

Prosody Detection improves Pretrained Automatic Speech Recognition

Anonymous ACL submission

Abstract

We show the performance of Automatic Speech Recognition (ASR) systems that use semi-supervised speech representations can be boosted by a complimentary prosody detection module, by introducing a joint ASR and prosody detection model. The prosody detection component of our model achieves a significant improvement on the state-of-the-art for the task, closing the gap in F1-score by 41%. Additionally, the ASR performance in joint training decreases WER by 28.3% on LibriSpeech, under limited resource fine-tuning. With these results, we show the importance of extending pretrained speech models to retain or relearn important prosodic cues.

1 Introduction

Models based on self-supervised speech representations have in recent years claimed state-of-the-art performance in ASR (Baeovski et al., 2020; Hsu et al., 2021). Moreover, they have permitted to bypass both a heavy speech-science informed featurisation component, as well a language dependent acoustic dictionary resource writing component. In doing so, such models have become seemingly less human reliant during development, provided adequate quantities of raw speech data and computational resources.

However, the training techniques for these self-supervised speech models do not reveal what exactly is deemed important by these models and later retained within their output speech representations. Subsequent studies have since introduced benchmarks and metrics to analyse the linguistic knowledge of these models at different levels, mostly from the point of view of assessing the existence of this linguistic knowledge. For example, the Zero-Resource Speech challenges¹ provide test beds to analyse phonetic, lexical, syntactic, and semantic

level book-keeping of these representations (Dunbar et al., 2017, 2020; Nguyen et al., 2021). More recently, ProsAudit was introduced to provide a similar book-keeping of the prosodic information retained in these speech representations (de Seyssel et al., 2023). However, these studies and corresponding benchmarks do not elucidate whether refocusing these SSL representations to retain more of the original linguistic signal could correspond to better performance downstream in basic but central speech tasks like ASR.

One main bottleneck to carrying out such a study is the sheer scarcity of datasets with prosodic annotations, and the previous expense in generating these annotations using trained linguists. One can, however, envisage a scheme for obtaining simple prosodic annotations for “important” words in an utterance from non-specialists, and even for non-written languages, by which annotators simply press a button when important segments of an utterance are heard. Still, currently, the only existing datasets for English (and any other language) are relatively small—the largest being the Boston University Radio News Corpus for English with 11 hours of data (Ostendorf et al., 1995). Is there any benefit for spoken language understanding tasks in extending current and/or developing new prosody datasets? To assess this question, one must investigate the role of prosody in these tasks.

Contributions. In this paper we study the role of prosodic information, specifically focusing on pitch accents, in ASR. Our contributions are as follows.

1. We streamline and significantly boost the performance of the current state-of-the-art model for pitch accent detection.
2. We present a multi-task model for integrating pitch accent detection into the ASR task, which improves the performance for ASR in limited resource settings.

¹<https://zerospeech.com>

- 079 3. We then automatically annotate the pitch ac- 129
080 cents of a small dataset using self-training, 130
081 and then apply it in our proposed joint model, 131
082 achieving even further ASR performance 132
083 boosts. 133

084 2 Related work 134

085 **Prosody Detection.** There is a long line of re- 135
086 search on automatic prosody detection (for exam- 136
087 ple, (Taylor, 1995; Rosenberg et al., 2015; Shahin 137
088 et al., 2016; Li et al., 2018; Stehwien et al., 2020; 138
089 Sabu et al., 2021)). With the advent of pretrained 139
090 speech models, and in particular, wav2vec (Schnei- 140
091 der et al., 2019) and wav2vec2 (Baevski et al., 141
092 2020), a new line of systems that builds on self- 142
093 supervised speech representations has achieved the 143
094 state-of-the-art in detecting prosodic boundaries in 144
095 Czech broadcast news recordings (Kunešová and 145
096 Řezáčková, 2022) and in pitch events and intona- 146
097 tion phrase boundaries in English broadcast news 147
098 (Zhai and Hasegawa-Johnson, 2023). This latter 148
099 model, called wav2TOBI, forms the point of depart- 149
100 ure for our multitask system presented here. The 150
101 question left open by these models, and others is 151
102 whether these self-supervised representations ade- 152
103 quately account for prosody, which has been shown 153
104 to aid in an array of linguistic tasks like SLU (Nöth 154
105 et al., 2002; Shriberg and Stolcke, 2004; Shriberg 155
106 et al., 1998; Rajaa, 2023; Wei et al., 2022), and 156
107 parsing (Tran et al., 2017; Gregory et al., 2004; 157
108 Kahn et al., 2005; Dreyer and Shafran, 2007; Kahn 158
109 and Ostendorf, 2012; Price et al., 1991; Beckman, 159
110 1996). Or whether a specialised module should 160
111 intervene and boost the prosodic signal for better 161
112 performance. 162

113 **Prosody with ASR.** In this paper, we are par- 163
114 ticularly interested in whether refocusing speech 164
115 pretrained models on prosody might aid in per- 165
116 formance for ASR. Prosody has been previously 166
117 shown to be of importance to ASR, both as engi- 167
118 neered features, as well as through learning from 168
119 prosody annotated datasets (Silverman et al., 1992; 169
120 Ostendorf et al., 2003; Hirose and Minematsu, 170
121 2004; Hirschberg et al., 2004; Hasegawa-Johnson 171
122 et al., 2005; Ananthakrishnan and Narayanan, 172
123 2007; Vicsi and Szaszák, 2010; Chen et al., 2012; 173
124 Kathania et al., 2020; Hasija et al., 2022; Coto- 174
125 Solano, 2021). However, to our knowledge, there is 175
126 no research that builds on pretrained speech models, 176
127 whose application to prosody detection and ASR 177
128 has resulted in the state-of-the-art performance. 178

State-of-the-art ASR A summary of state-of- 129
the-art performance in ASR over the Librispeech 130
dataset is given in the appendix (Table 3). For this 131
paper, for model comparability, we focus on the 132
wav2vec2 model, which is the pretrained model 133
fine-tuned by the wav2TOBI model for prosody de- 134
tection, and which is the model on which we base 135
our system presented here. 136

137 3 Modelling prosody and ASR 137

138 3.1 Datasets 138

139 Our research uses the BURNC (Ostendorf et al., 139
140 1995), Librispeech (Panayotov et al., 2015) and 140
141 Libri-light (Kahn et al., 2020) corpora. 141

142 The BURNC dataset is a broadcast news-style 142
143 read speech corpus which contains 11 hours of 143
144 speech, sourced from 7 different speakers (3 fe- 144
145 male and 4 male). It consists of audio snippets with 145
146 their transcriptions, phonetic alignments, parts-of- 146
147 speech tags and prosodic labels. We used 75% 147
148 of the data in this dataset for training, 15% for 148
149 development and 10% for testing. Because multi- 149
150 ple readers may have read the same news story in 150
151 BURNC, we ensure that no news stories appearing 151
152 in the test set also occur in the training set. 152

153 We represent pitch accent labels from the 153
154 BURNC following the binary labelling strategy pre- 154
155 sented in (Zhai and Hasegawa-Johnson, 2023). We 155
156 assign positive labels to time-frames correspond- 156
157 ing to audio segments labelled in the BURNC as 157
158 having pitch accents, and negative labels otherwise. 158

159 Following (Zhai and Hasegawa-Johnson, 2023), 159
160 we preprocess the BURNC audios by splitting them 160
161 into overlapping clips of 20s, at 10s intervals. 161

162 The Librispeech dataset consists of 1000 hours 162
163 of audio samples sourced from the LibriVox Project. 163
164 In our work, the dev-clean and test-clean data sub- 164
165 sets were used for model development and evalua- 165
166 tion, respectively. 166

167 The Libri-light dataset is made up of 60,000 167
168 hours of audio and, similarly to the Librispeech 168
169 corpus, was also sourced from the LibriVox Project. 169
170 We used the Libri-light limited resource training 170
171 data subsets, namely, train-1h (LS1), which con- 171
172 sists of 1 hour of labelled audio data. 172

173 3.2 A joint model for prosody and ASR 173

174 Our proposed system uses prosody annotations to 174
175 jointly learn pitch accent detection and automatic 175
176 speech recognition (cf Figure 1). 176

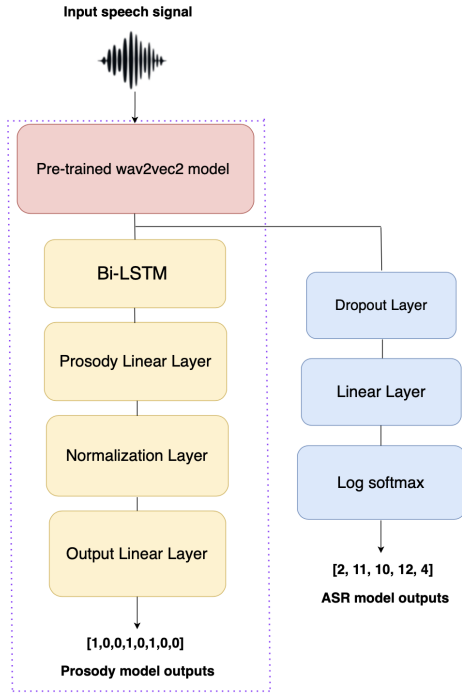


Figure 1: Joint Prosody-ASR Model

For both ASR and prosody detection, raw audio input is sent through the pretrained wav2vec2 model (Baevski et al., 2020) with a language modelling head on top for Connectionist Temporal Classification for the ASR task. For pitch accent detection, we built upon the prosodic event detection model proposed by Zhai and Hasegawa-Johnson (2023), wav2TOBI. In wav2TOBI, wav2vec2 timestep representations are concatenated with fundamental frequency features and fed through a BiLSTM, followed by a classification layer, with mean-squared loss. Our model streamlines wav2TOBI in the sense that we no longer require fundamental frequency features and make pure use of wav2vec2 output representations. On the other hand, we introduce an extra linear layer followed by layer normalisation before the classification layer.

We train our proposed joint model by minimising a joint loss function $\mathcal{L}_j = \mathcal{L}_{asr} + \mathcal{L}_{pad}$, which is the combination of the ASR model loss \mathcal{L}_{asr} and pitch accent detection model loss \mathcal{L}_{pad} .

Results for prosody detection. In the single task setting, for prosody detection, these simple changes result in significant improvements in pitch accent detection performance over wav2TOBI, even without recourse to the additional fundamental frequency features (Table 1).

Model	Tol	Prec	Rec	F1
wav2TOBI	0 ms	0.13	0.11	0.12
	40 ms	0.70	0.61	0.65
	80 ms	0.87	0.74	0.79
	100 ms	0.89	0.76	0.81
Ours Prosody	0 ms	0.37	0.36	0.36
	40 ms	0.82	0.8	0.81
	80 ms	0.89	0.86	0.87
	100 ms	0.9	0.87	0.88
Ours Semi-Sup	0 ms	0.49	0.48	0.48
	40 ms	0.85	0.82	0.83
	80 ms	0.92	0.88	0.9
	100 ms	0.93	0.89	0.9

Table 1: Prosody detection system performance at varying levels of error tolerance (Tol) in milliseconds. Our basic model (Ours Prosody) outperforms the state of the art wav2TOBI system in pitch accent detection. Our semi-supervised training approach (Ours Semi-Sup) further improves performance.

3.3 Semi-supervised Prosodic Event Detection

Our joint prosody-ASR modelling is limited to datasets where these prosodic labels are available. In order to address ASR performance for a dataset like LibriSpeech, where prosodic labels are unavailable, we resort to semi-supervision—specifically, self-training with model voting. In these experiments, we focus on the smaller Libri-light train-1h (LS1) dataset in order to minimise the possible extrapolation error of a larger dataset.

We partitioned the BURNC train set into three subsets as possible hold-outs. For each hold-out subset, we used the remaining 2/3s of the original train set to retrain a new model. We used each of the three models to obtain three predicted labels for each instance of the LS1 train set, and retained the majority class label of these for each instance. The prosody labelled version of LS1 was then added to the full BURNC train set, and then checked over the BURNC test set for performance gains. If there were gains, we repeated the process now with the prosody labelled LS1 as part of the partitioning step, replacing the labels of LS1 at each iteration. Otherwise the process halts. Our training process halted after 4 iterations.²

Results for prosody detection. The single task results for this approach on prosody detection are

²Note that we tried a number of different self-training techniques, but this simple voting technique worked the best.

given in Table 1. We observe that across all measures and error tolerances, this method improves performance and achieves, to our knowledge, the current state-of-the-art.

4 Experimental setup and results

We use the base-960h wav2vec2 pretrained model³. Our models all are trained for 30,000 steps, using default parameters.⁴ Results for word and character error rates (respectively WER and CER) are given in Table 2. All models were fine-tuned for ASR (resp. ASR and prosody jointly) on LS1, and thereafter possibly fine-tuned on BURNC (indicated by ft BURNC). In the fine-tuning process, following Baevski et al. (2020), the pretrained model remains frozen during the first 15K steps, after which the entire model is trained for the remaining 15K steps. The feature encoder remains frozen throughout fine-tuning. For the Prosody-ASR model, we use the prosody labelled version of LS1 outlined above.

Model	LibriSpeech		BURNC	
	WER	CER	WER	CER
ASR-only	6.0	1.0	23.0	7.0
ASR-only (ft BURNC)	4.9	1.0	12.0	4.0
Prosody-ASR	4.3	1.1	20.0	7.7
Prosody-ASR (ft BURNC)	4.3	1.1	20.0	7.0

Table 2: WER and CER results for ASR.

Results. We observe that while there is no great change to CER scores, the joint Prosody-ASR model improves both WER on LibriSpeech and BURNC test data by 28.3% and 13% respectively, showing that a refocus of the wav2vec2 representations on prosody helps to improve ASR performance over LibriSpeech. Interestingly, our semi-supervised approach yields worst WER for BURNC than bypassing prosodic labels for fine-tuning on both LS1 and BURNC. We posit that this may be due to the noisiness of the inferred prosodic labels in the LS1 dataset specifically for BURNC. We therefore tried fine-tuning the joint

³<https://huggingface.co/facebook/wav2vec2-base-960h>

⁴Default parameters are from https://huggingface.co/docs/transformers/en/model_doc/wav2vec2#transformers.Wav2Vec2ForCTC.

model solely on BURNC; however both WER and CER increased to 29.0 and 9.0 respectively. This may be due to the small size of the BURNC dataset in combination with the quantity of information that must be learned in the joint model.

5 Error analysis and discussion

We have shown above that pitch accent detection is useful for improving the performance of pre-trained speech models in ASR tasks within limited resource scenarios. However, even though we improve upon the WER in most of the experiments that we perform with our proposed joint model, we notice that experiments that involve the BURNC dataset tend to on average have higher CER scores. We list two reasons for this phenomenon below and discuss their impact.

Pre-processing mismatches and audio truncation. During the pre-processing of the transcriptions for the BURNC dataset, we transform the text into their uppercase representations and remove all punctuation marks that are not consequential in determining word meaning. For instance, given a word "CHIEF'S", we do not remove the apostrophe (') during pre-processing to form the word "CHIEFS" since doing so changes the inherent meaning of the word. Another example is "S.J.C's", we do not represent it as "SJCS". Even though this text pre-processing approach is well-warranted, it leads to higher CER scores during ASR since the BURNC dataset is filled with a plethora of acronyms, hyphenated and contracted words.

Following the data pre-processing approach utilised in (Zhai and Hasegawa-Johnson, 2023), we split our audios into overlapping clips of 20s, at intervals of 10s, for input to the wav2vec2 model. This however leads to the truncation of words in initial or final position. As a result of this, some of the words that are predicted by the model are incomplete and this leads to a higher CER score.

6 Conclusion

In this paper we have presented an approach for leveraging prosodic information to improve the performance of a pretrained speech model in a limited resource scenario. The results from our experiments demonstrate that re-focusing self-supervised speech models on supra-segmental speech cues such as prosody could lead to significant performance gains in downstream tasks.

7 Limitations

All experiments were carried out under the limited resource setting, with little fine-tuning data, due to the requirement of our method to use prosodic labels. More work is required to investigate the real impact when fine-tuning with larger ASR datasets.

Also, for prosodic cues, we only used pitch-accent, and with hard labels (0 or 1). It is not clear whether other aspects of prosody would also be important for ASR. This question remains open.

References

Sankaranarayanan Ananthakrishnan and Shrikanth Narayanan. 2007. Improved speech recognition using acoustic and lexical correlates of pitch accent in a n-best rescoring framework. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, volume 4, pages IV–873. IEEE.

Alexei Baevski, Michael Auli, and Abdelrahman Mohamed. 2019. Effectiveness of self-supervised pre-training for speech recognition. *arXiv preprint arXiv:1911.03912*.

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.

Mary E Beckman. 1996. The parsing of prosody. *Language and cognitive processes*, 11(1-2):17–68.

Sin-Horng Chen, Jyh-Her Yang, Chen-Yu Chiang, Ming-Chieh Liu, and Yih-Ru Wang. 2012. A new prosody-assisted mandarin asr system. *Ieee transactions on audio, speech, and language processing*, 20(6):1669–1684.

Rolando Coto-Solano. 2021. Explicit tone transcription improves asr performance in extremely low-resource languages: A case study in bribri. In *Proceedings of the first workshop on natural language processing for Indigenous languages of the Americas*, pages 173–184.

Maureen de Seyssel, Marvin Lavechin, Hadrien Titeux, Arthur Thomas, Gwendal Virlet, Andrea Santos Revilla, Guillaume Wisniewski, Bogdan Ludusan, and Emmanuel Dupoux. 2023. [ProsAudit, a prosodic benchmark for self-supervised speech models](#). In *Proceedings of INTERSPEECH 2023*, pages 2963–2967.

Markus Dreyer and Izhak Shafran. 2007. Exploiting prosody for pcfgs with latent annotations. In *INTER-SPEECH*, pages 450–453. Citeseer.

Ewan Dunbar, Xuan Nga Cao, Juan Benjumea, Julien Karadayi, Mathieu Bernard, Laurent Besacier, Xavier Anguera, and Emmanuel Dupoux. 2017. [The zero resource speech challenge 2017](#). In *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 323–330.

Ewan Dunbar, Julien Karadayi, Mathieu Bernard, Xuan-Nga Cao, Robin Algayres, Lucas Ondel, Laurent Besacier, Sakriani Sakti, and Emmanuel Dupoux. 2020. The zero resource speech challenge 2020: Discovering discrete subword and word units. In *NeurIPS Workshop on Self-Supervised Learning for Speech and Audio Processing*.

Michelle Gregory, Mark Johnson, and Eugene Charniak. 2004. Sentence-internal prosody does not help parsing the way punctuation does. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 81–88.

Mark Hasegawa-Johnson, Ken Chen, Jennifer Cole, Sarah Borys, Sung-Suk Kim, Aaron Cohen, Tong Zhang, Jeung-Yoon Choi, Heejin Kim, Taejin Yoon, et al. 2005. Simultaneous recognition of words and prosody in the boston university radio speech corpus. *Speech Communication*, 46(3-4):418–439.

Taniya Hasija, Virender Kadyan, Kalpna Guleria, Abdullah Alharbi, Hashem Alyami, and Nitin Goyal. 2022. Prosodic feature-based discriminatively trained low resource speech recognition system. *Sustainability*, 14(2):614.

Keikichi Hirose and Nobuaki Minematsu. 2004. Use of prosodic features for speech recognition. In *Eighth International Conference on Spoken Language Processing*.

Julia Hirschberg, Diane Litman, and Marc Swerts. 2004. Prosodic and other cues to speech recognition failures. *Speech communication*, 43(1-2):155–175.

Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460.

Jacob Kahn, Morgane Riviere, Weiyi Zheng, Evgeny Kharitonov, Qiantong Xu, Pierre-Emmanuel Mazaré, Julien Karadayi, Vitaliy Liptchinsky, Ronan Collobert, Christian Fuegen, et al. 2020. Libri-light: A benchmark for asr with limited or no supervision. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7669–7673. IEEE.

Jeremy G Kahn, Matthew Lease, Eugene Charniak, Mark Johnson, and Mari Ostendorf. 2005. Effective use of prosody in parsing conversational speech.

417			469
418			470
419			471
			472
			473
420	Jeremy G Kahn and Mari Ostendorf. 2012. Joint reranking of parsing and word recognition with automatic segmentation. <i>Computer Speech & Language</i> , 26(1):1–19.		474
421			475
422			476
423			477
424	Hemant Kathania, Mittul Singh, Tamás Grósz, and Mikko Kurimo. 2020. Data augmentation using prosody and false starts to recognize non-native children’s speech. <i>arXiv preprint arXiv:2008.12914</i> .		478
425			479
426			480
427			
428	Marie Kunešová and Markéta Řezáčková. 2022. Detection of prosodic boundaries in speech using wav2vec 2.0. In <i>International Conference on Text, Speech, and Dialogue</i> , pages 377–388. Springer.		481
429			482
430			483
431			484
432			485
433	Kun Li, Shaoguang Mao, Xu Li, Zhiyong Wu, and Helen Meng. 2018. Automatic lexical stress and pitch accent detection for l2 english speech using multi-distribution deep neural networks. <i>Speech Communication</i> , 96:28–36.		486
434			487
435			488
436			
437	Tatiana Likhomanenko, Qiantong Xu, Jacob Kahn, Gabriel Synnaeve, and Ronan Collobert. 2020. slimlpl: Language-model-free iterative pseudo-labeling. <i>arXiv preprint arXiv:2010.11524</i> .		489
438			490
439			491
440			492
441			493
442			494
443			
444			495
445			496
446			497
447			498
448			
449			499
450			500
451			501
452			502
453			
454			503
455			504
456			505
457			506
458			507
459			508
460			
461			509
462			510
463			511
464			512
465			513
466			514
467			515
468			
			516
			517
			518
			519
			520
			521

- 522 Trang Tran, Shubham Toshniwal, Mohit Bansal, Kevin
523 Gimpel, Karen Livescu, and Mari Ostendorf. 2017.
524 Parsing speech: a neural approach to integrating
525 lexical and acoustic-prosodic information. *arXiv*
526 *preprint arXiv:1704.07287*.
- 527 Klára Vicsi and György Szaszák. 2010. Using prosody
528 to improve automatic speech recognition. *Speech*
529 *Communication*, 52(5):413–426.
- 530 Kai Wei, Dillon Knox, Martin Radfar, Thanh Tran,
531 Markus Müller, Grant P Strimel, Nathan Susanj,
532 Athanasios Mouchtaris, and Maurizio Omologo.
533 2022. A neural prosody encoder for end-to-end
534 dialogue act classification. In *ICASSP 2022-2022*
535 *IEEE International Conference on Acoustics, Speech*
536 *and Signal Processing (ICASSP)*, pages 7047–7051.
537 IEEE.
- 538 Qiantong Xu, Tatiana Likhomanenko, Jacob Kahn,
539 Awni Hannun, Gabriel Synnaeve, and Ronan Col-
540 lobert. 2020. Iterative pseudo-labeling for speech
541 recognition. *arXiv preprint arXiv:2005.09267*.
- 542 Wanyue Zhai and Mark Hasegawa-Johnson. 2023. In
543 *Proceedings of the Annual Conference of the Inter-*
544 *national Speech Communication Association, IN-*
545 *TERSPEECH*, pages 2748–2752. Publisher Copy-
546 right: © 2023 International Speech Communica-
547 tion Association. All rights reserved.; 24th Inter-
548 national Speech Communication Association, Inter-
549 speech 2023 ; Conference date: 20-08-2023 Through
550 24-08-2023. [[link](#)].

8 Appendix

Model	Unlabeled Data	LM	dev-clean	dev-other	test-clean	test-other
10-min labeled						
DiscreteBERT (Baevski et al., 2019)	LS-960	4-gram	15.7	24.1	16.3	25.2
wav2vec 2.0 BASE (Baevski et al., 2020)	LS-960	4-gram	8.9	15.7	9.1	15.6
wav2vec 2.0 LARGE (Baevski et al., 2020)	LL-60k	4-gram	6.3	9.8	6.6	10.3
wav2vec 2.0 LARGE (Baevski et al., 2020)	LL-60k	Transformer	4.6	7.9	4.8	8.2
HUBERT BASE (Hsu et al., 2021)	LS-960	4-gram	9.1	15.0	9.7	15.3
HUBERT LARGE (Hsu et al., 2021)	LL-60k	4-gram	6.1	9.4	6.6	10.1
HUBERT LARGE (Hsu et al., 2021)	LL-60k	Transformer	4.3	7.0	4.7	7.6
HUBERT X-LARGE (Hsu et al., 2021)	LL-60k	Transformer	4.4	6.1	4.6	6.8
1-hour labeled						
DeCoAR 2.0 (Ling and Liu, 2020)	LS-960	4-gram	-	-	13.8	29.1
DiscreteBERT (Baevski et al., 2019)	LS-960	4-gram	8.5	16.4	9.0	17.6
wav2vec 2.0 BASE (Baevski et al., 2020)	LS-960	4-gram	5.0	10.8	5.5	11.3
wav2vec 2.0 LARGE (Baevski et al., 2020)	LL-60k	Transformer	2.9	5.4	2.9	5.8
HUBERT BASE (Hsu et al., 2021)	LS-960	4-gram	5.6	10.9	6.1	11.3
HUBERT LARGE (Hsu et al., 2021)	LL-60k	Transformer	2.6	4.9	2.9	5.4
HUBERT X-LARGE (Hsu et al., 2021)	LL-60k	Transformer	2.6	4.2	2.8	4.8
10-hour labeled						
SlimIPL (Likhomanenko et al., 2020)	LS-960	4-gram + Transformer	5.3	7.9	5.5	9.0
DeCoAR 2.0 (Ling and Liu, 2020)	LS-960	4-gram	-	-	5.4	13.3
DiscreteBERT (Baevski et al., 2019)	LS-960	4-gram	5.3	13.2	5.9	14.1
wav2vec 2.0 BASE (Baevski et al., 2020)	LS-960	4-gram	3.8	9.1	4.3	9.5
wav2vec 2.0 LARGE (Baevski et al., 2020)	LL-60k	Transformer	2.4	4.8	2.6	4.9
HUBERT BASE (Hsu et al., 2021)	LS-960	4-gram	3.9	9.0	4.3	9.4
HUBERT LARGE (Hsu et al., 2021)	LL-60k	Transformer	2.2	4.3	2.4	4.6
HUBERT X-LARGE (Hsu et al., 2021)	LL-60k	Transformer	2.1	3.6	2.3	4.0
100-hour labeled						
IPL (Xu et al., 2020)	LL-60k	4-gram + Transformer	3.19	6.14	3.72	7.11
SlimIPL (Likhomanenko et al., 2020)	LL-60k	4-gram + Transformer	2.2	4.6	2.7	5.2
Noisy Student (Park et al., 2020)	LS-860	LSTM	3.9	8.8	4.2	8.6
DeCoAR 2.0 (Ling and Liu, 2020)	LS-960	4-gram	-	-	5.0	12.1
DiscreteBERT (Baevski et al., 2019)	LS-960	4-gram	4.0	10.9	4.5	12.1
wav2vec 2.0 BASE (Baevski et al., 2020)	LS-960	4-gram	2.7	7.9	3.4	8.0
wav2vec 2.0 LARGE (Baevski et al., 2020)	LL-60k	Transformer	1.9	4.0	2.0	4.0
HUBERT BASE (Hsu et al., 2021)	LS-960	4-gram	2.7	7.8	3.4	8.1
SlimIPL (Likhomanenko et al., 2020)	LL-60k	Transformer	1.8	3.7	2.1	3.9
HUBERT X-LARGE (Hsu et al., 2021)	LL-60k	Transformer	1.7	3.0	1.9	3.5

Table 3: Comparison of ASR model performance on the Librispeech dataset (Hsu et al., 2021)