/23.1	9.4/12.5	21.2/21.4	13.6/14.7	13.1/8.6	7.3/5.7	3.7/7.6	1.6/4.5
hf-seamless-m4t-large	36.3/33.9	25.4/25.1	9.5/13.2	5.1/7.4	30.7/32.8	21.1/23.9	24.2/21.1	16.7/15.7	3.2/4.8	1.5/2.7
hf-seamless-m4t-medium	40.6/37.2	29.5/28.9	11.3/14.3	6.0/7.4	33.7/35.9	23.9/26.4	30.2/28.1	22.3/21.7	3.8/5.3	1.6/2.9
hubert-large-ls960-ft	31.1/33.6	15.2/19.8	24.1/28.8	10.6/13.6	37.6/34.4	20.6/16.3	19.3/18.3	8.1/7.8	2.1/3.7	0.6/1.6
hubert-xlarge-ls960-ft	31.1/34.3	15.0/20.0	24.1/28.7	10.5/13.9	37.3/34.9	20.4/15.9	18.1/17.4	7.3/7.6	2.0/3.8	0.6/1.7
mms-1b-all	37.0/36.2	19.1/20.8	22.5/27.5	8.9/12.5	34.1/30.6	19.6/15.1	19.4/16.9	8.3/7.6	4.2/6.2	1.3/2.7
mms-1b-ft02	75.4/73.3	35.1/33.9	42.6/45.3	17.8/19.9	50.6/52.3	24.2/26.5	37.2/35.7	15.7/15.2	15.8/17.3	5.1/5.9
moonshine-base	15.6/24.7	9.4/16.7	20.8/25.4	10.8/13.8	24.3/25.6	15.9/16.6	14.2/10.4	8.1/6.8	3.4/6.3	1.3/3.7
moonshine-tiny	21.3/25.3	12.8/16.7	26.7/31.7	14.4/17.3	31.2/32.7	19.7/20.2	16.6/14.1	9.1/8.6	4.5/7.2	1.8/4.2
parakeet-ctc-0.6b	17.0/23.1	10.0/16.3	10.7/21.1	5.1/11.2	24.7/19.1	16.9/11.5	12.0/8.6	6.1/5.2	2.0/5.1	0.7/2.5
parakeet-ctc-1.1b	15.7/21.4	9.0/15.3	10.5/20.1	5.2/11.0	24.0/17.7	16.6/10.7	12.2/7.9	6.2/5.0	1.8/5.4	0.5/2.6
parakeet-mnt-0.6b	18.8/24.0	11.7/17.9	8.5/19.9	4.2/10.8	25.2/18.7	17.5/11.5	11.7/9.0	6.2/5.4	1.8/5.5	0.6/3.2
parakeet-mnt-1.1b	18.6/23.5	11.7/17.2	6.7/19.6	3.4/10.5	25.7/17.9	18.4/11.4	11.3/8.4	6.0/5.0	1.5/5.0	0.5/3.3
parakeet-tdt-1.1b	17.1/23.5	10.2/16.9	7.2/19.6	3.4/10.6	24.5/16.6	17.1/10.0	10.2/7.8	4.9/4.7	1.3/6.0	0.4/2.9
parakeet-tdt_ctc-110m	18.5/18.8	10.7/13.6	12.7/17.7	6.9/10.1	22.2/14.8	15.7/9.2	12.6/8.2	6.2/5.0	2.6/6.7	0.9/3.8
seamless-m4t-v2-large	43.0/42.3	30.2/30.2	8.2/12.3	3.9/6.3	47.3/47.4	33.7/33.9	25.7/23.2	18.1/17.2	2.7/4.4	1.0/2.5
speechllm-1.5B	67.7/69.3	51.5/55.0	18.5/22.7	10.0/12.7	50.8/48.2	38.3/35.4	27.5/26.0	18.1/18.3	10.5/12.1	7.3/9.2
speechllm-2B	38.6/40.8	24.3/28.2	24.6/28.2	16.5/18.3	47.3/45.0	32.5/30.8	24.4/23.6	13.5/13.8	7.0/9.3	4.5/4.8
stt_en_conformer_ctc_large	15.3/19.9	7.9/13.4	10.4/15.4	4.7/8.0	24.8/20.0	16.4/10.7	13.2/10.6	5.9/5.6	2.2/3.7	0.7/2.4
stt_en_conformer_ctc_small	21.2/24.4	11.2/15.4	19.1/24.1	8.9/12.2	29.3/25.3	19.0/14.1	15.5/14.9	7.2/7.7	3.9/5.4	1.4/3.1
stt_en_fastconformer_ctc_large	20.3/24.0	11.7/15.3	9.5/19.3	4.6/10.2	27.3/21.5	18.3/13.0	14.5/14.7	7.2/8.2	1.9/5.2	0.7/2.8
stt_en_fastconformer_transducer_large	19.8/22.0	12.9/16.8	9.3/18.0	4.7/9.8	31.5/26.9	23.0/18.8	13.6/13.2	7.4/7.8	1.8/3.9	0.6/2.6
wav2vec2-base-960h	37.9/38.7	18.7/21.9	40.6/45.7	19.5/22.8	51.1/48.6	28.2/25.4	26.2/26.6	11.7/12.2	3.7/4.5	1.1/1.9
wav2vec2-conformer-rel-pos-large-960h-ft	35.0/38.7	18.5/24.1	23.7/28.0	10.7/13.7	38.4/36.2	21.7/17.6	18.5/17.2	8.5/7.9	1.6/3.3	0.5/1.5
wav2vec2-conformer-rope-large-960h-ft	34.3/36.4	18.0/22.7	23.6/28.5	11.6/15.0	36.9/33.9	21.4/16.7	17.9/17.7	7.3/7.6	1.8/3.8	0.5/1.6
wav2vec2-large-960h	34.0/36.4	16.4/20.2	34.1/38.6	16.2/19.4	46.4/43.4	25.4/21.7	20.6/20.5	8.6/9.1	2.9/4.3	0.8/2.1
wav2vec2-large-960h-lv60-self	29.1/31.5	15.5/19.5	23.1/28.8	11.0/15.0	36.7/32.5	20.8/15.7	17.6/17.2	7.5/8.0	1.7/3.5	0.5/1.9
wav2vec2-large-robust-ft-libri-960h	30.5/33.9	13.8/19.2	25.0/29.3	10.7/13.8	37.1/33.5	20.5/15.7	18.0/17.6	7.1/7.7	2.8/4.3	0.8/2.3
whisper-large	18.5/18.3	12.3/13.0	13.0/18.0	6.6/9.3	18.8/20.3	12.3/14.9	12.2/7.7	7.1/5.1	2.8/5.1	1.4/3.5
whisper-large-v2	18.6/17.1	12.1/11.8	11.3/15.5	5.7/8.0	19.0/21.5	13.0/15.8	12.5/7.1	7.3/4.9	2.8/5.1	1.5/3.2
whisper-large-v3	19.0/17.1	12.3/12.0	9.9/14.5	4.9/6.9	18.2/20.7	12.1/14.9	12.5/7.3	7.2/4.9	2.2/4.0	0.9/2.9
whisper-large-v3-turbo	19.0/17.5	12.3/11.9	12.6/16.1	6.3/8.1	18.8/21.1	12.9/15.6	12.2/6.8	7.1/4.6	2.4/4.4	1.1/2.5
whisper-medium.en	20.3/18.8	13.7/14.5	14.3/17.8	7.2/9.2	20.1/22.6	13.3/16.5	12.8/7.6	7.6/5.5	3.3/5.5	1.8/3.5
whisper-small.en	19.8/17.9	12.8/12.2	17.8/21.9	9.3/11.3	20.6/22.9	13.6/16.3	12.8/8.1	7.3/5.1	3.3/5.6	1.4/3.1
whisper-tiny	26.7/25.1	16.7/17.0	33.5/40.3	17.7/20.9	33.8/35.4	22.0/25.4	20.6/18.2	12.0/11.0	7.9/11.2	3.4/5.1

Table 8: Actual and approximated WER and CER, separated by a slash, across five standard datasets. The regression model is trained on nine datasets and tested on one, with this process repeated for all datasets, ensuring that the test data is always out-of-distribution.

Model	peoples_speech		slue_voxceleb		spgispeech_S		tedlium-dev-test		voxpupuli_en	
	WER/aWER	CER/aCER	WER/aWER	CER/aCER	WER/aWER	CER/aCER	WER/aWER	CER/aCER	WER/aWER	CER/aCER
asr-wav2vec2-librispeech	35.6/32.9	19.8/17.7	19.5/20.4	9.8/12.2	11.1/12.2	4.8/4.7	10.3/11.1	5.2/5.8	14.3/12.6	6.6/5.1
canary-1b	16.5/22.5	11.1/15.2	14.9/11.1	10.8/8.2	3.2/6.7	2.0/3.9	7.9/7.6	5.9/5.0	6.4/4.9	3.9/3.4
data2vec-audio-base-960h	43.4/38.6	24.4/20.8	26.1/27.6	13.0/15.5	19.2/19.8	8.2/7.9	13.6/14.2	6.3/6.4	18.9/17.5	8.5/7.1
data2vec-audio-large-960h	35.1/31.3	20.0/17.3	20.4/22.1	10.3/12.9	11.3/12.0	4.9/4.7	9.9/10.6	4.5/5.0	14.9/13.4	6.9/5.5
distil-large-v2	17.4/21.8	12.2/14.1	16.0/10.8	11.4/7.4	3.7/7.6	1.8/4.5	10.4/8.5	8.8/5.4	9.5/8.2	5.8/4.6
distil-large-v3	17.4/21.6	12.4/13.8	14.4/10.0	10.3/6.8	3.6/7.4	1.8/4.5	10.7/9.2	8.6/5.7	9.3/6.7	5.8/4.1
distil-small.en	19.0/22.5	13.3/14.3	15.9/11.4	11.3/7.8	4.0/7.9	1.9/4.7	10.8/8.8	9.1/5.6	10.2/7.4	6.4/4.3
hf-seamless-m4t-large	38.5/41.5	29.2/30.1	47.2/42.8	39.4/36.1	16.2/18.7	11.5/13.0	19.8/19.1	15.7/14.4	8.1/6.5	5.0/3.6
hf-seamless-m4t-medium	43.6/45.7	33.6/34.2	50.9/47.4	43.2/40.3	12.9/15.5	8.8/10.4	27.0/26.2	21.3/20.0	8.8/7.3	5.5/4.4
hubert-large-ls960-ft	34.1/31.3	18.8/17.7	20.8/22.0	10.1/12.3	11.6/12.4	4.9/4.6	11.0/12.0	5.3/5.4	15.0/13.6	6.9/5.4
hubert-xlarge-ls960-ft	35.5/31.5	19.5/16.0	20.3/22.1	9.9/12.3	11.9/12.3	4.8/4.5	10.1/11.2	4.2/5.0	14.5/12.8	6.7/5.3
mms-1b-all	32.2/36.0	16.8/18.7	27.6/26.3	14.6/14.8	10.0/12.5	3.8/4.9	13.5/13.3	7.3/6.5	8.9/7.6	4.4/3.1
mms-1b-ft02	52.4/52.7	26.0/25.5	51.7/48.7	26.1/23.2	19.1/22.9	5.9/8.6	29.7/30.3	12.9/12.9	22.6/20.5	9.3/7.9
moonshine-base	26.4/26.2	18.1/17.5	17.0/13.8	11.6/9.1	6.4/7.5	3.4/3.9	5.8/7.0	3.4/4.1	11.7/9.9	6.7/4.6
moonshine-tiny	31.8/30.8	20.5/19.1	20.1/17.2	13.4/11.8	9.1/9.9	4.8/5.3	9.8/9.5	6.7/5.8	14.9/12.9	8.2/7.2
parakeet-ctc-0.6b	24.2/20.0	16.5/12.6	13.1/11.0	8.7/7.8	6.4/7.5	3.6/3.8	4.3/7.2	2.6/4.4	7.0/7.2	4.1/3.7
parakeet-ctc-1.1b	20.7/18.1	13.9/11.8	13.0/11.6	8.7/8.3	6.4/7.1	3.7/3.6	5.0/7.4	3.0/4.4	6.7/6.5	3.9/3.3
parakeet-mnt-0.6b	21.9/17.9	15.4/12.2	14.5/11.6	10.0/8.6	4.9/7.3	2.9/3.7	5.0/6.9	3.0/4.0	6.4/6.7	3.8/3.7
parakeet-mnt-1.1b	23.3/17.2	16.6/11.8	14.1/10.9	9.8/8.9	4.5/7.7	2.7/4.4	5.2/7.7	3.5/4.9	5.6/6.0	3.4/3.2
parakeet-tdt-1.1b	24.5/17.9	16.6/12.6	13.5/10.6	9.1/7.9	5.4/7.7	3.2/4.2	4.4/7.1	2.8/4.4	5.5/6.0	3.3/3.1
parakeet-tdt-ctc-110m	16.6/21.4	11.5/15.1	15.3/11.5	10.7/8.3	3.8/7.3	2.2/4.0	5.2/6.6	3.3/4.1	7.5/6.2	4.6/3.1
seamless-m4t-v2-large	35.0/36.3	25.1/24.9	45.1/43.2	35.9/34.7	11.7/13.6	7.6/8.7	26.5/25.5	21.0/19.1	8.0/7.8	5.7/5.0
speechllm-1.5B	45.1/44.5	32.7/32.0	60.3/60.8	44.1/46.8	10.6/10.9	6.2/5.8	19.4/17.6	14.2/11.9	30.8/29.9	22.0/21.3
speechllm-2B	52.9/53.1	36.6/35.7	36.9/37.8	25.8/27.3	14.6/15.3	8.1/7.5	18.8/16.7	12.9/9.4	28.7/27.7	18.5/17.5
stt_en_conformer_ctc_large	24.2/21.5	15.2/13.0	12.6/14.7	7.6/9.4	7.9/7.3	4.1/3.4	5.9/7.7	3.3/4.4	6.9/5.3	3.9/2.7
stt_en_conformer_ctc_small	31.3/26.9	19.0/15.7	16.6/17.8	9.6/11.7	10.0/9.4	5.1/4.1	8.0/9.9	3.9/5.2	8.9/7.3	5.0/3.8
stt_en_fastconformer_ctc_large	26.9/20.5	18.3/12.9	15.4/14.0	9.9/9.7	6.9/8.4	3.7/4.1	5.7/7.8	3.1/4.5	6.3/6.0	3.8/3.3
stt_en_fastconformer_transducer_large	26.5/20.4	19.0/13.8	16.9/15.1	11.3/11.1	6.0/7.9	3.4/4.0	4.9/7.0	2.8/4.4	6.7/6.9	4.1/3.8
wav2vec2-base-960h	44.7/40.1	24.5/20.6	27.3/28.7	13.5/15.7	21.5/22.4	8.9/8.7	13.8/14.7	6.1/6.6	20.5/19.4	9.1/7.8
wav2vec2-conformer-rel-pos-large-960h-ft	37.3/34.7	21.2/19.7	20.3/22.1	10.6/13.4	12.0/12.2	5.2/4.6	11.7/12.3	6.6/6.3	14.8/13.1	6.9/5.5
wav2vec2-conformer-rope-large-960h-ft	35.3/32.1	20.4/19.3	20.6/22.1	10.4/12.9	11.7/12.5	5.1/4.8	10.9/11.8	5.5/6.0	14.5/13.4	6.9/5.4
wav2vec2-large-960h	38.9/35.7	21.6/19.2	23.2/25.1	11.5/13.7	16.3/17.1	6.9/6.6	12.2/13.2	5.6/6.1	18.1/16.7	8.2/7.0
wav2vec2-large-960h-lv60-self	32.5/29.4	18.7/15.7	20.2/20.8	10.7/12.9	10.4/12.2	4.3/4.7	9.5/10.5	4.2/4.9	13.5/12.5	6.4/5.2
wav2vec2-large-robust-ft-libri-960h	36.2/32.6	19.2/16.9	20.9/23.0	9.6/12.5	11.8/12.5	5.0/4.8	10.6/11.9	4.8/5.4	15.4/14.0	7.0/5.8
whisper-large	31.2/32.4	24.6/23.0	17.6/12.6	13.7/10.5	3.7/7.4	2.1/4.4	19.3/16.1	14.0/10.3	8.9/6.0	5.5/3.2
whisper-large-v2	18.8/25.2	14.2/18.1	18.7/15.1	14.8/12.1	4.1/7.7	2.4/4.8	28.3/25.3	19.4/14.4	8.7/6.8	5.5/3.8
whisper-large-v3	20.4/27.5	15.6/19.5	15.6/11.9	11.8/8.7	3.2/6.3	1.7/4.0	10.5/8.3	8.7/5.5	11.0/8.6	7.8/6.1
whisper-large-v3-turbo	16.0/23.7	12.0/16.5	15.3/11.9	11.5/9.0	3.1/6.3	1.7/3.9	9.9/8.1	8.5/5.2	13.3/11.1	9.8/7.9
whisper-medium.en	20.1/25.1	15.3/18.0	21.2/16.1	16.6/13.2	4.0/7.7	2.2/5.0	17.3/14.6	18.3/14.2	9.0/7.2	5.6/3.2
whisper-small.en	21.2/25.1	16.5/18.0	18.2/14.1	13.9/10.9	3.9/7.5	2.1/4.6	10.6/8.3	15.6/12.0	9.5/8.1	5.9/3.4
whisper-tiny	30.1/31.9	21.7/21.5	24.0/20.3	17.5/14.4	8.1/11.9	3.9/6.7	17.6/15.3	13.0/9.5	13.2/11.2	7.4/6.3

Table 9: Actual and approximated WER and CER, separated by a forward slash, across five standard datasets. The regression model is trained on nine datasets and tested on one, with this process repeated for all datasets, ensuring that the test data is always out-of-distribution.