



Noisy student-teacher training for robust keyword spotting

Hyun-Jin Park, Pai Zhu, Ignacio Lopez Moreno, Niranjan Subrahmanya¹

Google Inc., USA

{hjpark, paizhu, elnota}@google.com, niranjan.udupa@gmail.com

Abstract

We propose self-training with noisy student-teacher approach for streaming keyword spotting, that can utilize large-scale unlabeled data and aggressive data augmentation. The proposed method applies aggressive data augmentation (spectral augmentation) on the input of both student and teacher and utilize unlabeled data at scale, which significantly boosts the accuracy of student against challenging conditions. Such aggressive augmentation usually degrades model performance when used with supervised training with hard-labeled data. Experiments show that aggressive spec augmentation on baseline supervised training method degrades accuracy, while the proposed self-training with noisy student-teacher training improves accuracy of some difficult-conditioned test sets by as much as 60%.

Index Terms: keyword spotting, self-training, noisy student teacher, spec augmentation, semi-supervised

1. Introduction

Supervised learning has been the major approach in keyword spotting area [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11]. Although it has been successful, supervised learning requires high quality labeled data at scale, which often requires expensive human efforts or costs to obtain. Motivated from such difficulty, semi-supervised learning [12, 13, 14, 15, 16] and self-training [17, 18, 19, 20, 21] approaches were introduced recently. Those approaches utilize large unlabeled data in addition to smaller amounts of labeled data and achieve performance comparable to supervised models trained with large amounts of labeled data.

Semi-supervised learning approaches utilize unlabeled data to learn hidden layer activations that best predict neighboring data (temporally or spatially close sensory features), which is then used as the input feature to a classification network trained with small amount of labeled data (supervised training) [14, 13, 15, 16]. It is shown that semi-supervised learning with small amount of labeled data can achieve performance comparable to supervised learning with larger amount of data. On the other hand, Self-training approaches utilize unlabeled data by using a teacher network to generate soft-labels (pseudo-labels) which is then used to train student network [17, 18, 19, 20, 21]. Such student-teacher training step can be repeated as long as the performance improves, with the student being a teacher in the next step. Data augmentation is often used together during student-training step for further improvements [17, 18].

Data augmentation is another effective technique to boost model accuracy without requiring more training data. Augmenting data by adding reverberation or mixing with noise have been used in ASR (automatic speech recognition) and KWS (keyword spotting) [22] with some success. Recently introduced spectral augmentation [23, 24] is a new data augmentation technique shown to boost ASR accuracy significantly. In

a recent work, [25] showed that applying spectral augmentation on student's input can improve ASR accuracy in self-training setup.

In this paper, we explore an application of self-training with labeled and unlabeled data where aggressive data augmentation (spec augmentation) is applied to the input of both student and teacher. The proposed student-teacher training approach enables utilization of unlabeled (unsupervised) training data for KWS problem, and also helps in applying aggressive spectral augmentation to boost diversity of training data further.

Aggressive data augmentation can degrade accuracy in keyword spotting when used with supervised training with hard-labels ($\in \{0, 1\}$). If one applies very aggressive augmentation on a positive example, one can end up with an example that may actually be seen as negative but still has positive label. Although it may be not frequent, it can degrades the accuracy significantly by increasing false accept rate. With the proposed noisy student-teacher approach, teacher model generates soft-labels ($\in [0, 1]$) which dynamically reflects the degree of degradation in the input pattern. Supervised-training with predetermined hard-labels cannot reflect such changes of input pattern. With experiments, we show that the proposed noisy student-teacher training with spec augmentation boosts accuracy of the model in more challenging conditions (accented, noisy and far-field). Such benefits can be explained from the use of large scale unlabeled data and aggressive data augmentation.

We describe the proposed approach in Section 2. Then we show experimental setup in Section 3, and the results in Section 4. We conclude with discussions in Section 5.

2. Noisy student-teacher self-training with spec augmentation

2.1. Noisy student-teacher Self-training

We propose noisy student-teacher self-training approach which consists of two major stages. In the first stage, we train a teacher model (which is also a baseline model) using conventional supervised training method on labeled data (shown in Fig. 1 (a)). We use the same architecture and training method developed in previous work [1] for the first stage model. Also the same conventional data augmentation method (add reverberation and background noises)[22] is used in the first stage. The learned teacher model is passed to the second stage.

In the second stage, we train a student model using soft-labels generated from the teacher model trained in previous stage. Since the teacher provides soft-label, we can use additional unlabeled data for training student model. Also we can add more aggressive data augmentation on top of existing classical one to boost accuracy. Specifically we apply spectral augmentation which masks specific frequencies or time frames completely. Such strong modification might even change a positive pattern to a negative one, which will make an incorrect training example under supervised training method. But

¹This author contributed to this work while working at Google.

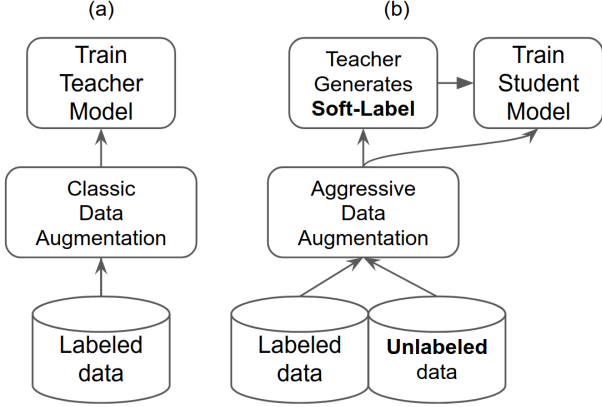


Figure 1: *Self-training with Noisy student-teacher* (a) *Teacher model is trained by labeled data.* (b) *Student model is trained by labeled+unlabeled data with teacher logit and data augmentation.*

in student-teacher approach, the teacher can compensate for such drastic changes by generating correspondingly lower confidence. To get such benefits, both the teacher and the student model is getting the same augmented data (Fig. 1(b)).

In the original self-training with noisy-student approaches [18, 25], a teacher model is provided with clean data and only the student is given noisy (augmented) data. This seems to be working well for multi-class classifications such as ImageNet (objects) or ASR (graphemes) problems. But we found that providing the same noisy input to the teacher and the student achieves better performance in Keyword Spotting problem. This seems to be due to the difference of the problem, where KWS is a binary classification task with highly unbalanced pattern space. In KWS, the space of positive pattern is much smaller than that of negative patterns. Thus augmenting a positive pattern can easily result in moving the pattern into the space of negative patterns. Also our approach is different from [18, 25] that both labeled and unlabeled data go through the teacher model to produce soft-labels used to train the student model. In previous works, labeled data is used for computing supervised loss on the student model while unlabeled data is used to generate soft-label. Also unlike [25], we don't have separate data selection for the second stage.

Algorithm 1 Self-training with noisy student-teacher

1. Train Teacher T_0 with labeled data L and classic augmentation.
 2. Train Student S_k by
 - (a) Apply aggressive augmentation on $L \cup D$
 $D = \text{Augment}(L \cup U)$
 - (b) Use teacher T_k to generate soft-label y^T by
 $y^{T_k}(i) = f^{T_k}(x(i))$ for $x(i) \in D$
 - (c) Student model trained using CE loss by
 $y^{S_k}(i) = f^{S_k}(x(i))$ for $x(i) \in D$
 $Loss = CE(y^{T_k}, y^{S_k})$
 3. Set $T_{k+1} = S_k$ and Repeat step 2
-

$$\text{Student_teacher_loss} = \alpha * \text{Loss}^E + \text{Loss}^D \quad (1)$$

$$\text{Loss}^D = \text{cross_entropy}(y_d^T, y_d^S) \quad (2)$$

$$\text{Loss}^E = \text{cross_entropy}(y_e^T, y_e^S) \quad (3)$$

$$y^T = [y_d^T, y_e^T] = f^T(\text{augment}(x)) \quad (4)$$

$$y^S = [y_d^S, y_e^S] = f^S(\text{augment}(x)) \quad (5)$$

The proposed method can be summarized by Algorithm 1. As shown, we can also have multiple iterations (indexed by k) of the second stage by using the student from previous iteration S_k as the teacher model T_{k+1} for next iteration. Losses for student-teacher training is computed by cross entropy (Eq. 1-5). Note that we compute two cross entropy's (for encoder and decoder labels), since our baseline model has both encoder and decoder as outputs [1, 2]. We combine two CE losses by weighted summation (Eq. 1).

2.2. Spec Augmentation

Data augmentation works by generating multiple variations of an original data example using various transforms [22, 23, 24], effectively multiplying number of training examples seen by the model. Classic approaches include adding reverberation or mixing with noise [22]. Recently proposed spectral augmentation method showed that one can boost ASR accuracy significantly by randomly masking blocks of frequency bins mostly [23, 24]. In this paper we explore the use of time and frequency masking, which is known to be most effective.

Spec augmentation is an aggressive data augmentation, in the aspect that it masks significant portion of input frequency bins or time frames in chunks. In ASR domain, such aggressive masking seems to help preventing over-fitting and facilitating the use of high level context. Also in ASR the target classes (phonemes or graphemes) are relatively well balanced in terms of prior. Meanwhile, KWS typically is a binary classification problem where positive pattern occupies only a small pattern space, while negative patterns span all the other spaces. One can easily transform a positive pattern to be a negative one by excessively masking chunks of frequency bins. In supervised learning with predetermined hard-label, those labels can simply be incorrect after some augmentation. To overcome such over-augmentation issue, we proposes to use spec augmentation with noisy student-teacher setup.

2.3. Model Architecture

For both the teacher (baseline) and the student model, we use the same two stage model architecture as in [2, 1]. The model consists of 7 simplified convolution layers and 3 projection layers, being organized into encoder and decoder sub-modules connected sequentially. Encoder module takes the input feature which is a 40-d vector of spectral frequency energies and generates encoder output of dimension N which learns to encode phoneme-like sound units. The decoder model takes the encoder output as input and generates binary output that predicts existence of a keyword in the input stream. For more details, please refer to [2, 1].

3. Experimental setup

3.1. Model setup

We implemented and compared the proposed model with baseline and some other variations as summarized in Table 1. In the table, Baseline_MP denotes the baseline model from our previous work [1] which uses supervised learning with labeled data (L) and max pooling loss. This model also becomes the first teacher model(T_0) in Algorithm 1. Model MP+sAug is a variant of the baseline model by simply adding spec augmentation on top of existing classic data augmentation. Model ST denotes a student-teacher trained model (using T_0) with classic augmentation and additional unlabeled data (U). Model ST+sAug is the same as model ST except that spectral augmentation is applied on top of classical augmentation. Model ST+sAug g2 denotes a second generation student-teacher trained model where the previous ST+sAug model was used as its teacher. ST+sAug NS(noisy student) is the same as ST+sAug model except that spec augmentation is applied only on student’s input similarly to [17, 25]. As shown in Table 1, all student-teacher trained models are trained using both labeled and unlabeled data, while supervised training models were trained using labeled data only.

Table 1: Summary of various models tested

Models	Loss	Training data
Baseline_MP	MaxPool CE Loss [1]	L
MP+sAug	MaxPool + sAug	L
ST	ST (Student-Teacher)	L + U
ST+sAug	ST + sAug	L + U
ST+sAug g2	ST + sAug 2nd gen.	L + U
ST+sAug NS	ST with noisy student + sAug	L + U

3.2. Training data set

We used both supervised (labeled) training data, and unsupervised (unlabeled) training data for experiments. Our supervised training data consists of 2.5 million anonymized utterances with the keywords (“Ok Google” or “Hey Google”). Supervised data is labeled by large ASR model similarly to [2, 1]. The unsupervised training data consists of 10 million anonymized utterances with the keywords and noises. The unsupervised data has relatively high noise level making it difficult for ASR model to generate reliable labels.

3.3. Evaluation data set

Evaluation is done with 6 positive data sets separate from training data, where each set represents a diverse environmental condition as summarized in Table 2. In the table, QLog, QLog low, QLog high are anonymous query logs from different time period and conditions. Especially QLog low has utterances that score relatively low confidences (difficult to detect), while Qlog high has utterances that score relatively high confidences (easy to detect).

Table 2: Summary of evaluation dataset

Eval set name	Description (# of utt’s)
Near Cl	Nearfield Clean (170K)
Near Cl Acc	Nearfield Clean Accented(16k)
Far Cl	Farfield Clean (1k)
Far Mus	Farfield w/ music noise (1k)
Far TV	Farfield w/ TV noise (1k)
QLog	Anonymous query logs (87k)
QLog low	Qlog’s with lower confidence (265k)
QLog high	Qlog’s with higher confidence (255k)

4. Results

Table 3: FR rate of models with various loss types at 0.1 FA/h

Models	Near Cl	Far Cl	Far Mus	Far TV
Baseline_MP	0.56%	1.83%	15.18%	27.94%
MP+sAug	1.17%	6.28%	32.24%	47.96%
ST	0.53%	1.57%	15.06%	27.58%
ST+sAug NS	0.74%	1.05%	17.65%	30.22%
ST+sAug	0.59%	0.96%	12.00%	24.58%
ST+sAug g2	0.53%	0.78%	13.65%	25.06%
Models	NearClAcc	QLog	QLog low	QLog high
Baseline_MP	1.64%	8.21%	11.07%	1.73%
MP+sAug	3.38%	18.74%	27.16%	8.48%
ST	1.58%	6.00%	8.63%	1.24%
ST+sAug NS	2.38%	7.38%	11.57%	2.14%
ST+sAug	1.71%	3.83%	7.28%	0.99%
ST+sAug g2	1.42%	3.12%	7.15%	0.91%

We evaluated 6 types of trained models on 8 evaluation sets and results are summarized by Tables and figures. Table 3 summarizes FR rates of the models at selected FA/h rate (0.1 FA per hour measured on 64K re-recorded TV noise set). Fig.2 and 3 shows the ROC (receiver operator characteristic) curves of tested models across different evaluation sets.

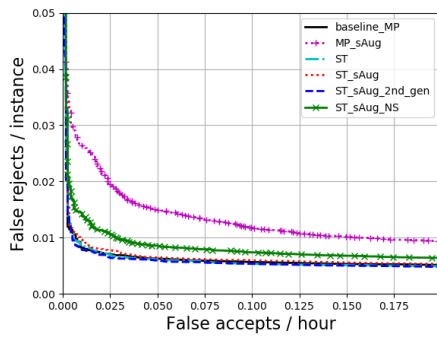
Results show that simply applying spec augmentation on top of classic augmentation with supervised training (MP+sAug model) doesn’t work well with our setup. This seems to be due to the risks of over augmentation (that transforms a positive labeled example into negative example with positive label).

The proposed model (ST+sAug and ST+sAug g2) showed significant improvements over baseline for difficult conditions such as Far-field and Query Logs. For example, Far-field Clean condition improved from 1.83% (baseline) to 0.78% (ST+sAug g2). Query Logs condition improved from 8.21% (baseline) to 3.12% (ST+sAug g2) (60% relative improvement). Clean Accented condition also improved from 1.64% to 1.42% (ST+sAug g2). ROC plots (Fig 2, 3) shows similar trends. Simple student-teacher training (ST model) shows small improvements over baseline (teacher model), assisted by extra unlabeled data. But the improvements are relatively minor compared to ST+sAug or ST+sAug g2.

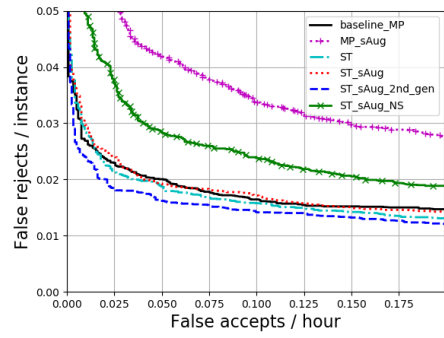
Direct adaptation of noisy student self-training method (Model ST+sAug NS) didn’t work well on our KWS problem setting. In this model, we apply spec augmentation on only the student model’s input, and not on the teacher model. Similarly to the MP+sAug model, the aggressive augmentation can transform positive example to a negative one, while the teacher model doesn’t see such changes resulting in incorrect soft-labels.

5. Conclusion

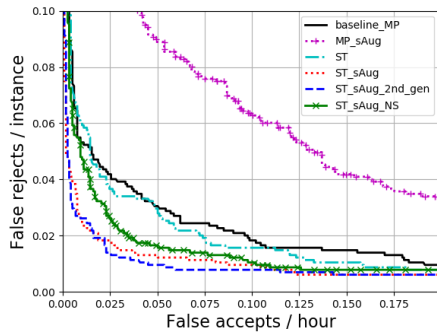
We presented self-training with noisy student-teacher for keyword spotting problem. The proposed approach enables the use of abundant unlabeled data and aggressive augmentation. Experimental results show that models with proposed approach significantly improves on evaluation set with difficult conditions. Experiments also show that applying aggressive augmentation directly in supervised learning approach doesn’t work well for keyword spotting problem, while semi-supervised training with noisy student-teacher can benefit from aggressive augmentation and unlabeled data. For future work, distillation from a larger teacher model or even multiple teachers can be explored for further improvements.



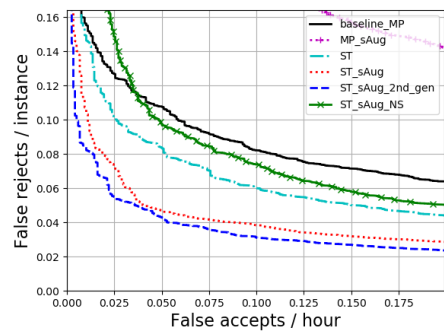
(a) Near-field Clean



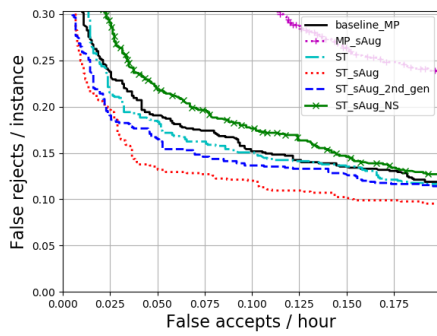
(e) Near-field Clean Accented



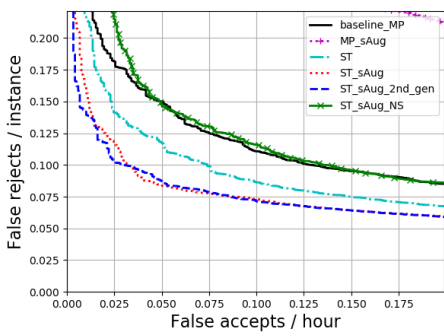
(b) Far-field Clean



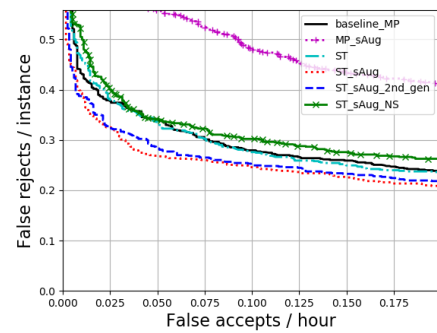
(f) Anonymous query Logs



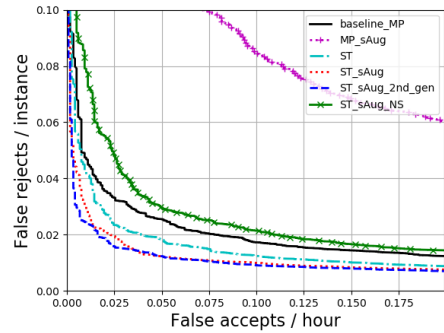
(c) Far-field with Music noise



(g) Anonymous query Logs with low confidence



(d) Far-field with TV noise



(h) Anonymous query logs with high confidence

Figure 2: ROC curves of models with various recipes and conditions

Figure 3: ROC curves of models with various recipes and conditions

6. References

- [1] H.-J. Park, P. Violette, and N. Subrahmanya, "Learning to detect keyword parts and whole by smoothed max pooling," *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7899–7903, 2020.
- [2] R. Alvarez and H. J. Park, "End-to-end Streaming Keyword Spotting," *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6336–6340, 2019.
- [3] Y. He, T. N. Sainath, R. Prabhavalkar, I. McGraw, R. Alvarez, D. Zhao, D. Rybach, A. Kannan, Y. Wu, R. Pang, Q. Liang, D. Bhatia, Y. Shangguan, B. Li, G. Pundak, K. C. Sim, T. Bagby, S.-Y. Chang, K. Rao, and A. Gruenstein, "Streaming end-to-end speech recognition for mobile devices," *ICASSP 2019*, pp. 6381–6385, 2018.
- [4] A. Gruenstein, R. Alvarez, C. Thornton, and M. Ghodrat, "A cascade architecture for keyword spotting on mobile devices," in *31st Conference on Neural Information Processing Systems (NIPS 2017)*, 2017. [Online]. Available: <https://arxiv.org/abs/1712.03603>
- [5] Y. He, R. Prabhavalkar, R. K., W. Li, A. Bakhtin, and I. McGraw, "Streaming small-footprint keyword spotting using sequence-to-sequence models," *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pp. 474–481, 2017.
- [6] T. Sainath and C. Parada, "Convolutional neural networks for small-footprint keyword spotting," in *Proceedings of Annual Conference of the International Speech Communication Association (Interspeech)*, 2015, pp. 1478–1482.
- [7] M. Wu, S. Panchapagesan, M. Sun, J. Gu, R. Thomas, S. Vitaladevuni, B. Hoffmeister, and A. Mandal, "Monophone-based background modeling for two-stage on-device wake word detection," in *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5494–5498.
- [8] J. Guo, K. Kumatani, M. Sun, M. Wu, A. Raju, N. Strom, and A. Mandal, "Time-delayed bottleneck highway networks using a DFT feature for keyword spotting," in *ICASSP - IEEE*, 2018, pp. 5489–5493.
- [9] M. Sun, D. Snyder, Y. Gao, V. Nagaraja, M. Rodehorst, S. Panchapagesan, N. Strom, S. Matsoukas, and S. Vitaladevuni, "Compressed time delay neural network for small-footprint keyword spotting," in *INTERSPEECH*, 2017.
- [10] S. Panchapagesan, M. Sun, A. Khare, S. Matsoukas, A. Mandal, B. Hoffmeister, and S. Vitaladevuni, "Multi-task learning and weighted cross-entropy for DNN-based keyword spotting," in *INTERSPEECH*, 2016.
- [11] S. Team, "Hey Siri: An On-device DNN-powered Voice Trigger for Apple's Personal Assistant," <https://machinelearning.apple.com/2017/10/01/hey-siri.html>, Apple Inc., 2017, accessed: 2018-10-06. [Online]. Available: <https://machinelearning.apple.com/2017/10/01/hey-siri.html>
- [12] A. van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint*, p. arXiv:1807.03748, 2018.
- [13] S. Löwe, P. O'Connor, and B. S. Veeling, "Putting an end to end-to-end: Gradient-isolated learning of representations," *NeurIPS*, 2019.
- [14] S. Schneider, A. Baevski, R. Collobert, and H. Auli, "wav2vec: Unsupervised pre-training for speech recognition," *Interspeech*, pp. 3465–3469, 2019.
- [15] O. H. Elibol, G. Keskin, and A. Thomas, "Semi-supervised and population based training for voice commands recognition," *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Brighton, United Kingdom*, pp. 6371–6375, 2019.
- [16] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *arXiv:2006.11477*, 2020.
- [17] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, "Self-training with noisy student improves imagenet classification," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10 687–10 698, 2020.
- [18] X. Zhai, A. Oliver, A. Kolesnikov, and L. Beyer, "S4: Self-supervised semi-supervised learning," *In Proceedings of the IEEE international conference on computer vision*, pp. 1476–1485, 2019.
- [19] Q. Sun, X. Li, Y. Liu, S. Zheng, T. Chua, and B. Schiele, "Learning to self-train for semi-supervised few-shot classification," *arXiv:1906.00562*, 2019.
- [20] A. R. Chowdhury, P. Chakrabarty, A. Singh, S. Jin, H. Jiang, L. Cao, and E. G. Learned-Miller, "Automatic adaptation of object detectors to new domains using self-training," *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 780–790, 2019.
- [21] J. He, J. Gu, J. Shen, and M. Ranzato, "Revisiting self-training for neural sequence generation," *arXiv:1909.13788*, 2019.
- [22] R. Prabhavalkar, R. Alvarez, C. Parada, P. Nakkiran, and T. Sainath, "Automatic gain control and multi-style training for robust small-footprint keyword spotting with deep neural networks," in *Proceedings of International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 4704–4708.
- [23] 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," *INTERSPEECH 2019*, pp. 2613–2617, 2019.
- [24] D. S. Park, Y. Zhang, C.-C. Chiu, Y. Chen, B. Li, W. Chan, Q. V. Le, and Y. Wu, "SpecAugment on large scale datasets," *ICASSP 2020*, pp. 6879–6883, 2020.
- [25] 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," *Submitted to Interspeech 2020*, pp. arxiv–2005.09 629, 2020.