

IMPROVING NOISY STUDENT TRAINING ON NON-TARGET DOMAIN DATA FOR AUTOMATIC SPEECH RECOGNITION

Yu Chen^{1,2}, Wen Ding^{2†} and Junjie Lai²

¹The University of Hong Kong, Hong Kong China

²NVIDIA, Shanghai China

u3603296@connect.hku.hk, {wend, julienl}@nvidia.com

ABSTRACT

Noisy Student Training (NST) has recently demonstrated extremely strong performance in Automatic Speech Recognition (ASR). In this paper, we propose a data selection strategy named *LM Filter* to improve the performance of NST on non-target domain data in ASR tasks. Hypotheses with and without a Language Model are generated and the CER differences between them are utilized as a filter threshold. Results reveal that significant improvements of 10.4% compared with no data filtering baselines. We can achieve 3.31% CER in AISHELL-1 test set, which is best result from our knowledge without any other supervised data. We also perform evaluations on the supervised 1000 hour AISHELL-2 dataset and competitive results of 4.73% CER can be achieved.

Index Terms— Data Selection Strategy, Noisy Student Training, Speech Recognition, Semi-supervised Learning

1. INTRODUCTION

In recent years, Semi-Supervised Learning (SSL) has attracted a lot of research interest in many fields of deep learning, such as Automatic Speech Recognition (ASR) [1, 2, 3], Computer Vision [4, 5, 6] and Natural Language Processing [7, 8, 9]. Among these methods, Noisy Student Training (NST) has recently demonstrated extremely strong performances in Image Classification [6] by introducing noise and randomness into traditional Teacher-student Learning [10, 11]. This method further demonstrates its robustness in the ASR field [12, 13, 14]. After combing with pre-train methods [15], NST is shown to be a vital component for achieving SOTA results on a number of datasets, e.g. Librispeech [16].

However, NST has not been widely investigated in ASR tasks when the domain of the supervised data does not match the unsupervised data. Noise and domain play an important role in ASR [17] and the abundant unsupervised data from social media may not always match the domain of the desired

task. Thus, proper data selection techniques are required to remove noise and select data that is close to the target domain [18]. The most common filter in ASR is the Confidence Score that selects the most trustworthy transcriptions based on confidence estimation and threshold [19, 20, 21]. However, this method is not always promising in scenarios with large amount of unlabelled data with domain mismatches. Another recent unsupervised data selection technique is investigated in [18], where a contrastive Language Model is applied as a data selector to better improve the target-domain ASR task.

In this paper, we propose a novel data selection strategy named *LM Filter* which can utilize model differences to filter more valuable non-target domain data to improve the performance of NST. We leverage concept of contrastive LM and data selection method in [22]. Our *LM Filter* is based on hypotheses from LM to gradually remove noisy data inside each iteration of NST method. The filter condition is relaxed through the NST iteration to make the model advance gradually in due order. This method has the following benefits:

- No additional data selection models are required. Model differences can be obtained from different decoding strategies (e.g. with/without LM).
- Label is not required to perform the data selection and it is totally unsupervised.
- Less time and resources are utilized to run the NST method and it can converge faster in fewer iterations.

Experiments on AISHELL-1 [23] as supervised data and WenetSpeech [24] as unsupervised data indicate a significant improvement of 10.4% comparing with no data filtering baselines. When combined AISHELL-2 [25] and WenetSpeech as unsupervised data, 3.31% character error rate(CER) is achieved on AISHELL-1 test set, which is the best result from our knowledge without any other supervised data on this test set. *LM Filter* further demonstrates its robustness in larger dataset such as AISHELL-2 (supervised) and WenetSpeech (unsupervised) to achieve promising result of 4.73% which has 13.6% improvement comparing with the baseline.

The rest of the paper is organized as follows. Section 2 briefly introduces the basic concepts and methods of NST in

Wen Ding[†] is the corresponding author. This work is done during Yu Chen’s internship at NVIDIA. Thanks to Yuekai Zhang and Hainan Xu for helpful suggestions. This work has been open-sourced into **WeNet toolkit**.

ASR. Our proposed data selection strategy LM Filter will be included in section 3. Experiment details are introduced in Section 4. Eventually we give our conclusions in section 5.

2. NOISY STUDENT TRAINING FOR ASR

Noisy Student Training [16] is an iterative self-training method evolved from Teacher Student Learning, the pipeline of which is illustrated in Fig 1. Initially a teacher model is

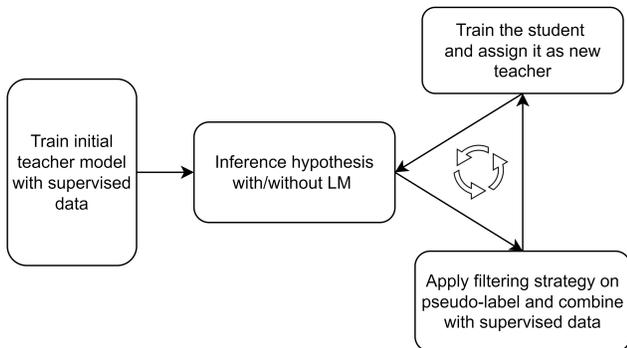


Fig. 1. Noisy Student Training pipeline for ASR

trained with supervised data and pseudo-labels are generated by the teacher. Then data augmentation methods such as SpecAug [26] and speed perturbation are applied during training and the student model is trained using both augmented supervised data and pseudo-data.

In our pipeline, we follow the design that the student model always has same parameters as the teacher, and adopts dropouts and stochastic depths so that the student could be more robust and general than the teacher when it is trained. After training finishes, the student model will be assigned as the new teacher and the whole pipeline will iterate. After several rounds of training, models trained to tolerant noises and augmentations will tend to have better performance generally.

3. DATA SELECTION STRATEGY

Data selection and filtering play a significant role in SSL especially in an out-of-domain situation. This circumstance occurs frequently in the industry when we have limited labelled data in the desired area or in low resource tasks.

Initially, standard NST is performed on AISHELL-1 as supervised data and AISHELL-2 as unsupervised data without any other filtering strategy. The generated pseudo labels have quite promising results of 8.38% CER but when unsupervised data is set to WenetSpeech which has different domain and recording settings, pseudo-label’s CER increases dramatically to 47.1% which is unacceptable for training. *LM Filter* is then proposed to improve the performances of NST when non-target domain data is provided.

Our hypothesis is that if a language model believes the sentence does not require any further modification, then this sentence has higher probability of being a correct pseudo-label. Here we introduce two definitions and examples to better understand how our *LM Filter* works.

- **CER-Hypo** is the CER between student model’s hypothesis with greedy decoding and student model’s hypothesis with Language model.
- **CER-Label** is the CER between student model’s hypothesis with Language model and the true label.

We evaluated our method on the Mandarin corpus using CER while the same definition can be applied to other languages e.g. English by replacing CER with WER. Two cases are listed in Fig 2. In case 1, the difference between the hypothesis with greedy decoding and the hypothesis decoding with LM is 1 character (eg. char “数” and char “诉”) so the CER-Hypo is 16.67% . The CER-Label is also 16.67% in this case, since it takes 1 substitution step to transfer the hypothesis to true label (eg. char “申” and char “胜”). In case 2, the sentence is more challenging than the first case for the initial student model. The student model learns partial acoustic features but the transcripts are mostly wrong. The LM tends to make more modifications due to the low probabilities of such sentence in the corpus. The CER-Hypo and CER-Label both are extremely high in this case.

Case 1

label: 完全能够申诉

Hyp: 完全能够胜数

Hyp(LM): 完全能够胜诉

CER-Label = 1/6 = 16.67% CER-Hypo = 16.67%

Case 2

label: 七十六号的电文纸是有数量的

Hyp: 其实掉的变温是有重到

Hyp(LM): 其治疗的重要

CER-Label = 12/13 = 92.31% CER-Hypo = 70.00%

Fig. 2. Examples of how to calculate the CER-Label and CER-hyp of sentences.

A large amount of cases suggest that **CER-Hypo** and **CER-Label** have strong positive correlations, sentences with lower CER-Hypo tend to have lower CER-Label. Our *LM Filter* uses CER-Hypo as a threshold (eg.10%) to filter out high CER-Label data. We also observe that unsupervised data with similar domain to supervised data are more likely to have lower CER-Hypo values. For unsupervised data from Youtube, similar topics in “readings”and “news” tend to have lower CER-Hypo and non-target domain such as “drama” and

“variety” are more likely to be removed by the *LM Filter*. We also propose a speaking rate filter for WenetSpeech dataset, which is the hypothesis length divided by audio time. The music and song audios that are common elements of drama and variety shows can be effectively removed by this filter.

4. EXPERIMENTS AND RESULTS

4.1. Datasets and domain description

We evaluate our proposed data selection strategy on the following three datasets: AISHELL-1, AISHELL-2 and WenetSpeech. AISHELL-1 is a 178-hour open-sourced Mandarin speech corpus, with strictly annotated and inspected transcriptions which mainly covers 5 topics of Finance, Technology, Sports, Entertainments and News. AISHELL-2 consists of 1k hours of Mandarin speech with the same device and recording environment settings as AISHELL-1. The major topics of these two datasets are similar, but the transcripts and audios of the test set are different. WenetSpeech has 10k hours of speech where transcripts are generated by OCR on video data from Youtube and Podcast, which lacks inspection and accuracy. Domains are diverse and mostly consists of Drama, Variety show and Audio books.

4.2. Experiment settings

First, we use AISHELL-1 as the supervised dataset and treat AISHELL-2 and WenetSpeech as unsupervised data. Initially 1k hours of WenetSpeech data are randomly selected to match the size of AISHELL-2. And then the size of WenetSpeech data is increased up to 4k hours to test the degree of saturation for unsupervised data. Eventually, we switch the supervised dataset to AISHELL-2 to evaluate the performances of our data selection strategy on industrial-level supervised datasets. The upper bound of data ratio for supervised and unsupervised data is set to 1:9.

The neural structures for both teacher and student models are the same, which is a 16-layer Conformer model [27]. Our language model is a 5-gram model with corpus contains training texts as well as extra wiki texts. All experiments are conducted in WeNet toolkit [28] and NVIDIA A100 GPUs. We perform 7 iterations of NST with and without data selection strategy on WenetSpeech and 5 iterations on AISHELL2.

4.3. Baselines

Supervised baseline using only AISHELL-1 data, which is the initial teacher of NST iterations is shown in Table 1. Then supervised training is done on AISHELL-1 data mixing with supervised AISHELL-2 and WenetSpeech. These two results are considered as ceilings of our model’s performance. Then standard NST experiments is conducted without data selection strategy using AISHELL-1 as supervised data and AISHELL-2 as unsupervised data, the results of which are

Table 1. CER for supervised baselines and standard NST first iteration with AISHELL-2 and 1k WenetSpeech dataset.

Supervised	Unsupervised	Test CER
AISHELL-1 Only	—	4.85
AISHELL-1 + WenetSpeech	—	3.54
AISHELL-1 + AISHELL-2	—	1.01
AISHELL-1	WenetSpeech	5.52
AISHELL-1	AISHELL-2	3.99

shown in Table 1. The 3.99% CER can be achieved after first NST iteration because these two datasets have similar domain and recording settings. The closer the topics and configs, the better performance the NST algorithm will have. In the case of the ideal data distribution, the filtering approach is not required. However, in the majority of recognition jobs, this condition is not typical. After first NST iteration with WenetSpeech pseudo-label, the CER increases to 5.52%, which is even higher than the supervised baseline using only AISHELL-1 data. To reduce CER of pseudo-labels and make training easier in early stages, an appropriate filter is required.

4.4. Data selection strategy performances

Performances of *LM Filter* of supervised AISHELL-1 data and unsupervised WenetSpeech data are shown in Table 2. The best 4.31% CER can be achieved after 7 iterations in the test set. There can be relatively 11.13% CER reduction compared with the supervised baseline. In addition to the test set’s CER, the following three metrics are used to assess the quality of the pseudo-label : *Pseudo CER* which is referred to the CER of pseudo-labels, *Filtered CER* and *Filtered hours* which are the CER and the duration of filtered unsupervised data.

Results in initial iteration indicate that the *LM Filter* can significantly decrease the Pseudo CER from 47.1% to 25.18% which makes the pseudo-label satisfactory for further training. The Pseudo CER and Filtered CER decrease as the number of iterations rises, and *LM Filter* permits more filtered data to be fed into the model. This suggests that *LM Filter* may gradually learn noisy information, and our student model could make even greater use of non-target domain data.

4.5. Discussions

Multiple NST iterations: Multiple NST iterations are conducted to show our proposed data selection strategy can achieve better performance and converge faster. Fig 3 displays all the results of AISHELL-1. When WenetSpeech is used as unsupervised data in section (a), after 7 iterations of NST without filter strategy, a negligible improvement is obtained. In contrast, *LM Filter* can yield a relative improvement of 10.4% with faster training time. With our *LM Filter*, CERs for test set and pseudo-label drop gradually and the

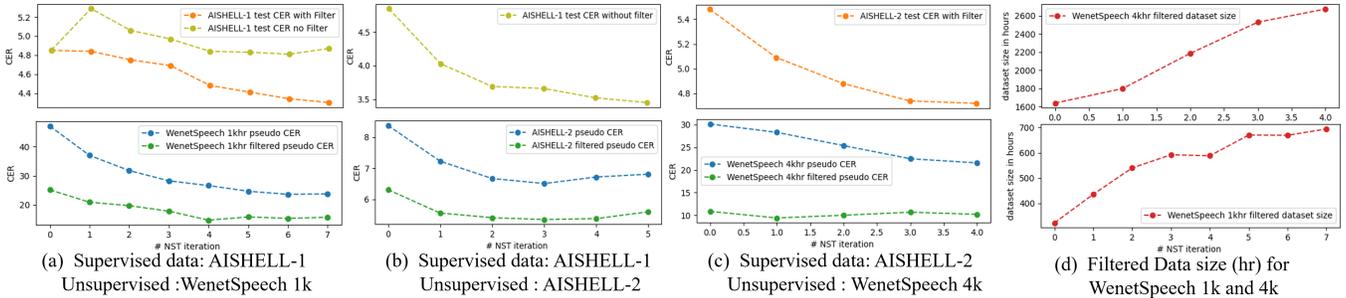


Fig. 3. This figure illustrates results of NST with our data selection strategy using different supervised and unsupervised data. CER performances in test sets of standard NST without *LM filter* are marked in yellow lines and with *LM filter* are in orange lines. Pseudo CERs are shown in blue lines and green lines gives the Filtered CER. Filtered hours of unsupervised data are shown in (d).

Table 2. Performances of *LM Filter* on supervised AISHELL-1 data and unsupervised 1k WenetSpeech data, including the CER of test set, Pseudo CER, Filtered CER and Filtered hours.

# NST Iter	AISHELL-1 test CER	Pseudo CER	Filtered CER	Filtered hours
0	4.85	47.10	25.18	323
1	4.86	37.02	20.93	436
2	4.75	31.81	19.74	540
3	4.69	28.27	17.85	592
4	4.48	26.64	14.76	588
5	4.41	24.70	15.86	670
6	4.34	23.64	15.40	669
7	4.31	23.79	15.75	694

filtered data size grows during each training iterations. In section (b), when using AISHELL-2 as unsupervised data, 3.45 % CER can be achieved after the 5 NST training without filter. The relatively small initial Pseudo CER of 8.38 % in AISHELL-2 indicates that unsupervised data with matched domain can generate effective pseudo-labels to acquire the requirement of NST training. Additionally, we perform extra iteration that combines all pseudo-labels that have been filtered by final NST models on both WenetSpeech and AISHELL-2, yielding the best CER result of 3.31%. According to our understanding, this is the best current result in AISHELL-1 test set without any further supervised data.

Impact of domains: Our experiments indicate that NST approach is very sensitive to the domain issue. Domain can have a significant impact on the effectiveness of the NST algorithm. The quality of pseudo-labels tends to rise if WenetSpeech samples are taken from tags that are more closely related to the AISHELL domain (such as Readings and News). In contrast, tags for drama and variety show that are not commonly used in AISHELL yielded inferior pseudo-labels. The pseudo-labels’ quality will further affect the filtered data size

Table 3. Results of AISHELL-2 test set when using supervised AISHELL-2 data and unsupervised 4k hr WenetSpeech data after applying *LM Filter*.

# NST Iter	AISHELL-2 test CER	Pseudo CER	Filtered CER	Filtered hours
0	5.48	30.10	11.73	1637
1	5.09	28.31	9.39	2016
2	4.88	25.38	9.99	2186
3	4.74	22.47	10.66	2528
4	4.73	22.23	10.43	2734

and NST iterations’ converging speed. Among all the topics, we also discover that sources like Audio books and Podcast most likely provide pseudo-labels with higher qualities.

Effectiveness on large dataset: To further demonstrate our *LM Filter*’s effectiveness on large supervised dataset, we conduct experiments using AISHELL-2 as supervised data. The CER results are shown in Table 3. In the AISHELL-2 test set, 13.6% relative improvement is achieved, which further demonstrates *LM Filter*’s scalability on larger supervised data under industrial scale. Detail performances is shown in plot (c) of Fig 3, it illustrates similar trends as AISHELL-1.

5. CONCLUSIONS

In this paper, a novel data selection strategy named *LM Filter* is proposed to improve the performances of NST in non-target domain data, which utilizes the model differences from decoding strategies. Results reveal significant improvements of 10.4% compared to baselines with no data filtering. we obtain 3.31% CER in AISHELL-1 test set, which is best result according to our knowledge without any further supervised data. In addition, we perform evaluations on 1k hour AISHELL-2 dataset and achieve 4.73% CER on test set, which further demonstrates the robustness of *LM Filter* with larger supervised data.

6. REFERENCES

- [1] Lori Lamel, Jean-Luc Gauvain, and Gilles Adda, “Lightly supervised and unsupervised acoustic model training,” *Comput. Speech Lang.*, vol. 16, no. 1, pp. 115–129, 2002.
- [2] Samuel Thomas, Michael L. Seltzer, Kenneth Church, and Hynek Hermansky, “Deep neural network features and semi-supervised training for low resource speech recognition,” in *ICASSP 2013*, 2013, pp. 6704–6708.
- [3] Gabriel Synnaeve, Qiantong Xu, Jacob Kahn, Tatiana Likhomanenko, Edouard Grave, Vineel Pratap, Anuroop Sriram, Vitaliy Liptchinsky, and Ronan Collobert, “End-to-end ASR: from supervised to semi-supervised learning with modern architectures,” in *ICML 2020 Workshop on Self-supervision in Audio and Speech*, 2020.
- [4] I. Zeki Yalniz, Hervé Jégou, Kan Chen, Manohar Paluri, and Dhruv Mahajan, “Billion-scale semi-supervised learning for image classification,” *arXiv e-prints*, p. arXiv:1905.00546, May 2019.
- [5] Yauhen Babakhin, Artsiom Sanakoyeu, and Hirotohi Kitamura, “Semi-supervised segmentation of salt bodies in seismic images using an ensemble of convolutional neural networks,” in *German Conference on Pattern Recognition*. Springer, 2019, pp. 218–231.
- [6] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, “Self-training with noisy student improves imagenet classification,” in *2020 CVPR*, 2020, pp. 10687–10698.
- [7] Junxian He, Jiatao Gu, Jiajun Shen, and Marc’Aurelio Ranzato, “Revisiting self-training for neural sequence generation,” in *2020 ICLR*, 2020.
- [8] David Yarowsky, “Unsupervised word sense disambiguation rivaling supervised methods,” in *1995 ACL*, 1995, pp. 189–196.
- [9] Roi Reichart and Ari Rappoport, “Self-training for enhancement and domain adaptation of statistical parsers trained on small datasets,” in *2007 ACL*, 2007, pp. 616–623.
- [10] H. Scudder, “Probability of error of some adaptive pattern-recognition machines,” *IEEE Transactions on Information Theory*, vol. 11, no. 3, pp. 363–371, 1965.
- [11] Ellen Riloff and Janyce Wiebe, “Learning extraction patterns for subjective expressions,” in *2003 EMNLP*, 2003, pp. 105–112.
- [12] 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,” *Proc. Interspeech 2020*, pp. 2817–2821, 2020.
- [13] Thibault Doutré, Wei Han, Min Ma, Zhiyun Lu, Chung-Cheng Chiu, Ruoming Pang, Arun Narayanan, Ananya Misra, Yu Zhang, and Liangliang Cao, “Improving streaming automatic speech recognition with non-streaming model distillation on unsupervised data,” in *2021 ICASSP*, 2021, pp. 6558–6562.
- [14] Haaris Mehmood, Agnieszka Dobrowolska, Karthikeyan Saravanan, and Mete Ozay, “FedNST: Federated Noisy Student Training for Automatic Speech Recognition,” in *Proc. Interspeech 2022*, 2022, pp. 1001–1005.
- [15] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 12449–12460, 2020.
- [16] Y. Zhang, J. Qin, D. S. Park, W. Han, C.-C. Chiu, R. Pang, Q. V. Le, and Y. Wu, “Pushing the limits of semi-supervised learning for automatic speech recognition,” *CoRR*, vol. abs/2010.10504, 2020.
- [17] Michael L. Seltzer, Dong Yu, and Yongqiang Wang, “An investigation of deep neural networks for noise robust speech recognition,” in *2013 ICASSP*, 2013, pp. 7398–7402.
- [18] Z. Lu., Yongqiang W., Y. Zhang, W. Han, Zhehuai Chen, and Parisa Haghani, “Unsupervised Data Selection via Discrete Speech Representation for ASR,” in *Proc. Interspeech 2022*, 2022, pp. 3393–3397.
- [19] George Zavalagkos and Thomas Colthurst, “Utilizing untranscribed training data to improve performance,” in *DARPA Broadcast News Transcription and Understanding Workshop, Landsdowne*. Citeseer, 1998.
- [20] D. Charlet, “Confidence-measure-driven unsupervised incremental adaptation for hmm-based speech recognition,” in *2001 ICASSP*, 2001, vol. 1, pp. 357–360 vol.1.
- [21] H.Y. Chan and P. Woodland, “Improving broadcast news transcription by lightly supervised discriminative training,” in *2004 ICASSP*, 2004, vol. 1, pp. I-737.
- [22] W. Zheng, A. Xiao, G. Keren, D. Le, F. Zhang, C. Fuegen, O. Kalinli, Y. Saraf, and A. Mohamed, “Scaling asr improves zero and few shot learning,” *Proc. Interspeech 2022*, pp. 5135–5139, 2022.
- [23] H. Bu, J. Du, X. Na, B. Wu, and H. Zheng, “Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline,” in *2017 O-COCOSDA*, 2017, pp. 1–5.
- [24] B. Zhang, H. Lv, P. Guo, Q. Shao, C. Yang, L. Xie, X. Xu, H. Bu, X. Chen, C. Zeng, et al., “Wenetspeech: A 10000+ hours multi-domain mandarin corpus for speech recognition,” in *2022 ICASSP*. IEEE, 2022, pp. 6182–6186.
- [25] Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu, “AISHELL-2: Transforming Mandarin ASR Research Into Industrial Scale,” *arXiv e-prints*, p. arXiv:1808.10583, Aug. 2018.
- [26] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, “SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition,” in *Proc. Interspeech 2019*, 2019, pp. 2613–2617.
- [27] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, et al., “Conformer: Convolution-augmented transformer for speech recognition,” *Proc. Interspeech 2020*, pp. 5036–5040, 2020.
- [28] Zhuoyuan Yao, Di Wu, Xiong Wang, Binbin Zhang, Fan Yu, Chao Yang, Zhendong Peng, Xiaoyu Chen, Lei Xie, and Xin Lei, “Wenet: Production oriented streaming and non-streaming end-to-end speech recognition toolkit,” in *Proc. Interspeech*, Brno, Czech Republic, 2021, IEEE.