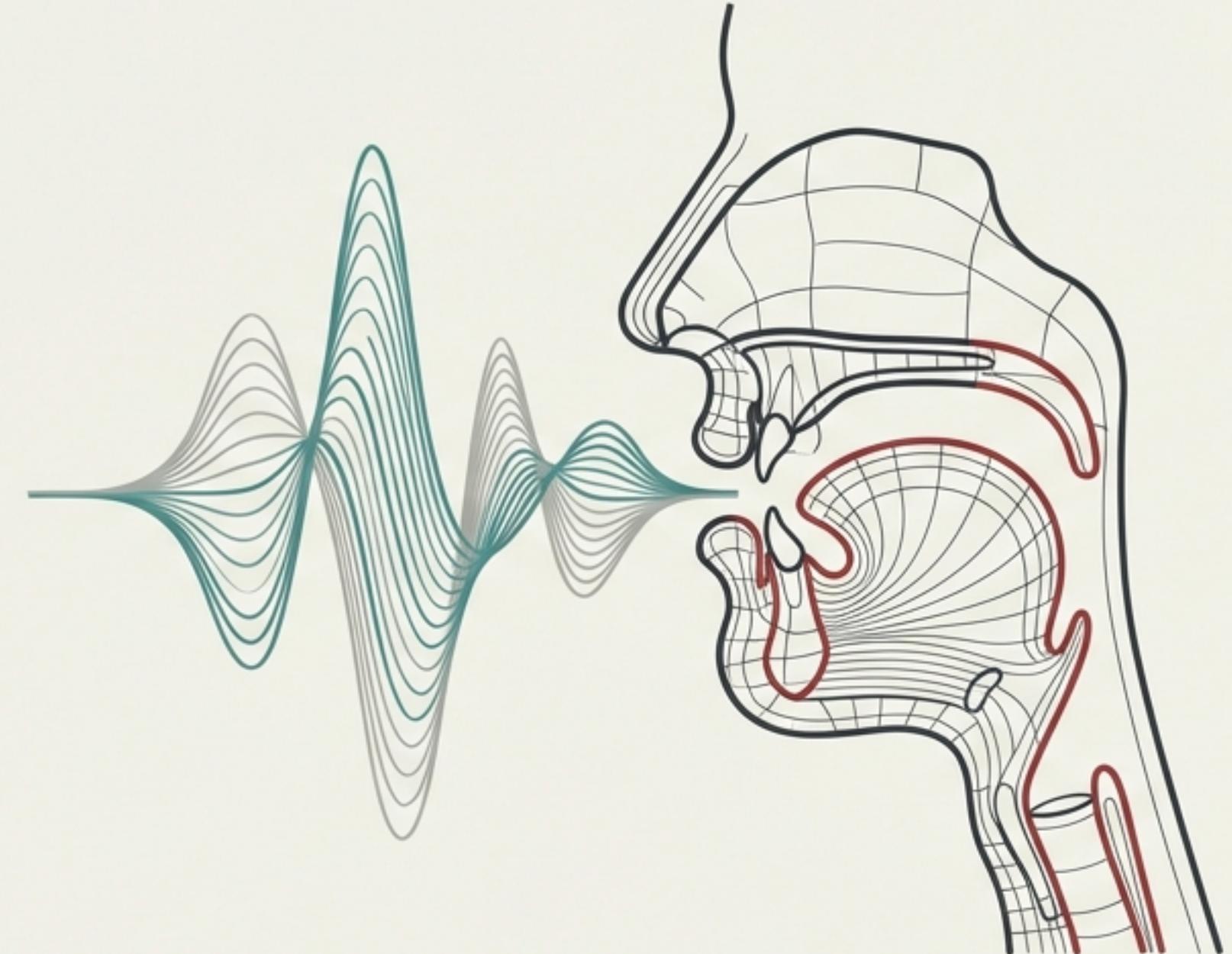


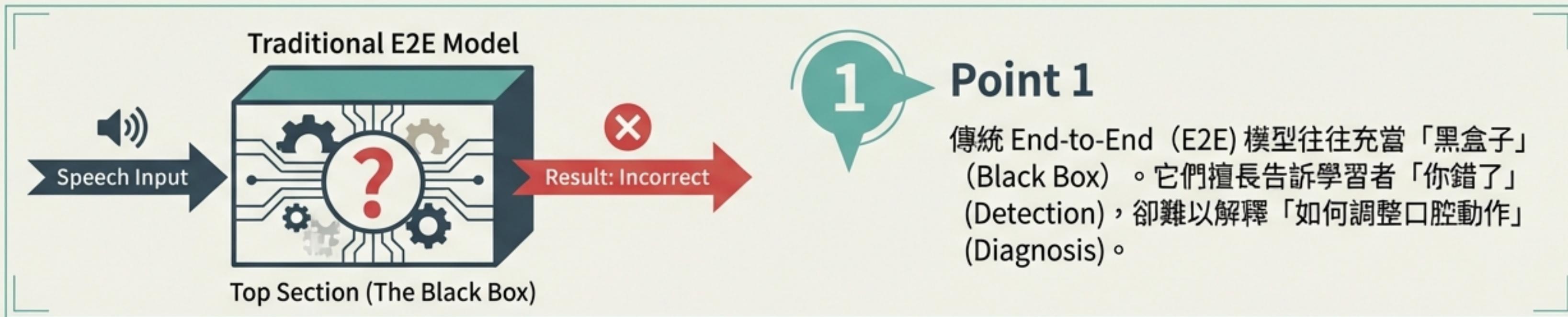
# 構音特徵增強型誤讀 檢測與診斷模型： 多維誤差分析

探索 End-to-End 架構下「檢測精度」  
與「診斷深度」的權衡與挑戰



Based on the study by Xing Wei, Catia Cucchiariini, Roeland van Hout, Helmer Strik  
Centre for Language and Speech Technology, Radboud University

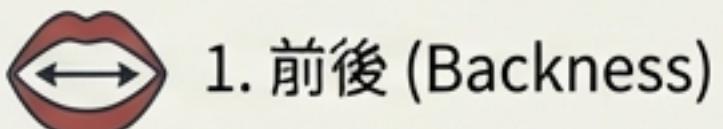
# 傳統模型的侷限：從「聽覺感知」到「構音指導」的缺口



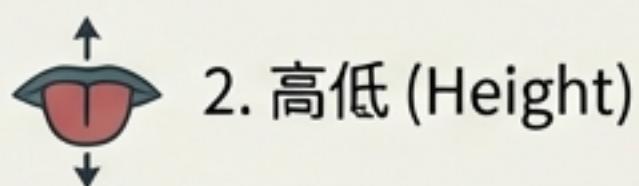
# 解決方案：引入構音特徵 (Articulatory Features, AFs)

**Concept:** AFs 將語音描述為聲道的物理配置，具有語言學基礎，對說話者變異具有魯棒性。

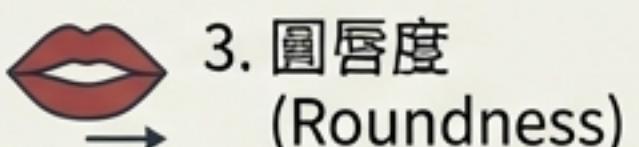
## Vowels (元音)



1. 前後 (Backness)

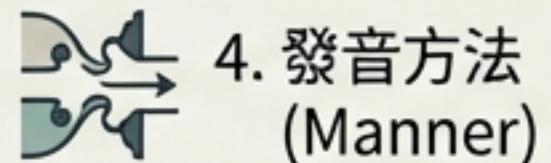


2. 高低 (Height)

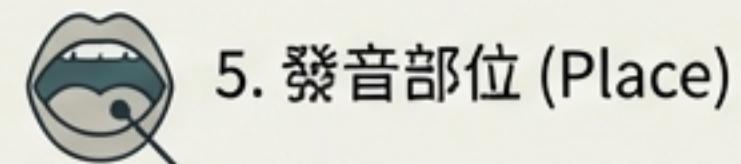


3. 圓唇度 (Roundness)

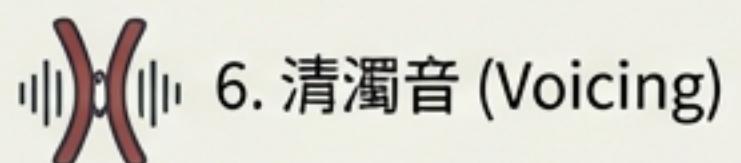
## Consonants (輔音)



4. 發音方法 (Manner)



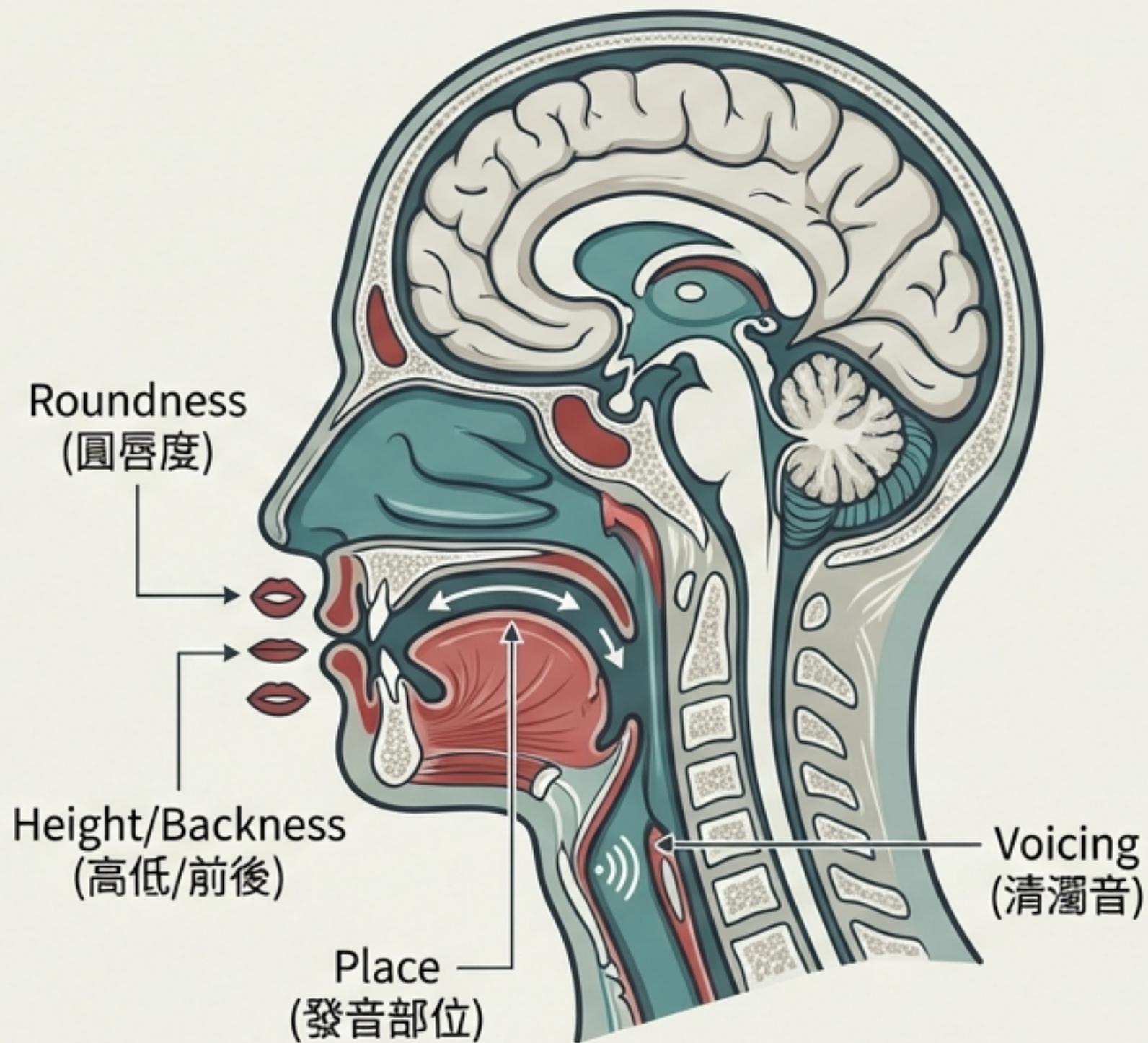
5. 發音部位 (Place)



6. 清濁音 (Voicing)

## Strategic Value Box

AFs 能夠提供細粒度的反饋 (Fine-grained feedback)，例如舌位高低或嘴唇圓展，這正是 L2 語音訓練的核心需求。

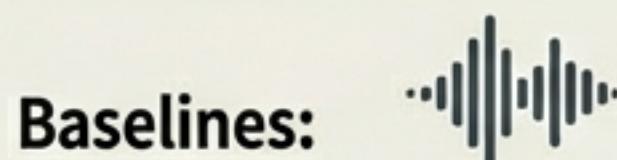


# 實驗設計：雙路徑建模與對照分析

## Method 1: Customized E2E (M1)



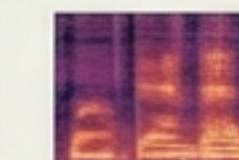
## Method 2: Wav2Vec 2.0 (M2)



Baselines:



Raw Speech (RS)



FBank Pitch (FP)



Fine-tuned only (FT)

Proposed:

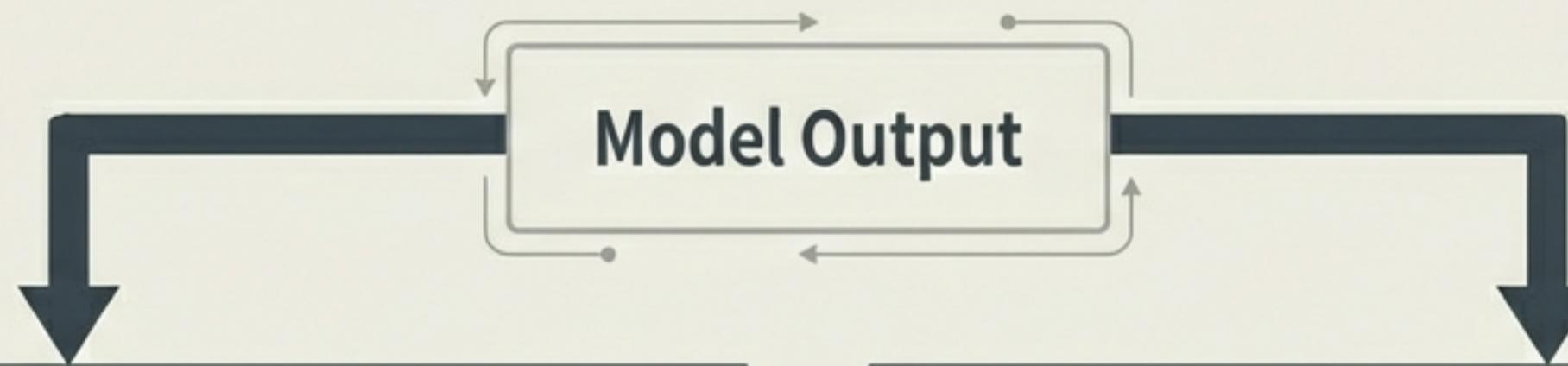


M1 (AF-Enhanced Conformer)



M2 (AF-Enhanced Wav2Vec)

# 關鍵變量：兩種輸出框架的定義 (PHN vs. ART)



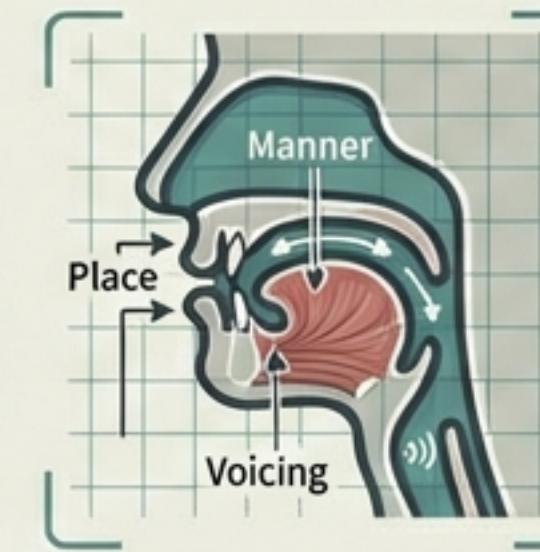
## Framework 1: Phoneme-based (PHN)

{ /t/ /d/ /æ/ }

模型直接預測音素。

- Strength: 標準化，適用於整體準確率評估。

## Framework 2: Articulatory-based (ART)



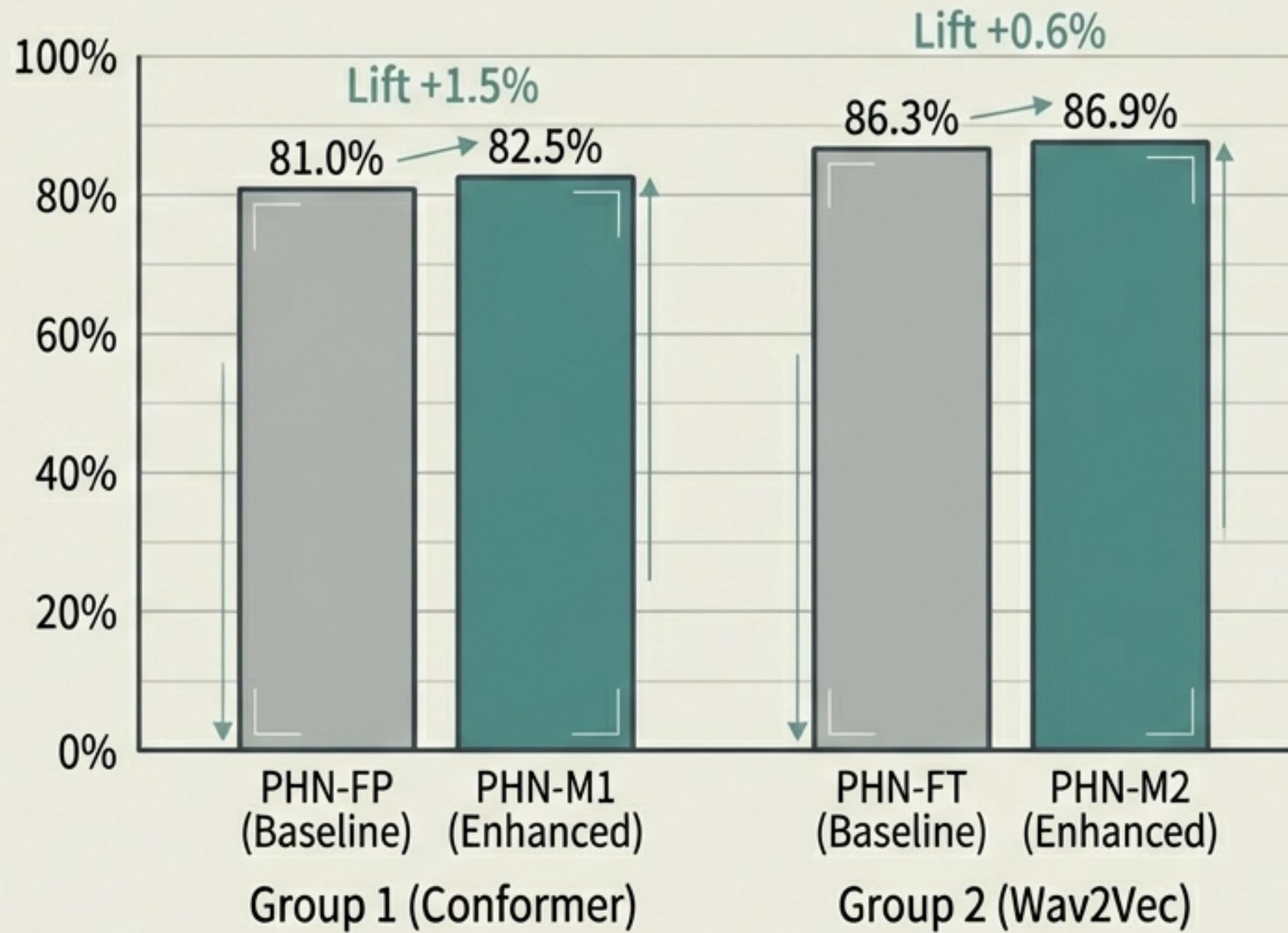
模型預測構音標籤  
(Articulatory Labels)。

- Strength: 能捕捉更細微的發音動作差異 (子音段分析)。

**Goal:** 比較哪種輸出形式能為 L2 學習者提供更精準的診斷。

# 整體效能驗證：構音特徵 (AF) 帶來的顯著提升

## Detection Accuracy (DA) Comparison



## Key Insights

**Key Finding:** AF 增強型模型 (M1, M2) 在所有指標上均優於其對應的基線模型 (FP, FT)。

### Data Detail Box



統計顯著性 (Statistical Significance):

- PHN-M1 vs PHN-FP 顯示出顯著差異 ( $p < 0.01$ )。
- 無論架構為何，引入構音知識都能幫助 AI 更準確地識別誤讀。

# 核心權衡：檢測精度 (Detection) vs. 診斷深度 (Diagnosis)

PHN (Generalist)



ART (Specialist)



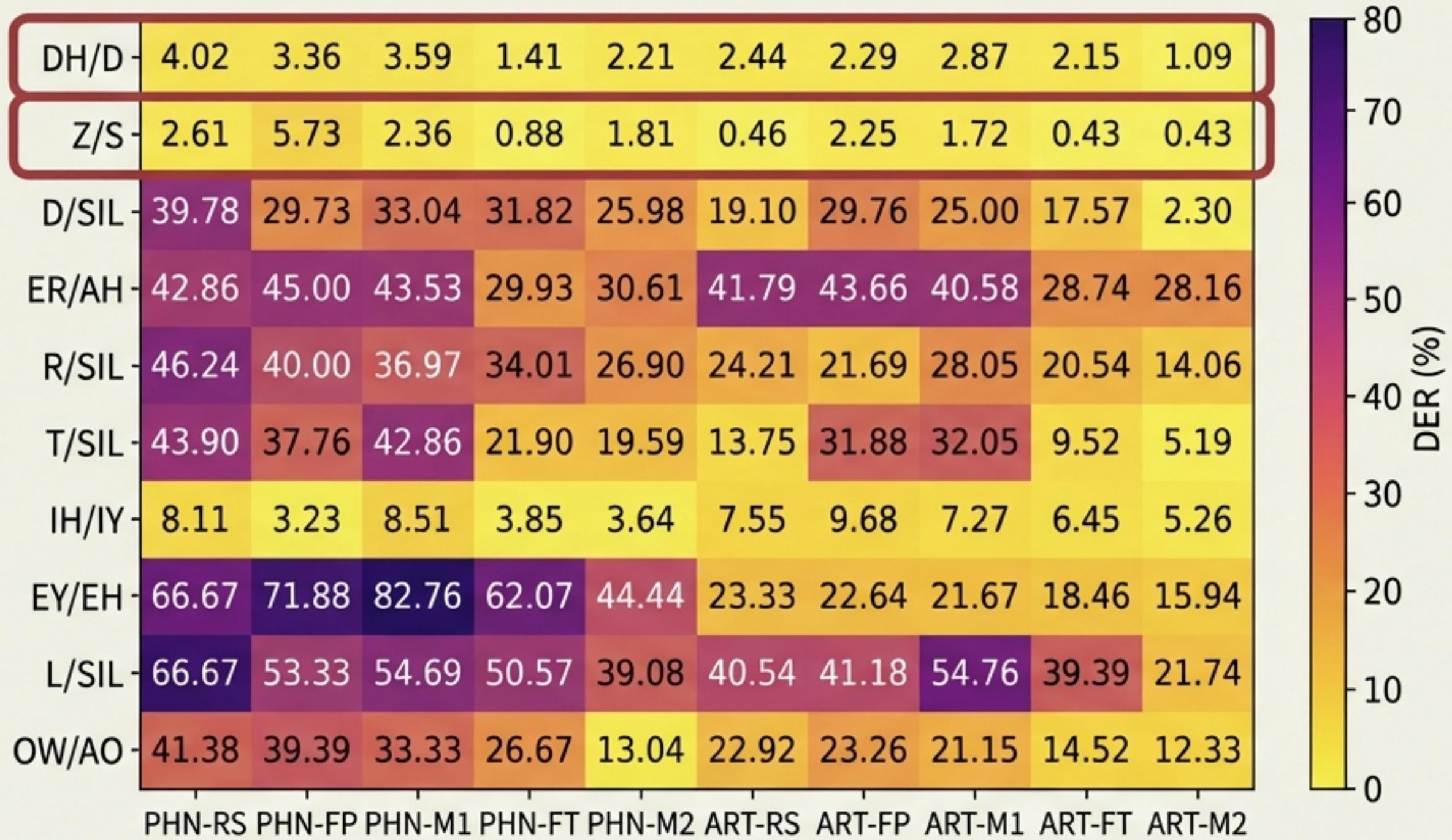
- Higher Detection Accuracy (DA)
- 能快速篩查問題 (Yes/No)

- Lower Diagnosis Error Rate (DER)
- 能精確診斷病因 (What/How)

這種權衡暗示了為了獲得更深刻的教學反饋 (ART)，我們可能需要接受些微的檢測精度損失 (PHN)。

PHN = 全科醫生 | ART = 專科醫生

# 診斷力可視化：高頻錯誤的精準識別



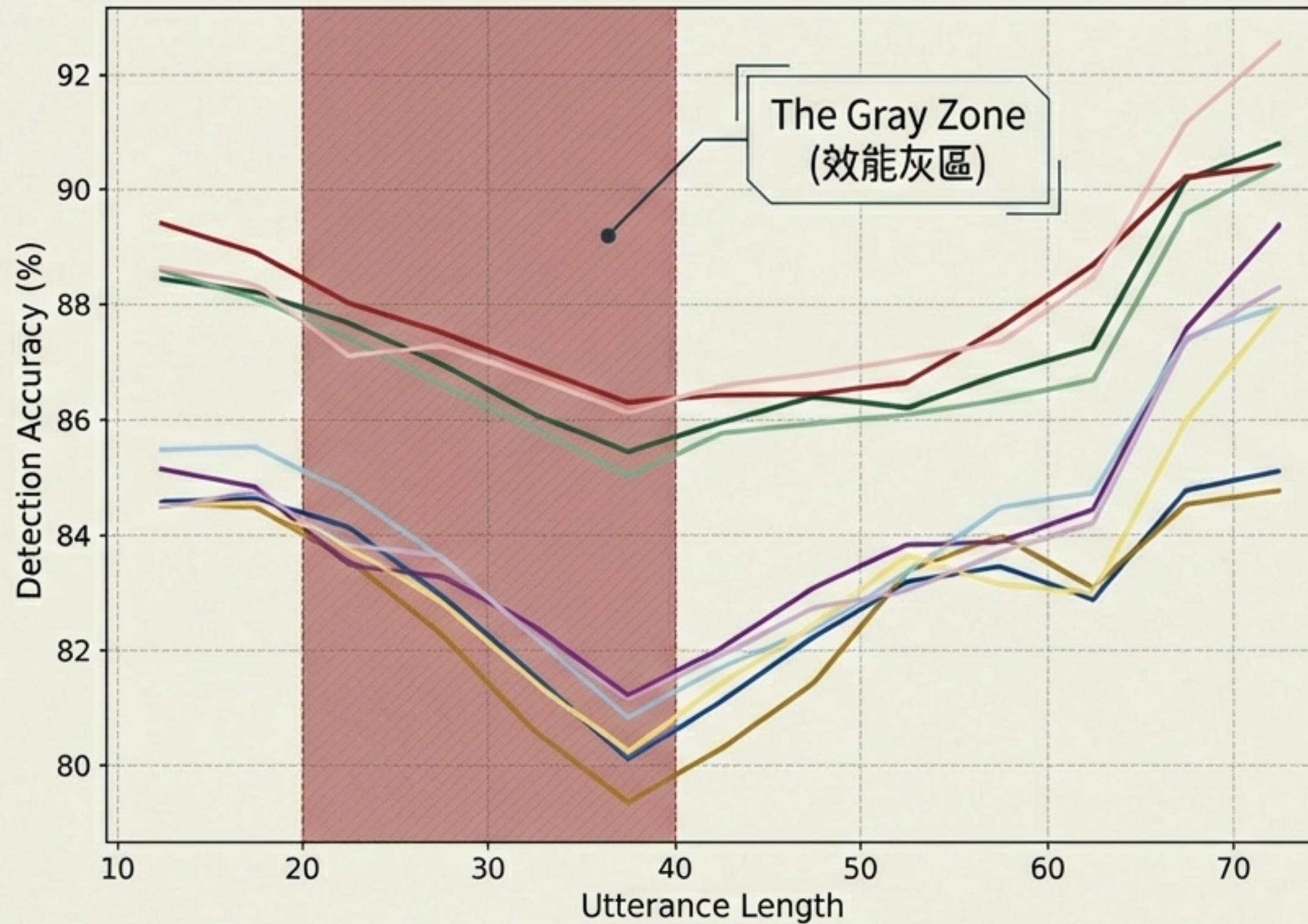
## Case Study: DH/D

將 /ð/ (咬舌音) 誤讀為 /d/ (塞音) 是最常見的 L2 錯誤。

**結果：**ART 模型（尤其是 ART-M2）在此類錯誤上的診斷錯誤率 (DER) 顯著降低。

**意義：**對於具體的構音動作錯誤，ART 架構提供了最具可操作性的反饋。

# 隱藏的陷阱：中等長度語句的「效能灰區」



## 為何發生？

- **Short (<20 Labels) :** 結構簡單，易於匹配。
- **Long (>40 Labels) :** 具備足夠的上下文冗餘 (Context Redundancy)。
- **Medium (20-40 Labels) :** 既缺乏足夠的上下文，又比短句複雜，導致錯誤接受率 (FAR) 與錯誤拒絕率 (FRR) 同時飆升。

# 說話者變異性：錯誤頻率 ≠ 可檢測性

Speaker THV

High Error Rate (27.00%)

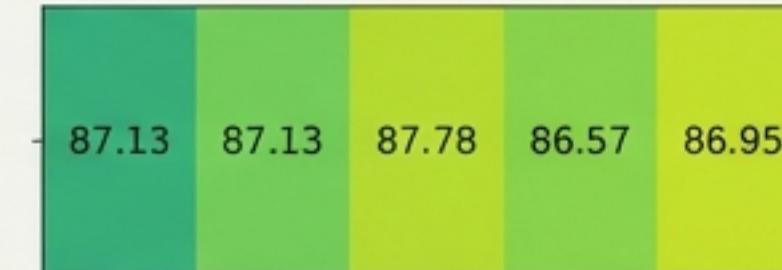
**Detection Accuracy: LOW**



Speaker TLV

High Error Rate (24.73%)

**Detection Accuracy: HIGH**



- Insight: 儘管兩者錯誤頻率相當，TLV 的錯誤更具「聲學顯著性」(Acoustically Distinct)。

結論：單純統計錯誤數量不足以預測系統效能，聲學特徵的清晰度 (Distinctiveness) 至關重要。

# 錯誤類型分佈與模型偏好

## Top Mispronunciations

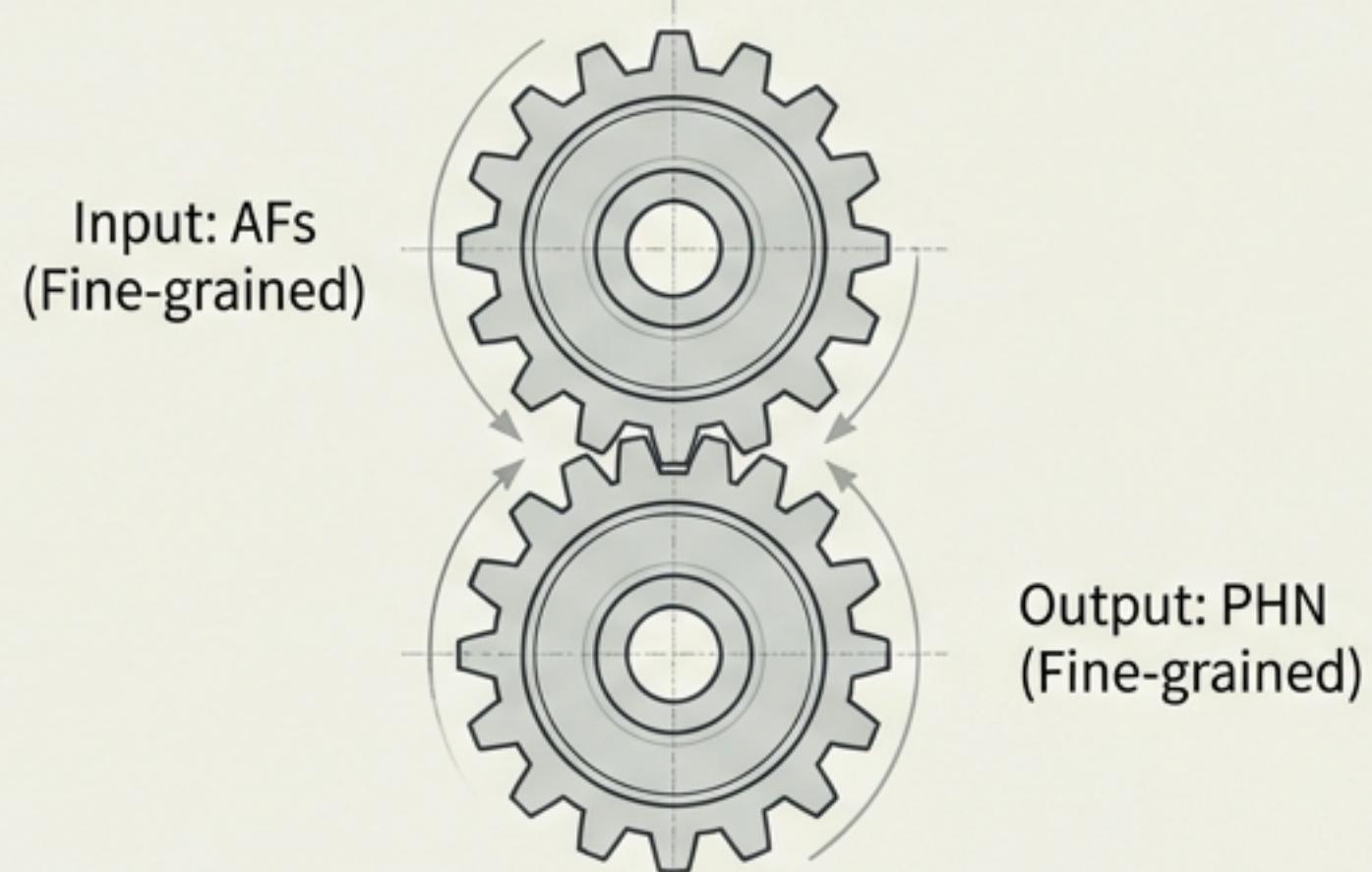
Error	Count	Type	Status
DH/D (Substitution)	348	Fricative → Stop	✓ High Detection
Z/S (Voicing Error)	296	Voiced → Unvoiced	✓ High Detection
L/SIL (Deletion)	126	Missing Sound	⚠ Persistent Challenge

## Analysis

- **ART 模型優勢區**：特徵替換類錯誤 (如清濁音、發音方法改變)。
- **ART 模型劣勢區**：缺失/插入類錯誤 (如 L/SIL)，即使在最佳模型中仍保持高 DER。

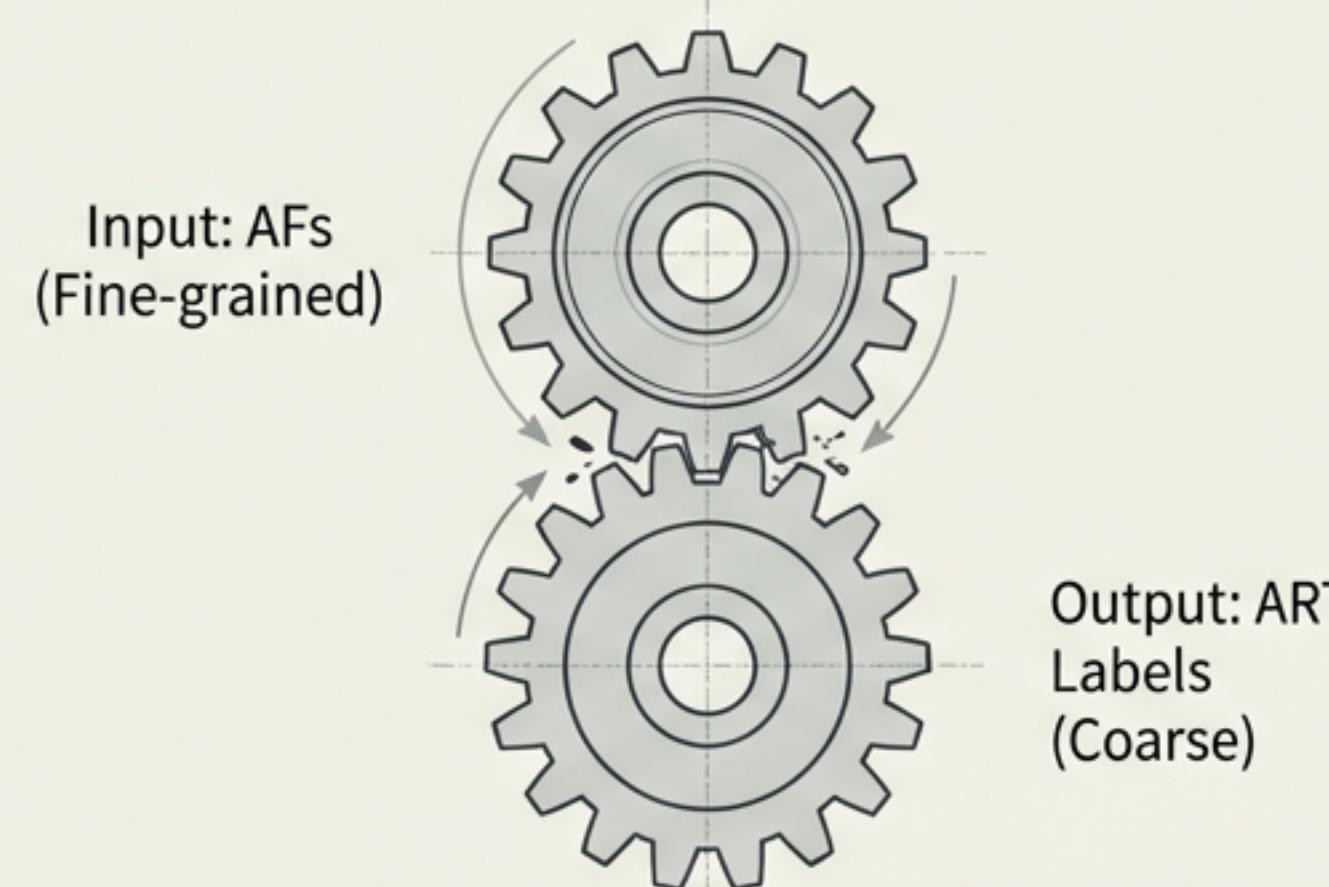
# 深度解析：粒度匹配 (Granularity Matching) 假說

The Match



High Accuracy (M1/M2 PHN)

The Mismatch

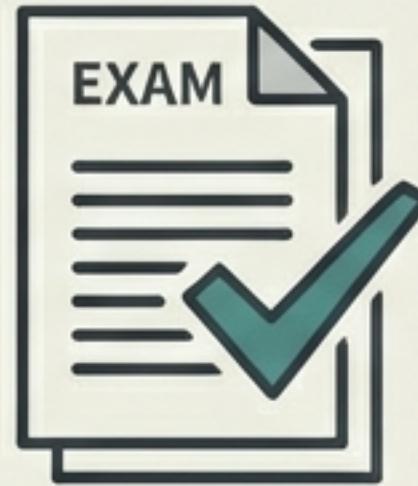


Information Loss

## 為何 PHN 準確率有時更高？

當輸入特徵的細緻度與輸出標籤的細緻度高度一致時，AF 的效益最大化。ART 標籤相對粗糙，可能導致信息丟失。

# 針對教育科技 (EdTech) 的實務建議



## 評分與分級 (Scoring)

建議使用 **PHN** 輸出架構。  
優先考慮整體檢測準確率  
(DA)，以確保考試公平性。



## 輔導與糾音 (Tutoring)

建議使用 **ART** 輸出架構。  
優先考慮診斷精度 (Low DER)，提供具體的口腔動作指導。



## 數據增強 (Data)

針對**中等長度語句 (20-40 labels)**進行專項訓練數據擴充，填補上下文建模短板。

# 綜合效能比較總表

Feature	PHN Framework	ART Framework
<b>Detection (DA)</b>	Winner (Slight Advantage) 	Good
<b>Diagnosis (DER)</b>	Average	 Winner (Significant for specific errors)
<b>Best Model</b>	Wav2Vec 2.0 (M2)	Wav2Vec 2.0 (M2)
<b>Primary Use Case</b>	Testing / Scoring	Learning / Feedback
<b>Weakness</b> 	Less actionable feedback	Coarser granularity

Robustness Issue: All models struggle with 20-40 label length; Speaker variance remains high.

# 結論：從「糾錯」邁向「教學」



整合構音特徵 (AF) 使系統不僅能聽出「什麼」 (What)，更能理解「如何」 (How)

## Future Directions List

1. 擴大數據集多樣性 (L1 Backgrounds).
2. 驗證「可檢測性」假說.
3. 解決中等長度語句的上下文建模問題.

## Final Thought

未來的 MDD 系統將不再是單純的評判者，而是具備生理學知識的虛擬教練。