

# Leveraging Allophony in Self-Supervised Speech Models for Atypical Pronunciation Assessment Study Guide

---

音位變體於自監督語音模型之非典型發音評估研究：學習指南

## 一、核心概念複習

什麼是音位變體 (Allophony)？請解釋音位變體的定義，並舉出英語中 /t/ 的至少兩種不同的音位變體及其出現的語音環境。

非典型發音評估的重要性為何？說明準確區分典型與非典型發音（例如語音障礙者、非母語者）的意義與挑戰。

傳統基於音位分類器的發音評估方法有哪些局限？闡述這些方法在建模音位變體以及處理與訓練資料分布不同的異常發音時所遇到的問題。

自監督語音模型 (S3M) 在聲學建模方面展現了什麼優勢？簡述 S3M 的特性及其在語音處理領域的重要性。

什麼是 Goodness of Pronunciation (GoP)？解釋 GoP 的基本概念以及其在傳統發音評估中的作用。

## 二、MixGoP 方法理解

MixGoP 方法的核心思想是什麼？說明 MixGoP 如何克服傳統方法的局限，以提升非典型發音評估的效能。

MixGoP 如何利用高斯混合模型 (GMM) 建模音位變體？解釋 GMM 在 MixGoP 中扮演的角色，以及其如何表示一個音位的多個聲學實現。

MixGoP 如何結合自監督語音模型 (S3M) 的特徵？說明 MixGoP 如何利用 S3M 提取的語音特徵進行音位變體的建模和發音評估。

MixGoP 的發音異常程度評估方式是什麼？解釋 MixGoP 如何計算語音片段屬於特定音位的可能性，並以此判斷發音的典型性。

MixGoP 相較於傳統 GoP 方法的主要不同之處在哪裡？強調 MixGoP 在處理音位變體和異常發音假設上的創新。

## 三、實驗設計與結果

本研究使用了哪些資料集？這些資料集分別代表了什麼類型的語音？列舉研究中使用的五個資料集，並說明其在語音障礙和非母語語音方面的特性。

研究中比較了哪些不同的語音特徵？結果顯示 S3M 特徵的表現如何？列出研究中使用的 MFCC、Mel spectrogram、TDNN-F 和 S3M 等特徵，並總結 S3M 特徵相對於傳統特徵的優勢。

研究中與 MixGoP 進行比較的基線方法有哪些？MixGoP 在實驗結果中表現如何？列出基於音位分類器和基於 OOD 偵測的基線方法，並說明 MixGoP 在不同資料集上的效能表現。

研究如何分析 S3M 特徵捕捉音位變體的能力？使用了哪些方法和指標？說明研究中利用 UMAP 可視化和歸一化互信息 (NMI) 量化 S3M 特徵中音位變體資訊的方法。

研究針對 MixGoP 的樣本效率和子分布數量敏感性進行了哪些分析？結果如何？簡述研究中關於訓練樣本數量和 GMM 子分布數量對 MixGoP 效能影響的實驗結果。

#### 四、主要數學公式

傳統 GoP 分數：

$$GoP_p(s) = \log P_\theta(p|s)$$

解釋公式中 (p)、(s) 和  $P_\theta(p|s)$  的含義。

MixGoP 音位可能性：

$$P_\theta(s|p) = \sum_{c=1}^C \pi_{c|p} \mathcal{N}(Enc(s) | \mu_{c|p}, \Sigma_{c|p})$$

解釋公式中 (C)、 $\pi_{c|p}$ 、 $\mathcal{N}$ 、 $Enc(s)$ 、 $\mu_{c|p}$  和  $\Sigma_{c|p}$  的含義，並說明這個公式如何表示一個音位的多個子分布。

MixGoP 分數：

$$MixGoP_p(s) = \log P_\theta(s|p)$$

解釋這個分數如何用於評估發音的異常程度。

全句發音分數：

$$Pronunciation(x) = \frac{1}{N} \sum_{i=1}^N GoP_{p_i}(s_i)$$

解釋這個公式如何從單個音位的 GoP 分數得到整句話的發音評估。

#### 五、重要術語

Allophony (音位變體): 一個音位的不同語音實現，它們在特定的語音環境中出現，並且不會改變詞語的意義。

Atypical Pronunciation Assessment (非典型發音評估): 評估與典型發音存在偏差的語音，例如語音

障礙者或非母語者的發音。

Self-Supervised Speech Model (S3M) (自監督語音模型): 通過在大量未標註語音資料上進行預訓練，學習語音的通用表示的模型。

Gaussian Mixture Model (GMM) (高斯混合模型): 一種概率模型，用若干個高斯分布的加權和來表示一個複雜的分布。

Goodness of Pronunciation (GoP): 一種衡量語音片段與目標音位相符程度的評估分數。

Phoneme Classifier (音位分類器): 一種將語音片段分類到不同音位的模型。

Out-of-Distribution (OOD) (分布外): 指測試資料的分布與訓練資料的分布存在顯著差異。

Log Likelihood (對數似然): 在給定模型參數下，觀察到實際資料的可能性（似然性）的對數。

Mahalanobis Distance (馬氏距離): 一種考慮資料的協方差結構的多維空間中點之間的距離度量。

Kendall-tau Correlation Coefficient (肯德爾 tau 相關係數): 一種衡量兩個排序變數之間相關性的非參數指標。

測驗

(每題請以 2-3 句話回答)

請簡述音位變體現象對於自動發音評估帶來的挑戰。

傳統的基於單一分布假設的發音評估方法在處理非母語者的發音時可能遇到什麼困難？

自監督語音模型相較於傳統的聲學模型，在捕捉語音特徵方面的主要優勢是什麼？

MixGoP 方法如何利用高斯混合模型來解決傳統方法無法有效建模音位變體的問題？

在 MixGoP 方法中，S3M 特徵扮演了什麼關鍵角色？它們是如何被應用於發音評估的？

MixGoP 的發音異常程度評估是基於什麼原理？分數越高或越低代表什麼意義？

本研究的實驗結果顯示，MixGoP 在哪些類型的語音資料集上表現出了最顯著的優勢？

研究人員是如何驗證自監督語音模型特徵能夠有效捕捉音位變體資訊的？請簡述其方法。

根據研究結果，GMM 中子分布的數量對 MixGoP 的效能有何影響？是否存在一個最佳的子分布數量？

請簡述 MixGoP 方法相較於傳統 GoP 方法，在處理異常發音時的核心優勢。

測驗答案

音位變體是指同一個音位在不同語音環境下會有多種不同的發音方式，這使得將每個音位視為單一聲學目標的傳統方法難以準確評估發音是否符合語境。

非母語者的發音往往會受到母語語音系統的影響，產生許多在母語中存在的音位變體但在目標語言中不常見或被視為錯誤的發音，傳統基於單一分布假設的方法難以將這些發音視為可接受的變異。

自監督語音模型通過在大規模無標註資料上進行預訓練，能夠學習到更豐富、更底層的語音聲學特徵，這些特徵更能捕捉語音的細微變化，而無需人工標註的音位資訊。

MixGoP 方法為每個音位建立一個高斯混合模型，模型中的每個高斯分布都代表該音位的一種可能的音位變體，從而能夠捕捉一個音位內部的多個聲學子分布。

在 MixGoP 中，S3M 特徵作為高斯混合模型的輸入，為每個語音片段提供豐富的聲學表示。模型利用這些特徵來學習每個音位的不同音位變體的聲學特性。

MixGoP 基於測試語音片段在訓練好的 GMM 中屬於目標音位的對數似然分數來評估異常程度。分數越低表示該語音片段越不像模型所學習到的該音位的典型發音，因此被認為越異常。

本研究的實驗結果顯示，MixGoP 在語音障礙（dysarthria）資料集上相較於其他基線方法表現出了更為顯著的優勢，這可能與語音障礙者的發音具有較大的分布外特性有關。

研究人員首先使用 UMAP 等降維技術將 S3M 特徵在二維空間中可視化，觀察同一音位是否形成多個不同的聚類。然後，他們通過計算聚類索引與周圍語音環境的歸一化互信息來量化 S3M 特徵中包含的音位變體資訊量。

研究發現，增加 GMM 中子分布的數量通常可以帶來一定的效能提升，因為更細緻的子分布可以更好地捕捉音位變體的細微差別。然而，當子分布數量過多時，效能提升會趨於平緩甚至可能下降。

MixGoP 的核心優勢在於它顯式地建模了每個音位的多種可能的發音變體（音位變體），並且沒有像傳統方法那樣假設測試的異常發音與訓練的典型發音具有相同的分布，因此更能有效地識別和評估非典型發音。

## 申論題

探討自監督語音模型 (S3M) 在非典型發音評估領域的潛在應用與挑戰。除了本研究中使用的 MixGoP 方法外，你認為還可以如何利用 S3M 的特性來改進發音評估？

本研究強調了建模音位變體對於提升非典型發音評估準確性的重要性。請進一步闡述在設計自動發音評估系統時，考慮語音的上下文環境和音位變體的重要性。你認為應該如何有效地融入這些資訊？

MixGoP 方法使用了高斯混合模型 (GMM) 來建模音位變體。請比較 GMM 與其他可能的生成式模型（例如隱馬可夫模型 HMM、變分自編碼器 VAE）在音位變體建模方面的優勢與劣勢，並說明 MixGoP 選擇 GMM 的可能原因。

研究結果顯示，MixGoP 在語音障礙資料集上的表現優於非母語語音資料集。請分析可能導致這種差異的原因，並思考如何針對不同類型的非典型發音進一步優化 MixGoP 或設計更專用的評估方法。

本研究的局限性提到了語言泛化性、時間對齊品質以及評估指標的改進空間。請選擇其中一個局限性進行深入探討，提出可能的解決方案或未來研究方向，以克服這些限制並提升研究的價值。