第一章 導論

1.1 研究動機

語言是人與人彼此交流最主要的橋樑,而人們互相溝通最自然的方式便是透過說話的語音(speech)達成。人類往往是自幼就牙牙學語開始說話,直到已屆學齡左右才開始學習認字與書寫。雖然在這個資訊爆炸的時代,人們已經習慣以文字呈現的語言作為獲取資訊的主要媒介,但不論如何,各種書寫系統其背後承載的語言必定有語音的形式作為對應。更何況世界上現存大約七千多種[1]語言中,絕大多數不見得存在成熟且普及的文字系統,卻無礙於這些語言被人們所熟悉和使用。因此,「語音」作為語言不可或缺的存在方式,了解它和研究它的價值自然不言而喻。

然而,相對於穩定、易於處理和保存的文字文本,語音訊號的變化萬千,蘊藏了大量從語者風格、表達內容到抑揚頓挫(韻律,prosody)等不同層次的訊息,使得對它的處理、研究相比之下複雜度與難度劇增。由於語音的這種特性,過往對於語言最有興趣的語言學家們,即便明白語音作為多數語言主體的事實,也不得不藉文字符號為依託進行探索。進入資訊化時代後,藉助電腦硬體等計算設備的幫助,從語料庫、計算語言學到自然語言處理等透過科技的力量發展語言處理技術的領域,頗長一段時間也是專注於文字的處理與分析。而嘗試結合訊號處理發展的語音技術領域,當時則是透過語言學家對語言的領域知識,例如從phoneme、morphology、語法等等用以刻劃人類語音和語言特性的概念,將之結合機器學習建立模型,開發技術以方便人們能以語音這種更靈活的媒介,更好的讓電腦、手機等科技工具可以更接近「直接溝通」的使用方式,便利人們的日常生活。

近年來,由於圖形處理器(graphics processing unit,GPU)等硬體平行運算技術的進步,深層學習(deep learning)快速崛起成為人工智慧的主流,有了此項機器學習的技術,模型的彈性能夠更好的萃取資料、更貼近的尋找資料背後的機制並進行預測,使得人們不再非得依賴大量費時費工的人類標注過程,進而使得利用大量語料庫發展語言技術,進一步推進語言科技發展成為可能。尤其在自監督學習(self-supervised learning)技術出現之後,深層學習模型可以依照人們給定的方向,更細緻的從大量未標注、相較容易取得的語音或文字的語料,找出其中的語音、語法及語義等等結構,形成帶有對人類語言有前所未見表現的基石模型(foundation model),是這個領域的一大里程碑。尤其在以處理文字為主體的自然語言處理領域,甚至出現了幾乎使人類真偽難辨的生成式模型,改變了人們生活的方方面面。

借鏡文字方面的成功經驗,語音處理領域的研究者們也開始嘗試將語言模型 (language model)的概念套用於變化莫測的語音訊號之上,原先人們藉助訊號處理知識一直使用的各種語音訊號特徵 (feature)也在自監督學習的架構之下,出現了許多模型從大量語音資料中得到的「語音表徵 (speech representation)」,作為精煉語音資訊的另外一種新選擇開始廣泛被採用。然而,相比於文字符號的穩定與單純,語音的複雜性使得它處理起來會需要更大量的資料和運算資源來擷取其中不同層次的細節,而且作為物理訊號,語音還必須處理掉環境中的雜訊等干擾。為了從紛亂的聲音中提取出最重要的訊息,向量量化 (vector quantization)的技巧因而經常被使用在語音 [2,3,4]或影像的領域中。爾後,[5]基於模仿人類學習語言的過程,藉助諸如 CPC ([6])、HuBERT ([7])、wav2vec 2.0 ([8])等自監督學習模型的幫助,引入向量量化的技術,提出了「無文字 (textless)」的學習架構,轉而以語音表徵量化後的「離散單元 (discrete unit)」作為操作對象,企圖以

單純大量的語音資料中訓練出一個不依賴文字的語言模型。此種學習架構的優勢在於在能保有利用大量未標注文字轉寫語音資料的同時,與連續表徵相比資訊的位元率(bit rate)利用更有效率、容易儲存、處理與傳輸,以及形式上更像文字的特性。自此之後有一系列(((補足 unit 相關的續作,例如英翻臺)))使用 discrete units 進行任務訓練的研究,一定程度的印證了 unit 捕捉語音內容的有效性。

Unit 固然在編碼語音之上有不錯的效果,因此 empirically 被領域內的不少學者直接視為「機器自己學習出來的文字」,可是究竟這項 claim 本身有多符合實際狀況,有學者便嘗試重新從語言學的概念汲取靈感,對其進行語音學層面的分析析析 (((放唐顥等等的語音學研究)))

[[[確認一下 acoustic piece 那邊的文線是不是也是如我那樣說的]]]

[[[確認研究動機要寫到哪裡]]]

1.2 研究方向

本論文旨在探討自監督學習模型的離散語音表徵與音位(phoneme)之間的關係。

1.3 主要貢獻

1.4 章節安排

本論文之章節安排如下:

參考文獻

- [1] D. M. Eberhard, G. F. Simons, and C. D. Fennig, *Ethnologue: Languages of the World*, 27th ed. Dallas, Texas: SIL International, 2024. [Online]. Available: https://www.ethnologue.com
- [2] J. Chorowski, R. J. Weiss, S. Bengio, and A. Van Den Oord, "Unsupervised speech representation learning using wavenet autoencoders," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 27, no. 12, pp. 2041–2053, 2019.
- [3] L.-W. Chen, S. Watanabe, and A. Rudnicky, "A vector quantized approach for text to speech synthesis on real-world spontaneous speech," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 11, 2023, pp. 12644–12652.
- [4] X. Zhao, Q. Zhu, J. Zhang, Y. Zhou, and P. Liu, "Speech enhancement with multi-granularity vector quantization," in 2023 Asia Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC). IEEE, 2023, pp. 1937–1942.
- [5] K. Lakhotia, E. Kharitonov, W.-N. Hsu, Y. Adi, A. Polyak, B. Bolte, T.-A. Nguyen, J. Copet, A. Baevski, A. Mohamed *et al.*, "On generative spoken language modeling from raw audio," *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 1336–1354, 2021.
- [6] A. van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," 2019.

- [7] W.-N. Hsu, Y.-H. H. Tsai, B. Bolte, R. Salakhutdinov, and A. Mohamed, "Hubert: How much can a bad teacher benefit asr pre-training?" in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6533–6537.
- [8] A. Baevski, Y. Zhou, A. Mohamed, and M. 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.