
ITERATIVE PSEUDO-LABELING FOR SPEECH RECOGNITION

A PREPRINT

Qiantong Xu, Tatiana Likhomanenko, Jacob Kahn, Awni Hannun, Gabriel Synnaeve, Ronan Collobert
Facebook AI Research, Menlo Park & New York, USA
{qiantong, antares, jacobkahn, awni, gab, locronan}@fb.com

August 28, 2020

ABSTRACT

Pseudo-labeling has recently shown promise in end-to-end automatic speech recognition (ASR). We study Iterative Pseudo-Labeling (IPL), a semi-supervised algorithm which efficiently performs multiple iterations of pseudo-labeling on unlabeled data as the acoustic model evolves. In particular, IPL fine tunes an existing model at each iteration using both labeled data and a subset of unlabeled data. We study the main components of IPL: decoding with a language model and data augmentation. We then demonstrate the effectiveness of IPL by achieving state-of-the-art word-error rate on the LIBRISPEECH test sets in both standard and low-resource setting. We also study the effect of language models trained on different corpora to show IPL can effectively utilize additional text. Finally, we release a new large in-domain text corpus which does not overlap with the LIBRISPEECH training transcriptions to foster research in low-resource, semi-supervised ASR.

Index Terms: speech recognition, language modeling, pseudo-labeling, semi-supervised learning, deep learning

1 Introduction

Recent advances in end-to-end speech recognition are largely due to acoustic model (AM) architecture improvements. Some of the most promising are from the Transformer family [1, 2, 3, 4, 5], which give state-of-the-art results on many ASR benchmarks and close the gap between end-to-end and hybrid systems. Given the performance gain from new architectures, research has shifted focus towards leveraging self- and semi-supervised techniques to better utilize unlabeled data. For example, pseudo-labeling successfully boosts the performance on LIBRISPEECH [6] baselines by a large margin [1]. Many algorithms exist which incorporate unlabelled data to improve ASR in the low-resource setting, including representation learning [7, 8, 9], pseudo-labeling [10], local prior matching [11], pseudo-label augmentation [12], adversarial training [13] and back translation [14]. While many of these methods outperform a supervised baseline with limited resources, a large gap to fully-supervised training remains. Furthermore, not all approaches scale easily to large amounts of data, such as that recently used in the LIBRILIGHT benchmark [7].

In this work, we study iterative pseudo-labeling (IPL), a straightforward method that can easily scale to large unlabeled datasets and further boost the performance in both standard- and low-resource settings. IPL is motivated by the simplicity and effectiveness of pseudo-labeling (PL) [10, 1]. A simple extension to [1] involves conducting more iterations of PL as the model trains so as to continuously refine and improve the quality of generated pseudo-labels. That said, training a model from scratch after each round of pseudo-labeling and relabeling a large collection of unlabeled data is expensive. IPL mitigates these challenges by 1) labeling only a subset of the unlabeled data in each iteration, and 2) fine tuning the existing model on this subset, rather than training from scratch. An intuitive motivation for this is shown in Figure 1, where the same acoustic model is trained to convergence with a fixed learning rate; both settings reach a similar word-error rate (WER). Training from scratch with PL is also shown to be roughly equivalent to iterating in a machine translation [15] and a image classification [16] setting.

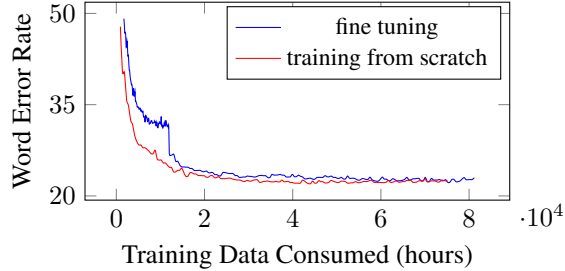


Figure 1: WER on dev-other for two different strategies of using the unlabeled data. In fine tuning, the model is first trained on train-clean-100 only for 160 epochs and then used to generate pseudo-labels on train-clean-360. Both train-clean-100 and train-clean-360 are then used to fine tune the same model. In "from scratch training", same pseudo-labels on train-clean-360 together with the true labels on train-clean-100 are used to train a new model with the same architecture from scratch.

2 Related Work

Semi-supervision in ASR is well-studied [17, 18, 19, 20, 21, 22]; our work builds primarily on recent work with end-to-end systems, especially PL [10, 1]. In [10], PL is shown to be effective with only 100h of labeled audio. The model uses a sequence-to-sequence loss and requires additional pseudo-label filtering to achieve the best results. To mitigate the instability found in sequence-to-sequence decoding, as in [1] all pseudo-labels are generated with models trained with Connectionist Temporal Classification (CTC) loss [23]. Our work extends [1] by conducting more rounds of PL and fine tuning the existing model and demonstrates the effectiveness of the IPL approach in settings with both 960h and 100h of labeled audio. Still other work [8] learns discrete audio feature representations directly from the waveform and works quite well with limited data, even in settings with under 100h of labeled audio. In this setting, learned acoustic features are presumably more competitive than an acoustic model trained end-to-end with MFCC or log-mel filterbank features. Other work including [11], those with CPC baselines [7], and those using adversarial training [13] and back-translation-style techniques [14] also provide promising end-to-end semi-supervised approaches, but results are not comparable as newer end-to-end approaches outperform these works' baselines.

3 Method

In this section, we first introduce the iterative pseudo-labeling algorithm (IPL). We then give theoretical justifications why IPL facilitates effective training. Finally, we perform analysis and experiments on a small-scale labeled dataset.

3.1 Iterative Pseudo-Labeling

Algorithm 1: Iterative pseudo-labeling

Data: Labeled data $L = \{x_i, y_i\}_{i=1}^l$, Unlabeled data $U = \{x'_j\}_{j=1}^u$

Result: Acoustic model p_θ

Initialize p_θ by training on only labeled data L ;

repeat

1. Draw a subset of unpaired data $\tilde{U} \in U$;
2. Apply p_θ and decoding with LM to the subset \tilde{U} to generate $\hat{U} = \{(x, \hat{y}) | x \in \tilde{U}\}$;
3. Fine tune p_θ on $L \cup \hat{U}$ with data augmentation;

until convergence or maximum iterations are reached;

As listed in Algorithm 1, IPL utilizes both labeled and unlabeled data as in the conventional semi-supervised learning. The model minimizes the following loss function:

$$\mathcal{L} = \mathcal{L}_L + \lambda \mathcal{L}_U \tag{1}$$

where \mathcal{L}_L and \mathcal{L}_U denote the parts of the loss function on labeled and unlabeled data accordingly:

$$\mathcal{L}_L = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})} \log(p_\theta(\mathbf{y} | \mathbf{x})) \tag{2}$$

$$\mathcal{L}_U = -\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\hat{\mathbf{y}} \sim p_\theta(\mathbf{y} | \mathbf{x})} \log(p_\theta(\hat{\mathbf{y}} | \mathbf{x})). \tag{3}$$

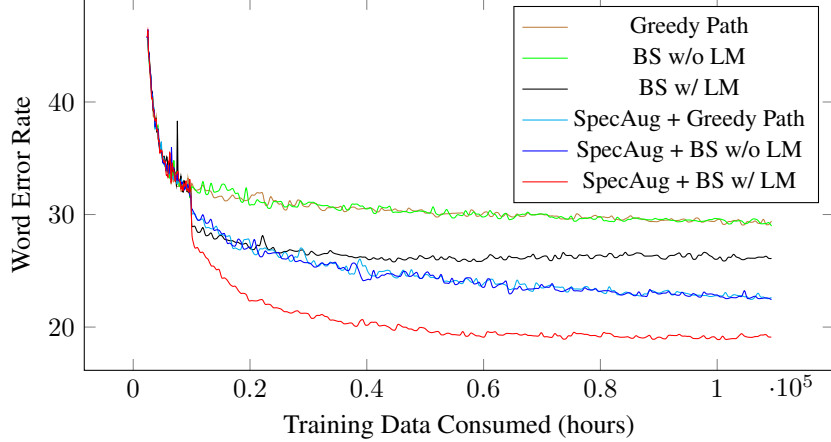


Figure 2: WER on dev-other with different IPL training strategies. Beam-search (BS) decoding is performed with an 4-gram LM.

Note that in ASR, instead of sampling from $p_\theta(\mathbf{y}|\mathbf{x})$, the transcriptions as well as the pseudo-labels are usually selected from the greedy path:

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} p_\theta(\mathbf{y}|\mathbf{x}). \quad (4)$$

3.2 Avoidance of Local Minima

As discussed in [15], one bane of loss (1) optimization in fine tuning is that it tends to get stuck at a local minima after each round of training with existing pseudo-labels; the conditional log likelihood (3) is already maximized when $p_\theta(\mathbf{y}|\mathbf{x})$ matches the underlying data distribution $p_{\theta^*}(\mathbf{y}|\mathbf{x})$, so that $\nabla_{\theta} \mathcal{L}|_{\theta=\theta^*} = 0$. The IPL algorithm has two distinct components: one with respect to the target (\mathbf{y}) and the other with respect to the data (\mathbf{x}), that we found to be effective to overcome this behaviour.

3.2.1 External Language Model

In modern ASR systems, in addition to the acoustic model p_θ , a decoding procedure (typically either WFST-based [24] or beam-search-based (BS) [25, 26] as well as lattice/beam rescoring) is always used to integrate an external language model (LM). Thus, instead of using the greedy path (4) as transcriptions, we consider:

$$\hat{\mathbf{y}}' = \underset{\mathbf{y}}{\operatorname{argmax}} \log p_\theta(\mathbf{y}|\mathbf{x}) + \alpha \log p_{\text{LM}}(\mathbf{y}) + \beta |\mathbf{y}|, \quad (5)$$

where α and β are hyper-parameters [1] usually optimized on validation set. This differs $\hat{\mathbf{y}}'$ from $\hat{\mathbf{y}}$ by introducing extra LM knowledge into transcriptions so that the learned weights θ are no longer optimal given the new labels $\hat{\mathbf{y}}'$, and the model will keep training with $\nabla_{\theta} \mathcal{L}|_{\theta \neq \theta^*} \neq 0$. This is also observed in machine translation [15], where the gain from using greedy-path decoding is limited in self-training with PL.

3.2.2 Data Augmentation

With respect to data, when data augmentation is introduced, the log likelihood the model optimizes also changes. We can rewrite (3) as

$$\mathcal{L}_U = -\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), \mathbf{x}' \sim q(\mathbf{x}'|\mathbf{x})} \mathbb{E}_{\hat{\mathbf{y}} \sim p_\theta(\mathbf{y}|\mathbf{x})} \log(p_\theta(\hat{\mathbf{y}}|\mathbf{x}')), \quad (6)$$

where $q(\cdot)$ is the data augmentation function, which is SpecAugment [27] in our experiments. The model weights optimized before could be no longer optimal given the new augmented input; the model keeps updating with $\nabla_{\theta} \mathcal{L}|_{\theta \neq \theta^*} \neq 0$.

3.2.3 Empirical Study

We use train-clean-100 and train-clean-360 in LIBRISPEECH as labeled and unlabeled training data and dev-other as a validation set. The AMs detailed in Section 4.5, are first trained on labeled data for 100 epochs ($\approx 10^4$ hours) and then continue with IPL. Pseudo-labeling is performed every 10 epochs with all unlabeled data without down-sampling. The learning rate is fixed throughout training for a fair comparison. As shown in Figure 7, if both data

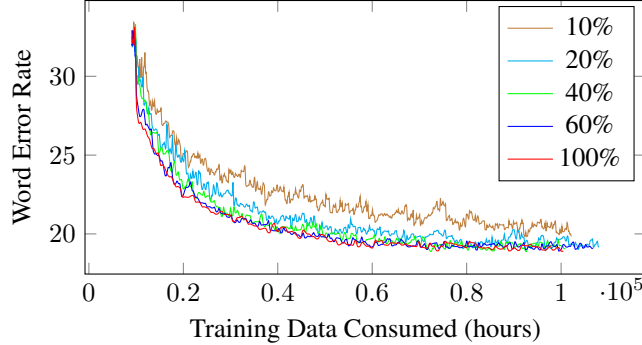


Figure 3: WER on dev-other with different amount of unlabeled data used in each iteration of IPL training (data augmentation and BS decoding with 4-gram LM are applied).

augmentation and LM decoding are used, there is a dramatic WER drop when unlabeled data is first in use, and the WER keeps decreasing as IPL progresses. If either data augmentation or LM decoding is removed, convergence degrades noticeably. Further, if both are removed, adding unlabeled data provides no benefit to IPL, which is in consistent with hitting local minima. One should note that the model is not fully converged at epoch 100; the WER continues decreasing. The contribution of the beam search alone is also limited.

3.3 Dataset Distribution Approximation

Usually, unlabeled dataset U is much larger than the labeled one L ; it is very time-consuming to label the entire U in each round of PL. If only insufficient data is selected in each round, however, the $p(\mathbf{x})$ in (3) will be poorly estimated, which will increase the chance of model overfitting. To balance the trade-off between the accuracy of $p(\mathbf{x})$ approximation and the PL efficiency, with the same setup as in Section 3.2, we conduct an empirical ablation where the PL is performed only on a randomly sampled subset of U . Note that we treat samples from L and U equally following [1], so that the λ in (1) is implicitly set to $|U|/|L|$, the ratio between the number of samples in the two sets.

As shown in Figure 3, even though there is an up to 5-time gap in λ , using 20% to 40% of U can reach the same WER as using 100%. This motivates us to 1) pay less attention to λ tuning and 2) down sample U in PL, so as to perform more rounds of PL in total to better utilize large unlabeled set. On the other hand, using only 10% from U significantly hurts the convergence, which shows the importance of $p(\mathbf{x})$ approximation and sets up a lower bound of down-sample rate.

4 Experiments

4.1 Audio Data

Audio data for our experiments comes from two sources: LIBRISPEECH, containing 960h of audio and paired transcriptions, and audio from LIBRIVOX (54K hours of audio) extracted following [7]. Three setups of labeled data are used: 1) full LIBRISPEECH (960h), 2) train-clean-100 subset (100h) from LIBRISPEECH and 3) the 10-hour training subset from LIBRISPEECH prepared in [7], LibriLight-10. We use the standard development (for all hyper-parameters optimization) and test (for final evaluation only) sets from LIBRISPEECH.

4.2 Gutenberg Text Corpus

As discussed in [10], it is important to remove components of the LM training corpus that overlap with unlabeled audio to ensure the LM has no information about ground truth transcriptions from the unlabeled audio. We study the contribution of the LM to IPL and conduct rigorous experiments when a subset of LIBRISPEECH is used as unlabeled audio. For LM training, we prepare a larger in-domain text corpus using books from Project Gutenberg [28]. To prepare the corpus, we first start with a large subset of English books from Project Gutenberg (which includes some of the 14.5k books present in the LIBRISPEECH LM corpus [6] with 0.8B words) and filter out all books present in LIBRIVOX audio data. We perform the same procedure as in [1] along with a manual matching step to find exact or similar titles (after normalization) in LIBRIVOX (α , β are the same as in [1]), filtering out the resulting books. Similarly, we remove from the corpus books present in the LIBRISPEECH validation and test sets. The resulting filtered corpus is normalized in the same way as in [1] which mimics the normalization in the LIBRISPEECH corpus, but has additional mappings of some

abbreviations and does not split text mid-sentence. We denote this final corpus as $GB \setminus LV$ (2.16B words from 34k books). Further, with the same procedure, we filter out books containing LIBRISPEECH training transcriptions (960h) and form a new corpus denoted as $GB \setminus LV \setminus LS$ (2.11B words from 33.4k books).

4.3 Models

4.3.1 Acoustic Model

We use the best-performing Transformer architecture on LIBRISPEECH and LIBRIVOX with 322M parameters from [1] in our experiments. In particular, there is a convolutional front-end containing 6 layers of 1-D convolutions with kernel-width 3 followed by 36 4-head Transformer blocks [29] with self-attention dimension $D_{tr} = 768$. The 2nd, 4th and the final convolutions in the front-end have stride 2, so the overall sub-sampling rate of the model is 8. The AMs take 80-channel log-mel filterbanks as input and are trained end-to-end with CTC loss.

4.3.2 Language Model

For fair comparison with existing works, IPL experiments in Table 2 use BS decoding with the 4-gram LM (200k top words) used in [1], which is trained on the official LIBRISPEECH LM corpus with transcriptions in LIBRIVOX excluded, denoted as $LS \setminus LV$. The same LM is used for the final beam-search decoding of trained models. Also we are using Transformer LM from [1] for the final beams rescoring. For IPL experiments, where beam rescoring is used for PL generation in addition to the n -gram BS decoding, transformer LM (with the same architecture as transformer LM from [1]) is trained on $LS \setminus LV$. As an ablation study, we train 5-gram LMs on $GB \setminus LV \setminus LS$ and $GB \setminus LV$ with top 200k words in each and without pruning using the KenLM toolkit [30]. The LMs perplexities are listed in Table 1.

Table 1: Perplexities of language models.

Data	$LS \setminus LV$	$GB \setminus LV \setminus LS$	$GB \setminus LV$	Transf. $LS \setminus LV$	Transf.
dev-clean	161.7	101.6	99.7	58.3	48.2
dev-other	152.5	112.9	110.5	59.3	50.2

4.4 Model Training

We use word pieces (WP) [31] as modeling units in our experiments. Following [1], we use the same 10k WP estimated from the training transcriptions, if full LIBRISPEECH training set is used. If `train-clean-100` is in use, we switch to the 5k WP estimated on `train-clean-100` transcriptions as in [11]. We use a lexicon, including words only in training and validation sets, to limit the search space of the BS decoding in IPL. For LibriLight-10 setup, the same units and lexicon as `train-clean-100` are used.

Dropout [32] and layer drop [33] are tuned and used to regularize each model. For models that either use LibriLight-10 as labeled data or trained without LIBRIVOX, we set both dropout and layer drop to 0.3; for models trained on LIBRIVOX, layer drop is set to 0.2 while dropout is 0.2 and 0.15 for models using 100 and 960 hours labeled data, respectively. All models are trained on 64 GPUs with a batch size of 4 per GPU if using LIBRIVOX and 6 otherwise. We use the Adagrad [34] optimizer; the learning rate is initialized to 0.03 and is never decreased for models trained on LIBRIVOX but is halved once at epoch 800 for LIBRISPEECH-only models.

In terms of IPL training, we implemented the automated pipeline in `wav2letter++` [25]. For models trained with LIBRIVOX data, we use only BS with n -grams LM in decoding; while for models trained on LIBRISPEECH only, we apply two stage decoding in PL generating: 1) BS decoding with n -grams LM and 2) beam rescoring with Transformer LM. We use random search with 256 jobs to optimize the hyper-parameters in decoding [26] on `dev-other` and use the optimal values in the subsequent pseudo-labeling. As mentioned in Section 3.3, we only select 20% to 40% of the data in each round of PL if LIBRIVOX is the unlabeled dataset. Otherwise, if the rest of LIBRISPEECH is the unlabeled set, the entire unlabeled set is pseudo-labeled. Pseudo-labels are regenerated every 10 epochs. Note that there is no filtering applied to the pseudo-labels generated, i.e. the whole unlabeled set will be used in training.

4.5 Results

In this section we compare our results with other recent work in semi-supervised learning. All results are listed in Table 2. Given LibriLight-10 as labeled data, our method reaches 26.02% and 19.92% on `test-other` with the full LIBRISPEECH or LIBRIVOX as unlabeled data, respectively; while with `train-clean-100`, we further get 8.95%

Table 2: Comparison of WER with other semi-supervised methods on LIBRISPEECH (LS) and LIBRIVOX (LV) data. The beam-search decoding with 4-gram $LS \setminus LV$ LM is used in IPL training for pseudo-labels generation (for models used LIBRISPEECH as unlabeled data beam rescoring with transformer $LS \setminus LV$ LM is applied for pseudo-labels generation). LM column refers either to the greedy path ("") or to the final decoding and rescoring ("*") with external LMs.

Method	Data (hours)		LM	Dev WER		Test WER	
	Labeled	Unlabeled		clean	other	clean	other
Semi-supervision with PL [10]	LS-100	LS-360	GCNN	5.37	22.13	5.93	24.07
Local Prior Matching [11]	LS-100	LS-860	GCNN	4.87	13.84	4.88	15.28
DeCoAR [9]	LS-100	LS-860	4-gram	-	-	4.74	12.20
vq-wav2vec + BERT [8]	LS-100	LS-860	4-gram	4.0	10.9	4.5	12.1
Semi-supervision with PL, CTC [1]	LS-960	LV-54K	GCNN + Transf.*	2.01	3.95	2.31	4.54
Semi-supervision with PL, S2S [1]			GCNN WP + Transf.*	2.00	3.65	2.09	4.11
Ours, IPL	LL-10	LS-960	-	23.84	25.70	24.58	26.44
			4-gram + Transf.*	23.51	25.48	24.37	26.02
		LV-54K	-	19.76	21.67	20.63	22.38
			4-gram + Transf.*	17.00	19.34	18.03	19.92
	LS-100	LS-860	-	5.41	9.32	5.95	10.28
			4-gram + Transf.*	4.98	7.97	5.59	8.95
		LV-54K	-	4.35	7.90	5.07	8.84
			4-gram + Transf.*	3.19	6.14	3.72	7.11
LS-960	LV-54K	-	2.05	4.12	2.21	4.71	
		4-gram + Transf.*	1.85	3.26	2.10	4.01	

and 7.11% on test-other with the rest of LIBRISPEECH or LIBRIVOX as unlabeled data, respectively. If all of LIBRISPEECH is used as labeled data, we achieve 4.01% on test-other. Our result achieves a clear state-of-the-art in all the three semi-supervised learning setups.¹ The final decoding and rescoring strategies are the same as in [1].

4.6 Analysis

4.6.1 Effectiveness

We conduct 3 rounds of pseudo-labeling on the entire unlabeled set and retraining a new AM from scratch. All the models in this section are trained with n -grams BS decoding in pseudo-labeling. As shown in Table 3, WER on dev-other decreases with better PL generated, but the marginal gain diminishes as iterations continue. IPL, however, clearly outperforms the 3 rounds PL baseline, indicating it is effective to accumulate gains through training with more (up to 80) rounds of PL updates.

Table 3: WER of greedy path on dev-other for IPL and training from scratch for multiple rounds. 4-gram $LS \setminus LV$ LM is used for pseudo-labels generation.

Data		# Rounds of PL				IPL
Labeled	Unlabeled	0	1	2	3	
LS-100	LS-860	27.76	17.1	15.8	15.09	10.69
LS-100	LS + LV	27.76	16.3	12.9	10.95	7.90
LS-960	LV-54K	7.31	5.00	4.69	4.57	4.12

¹The results of this work were first published in May 2020 and the works that we are comparing to here are selected by then. Comparison with the latest methods is shown in Appendix.

4.6.2 Efficiency

Given the same amount of time for 3 rounds of PL training in Table 3 to finish, which is 4, 11 and 17 days from top to bottom, IPL achieves WER 10.69%, 8.50% and 4.14%, respectively. To achieve the same WER after round 3, however, IPL takes only 0.7, 3.3 and 8 days. This efficiency derives directly from the two proposed changes in IPL: 1) fine tuning the existing model with new labels to save computation in re-bootstrapping and 2) down sampling the unlabeled set to shorten the PL time, i.e. 20% down-sample rate leads to 5 time speed up in labeling.

4.6.3 LM Study

Comparison of IPL with different LMs is shown in Table 4 with LMs perplexity in Table 1. All the models in this section are trained with n -grams BS decoding in pseudo-labeling. Given a better LM, IPL better transfers LM knowledge into AM and achieves better performance. IPL can thus effectively leverage large amount of unpaired text, in addition to unpaired audio. However, although the difference in perplexity between the $GB \setminus LV$ and $GB \setminus LV \setminus LS$ LMs is small, there is still a large gap in WER. This is because the LM implicitly leaks the labels of unlabeled audio. Fortunately, comparing WER between $LS \setminus LV$ and $GB \setminus LV \setminus LS$, when the transcription leaking is completely removed, it is still possible to reach similar (or even better) WER by utilizing more in-domain text.

Table 4: IPL training with different LMs in BS decoding. WER on dev-other (test-other) is reported for the greedy path (top), an extra BS decoding with the same n -grams LM used in IPL training (middle) and Transformer LM rescoring (bottom).

Decoding	Data		Language Model		
	Labeled	Unlabeled	$LS \setminus LV$	$GB \setminus LV \setminus LS$	$GB \setminus LV$
None	LS-100	LS-860	10.69 (11.48)	10.80 (11.61)	10.19 (11.09)
	LS-100	LS + LV	7.90 (8.84)	7.21 (8.28)	6.82 (7.89)
	LS-960	LV-54K	4.12 (4.71)	-	4.02 (4.42)
4-gram	LS-100	LS-860	10.05 (10.90)	9.82 (10.49)	9.09 (9.82)
	LS-100	LS + LV	7.21 (8.19)	6.70 (7.74)	6.15 (7.35)
	LS-960	LV-54K	3.67 (4.33)	-	3.65 (4.10)
Trans.	LS-100	LS-860	8.72 (9.51)	8.90 (9.67)	8.25 (9.11)
	LS-100	LS + LV	6.14 (7.11)	5.96 (6.99)	5.56 (6.71)
	LS-960	LV-54K	3.26 (4.01)	-	3.42 (3.83)

4.6.4 Decoding Parameters Mismatch

As shown in Table 4, there is still an observable improvement in WER across greedy path and the BS decoding with n -grams LM. One inhibitor to transferring LM knowledge into the AM is a mismatch in decoding parameters: the parameters optimized on dev-other may not be optimal for decoding unlabeled audio. Thus, in the final stage of IPL training, the marginal WER improvement on dev-other is not reflected similarly on the one on unlabeled audio, so as to prevent the AM from improving.

5 Conclusion

We have shown that iterative pseudo-labeling can give superior results in both standard and low-resource settings and provides an efficient algorithm with which to train as compared to conventional pseudo-labeling approaches. Iterative pseudo-labeling benefits from beam-search decoding with a language model and data augmentation along with dataset sub-sampling which also improves efficiency. With our Transformer acoustic model, IPL achieves the state-of-the-art results on LIBRISPEECH test sets.

References

- [1] G. Synnaeve, Q. Xu, J. Kahn, T. Likhomanenko, E. Grave, V. Prapat, A. Sriram, V. Liptchinsky, and R. Collobert, “End-to-end asr: from supervised to semi-supervised learning with modern architectures,” *arXiv preprint arXiv:1911.08460*, 2019.
- [2] A. Mohamed, D. Okhonko, and L. Zettlemoyer, “Transformers with convolutional context for asr,” *arXiv preprint arXiv:1904.11660*, 2019.
- [3] S. Karita, N. Chen, T. Hayashi, T. Hori, H. Inaguma, Z. Jiang, M. Someki, N. E. Y. Soplin, R. Yamamoto, X. Wang *et al.*, “A comparative study on transformer vs rnn in speech applications,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 449–456.
- [4] K. J. Han, R. Prieto, and T. Ma, “State-of-the-art speech recognition using multi-stream self-attention with dilated 1d convolutions,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 54–61.
- [5] Y. Wang, A. Mohamed, D. Le, C. Liu, A. Xiao, J. Mahadeokar, H. Huang, A. Tjandra, X. Zhang, F. Zhang *et al.*, “Transformer-based acoustic modeling for hybrid speech recognition,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6874–6878.
- [6] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 5206–5210.
- [7] J. Kahn, M. Rivière, W. Zheng, E. Kharitonov, Q. Xu, P.-E. Mazaré, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen *et al.*, “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)*. IEEE, 2020, pp. 7669–7673.
- [8] A. Baeovski, M. Auli, and A. Mohamed, “Effectiveness of self-supervised pre-training for speech recognition,” *arXiv preprint arXiv:1911.03912*, 2019.
- [9] S. Ling, Y. Liu, J. Salazar, and K. Kirchhoff, “Deep contextualized acoustic representations for semi-supervised speech recognition,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6429–6433.
- [10] J. Kahn, A. Lee, and A. Hannun, “Self-training for end-to-end speech recognition,” in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7084–7088.
- [11] W.-N. Hsu, A. Lee, G. Synnaeve, and A. Hannun, “Semi-supervised speech recognition via local prior matching,” *arXiv preprint arXiv:2002.10336*, 2020.
- [12] Y. Chen, W. Wang, and C. Wang, “Semi-supervised asr by end-to-end self-training,” *arXiv preprint arXiv:2001.09128*, 2020.
- [13] A. H. Liu, H.-y. Lee, and L.-s. Lee, “Adversarial training of end-to-end speech recognition using a criticizing language model,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6176–6180.
- [14] M. K. Baskar, S. Watanabe, R. Astudillo, T. Hori, L. Burget, and J. Černocký, “Self-supervised sequence-to-sequence asr using unpaired speech and text,” *arXiv preprint arXiv:1905.01152*, 2019.
- [15] J. He, J. Gu, J. Shen, and M. Ranzato, “Revisiting self-training for neural sequence generation,” in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=SJgdnAVKDH>
- [16] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, “Self-training with noisy student improves imagenet classification,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 687–10 698.
- [17] F. Wessel and H. Ney, “Unsupervised training of acoustic models for large vocabulary continuous speech recognition,” *IEEE Transactions on Speech and Audio Processing*, vol. 13, no. 1, pp. 23–31, 2004.
- [18] L. Lamel, J.-L. Gauvain, and G. Adda, “Lightly supervised and unsupervised acoustic model training,” *Computer Speech & Language*, vol. 16, no. 1, pp. 115–129, 2002.
- [19] J. Ma and R. Schwartz, “Unsupervised versus supervised training of acoustic models,” in *Ninth Annual Conference of the International Speech Communication Association*, 2008.
- [20] K. Yu, M. Gales, L. Wang, and P. C. Woodland, “Unsupervised training and directed manual transcription for lvcsr,” *Speech Communication*, vol. 52, no. 7-8, pp. 652–663, 2010.

- [21] V. Manohar, D. Povey, and S. Khudanpur, “Semi-supervised maximum mutual information training of deep neural network acoustic models,” in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [22] H. Liao, E. McDermott, and A. Senior, “Large scale deep neural network acoustic modeling with semi-supervised training data for youtube video transcription,” in *2013 IEEE Workshop on Automatic Speech Recognition and Understanding*. IEEE, 2013, pp. 368–373.
- [23] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [24] M. Mohri, F. Pereira, and M. Riley, “Weighted finite-state transducers in speech recognition,” *Computer Speech & Language*, vol. 16, no. 1, pp. 69–88, 2002.
- [25] V. Pratap, A. Hannun, Q. Xu, J. Cai, J. Kahn, G. Synnaeve, V. Liptchinsky, and R. Collobert, “wav2letter++: The fastest open-source speech recognition system,” *arXiv preprint arXiv:1812.07625*, 2018.
- [26] N. Zeghidour, Q. Xu, V. Liptchinsky, N. Usunier, G. Synnaeve, and R. Collobert, “Fully convolutional speech recognition,” *arXiv preprint arXiv:1812.06864*, 2018.
- [27] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” *Proc. Interspeech 2019*, pp. 2613–2617, 2019.
- [28] “Project gutenber,” <https://www.gutenberg.org>.
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [30] K. Heafield, “Kenlm: Faster and smaller language model queries,” in *Proceedings of the sixth workshop on statistical machine translation*. Association for Computational Linguistics, 2011, pp. 187–197.
- [31] T. Kudo and J. Richardson, “Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2018, pp. 66–71.
- [32] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [33] A. Fan, E. Grave, and A. Joulin, “Reducing transformer depth on demand with structured dropout,” in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=SylO2yStDr>
- [34] J. Duchi, E. Hazan, and Y. Singer, “Adaptive subgradient methods for online learning and stochastic optimization,” *Journal of machine learning research*, vol. 12, no. Jul, pp. 2121–2159, 2011.
- [35] D. S. Park, Y. Zhang, Y. Jia, W. Han, C.-C. Chiu, B. Li, Y. Wu, and Q. V. Le, “Improved noisy student training for automatic speech recognition,” *arXiv preprint arXiv:2005.09629*, 2020.
- [36] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *arXiv preprint arXiv:2006.11477*, 2020.

A Learning Curves

A.1 Typical IPL curves

We provide training curves for three scenarios with different amount of labeled data. There is an obvious inflection point when IPL starts and unlabeled data is in use. In the first several rounds of IPL, WER decreases dramatically each time when new data and their PLs are generated and consumed. IPL is also able to keep accumulating marginal gains until the end of training. Surprisingly, IPL is able to bootstrap from an immature model, even with about 80% WER. Key to the success here is the LM and lexicon that limits the search space of words.

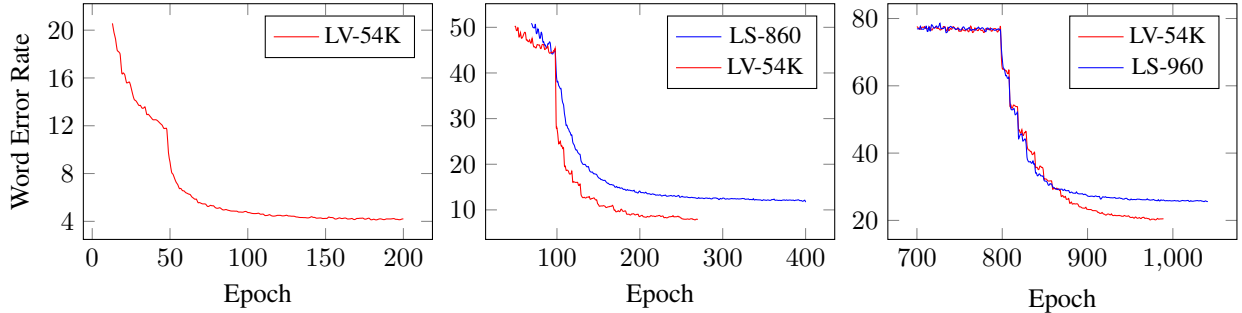


Figure 4: Typical IPL training curves on dev-other. The labeled training data used in each sub-plot is (left) full LIBRISPEECH, (middle) train-clean-100, (right) LibriLight-10. LIBRISPEECH and LIBRIVOX are used as the unlabeled data. Each epoch may contain different amount of data depending on the subset size selected from unlabeled data.

A.2 IPL v.s. Training From Scratch

If the AM is trained from scratch each time as the PL evolves, we can see the gain is diminishing. In addition, both training AM from scratch and label the whole 54K hours dataset are time consuming. Ignoring the labeling time and given the same amount of training time for an AM to converge from scratch, IPL can always outperform and reach better WER.

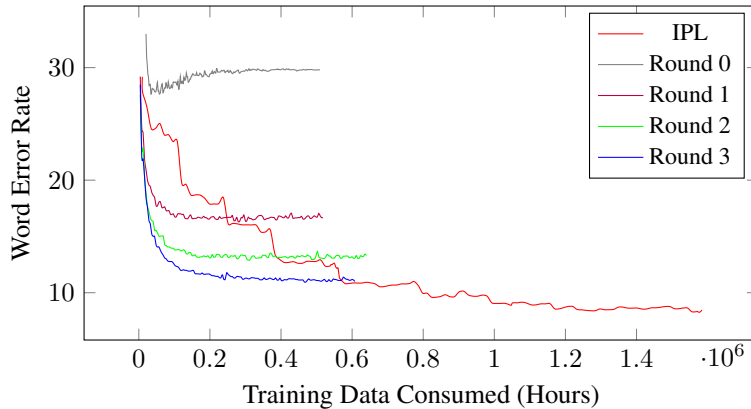


Figure 5: Comparison between IPL and training from scratch with pseudo-labels. Label and unlabeled training data used in these experiments are train-clean-100 and LIBRIVOX plus the rest of LIBRISPEECH, respectively.

A.3 IPL v.s. Retraining

We show in Section 1 that fine tuning and training from scratch converge to the same ball park in the early stage of IPL. This experiment is designed to verify if this still holds when the IPL model is converged. As shown in Figure 6, if we train a new AM from scratch using the pseudo-labels generated from a fully converged IPL model, it will not outperform the IPL model. Even with beam search decoding, both models can still reach similar results. For the model

trained on train-clean-100, WER are 8.83 and 8.87 for IPL and retraining models, respectively; while for models trained on LibriLight-10, WERs are 25.69 and 25.62.

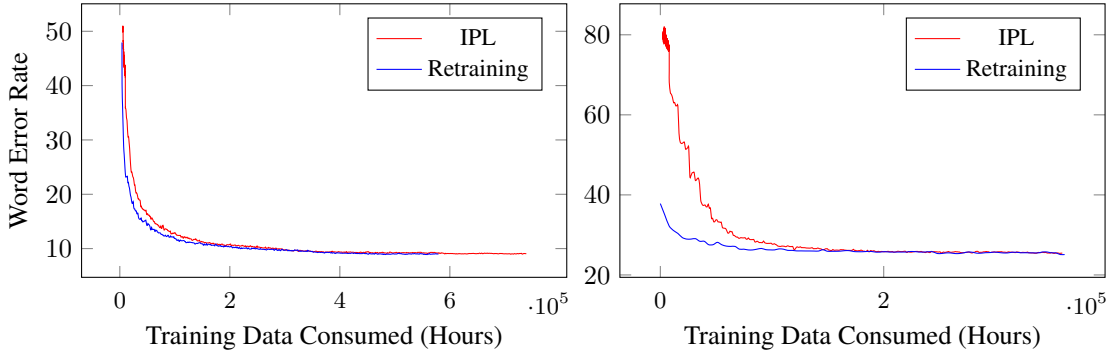


Figure 6: Retraining a final model from scratch with the pseudo-labels generated from a fully converged IPL model. Labeled training data used in the subplots are (left) train-clean-100 and (right) LibriLight-10, while the rest of LIBRISPEECH is used as unlabeled data in both.

B Decoding Parameters

B.1 LM weights for different LMs

We conduct grid search over LM weight used in beam-search decoding and rescoring for an IPL model. Specifically, we use a 4-gram LM in the beam-search decoding to generate a beam of candidates and rescore the beam using an NN-LM. The final WER is not sensitive to the LM weight used in the beam-search decoding, which means as long as the LM weight is reasonable (e.g. between 0 to 2) the good candidates can always be preserved in the beam and later selected out by rescoring. This enables us to pay less attention to the beam-search decoder parameter tuning and focus on rescoring.

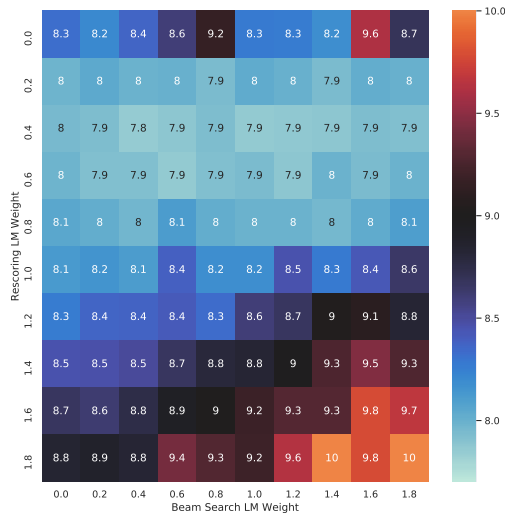


Figure 7: WER with different LM weights used in beam-search decoding and rescoring.

B.2 LM weights on different dataset

We conduct parameter sweep on both development and training set to see if there is a mismatch between the optimal decoding/rescoring parameters on them. As shown in Figure 8, the optimal LM weight is aligned between development and train for both immature and converged model. This shows that the development set can be used as a good proxy of decoder parameter tuning.

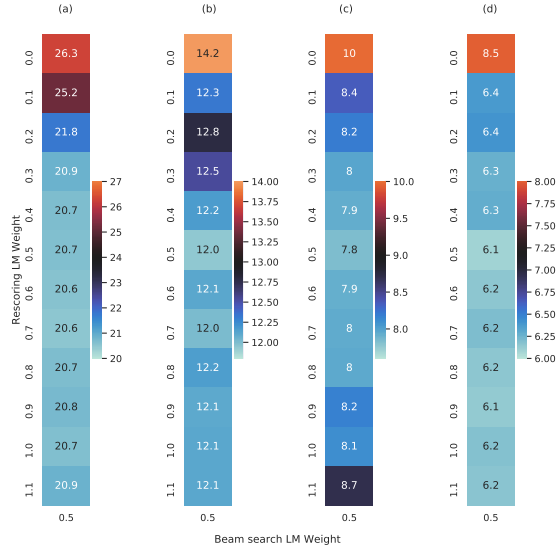


Figure 8: Comparison of the WER heatmap of rescoreing LM weight on development and training set. Left: decoding a checkpoint of an AM that is not fully converged, with (a) WER on dev-other and (b) WER on train-clean-100 and train-other-500. Right: decoding a fully converged AM, with (c) WER on dev-other and (d) WER on train-clean-100 and train-other-500.

C Comparison With The Latest Methods

In Table 5, we mainly compare with two latest works: one with multi-rounds of pseudo-labeling [35] and the other with unsupervised pre-training [36]. Both works achieve better results than us but with different approach. [35] uses the same pseudo-labeling, but is equipped with 1) Listen, Attend and Spell (LAS) network as an acoustic model reaching better WER on labeled data only and 2) filtering mechanism on the pseudo-labels. [36] learns pre-trained speech features on large scale clean speech dataset and further fine tunes the model on labeled transcriptions only.

Table 5: Comparison of WER with latest semi-supervised methods on LIBRISPEECH (LS) and LIBRIVOX (LV) data.

Method	Data (hours)		LM	Dev WER		Test WER	
	Labeled	Unlabeled		clean	other	clean	other
Improved T/S [35]	LS-100	LS-860	LSTM	3.9	8.8	4.2	8.6
	LS-960	LV-54K	LSTM	1.6	3.4	1.7	3.4
wav2vec 2.0 [36]	LL-10	LS-860	Transf.	2.9	5.7	3.2	6.1
	LL-10	LV-54K	Transf.	2.5	5.2	2.6	5.2
	LS-100	LS-860	Transf.	2.1	4.8	2.3	5.0
	LS-100	LV-54K	Transf.	2.0	4.1	2.1	4.4
	LS-960	LV-54K	Transf.	1.6	3.2	1.9	3.5
Ours, IPL	LL-10	LS-960	-	23.84	25.70	24.58	26.44
		4-gram + Transf.*	23.51	25.48	24.37	26.02	
	LV-54K	-	19.76	21.67	20.63	22.38	
		4-gram + Transf.*	17.00	19.34	18.03	19.92	
	LS-100	LS-860	-	5.48	9.32	5.95	10.31
		4-gram + Transf.*	4.98	7.97	5.59	8.95	
LV-54K	-	4.35	7.90	5.07	8.84		
	4-gram + Transf.*	3.19	6.14	3.72	7.11		
LS-960	LV-54K	-	2.05	4.12	2.21	4.71	
		4-gram + Transf.*	1.85	3.26	2.10	4.01	

IPL, however, can easily absorb and combine the advantages of both methods by leveraging a better acoustic model and pre-trained features. This will be a future work to further push the performance of semi-supervised ASR.