Fine-tuning TitaNet-Large Model for Speaker Anonymization Attacker Systems

Candy Olivia Mawalim, Aulia Adila, Masashi Unoki

Graduate School of Advanced Science and Technology Japan Advanced Institute of Science and Technology, Ishikawa, Japan {candylim, adila, unoki}@jaist.ac.jp

Abstract-Speaker anonymization techniques are crucial for safeguarding user privacy in voice-based applications. However, these methods are susceptible to adversarial attacks that can compromise their effectiveness. This paper proposes attacker systems that leverage the power of fine-tuned TitaNet-Large and ECAPA-TDNN models to identify the original speaker from anonymized speech generated by various anonymization methods. Both pre-trained models are renowned for their stateof-the-art ability to extract robust speaker embeddings. Finetuning these models with anonymized speech enables them to identify underlying patterns in anonymized speech. We evaluated the proposed attacker systems against multiple anonymization techniques that performed effectively in a series of voice privacy challenges. Our experimental results underscore the effectiveness of the fine-tuned TitaNet-Large model in breaking through these anonymization methods, as indicated by the reduced equal error rate (EER). This highlights the importance of robust and adaptive anonymization strategies to counter such emerging semiinformed threats.

Index Terms—speaker anonymization, attacker model, semiinformed, TitaNet-Large, voice privacy

I. INTRODUCTION

Speaker anonymization techniques are essential for safeguarding user privacy in voice-based applications [1], [2]. However, the effectiveness of these techniques is constantly challenged by the emergence of sophisticated adversarial attacks [3]. The First Attacker Challenge has highlighted the growing threat to speaker anonymization systems [4]. This challenge focuses on developing attacker systems in the form of automatic speaker verification (ASV) systems, which are capable of identifying speakers even after anonymization.

This paper introduces an attack strategy that compares the power of fine-tuned TitaNet-Large (TitaNet-L) [5] and ECAPA-TDNN models [6]. Both models are renowned for their state-of-the-art performance in the speaker verification task. By fine-tuning these models on a diverse dataset of anonymized and original speech samples, we enable them to identify underlying patterns in anonymized speech, thereby accurately recognizing the original speaker.

Our work focuses on the semi-informed attack scenario, where the attacker has access to the speaker anonymization system and anonymized speech samples. We evaluated our proposed systems against the state-of-the-art anonymization techniques submitted to the VoicePrivacy 2024 Challenge [7]. Our experimental results underscore the effectiveness of the fine-tuned TitaNet-L models in breaking through these defenses, as indicated by the reduced equal error rate (EER). These results highlight the limitations of current anonymization techniques to counter emerging semi-informed threats.



Fig. 1. Fine-tuning phase (Top) and Inference phase (Bottom) of our proposed attacker systems based on fine-tuned TitaNet-L model.

II. PROPOSED ATTACKER SYSTEMS

Figure 1 shows our proposed system which is based on fine-tuning the TitaNet-Large (TitaNet-L) model [5]. TitaNet-L is the largest variant of the TitaNet architecture with 25.3 million parameters. TitaNet has an encoder-decoder structure. The ConvASR encoder acts as a high-level feature extractor, processing input audio from normalized mel spectrograms. It combines local features extracted through 1D depth-wise separable convolutions with global context information obtained via global average pooling in Squeeze-and-Excitation (SE) layers. This process improves the ability of the network to distinguish speaker-specific characteristics from input audio.

The speaker decoder incorporates an attentive statistics pooling layer to compute attention features across channel dimensions, creating time-independent, utterance-level speaker representations. This layer calculates weighted statistics, allowing the network to focus on the most relevant information for speaker verification. The features are then passed through linear layers to reduce dimensionality and map the resulting 192-dimensional features to the final number of classes, representing the different speakers in the training set.

The model generates fixed-length speaker embeddings, known as t-vectors, as its final output from the decoder. These embeddings encapsulate speaker-specific information and are used in speaker verification tasks, where the cosine similarity between t-vectors serves as the scoring backend.

To fine-tune the TitaNet-L model, we used a 9:1 trainingvalidation ratio. We only fine-tuned the final decoder layer of

 TABLE I

 EER (%) COMPARISON OF TITANET-LARGE AND ECAPA-TDNN ATTACKER MODELS ACROSS VARIOUS EXPERIMENTAL CONDITIONS

Method	Training data	Gender	Development								Test						Global avg.
			B3	B4	B5	T8-5	T10-2	T12-5	T25-1	B3	B4	B5	T8-5	T10-2	T12-5	T25-1	
ECAPA-TDNN	Original	F M Avg.	28.43 22.04 25.23	34.37 31.06 32.71	35.82 32.92 34.37	39.63 40.84 40.23	43.63 40.04 41.84	43.32 44.10 43.71	42.65 40.06 41.36	27.92 26.72 27.32	29.37 31.16 30.27	33.95 34.73 34.34	42.50 40.05 41.27	41.97 38.75 40.36	43.61 41.88 42.75	42.34 41.92 42.13	37.82 36.16 36.99
Fine-tuned ECAPA-TDNN	Anonymized (B5 and T12-5)	F M Avg.	41.51 38.08 39.80	42.80 42.16 42.48	37.97 36.34 37.15	46.57 47.06 46.82	39.48 34.66 37.07	39.18 36.15 37.66	42.57 41.59 42.08	38.47 37.98 38.23	41.04 42.16 41.60	35.37 35.80 35.58	46.97 45.94 46.45	34.61 31.91 33.26	35.73 36.27 36.00	41.13 40.82 40.97	40.24 39.07 39.65
TitaNet-L	VoxCeleb 1,2; SRE; LibriSpeech; RIR noise; Fisher; Switchboard	F M Avg.	43.18 37.60 40.39	46.59 44.38 45.48	48.27 49.35 48.81	45.20 42.24 43.72	33.65 30.43 32.04	49.29 49.40 49.34	48.46 47.20 47.83	42.34 40.32 41.33	44.53 43.43 43.98	46.70 47.66 47.18	45.44 44.77 45.10	27.14 29.40 28.27	46.21 50.78 48.49	48.87 46.55 47.71	43.99 43.11 43.55
Fine-tuned TitaNet-L*	Anonymized (B5)	F M Avg.	36.09 34.01 35.05	32.67 33.23 32.95	34.66 34.00 34.33	43.18 44.24 43.71	33.40 31.95 32.68	36.19 33.51 34.85	38.07 35.44 36.75	34.13 33.63 33.88	34.13 32.50 33.31	33.21 31.64 32.42	43.79 43.43 43.61	33.03 33.85 33.44	32.71 32.25 32.48	36.32 34.08 35.20	35.83 34.84 35.33
Fine-tuned TitaNet-L**	Anonymized (B5 and T12-5)	F M Avg.	34.66 33.56 34.11	34.23 32.92 33.57	31.96 27.02 29.49	41.73 41.93 41.83	33.10 27.64 30.37	31.25 28.42 29.83	36.54 36.37 36.46	32.15 30.72 31.44	32.12 33.80 32.96	26.24 28.06 27.15	41.42 40.92 41.17	26.66 28.29 27.47	26.82 29.16 27.99	33.75 36.07 34.91	33.04 32.49 32.77
Fine-tuned TitaNet-L	Anonymized (all)	F M Avg.	45.57 49.84 47.71	46.57 51.41 48.99	49.00 50.93 49.97	48.30 51.54 49.92	48.72 51.06 49.89	49.15 51.53 50.34	48.44 51.09 49.77	46.35 47.89 47.12	49.60 48.78 49.19	49.45 49.89 49.67	51.82 45.88 48.85	49.45 49.00 49.23	49.79 49.95 49.87	49.77 48.56 49.17	48.71 49.81 49.26

The abbreviations F, M, and Avg. under Gender correspond to female, male, and average, respectively. "Original" training data refers to the original, unanonymized training data used for building speaker anonymization models. In contrast, "Anonymized" training data consists of speech samples that have been processed by a specific speaker anonymization model. * and ** indicate the submitted system 1 and 2. Bold font represents the EER obtained by the most effective attacker system in each category.

the model for a maximum of 10 epochs with a batch size of 8. We also performed speed perturbation as a data augmentation technique to enhance the model's robustness to variations in speech speed. We used the angular softmax loss function to improve the discriminative power of embeddings by increasing the angular margin between different classes.

The fine-tuned model was optimized using the AdamW optimizer with a learning rate of 0.0001 and a weight decay of 0.0002. A cosine annealing learning rate scheduler with warmup was used to gradually decrease the learning rate over time, improving convergence and generalization.

III. EXPERIMENTS

We utilized fine-tuned TitaNet-L models¹ [5] as our primary attacker models. We fine-tuned our attacker models on the anonymized speech dataset. Specifically, we focused on anonymized speech generated by the B5 and T12-5 systems, which exhibited high EER scores and shared similar feature extraction and modification techniques.

We compared our proposed attacker model with a baseline ECAPA-TDNN model² [4] and its fine-tuned version trained on B5 and T12-5 anonymized speech. Our results, summarized in Table I, demonstrate that our fine-tuned TitaNet-L model outperforms the baseline model overall on anonymized B5 and T12-5 data. Fine-tuning an ASV system on anonymized speech generated by similar methods proves to be an effective strategy for creating a robust attacker model for those particular methods. This approach leads to a significant reduction in the EER, exceeding 10% in both attacks on B5 and T12-5 speaker anonymization methods. However, we also observed a reduction in performance on systems with significantly different anonymization techniques, such as B3, B4, and T8-5.

As a straightforward approach, we also fine-tuned the TitaNet-L model on all anonymized data. However, this approach was ineffective and in fact led to a decrease in performance as the model became confused by the diverse and potentially misleading information present in the anonymized data. This is evident in the increase in EER observed when

training on all anonymized data (last row of Table I) compared to the original pretrained model.

IV. CONCLUSION

This paper presents attacker systems based on fine-tuning the TitaNet-L model to compromise speaker anonymization systems. Our results demonstrate the effectiveness of this approach, particularly against similar techniques such as B5 and T12-5. However, its performance degrades against fundamentally different techniques. The ineffectiveness of directly fine-tuning TitaNet-L on all anonymized data underscores the importance of researching training data quality, as well as the sources and anonymization methods applied to datasets, for developing robust attacker models.

ACKNOWLEDGMENT

This work was supported by JSPS KAKENHI (20KK0233, 23K18491) and the SCAT Foundation.

REFERENCES

- [1] N. A. Tomashenko, X. Wang, E. Vincent, J. Patino, B. M. L. Srivastava, P. Noé, A. Nautsch, N. W. D. Evans, J. Yamagishi, B. O'Brien, A. Chanclu, J. Bonastre, M. Todisco, and M. Maouche, "The voiceprivacy 2020 challenge: Results and findings," *Comput. Speech Lang.*, vol. 74, p. 101362, 2022.
- [2] M. Panariello, N. A. Tomashenko, X. Wang, X. Miao, P. Champion, H. Nourtel, M. Todisco, N. W. D. Evans, E. Vincent, and J. Yamagishi, "The VoicePrivacy 2022 Challenge: Progress and Perspectives in Voice Anonymisation," *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 32, pp. 3477–3491, 2024.
- [3] T. Yang, L. S. Cang, M. Iqbal, and D. J. Almakhles, "Attack risk analysis in data anonymization in internet of things," *IEEE Trans. Comput. Soc. Syst.*, vol. 11, no. 4, pp. 4986–4993, 2024.
 [4] N. Tomashenko, X. Miao, E. Vincent, and J. Yamagishi, "The
- [4] N. Tomashenko, X. Miao, E. Vincent, and J. Yamagishi, "The First VoicePrivacy Attacker Challenge Evaluation Plan," 2024. [Online]. Available: https://www.voiceprivacychallenge.org/attacker/docs/Attacker_ Challenge_Eval_Plan_v2.2.pdf
- [5] N. R. Koluguri, T. Park, and B. Ginsburg, "TitaNet: Neural Model for Speaker Representation with 1D Depth-Wise Separable Convolutions and Global Context," in *Proc. ICASSP 2022.* IEEE, 2022, pp. 8102–8106.
 [6] B. Desplanques, J. Thienpondt, and K. Demuynck, "ECAPA-TDNN:
- [6] B. Desplanques, J. Thienpondt, and K. Demuynck, "ECAPA-TDNN: emphasized channel attention, propagation and aggregation in TDNN based speaker verification," in *Proc. of Interspeech 2020*. ISCA, 2020, pp. 3830–3834.
- [7] N. A. Tomashenko, X. Miao, P. Champion, S. Meyer, X. Wang, E. Vincent, M. Panariello, N. W. D. Evans, J. Yamagishi, and M. Todisco, "The VoicePrivacy 2024 Challenge Evaluation Plan," *CoRR*, vol. abs/2404.02677, 2024.

¹https://huggingface.co/nvidia/speakerverification_en_titanet_large

²https://huggingface.co/speechbrain/spkrec-ecapa-voxceleb