

## 論文分析報告

### 主要目標

本論文旨在提升自動發音評估（特別是異常發音評估）的準確性與魯棒性，核心目標為**有效建模音位的異音現象（allophony）**，並藉由結合自監督語音模型（Self-Supervised Speech Models, S3M）與高斯混合模型（Gaussian Mixture Models, GMM），提出一種新方法「MixGoP」，以更精確地區分典型與非典型（如語音障礙、非母語者）發音<sup>[1]</sup>。

### 相關研究

- **傳統方法**：早期自動語音評估多以隱馬可夫模型（HMM）或深度神經網路（DNN）為基礎，將每個音位視為單一聲學分布，忽略了異音的多樣性<sup>[1]</sup>。
- **GoP（Goodness of Pronunciation）**：Witt & Young (2000) 提出GoP分數，衡量語音片段與目標音位的相符程度。後續研究將GoP應用於DNN、S3M等新型模型<sup>[1]</sup>。
- **自監督語音模型（S3M）**：如wav2vec 2.0、XLS-R、WavLM等，證明能捕捉豐富的語音特徵，並在各種語音任務上表現優異，但對於異音的建模分析仍有限<sup>[1]</sup>。
- **異常語音偵測**：Yeo et al. (2023a)等將異常發音評估視為分布外（OOD）偵測問題，嘗試跳脫傳統分類器的限制<sup>[1]</sup>。

## 研究方法與流程

### 1. 問題診斷

- 傳統GoP方法假設每個音位只有單一分布（單峰），無法反映異音的多樣性。
- 傳統分類器依賴softmax，假設測試語音與訓練語音分布一致，對於異常語音（如語音障礙、非母語）不適用<sup>[1]</sup>。

### 2. MixGoP方法設計

- **核心理念**：用高斯混合模型（GMM）為每個音位建模，允許每個音位包含多個子分布（對應不同異音），並用S3M萃取的特徵作為輸入。
- **流程簡述**：
  - 用S3M（如WavLM、XLS-R）萃取語音特徵。
  - 依據已標註的音位起訖，將語音切分為音位片段。
  - 對每個音位片段進行特徵聚類（k-means初始化），再以GMM建模，捕捉異音分布。
  - 用訓練好的GMM計算測試語音片段對應音位的對數似然分數（MixGoP分數）。
  - 以分數評估發音的典型性，並與人工標註的流利度/錯誤率做相關性分析<sup>[1]</sup>。

### 3. 實驗設計

- **資料集**：涵蓋三個語音障礙（dysarthria）資料集與兩個非母語者資料集，訓練用健康/母語語者語音，測試用異常語音<sup>[1]</sup>。
- **特徵比較**：比較MFCC、Mel spectrogram、TDNN-F、S3M等多種特徵。
- **基線方法**：包括傳統GoP（GMM-GoP、NN-GoP、DNN-GoP、MaxLogit-GoP）、kNN、單類SVM等<sup>[1]</sup>。
- **評估指標**：Kendall-tau相關係數（發音分數與人工標註的流利度/錯誤率之相關）<sup>[1]</sup>。

### 4. 分析與驗證

- 驗證S3M特徵是否能有效捕捉異音（以聚類與互資訊量化）。
- 探討模型在不同資料集、特徵、子分布數量下的表現。
- 分析MixGoP在樣本數不足時的表現穩定性<sup>[1]</sup>。

## 主要數學公式說明（含隱喻與變數解釋）

### 1. 傳統GoP分數

$$\text{GoP}_p(\mathbf{s}) = \log P_{\theta}(p \mid \mathbf{s})$$

- **隱喻**：想像你是語音老師，聽到一段學生的發音（ $\mathbf{s}$ ），你根據經驗判斷這是不是你預期的音位（ $p$ ），這個分數就是你信心的「對數」。
- **變數說明**：
  - $p$ ：目標音位
  - $\mathbf{s}$ ：語音片段
  - $P_{\theta}(p \mid \mathbf{s})$ ：給定語音片段屬於音位 $p$ 的機率，由參數 $\theta$ 決定<sup>[1]</sup>。

### 2. 傳統分類器的logits計算

$$f_{\theta}(\mathbf{s}) = \mathbf{W} \cdot \text{Enc}(\mathbf{s})$$

- **隱喻**：每個音位都有一個「代表座標」（ $\mathbf{W}$ 中的每一列），你把語音片段轉成特徵（ $\text{Enc}(\mathbf{s})$ ），再看它跟各音位代表座標的接近程度。
- **變數說明**：
  - $\mathbf{W}$ ：權重矩陣，每一列代表一個音位的中心
  - $\text{Enc}(\mathbf{s})$ ：語音特徵向量
  - $f_{\theta}(\mathbf{s})$ ：每個音位的打分<sup>[1]</sup>。

### 3. MixGoP的GMM建模公式

$$P_{\theta}(\mathbf{s} \mid p) = \sum_{c=1}^C \pi_p^c \mathcal{N}(\text{Enc}(\mathbf{s}) \mid \boldsymbol{\mu}_p^c, \boldsymbol{\Sigma}_p^c)$$

- **隱喻**：想像每個音位像一個水果籃，裡面有不同品種的蘋果（異音），每個品種有不同的形狀（均值 $\boldsymbol{\mu}_p^c$ ）、大小（協方差 $\boldsymbol{\Sigma}_p^c$ ），出現機率（ $\pi_p^c$ ）。你拿到一顆蘋果（語音片段），要判斷它屬於這個籃子的機率，就是看它跟每個品種的相似度加總。
- **變數說明**：
  - $C$ ：每個音位的子分布（異音）數量
  - $\pi_p^c$ ：第 $c$ 個子分布的權重（機率），所有子分布加總為1
  - $\mathcal{N}(\cdot \mid \boldsymbol{\mu}_p^c, \boldsymbol{\Sigma}_p^c)$ ：多維高斯分布，均值 $\boldsymbol{\mu}_p^c$ 、協方差 $\boldsymbol{\Sigma}_p^c$
  - $\text{Enc}(\mathbf{s})$ ：語音特徵
  - $\theta = \{\boldsymbol{\mu}_p^c, \boldsymbol{\Sigma}_p^c, \pi_p^c\}$ ：所有GMM參數<sup>[1]</sup>。

### 4. MixGoP分數

$$\text{MixGoP}_p(s) = \log P_{\theta}(\mathbf{s} \mid p)$$

- **隱喻**：你根據蘋果（語音片段）屬於這個蘋果籃（音位）的機率，給它一個對數分數，分數越低代表越不像這個籃子的蘋果（異常發音）<sup>[1]</sup>。

### 5. Mahalanobis距離（GMM中的核心）

$$-\frac{1}{2}(\text{Enc}(\mathbf{s}) - \boldsymbol{\mu}_p^c)^T (\boldsymbol{\Sigma}_p^c)^{-1} (\text{Enc}(\mathbf{s}) - \boldsymbol{\mu}_p^c)$$

- **隱喻**：這就像你量測蘋果到每個品種中心的距離，不只是單純的直線距離，還考慮每個方向的變異（有的品種蘋果形狀比較分散，有的很集中），這種距離能更準確反映「像不像」<sup>[1]</sup>。
- **變數說明**：
  - $\text{Enc}(\mathbf{s})$ ：語音特徵
  - $\boldsymbol{\mu}_p^c$ ：第 $c$ 個子分布的中心
  - $\boldsymbol{\Sigma}_p^c$ ：第 $c$ 個子分布的協方差矩陣<sup>[1]</sup>。

### 6. 全句發音分數

$$\text{Pronunciation}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \text{GoP}_p(s)$$

- **隱喻**：把一句話的每個音位分數平均，得到整句發音的整體表現<sup>[1]</sup>。
- **變數說明**：
  - $\mathbf{x}$ ：整句語音
  - $N$ ：音位數

- $\text{GoP}_p(s)$ ：每個音位的發音分數<sup>[1]</sup>。

## 總結

- **MixGoP方法**：以GMM建模音位異音，結合S3M特徵，不僅克服傳統單峰假設，也不再假設異常語音與訓練分布一致，對異常發音評估更為精準。
- **實驗結果**：在五個語音障礙與非母語資料集上，MixGoP在四個資料集取得最先進表現，證明其有效性<sup>[1]</sup>。
- **S3M特徵分析**：S3M比傳統特徵（MFCC、Mel spectrogram）更能捕捉異音資訊，且異音建模能力與下游評估表現呈正相關<sup>[1]</sup>。
- **數學公式易懂化**：MixGoP的核心在於用多個「子分布」來描述一個音位的多種發音實現，並以機率和距離的方式量化發音的「正常」或「異常」程度。

**結論**：本論文提出的MixGoP方法，結合S3M與GMM，有效捕捉音位異音變異，顯著提升異常發音自動評估的準確性與泛化能力，並為未來語音評估與語音模型特徵分析提供了新視角<sup>[1]</sup>。



1. [https://ppl-ai-file-upload.s3.amazonaws.com/web/direct-files/31472917/d431f090-fb56-4a89-b8e6-3744029796d1/Leveraging-Allophony-in-Self-Supervised-Speech-Models-for-Atypical-Pronunciation-Assessment\\_v2.pdf](https://ppl-ai-file-upload.s3.amazonaws.com/web/direct-files/31472917/d431f090-fb56-4a89-b8e6-3744029796d1/Leveraging-Allophony-in-Self-Supervised-Speech-Models-for-Atypical-Pronunciation-Assessment_v2.pdf)