

114-1

期末報告

注射式喉成型術前聲音檢測資料 可否作為聲帶治療成效的預測工具 -應用人工智慧方法建模與分析

報
合作醫師

醫師

報告日期：2025. 01. 12

OUTLINE

01 研究動機

02 研究問題

03 研究目的

04 原始資料敘述

05 資料前處理

- 音檔切割
- 前處理後的資料敘述

06 模型設計

07 模型評估與解釋

08 結果與討論

附錄 01 補充文字的連結

附錄 02 資料的連結

01 研究動機

- 閉合功能健全的聲帶才能使我們發出好的聲音狀況，而閉鎖功能不全的病人(如聲帶麻痺及聲帶萎縮)就會造成發聲的困擾，更會使發出聲音吃力更甚至沙啞的情形。
- 玻尿酸注射性喉成形手術是可以大幅改善閉合功能不全的問題。

-Journal of Voice. 2024. Jennifer A. Silver

- 越來越多相關文章提出運用AI方式學習讓機器判讀人類聲音，對其評估出一個標準的模式，也讓疾病初步的判別更加方便簡單。

02 研究問題

- 內視鏡檢查我們僅知道病人需要進行手術幫忙，但病人總會問說“效果怎麼樣?”、“我聲音會進步多少?”、“我需要使用多少玻尿酸?”。
- ☒ 但我們僅給出”會進步”這樣的結果

03 研究目的

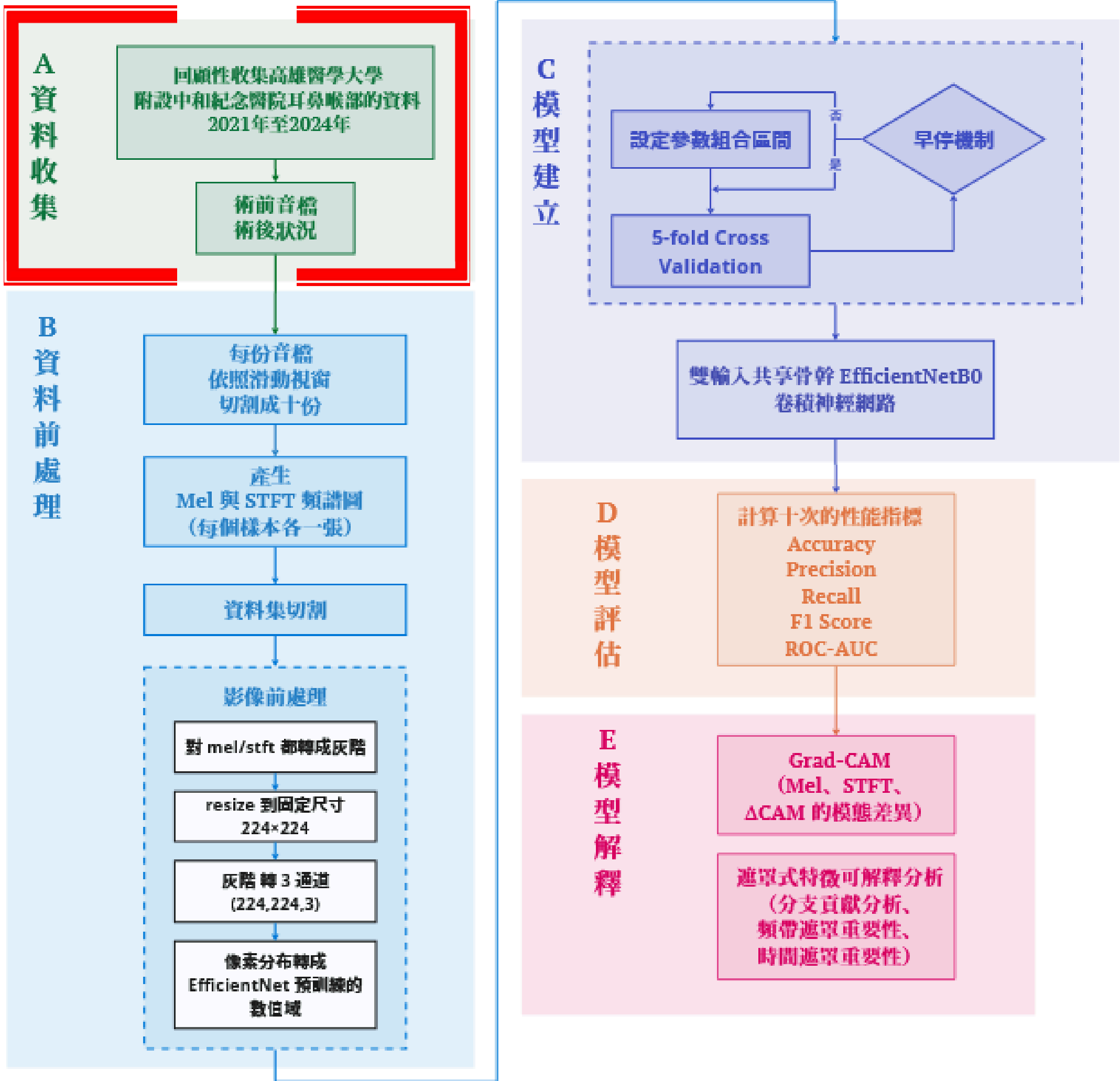
機器學習（於上學期完成）

- 透過人工智慧方式建模由手術前的聲音狀況，預測病人手術後發聲的情形

深度學習

- 以術前 a 母音音檔轉梅爾頻譜圖與頻譜圖訓練 CNN，預測病人預後（好/不好）

資料收集



04 原始資料敘述

- 資料集來自 [REDACTED]
- 每位病人在手術前 (S) 與手術後 (A) 各錄製 4 個母音音檔 (A, I, U, E)
- 原始每人有 4 (母音) \times 2 (術前/後) = 8 筆錄音
- 原始錄音筆數 = 43 位 \times 4 母音 \times 2 (術前術後) = 344 筆

for AI assessment-20250524T061350Z-1-001 > for AI assessment > 16422748well-2022 > pre				
<input type="checkbox"/>	名稱	修改日期	類型	大小
	16422748S1.txt	2022/4/15 下午 03:40	文字文件	2 KB
	16422748S1.wav	2022/4/15 下午 02:43	WAV 檔案	76 KB
	16422748S2.txt	2022/4/15 下午 02:44	文字文件	2 KB
	16422748S2.wav	2022/4/15 下午 02:44	WAV 檔案	112 KB
	16422748S3.txt	2022/4/15 下午 02:44	文字文件	2 KB
	16422748S3.wav	2022/4/15 下午 02:44	WAV 檔案	55 KB
	16422748S4.txt	2022/4/15 下午 02:45	文字文件	2 KB
	16422748S4.wav	2022/4/15 下午 02:45	WAV 檔案	68 KB

04 原始資料敘述

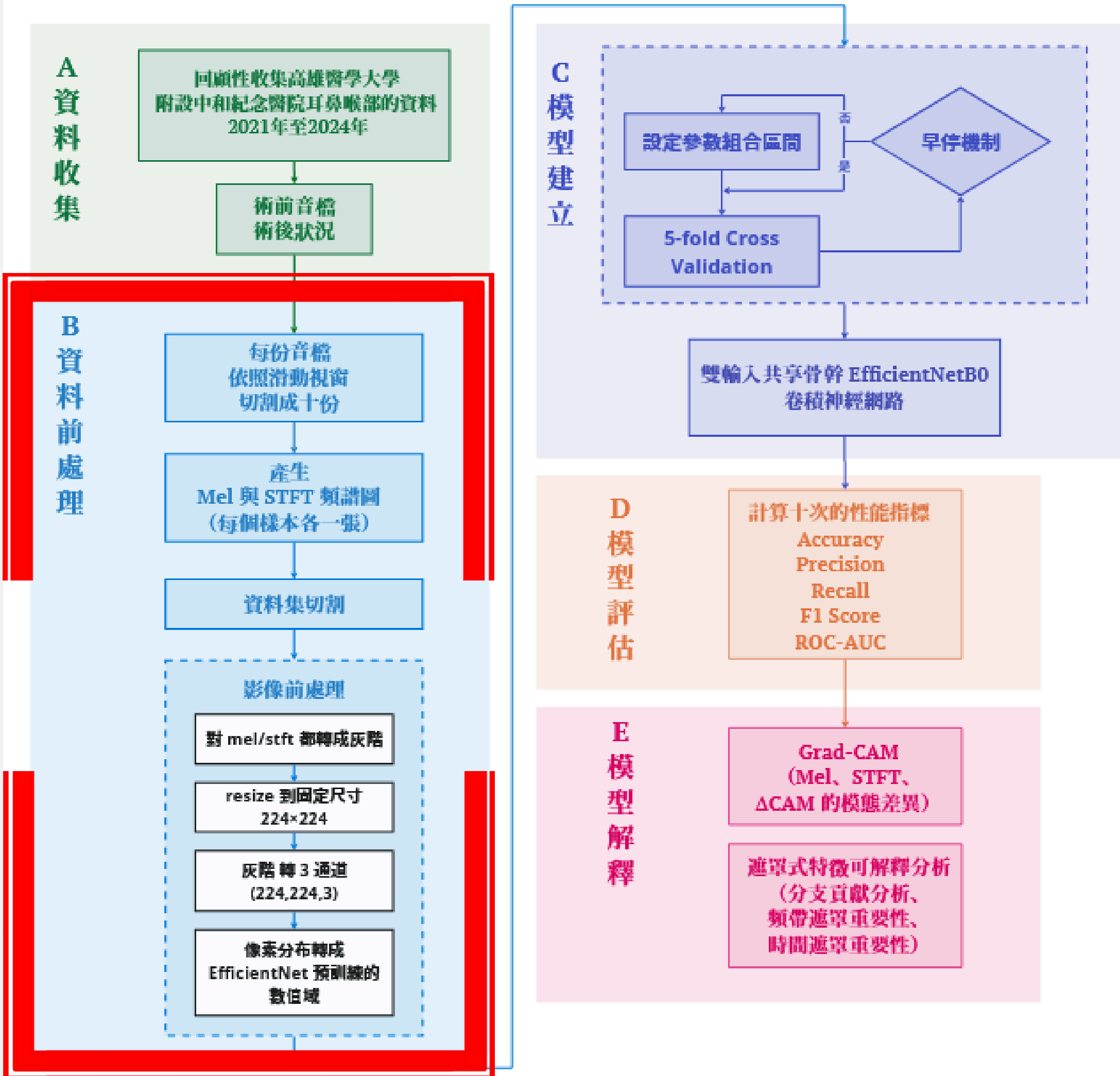
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	type	Date	院區	姓名	ID	病歷號	Sex	Age	總次數	次數	側性	Cause	Injection Tech	Sub/Tran	Amount R	Inj point R	Amount L	Inj point L	
2	1									1	1	L	VP	CTm	tran			0.8	1
3	1									1	1	R	VP	CTm	tran	1(0.3+0.7)	2		
4	1									1	1	R	VP	CTm	tran	1(0.4+0.6)	2		
5	1									1	1	L	VP	CTm	tran			0.8(0.3+0.5)	2
6	1									1	1	L	VP	CTm	tran			0.5	2
7	1									1	1	L	VP	CTm	tran			1	1
8	1									1	1	L	VP	CTm	tran			0.7(0.2+0.5)	2
9	1									1	1	L	VP	CTm	tran			1	2
10	1									1	1	L	VP	CTm	tran			0.3(0.1+0.2)	2
11	1									1	1	left	palsy	CT	Sub			1	2
12	1									1	1	right	palsy	CT	Tran	1	2		
13	1									1	1	left	palsy	CT	Tran			1	2
14	1									1	1	left	palsy	CT	Tran			1	2
15	1									1	1	left	palsy	CT	Tran			1	3
16	1									1	1	right	palsy	CT	Tran	0.9	2		
17	1									1	1	left	palsy	CT	Tran			1	2
18	1									1	1	left	palsy	CT	Tran			0.8	1
19	1									1	1	right	palsy	CT	Tran	1	2		
20	1									1	1	right	palsy	CT	Tran	1	2		
21	1									1	1	left	palsy	CT	Tran			1	3
22	1									1	1	left	palsy	CT	Tran			1	2
23	1									1	1	left	palsy	CT	Tran			1	8

04 原始資料敘述

- 音檔總數：86
- 最短音檔長度：0.33 秒
- 最長音檔長度：14.65 秒
- 音檔長度不一

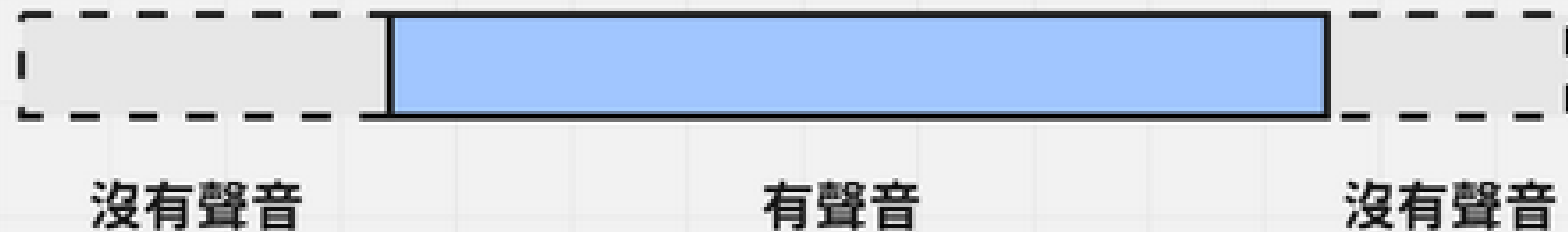
資料集	最小值(ms)	最大值(ms)	平均(ms)
palsy 手術前 S	323	14653	4347.22
palsy 手術後 A	702	10750	6196.13
palsy total	323	14653	5271.67
atro 手術前 S	3405	10246	8799.07
atro 手術後 A	3262	10256	8805.4
atro total	3262	10256	8802.24
total	323	14653	7036.95

資料前處理



05 資料前處理：音檔切割

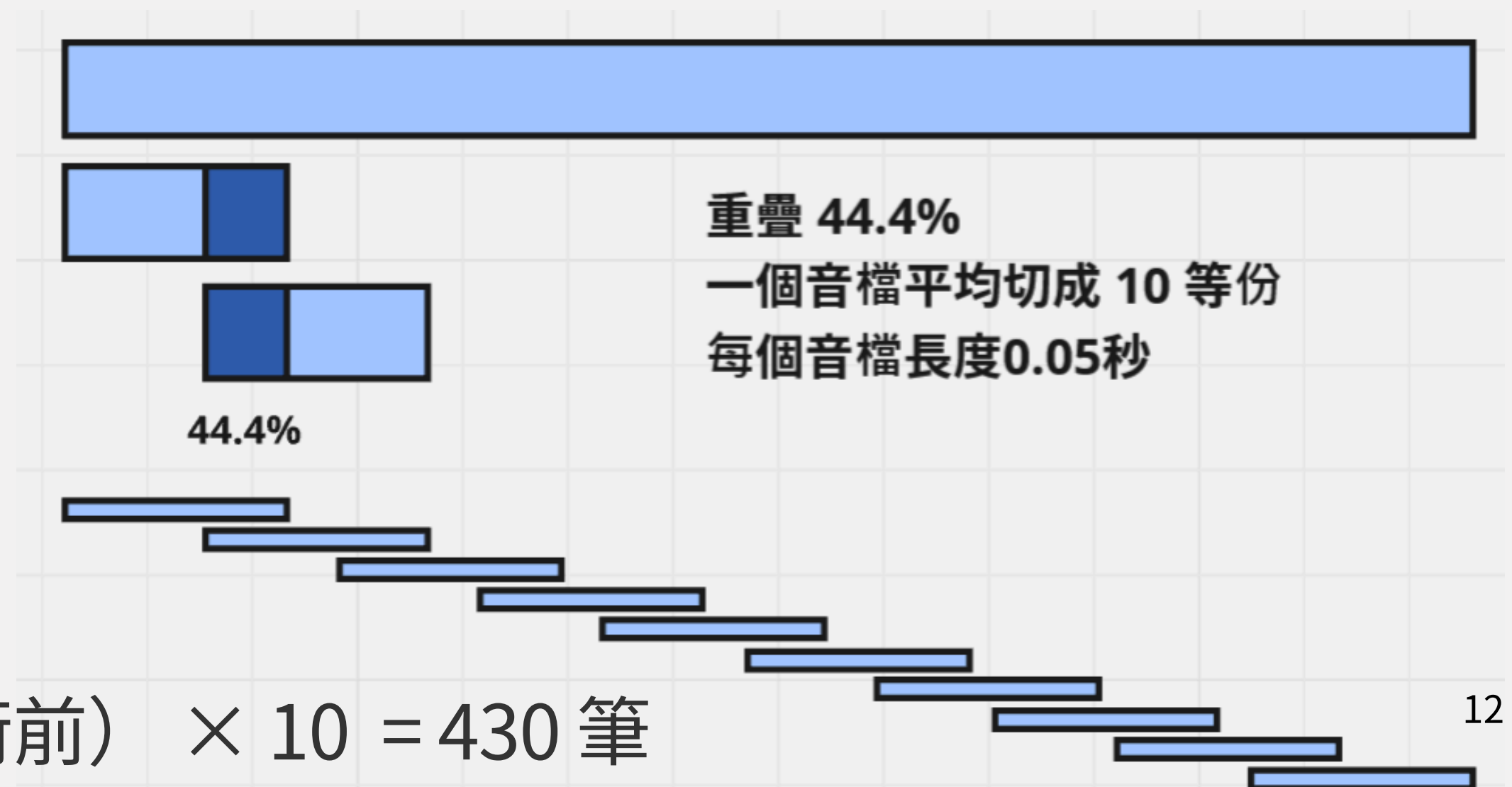
- 問題：資料筆數過少
- 解決方式：
 - 只選用母音A的術前音檔進行模型建構
 - 設定能量門檻（ENERGY THRESHOLD）
 - 將每分音檔切成十等份進行聲學特徵分析
- 原始錄音筆數 = 43 位 \times 1（術前） = 43 筆



秒數	檔名	切空白前	切空白後
min (ms)	20_S1	323	320
max (ms)	22_S1	14650	14390

05 資料前處理：音檔切割

- 滑動視窗 (Sliding Window)
 - 總保留音訊長度為 0.3 秒
 - 每段視窗長度為 0.05 秒
 - 總共要切出 10 段音訊
 - 重疊百分比44.4%



- 資料筆數 = 43 位 \times 1 (術前) \times 10 = 430 筆

05 資料前處理：前處理後的資料敘述

- 每個樣本各轉一張Mel 與 STFT 頻譜圖
- 對照CSV 包含每筆樣本的音檔檔名 (sound.files) 與標籤 (outcome)

輸出變項	Outcome	術後狀況	類別型	0	well	170 (44.18%)
				1	poor	260 (55.81%)

- 輸入：同一段語音的兩種時頻表示 (兩張影像)
- 輸出： $p(y=1)$ 的機率 (sigmoid)

05 資料前處理：前處理後的資料敘述

- 目標：類別型變數 —— 術後狀況
- 音檔層級

訓練集

Fold	正樣本 (pos)	負樣本 (neg)	總數	正樣本比例
1	209	126	335	0.6239
2	220	116	336	0.6548
3	189	150	339	0.5575
4	219	126	345	0.6348
5	199	146	345	0.5768

測試集

Fold	正樣本 (pos)	負樣本 (neg)	總數	正樣本比例
1	50	40	90	0.5556
2	39	50	89	0.4382
3	70	16	86	0.8140
4	40	40	80	0.5000
5	60	20	80	0.7500

05 資料前處理：前處理後的資料敘述

- 目標：類別型變數 —— 術後狀況
- 病患層級

訓練集

Fold	正樣本 (pos)	負樣本 (neg)	總人數	正樣本比例
1	21	13	34	0.6176
2	22	12	34	0.6471
3	19	15	34	0.5588
4	22	13	35	0.6286
5	20	15	35	0.5714

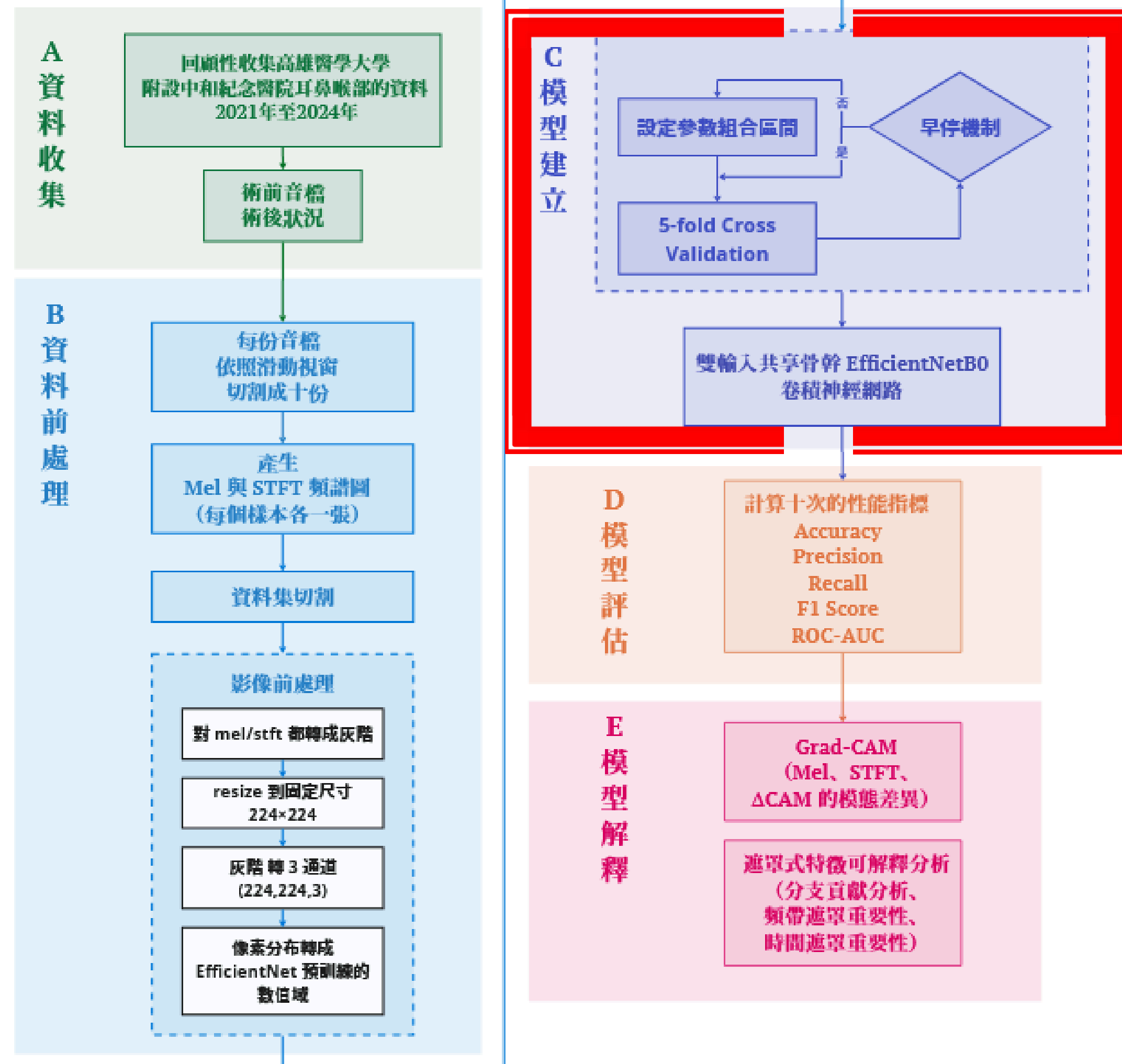
測試集

Fold	正樣本 (pos)	負樣本 (neg)	總人數	正樣本比例
1	5	4	9	0.5556
2	4	5	9	0.4444
3	7	2	9	0.7778
4	4	4	8	0.5000
5	6	2	8	0.7500

05 進模型前的影像前處理

- 對 mel/stft 都轉成灰階
 - 因為實驗過程中會有音檔過短導致整張變成紫色單色圖的問題
- resize 到固定尺寸：224×224
- 灰階 → 3 通道(224,224,3)
- 把像素分布轉成 EfficientNet 預訓練的數值域

06 模型建立



06 模型建立

- 同時輸入 Mel 頻譜圖 與 STFT 頻譜圖
- 兩個輸入共用同一個 EfficientNetB0 特徵擷取器（非兩個獨立 CNN）

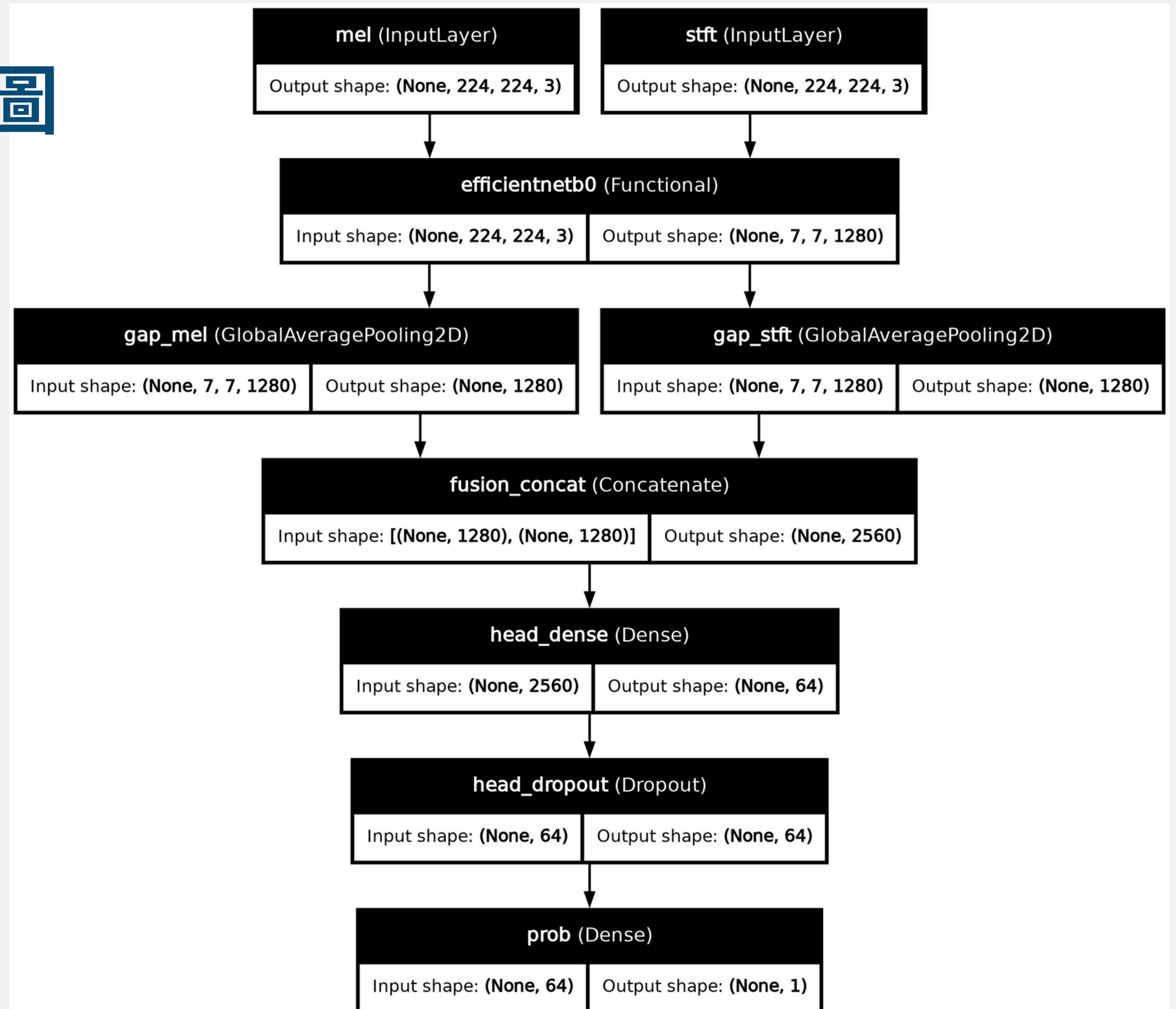
兩階段訓練（freeze → finetune）

1：Warmup（特徵抽取器凍結）

- backbone：全部凍結
- 只訓練 head（Dense+Dropout+Sigmoid）

2：Finetune（只解凍 backbone 最後 20%）

06 模型架構圖

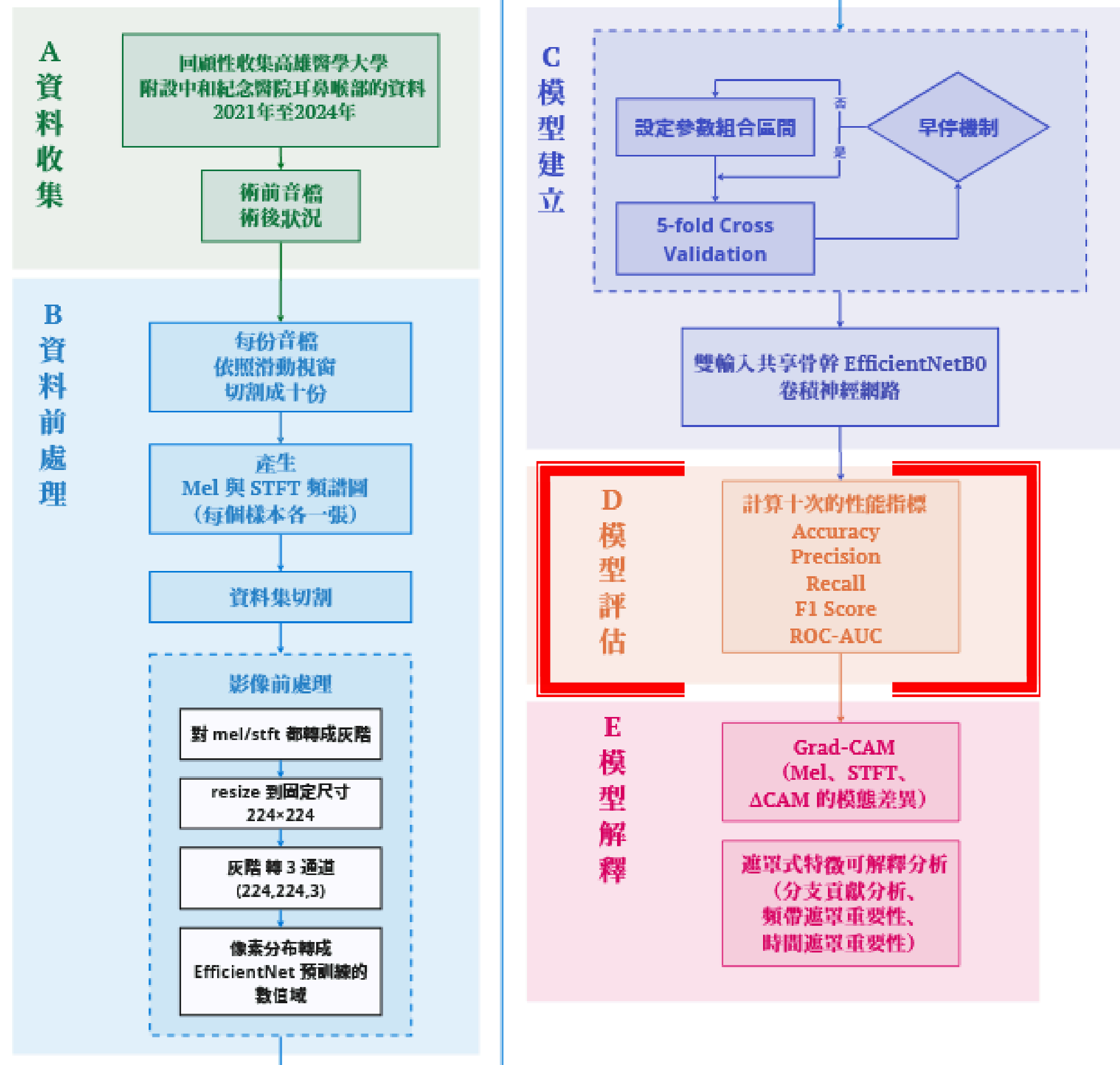


06 模型超參數設定

超參數名稱	意義（用途說明）	設定的範圍 / 數值
IMG_H, IMG_W	輸入頻譜圖影像尺寸 （統一 CNN 輸入規格）	224 × 224
N_FOLDS	外層交叉驗證折數 （評估泛化能力）	5
SEED	隨機種子	42
BATCH_SIZE	每次梯度更新的樣本數	32
VAL_RATIO_IN_TRAIN	在外層 train 中再切出內層 validation 的比例	0.15
EPOCHS_WARMUP	第一階段 warmup 訓練 epoch 數（凍結 backbone）	20
EPOCHS_FINETUNE	第二階段 fine-tune epoch 數 （部分解凍 backbone）	60
LR_WARMUP	warmup 階段學習率	3e-4
LR_FINETUNE	fine-tune 階段學習率（小步調微調）	1e-5

超參數名稱	意義（用途說明）	設定的範圍 / 數值
WEIGHT_DECAY_L2	Dense 層的 L2 正則化係數（抑制過擬合）	1e-4
HEAD_UNITS	分類頭 Dense 層神經元數	64
DROPOUT	分類頭 Dropout 比例（防止過擬合）	0.35
PATIENCE	EarlyStopping 容忍無進步的 epoch 數	15
unfreeze_ratio	fine-tune 時解凍 backbone 的比例	0.2（最後 20% 層）
label_smoothing	Binary Crossentropy 的標籤平滑係數	0.05
class_weight[0]	類別 0 的 loss 權重 （補償少數類）	tr_pos / tr_neg（動態計算）
THR_SWEEP_STEPS	驗證集上搜尋最佳 threshold 的切割數	201
best_monitor	選擇最佳 epoch 的主要指標	val_auroc_prob
Optimizer	參數更新方法	Adam
Loss function	訓練損失函數	Binary Crossentropy

07 模型評估



07 模型評估(1/2) - 評估指標

- 混淆矩陣
- Accuracy
- Precision
- Recall
- F1 score
- ROC / AUC

07 模型評估(2/2) - 模型可解釋種類

- **Grad-CAM：模型在看哪裡、怎麼判斷**
 - 用梯度找出 **單一筆樣本** 模型在影像上 **最敏感、最依賴** 的區域
 - 用熱區圖將決策證據視覺化
 - 另外加上 Δ CAM (mel - stft) 的模態差異診斷，可以看出同一筆資料模型證據主要來自哪個分支
- **Occlusion / Masking (遮罩測試)**
 - 把影像切成 **頻帶區塊 / 時間區塊** 逐段遮掉，看輸出機率掉多少
 - 頻帶遮罩 (Frequency Importance) 評估模型對不同頻率區段的依賴程度
 - 時間遮罩 (Temporal Importance) 分析模型是否依賴特定時間區段或整體時間結構

我選擇結合這兩種方法，希望可以看到 局部 + 全域總結 的價值

07 模型評估

訓練集

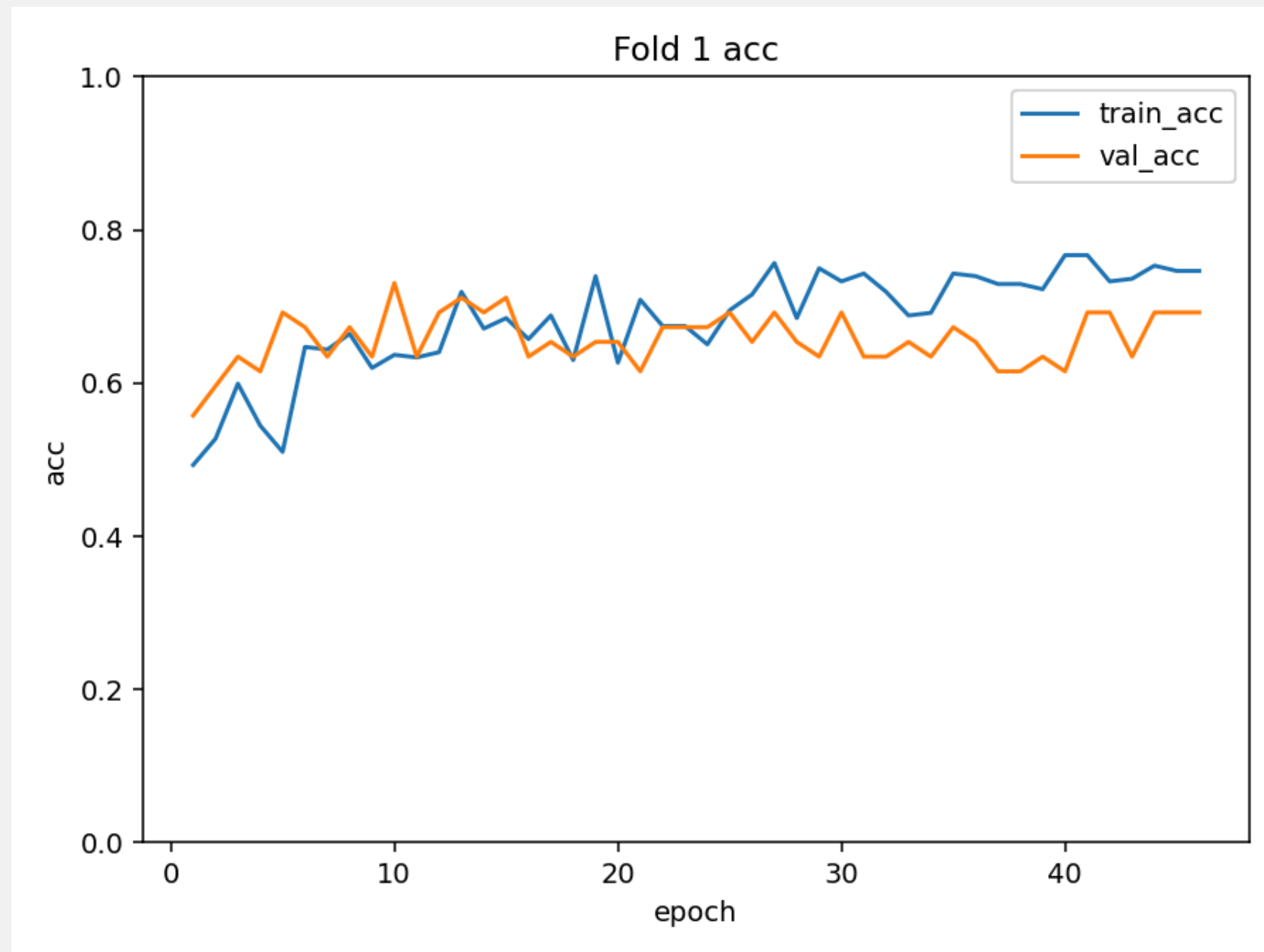
Fold	ACC (%)	Precision (%)	Recall (%)	F1	AUROC (%)	Best Epoch
1	68.15%	65.67%	99.44%	79.10	85.57%	31
2	91.10%	90.37%	95.48%	92.86	97.77%	79
3	91.78%	95.27%	90.96%	93.06	97.25%	80
4	90.75%	96.30%	88.14%	92.04	97.16%	80
5	90.07%	90.66%	93.22%	91.92	96.99%	80
Mean ± Std	86.37% ± 9.13%	87.65% ± 11.25%	93.45% ± 3.86%	89.80 ± 5.37%	94.99% ± 5.10%	70.0 ± 19.5

測試集

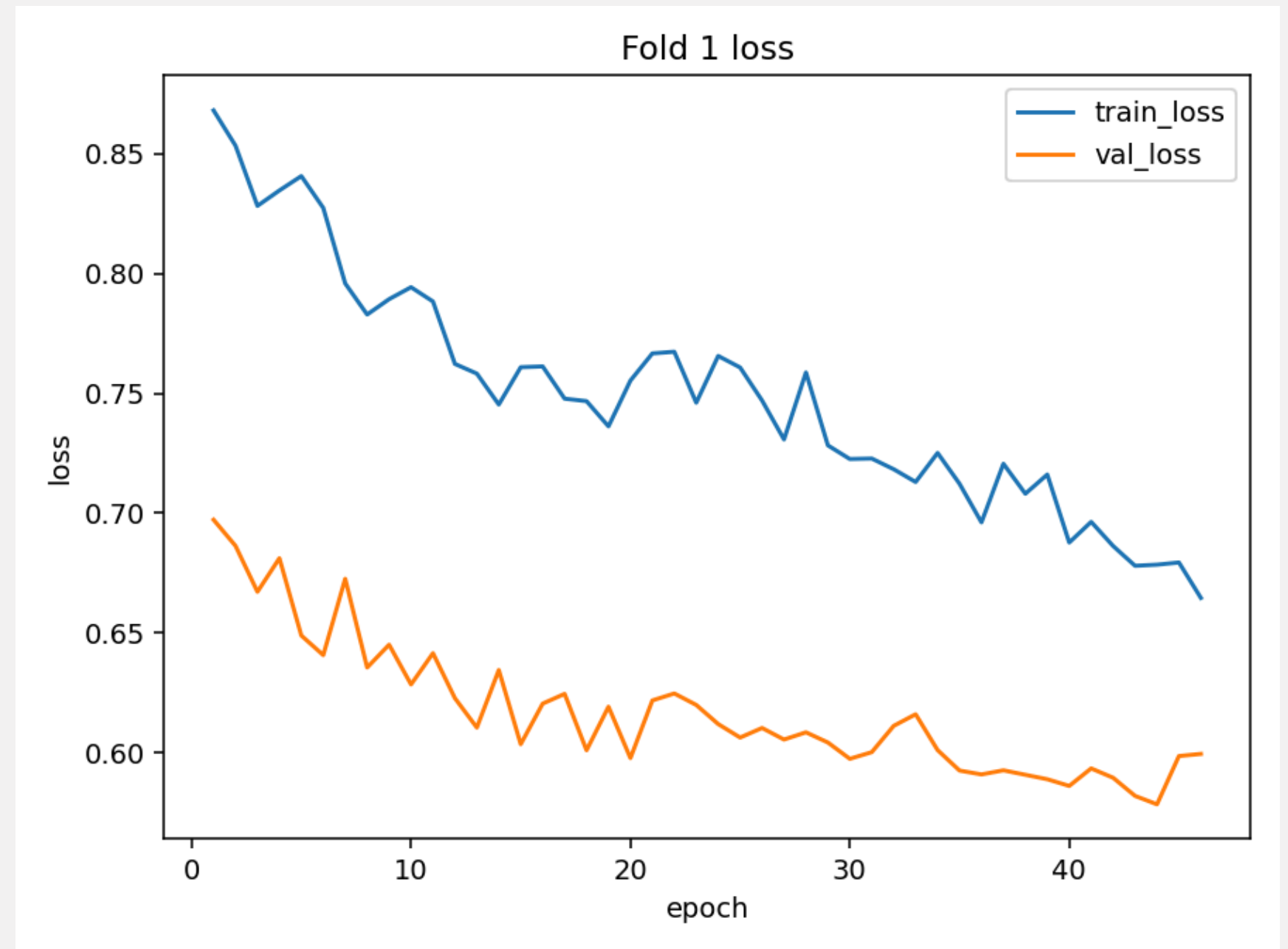
Fold	ACC (%)	Precision (%)	Recall (%)	F1	AUROC (%)
1	68.60%	66.23%	98.08%	79.07	75.48%
2	86.05%	85.71%	92.31%	88.89	88.86%
3	76.74%	79.63%	82.69%	81.13	83.20%
4	73.26%	85.37%	67.31%	75.27	87.44%
5	83.72%	85.19%	88.46%	86.79	90.84%
Mean ± Std	77.67% ± 6.47%	80.43% ± 7.44%	85.77% ± 10.51%	82.23 ± 5.00	85.16% ± 5.45%

07 模型評估 - fold 1

準確度曲線

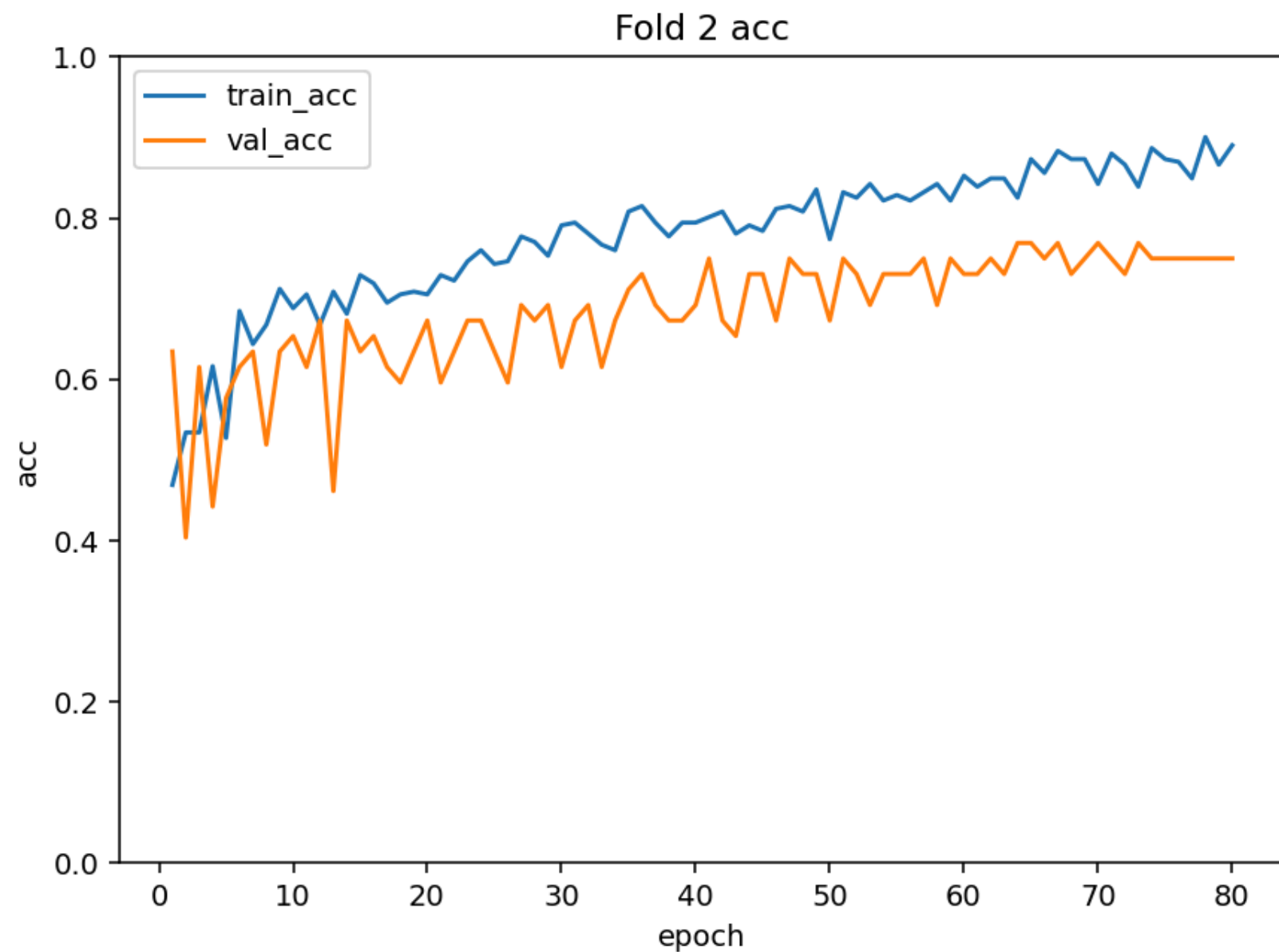


損失曲線圖

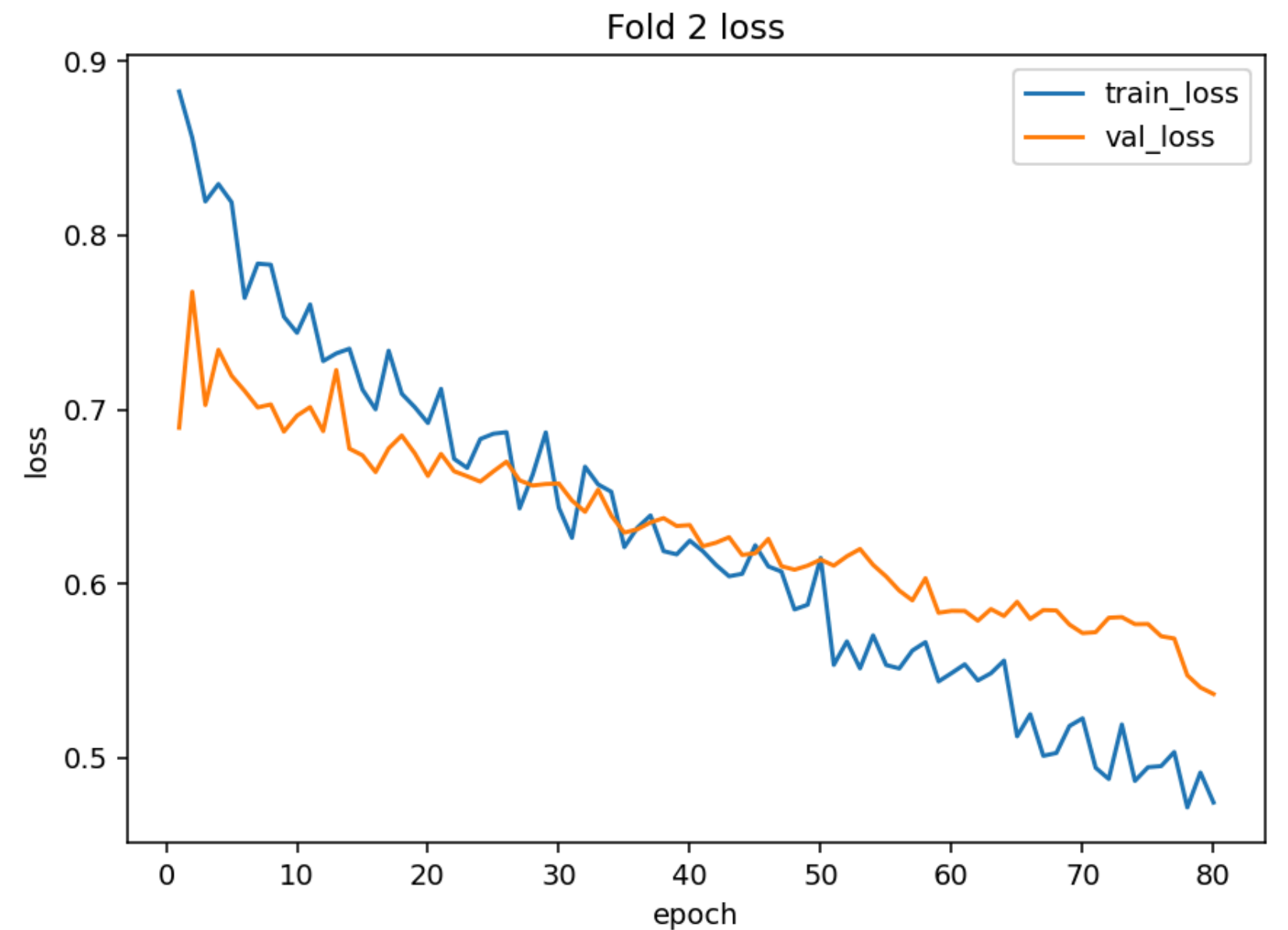


07 模型評估 - fold 2

準確度曲線

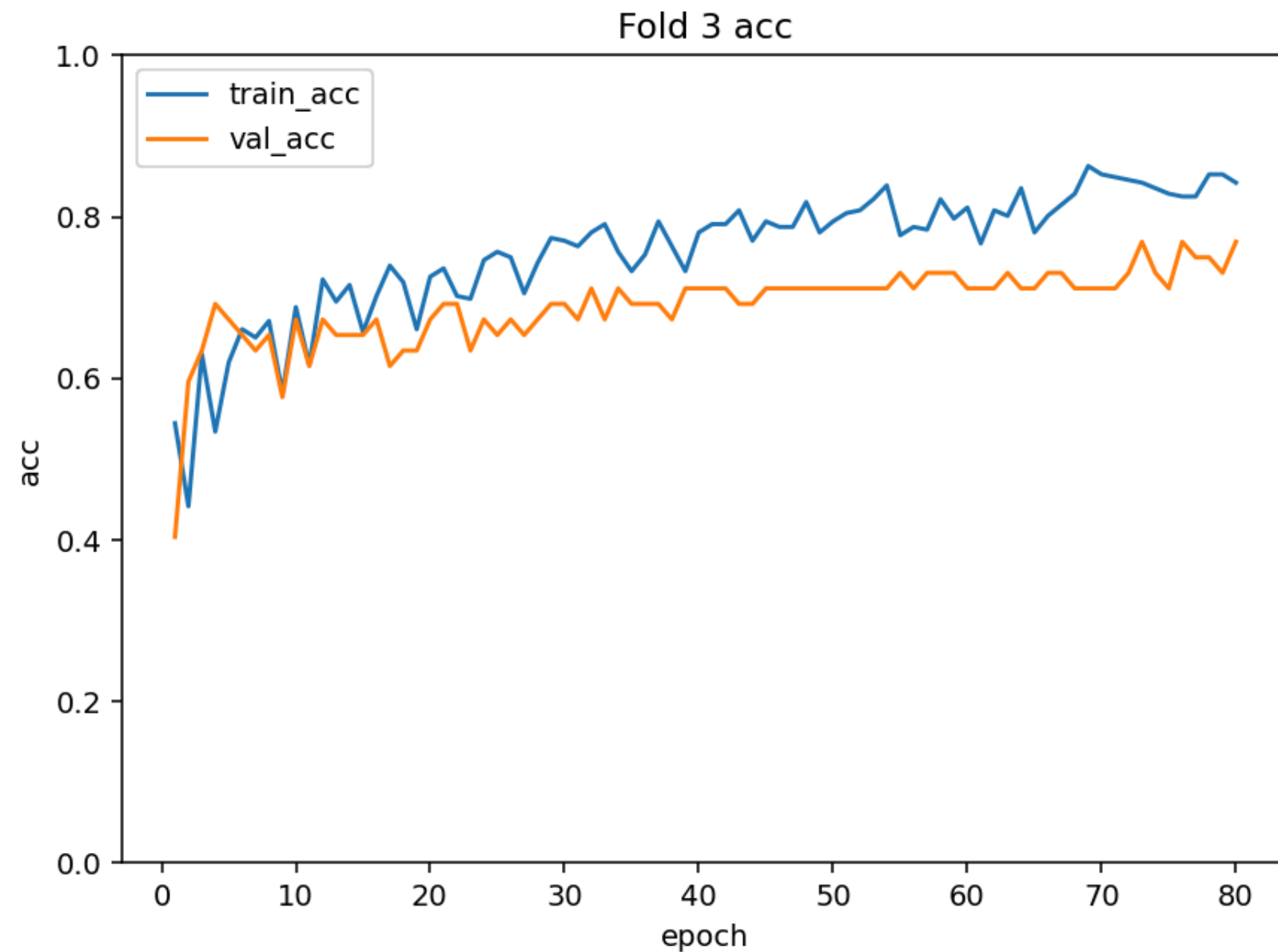


損失曲線圖

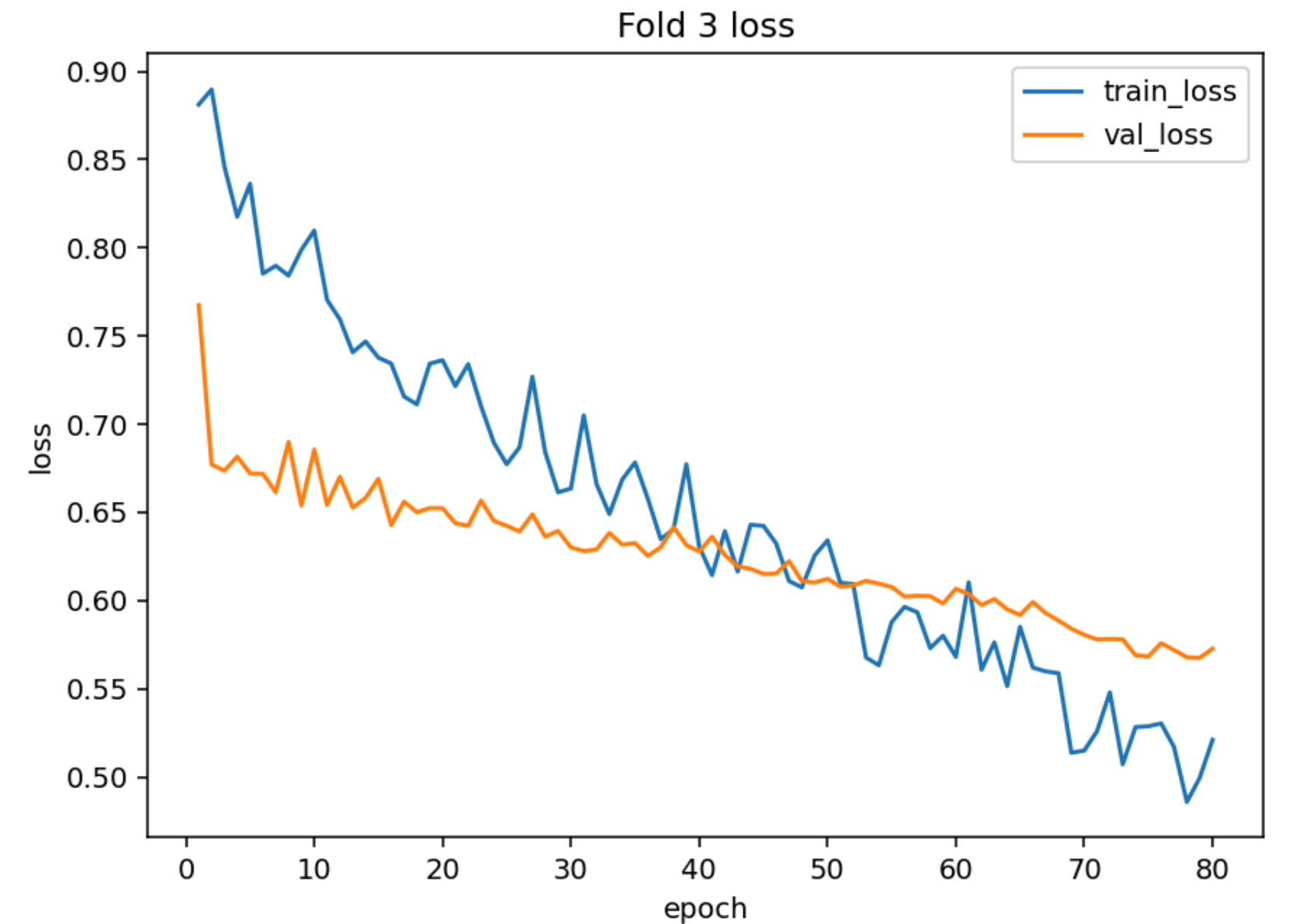


07 模型評估 - fold 3

準確度曲線

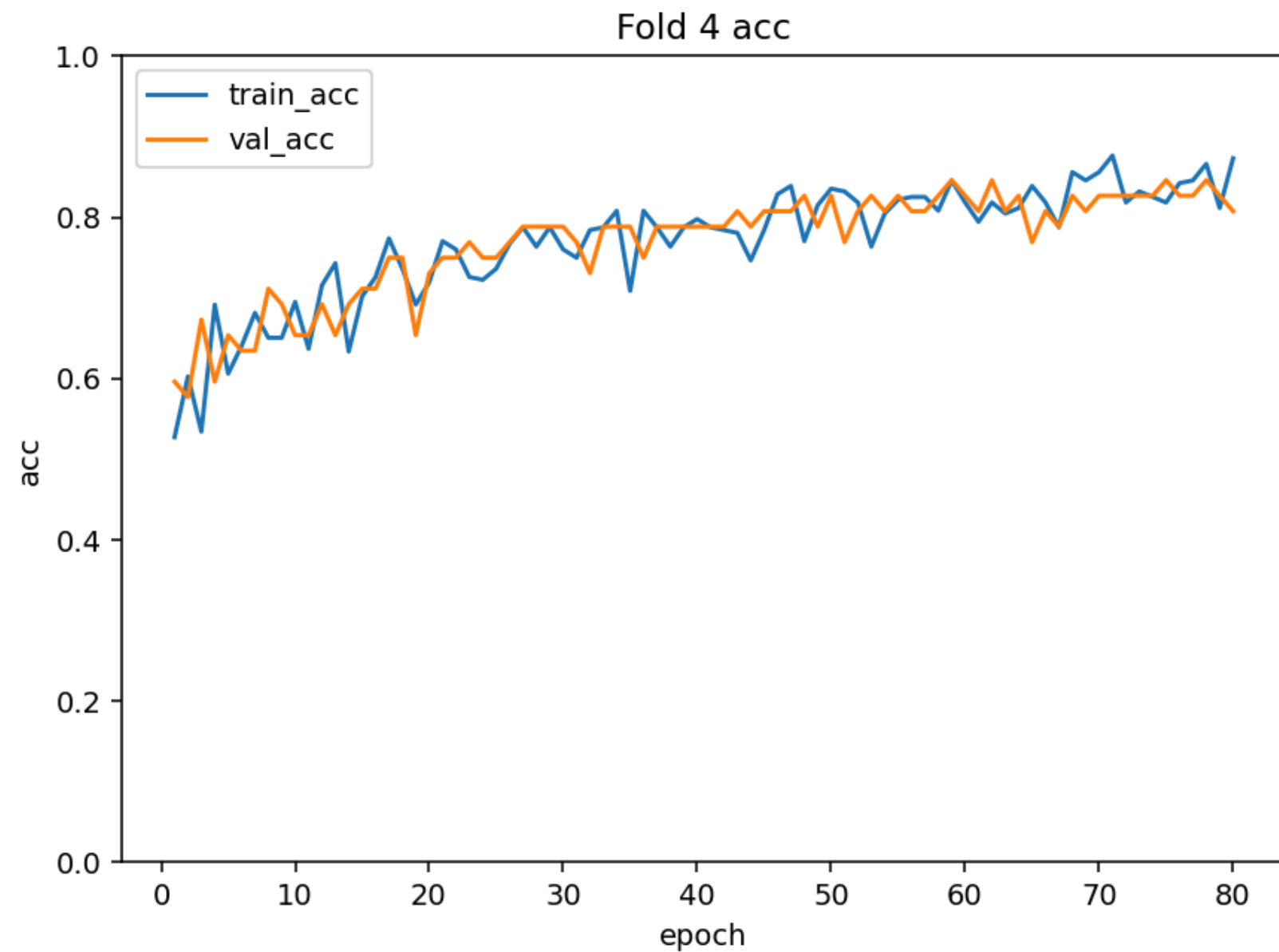


損失曲線圖

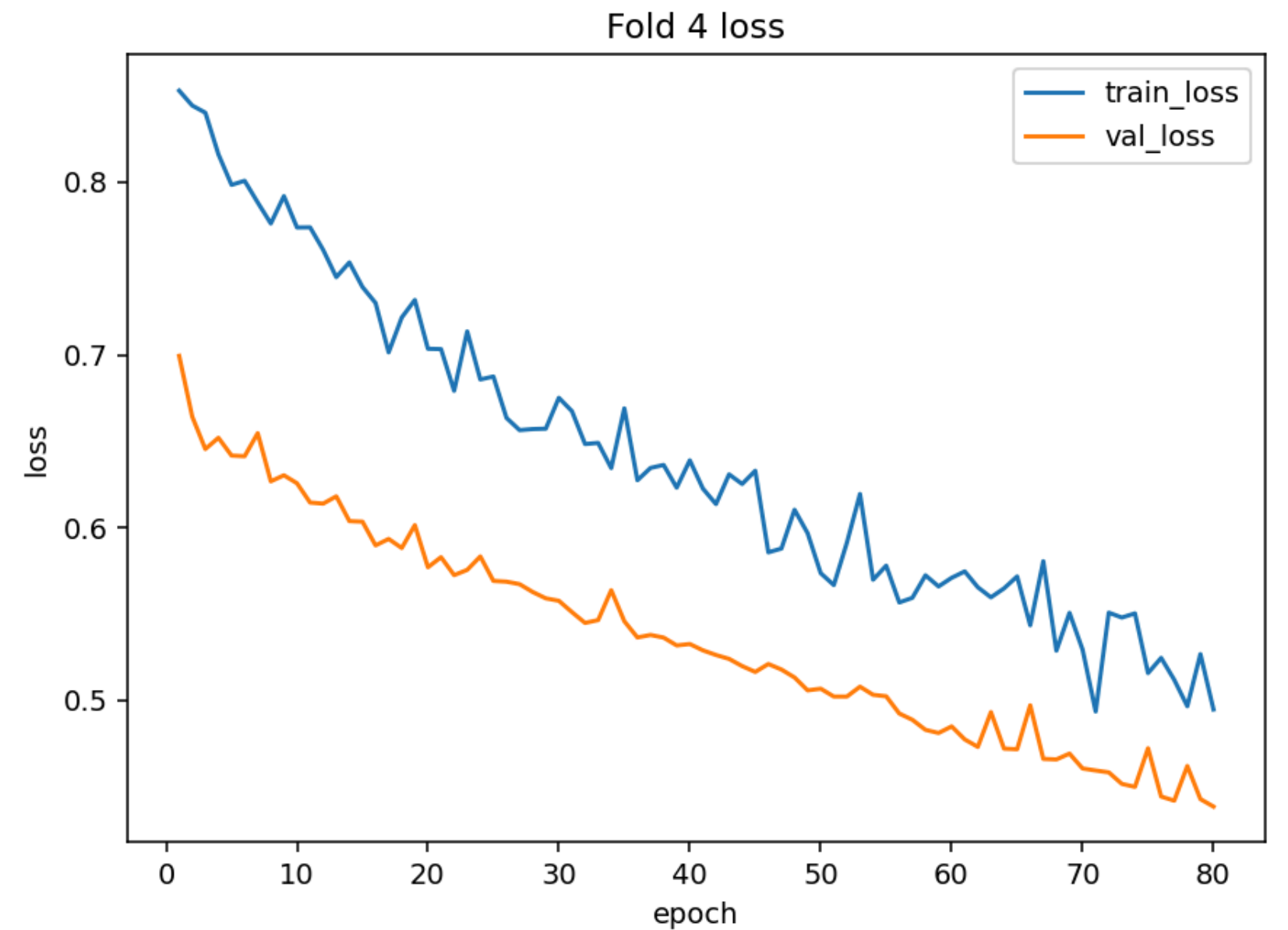


07 模型評估 - fold 4

準確度曲線

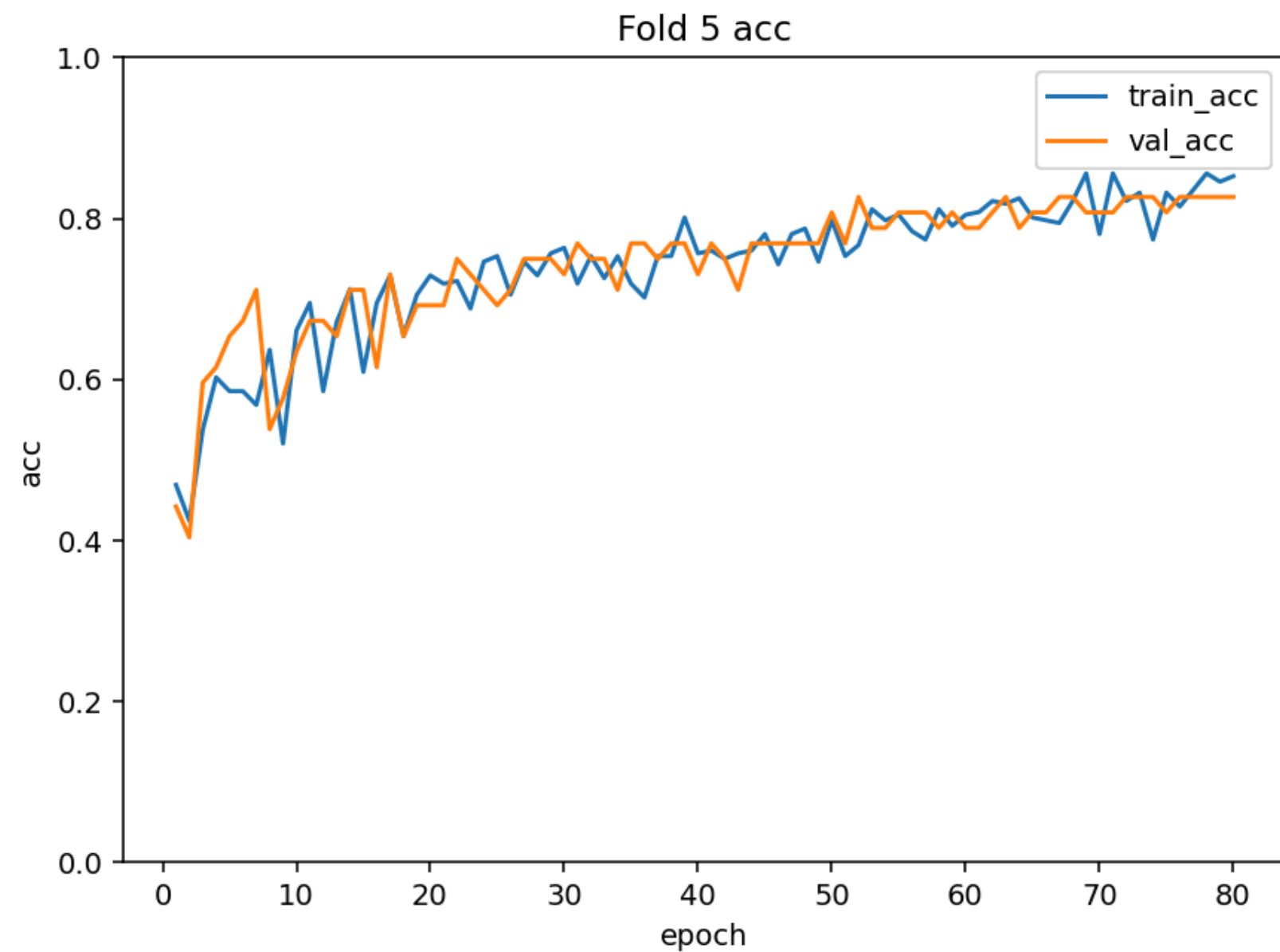


損失曲線圖

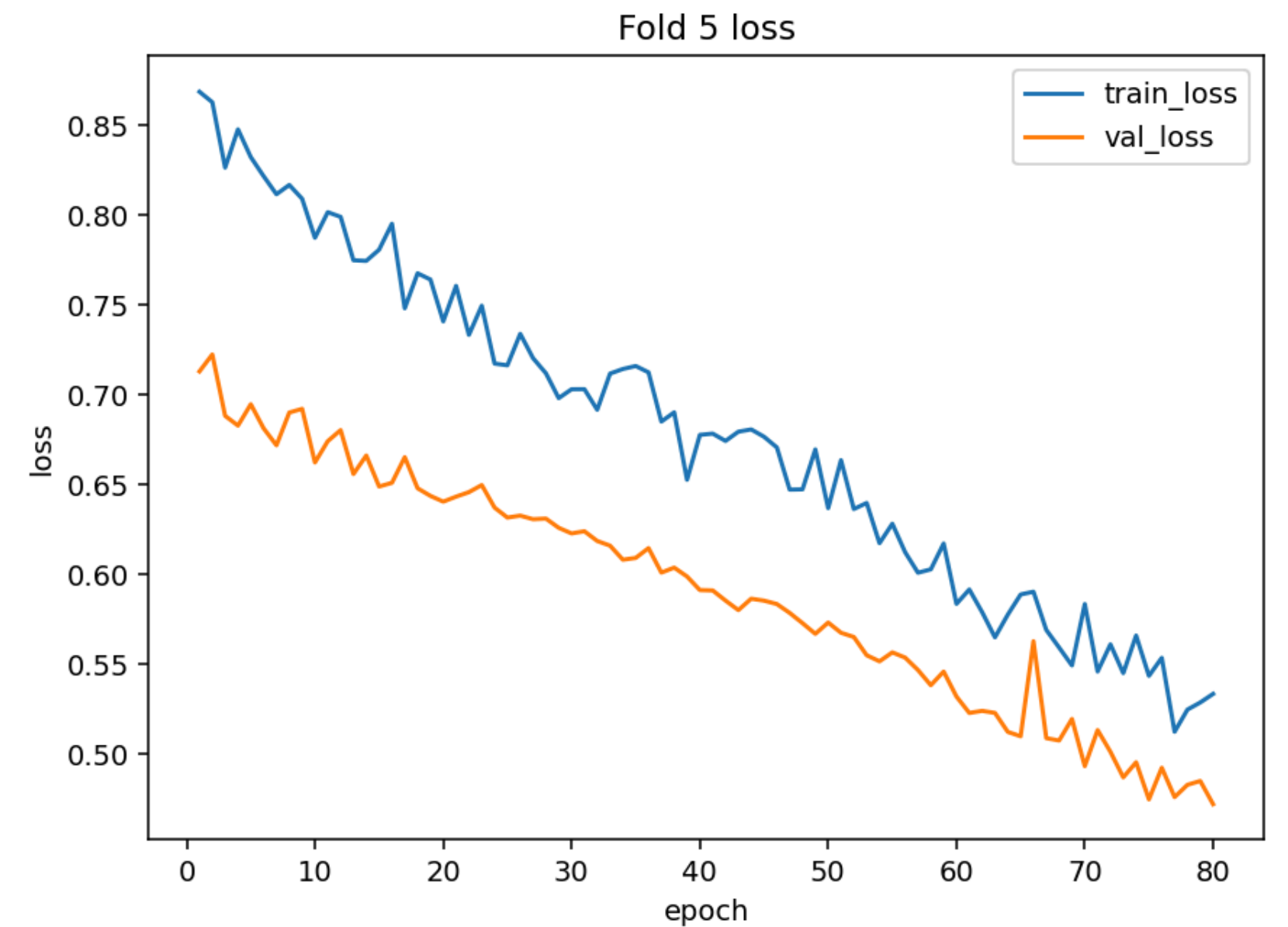


07 模型評估 - fold 5

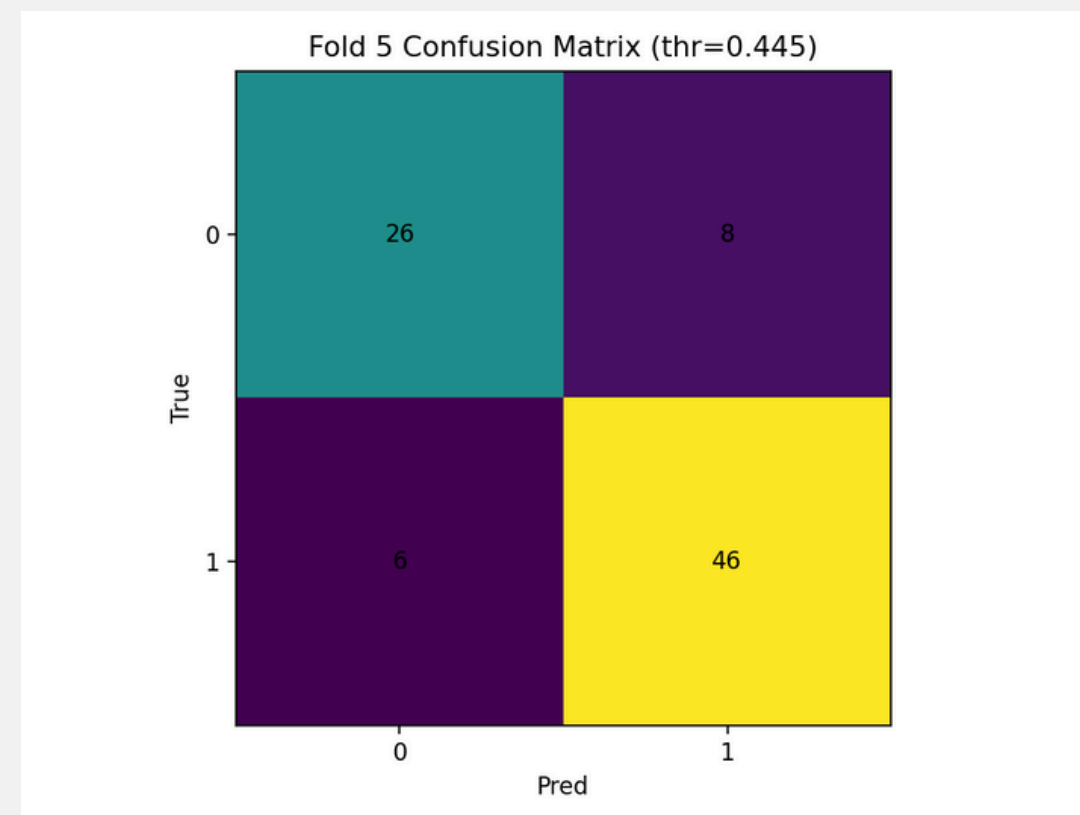
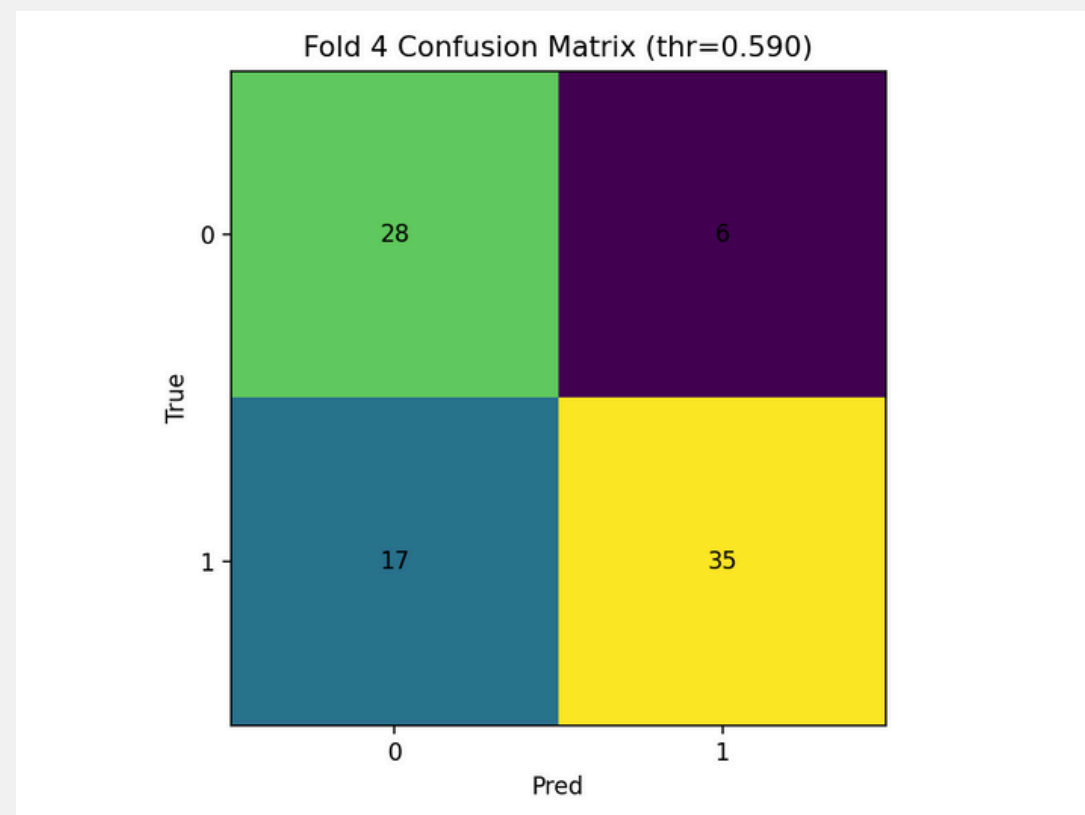
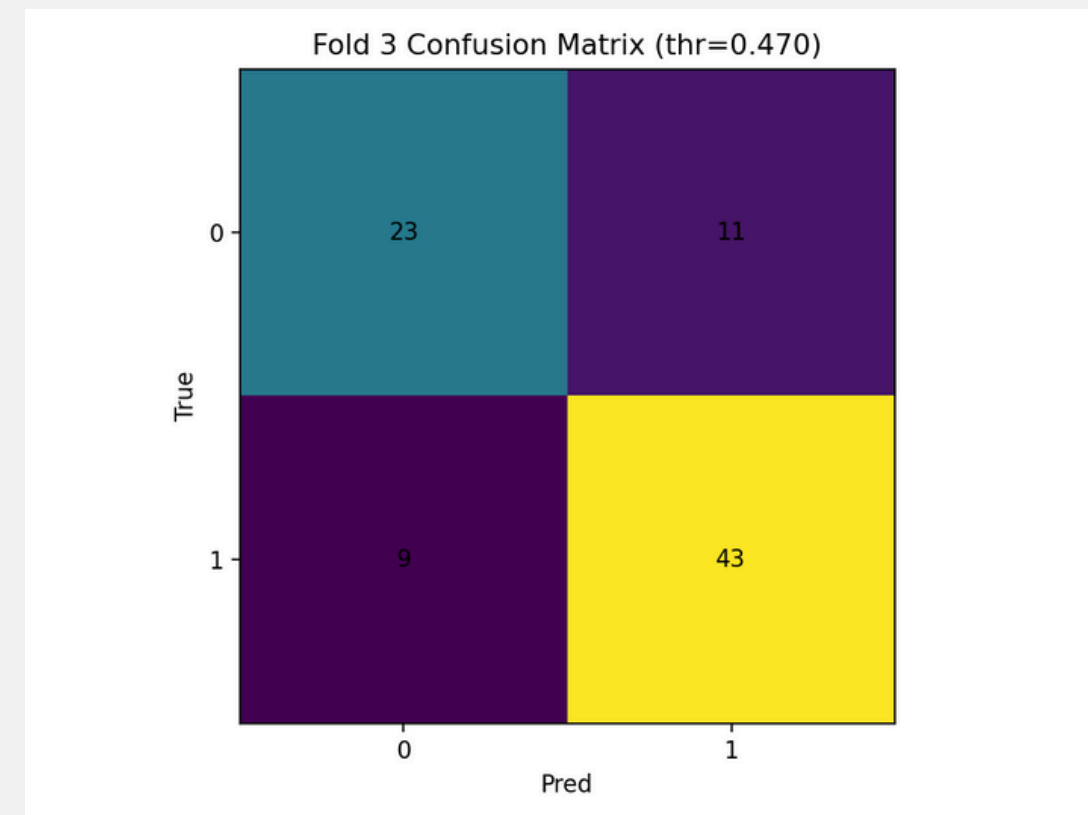
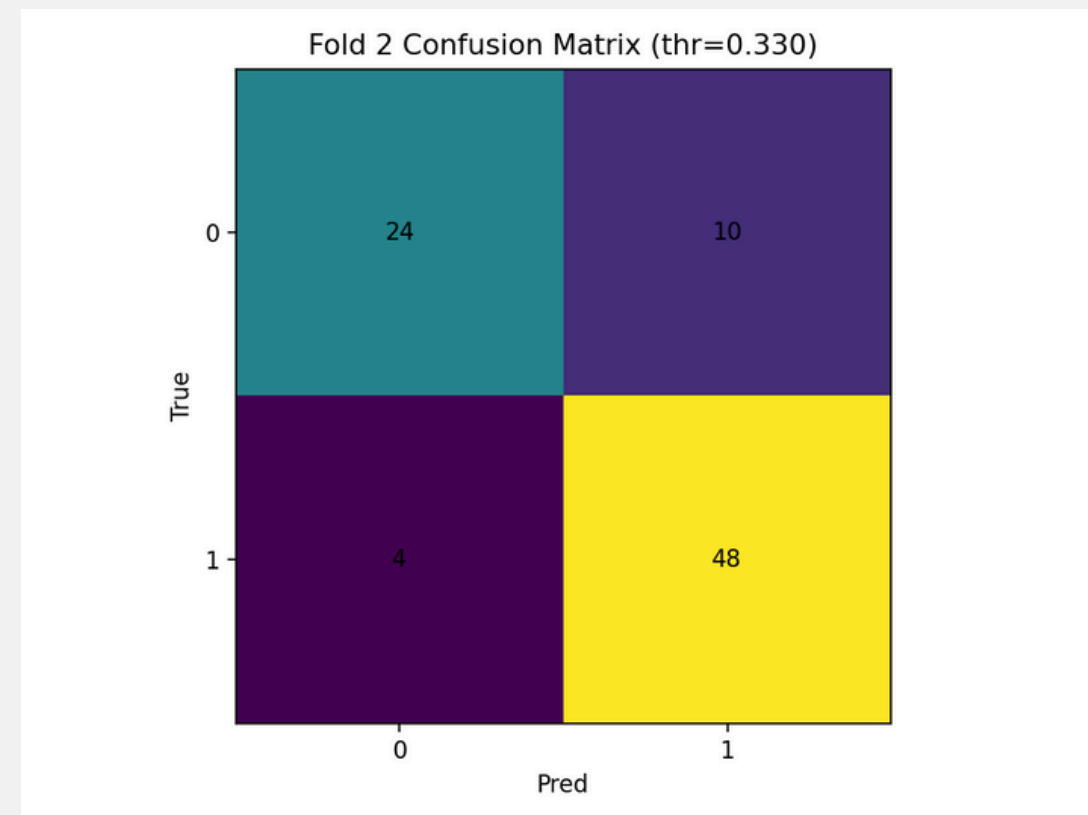
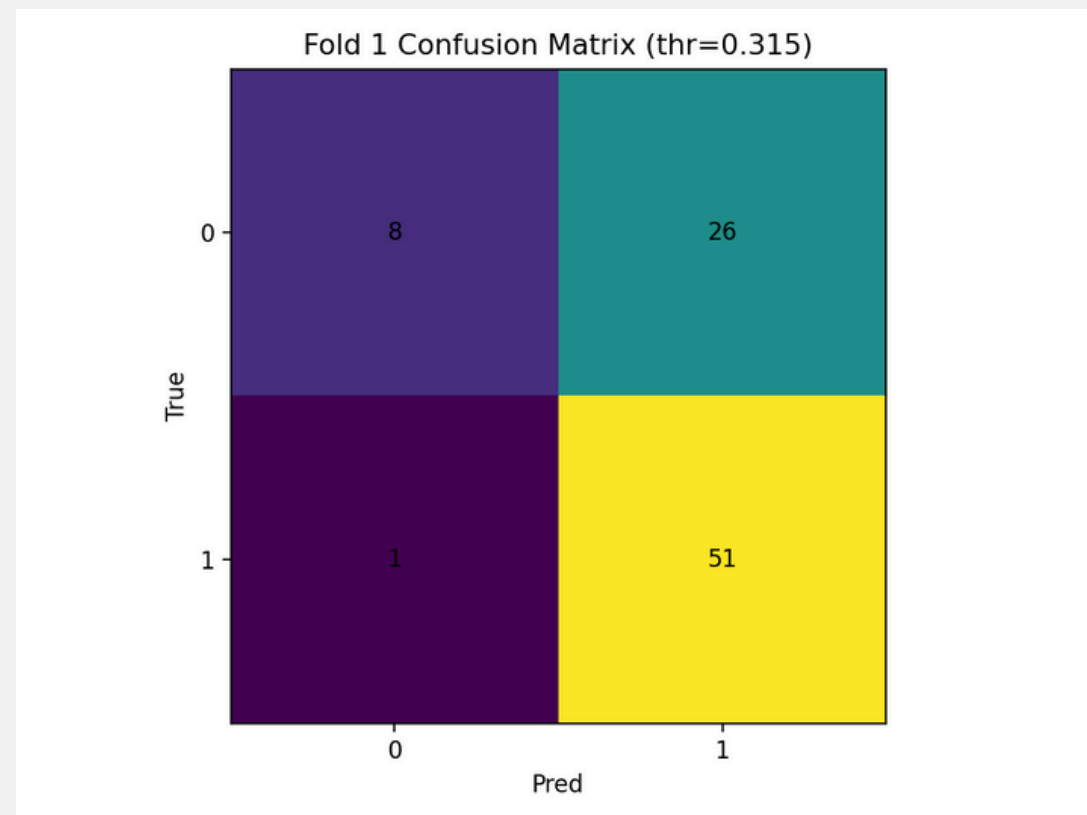
準確度曲線



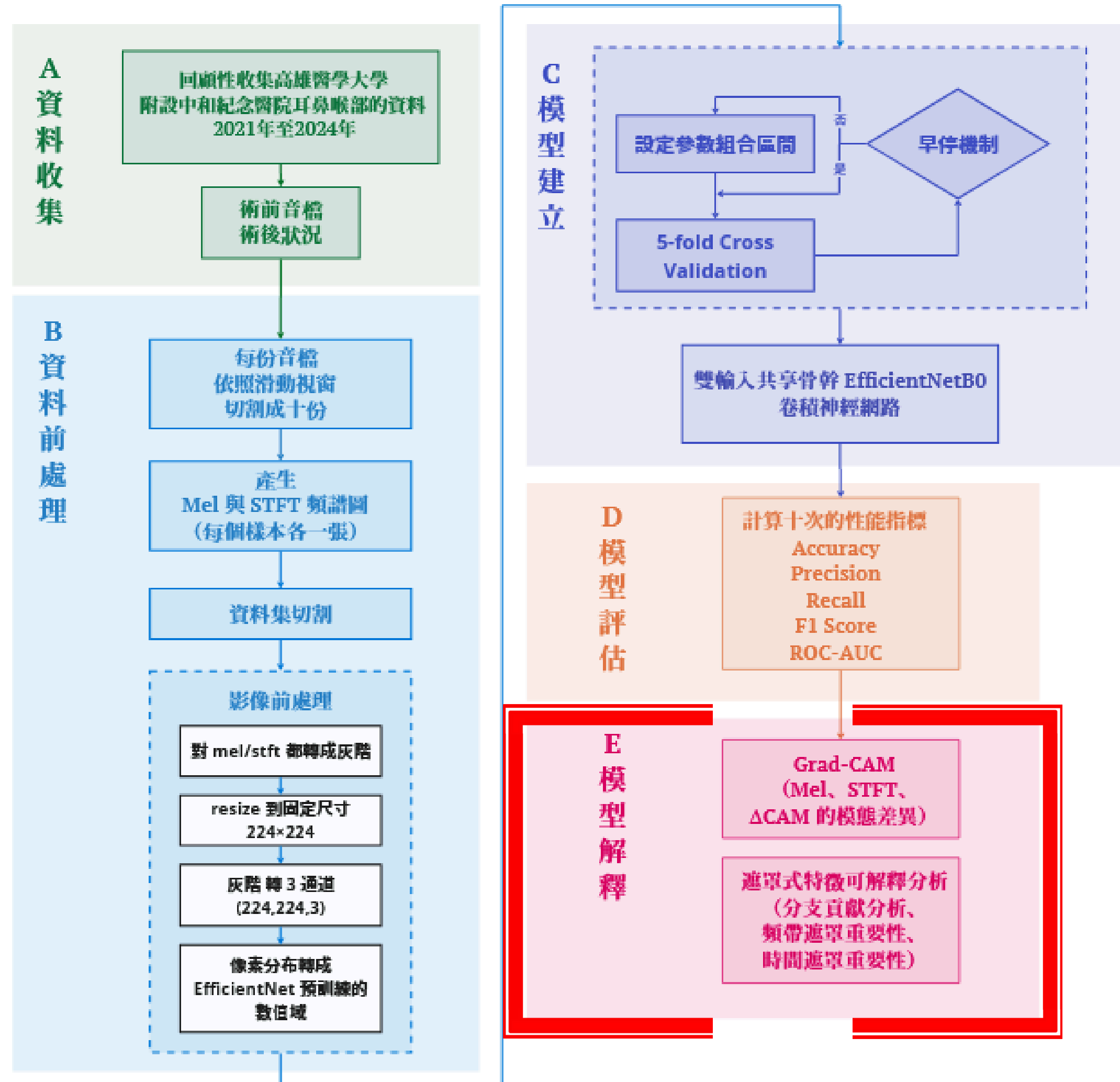
損失曲線圖



07 模型評估 混淆矩陣

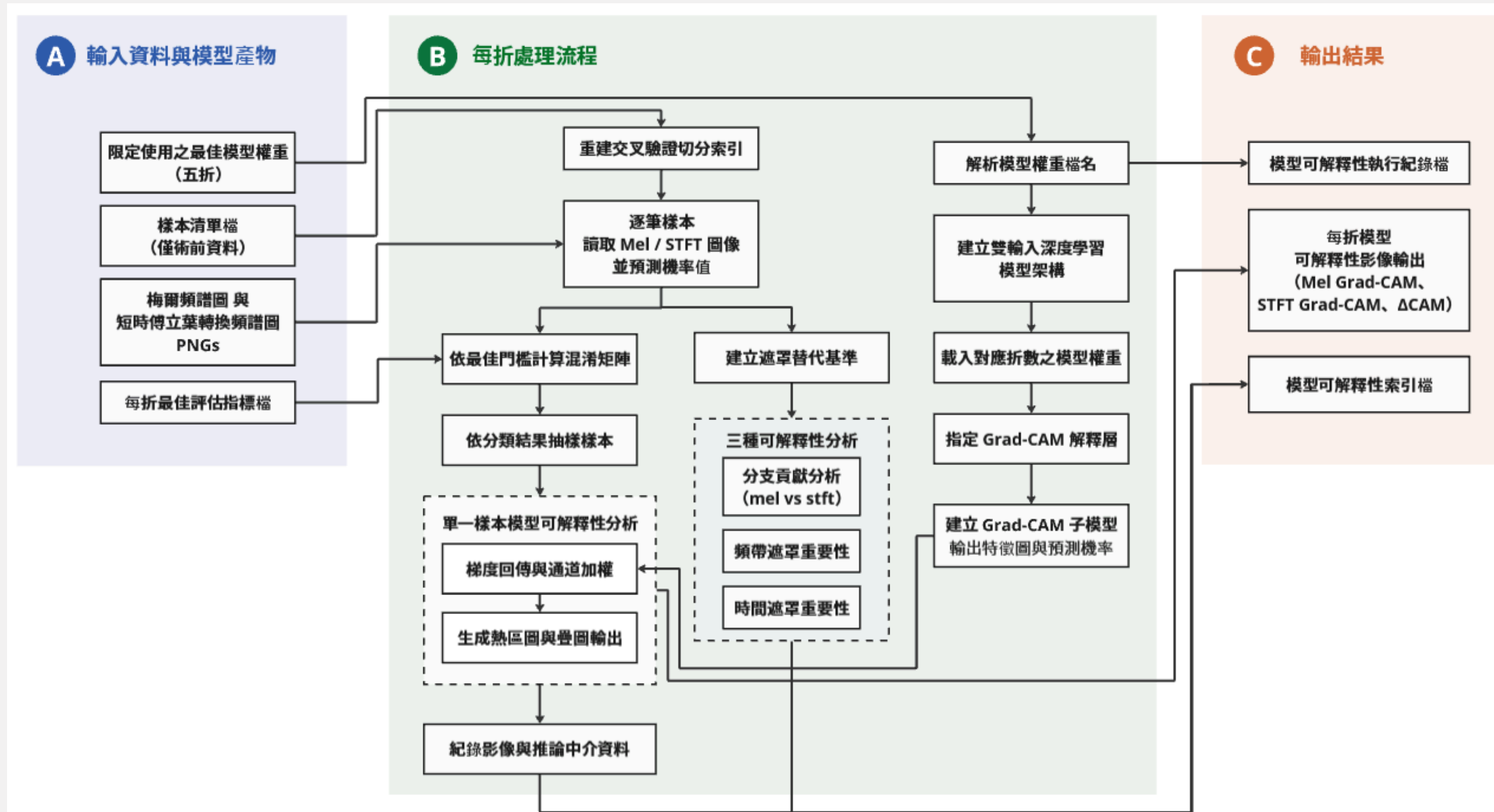


07 模型解釋



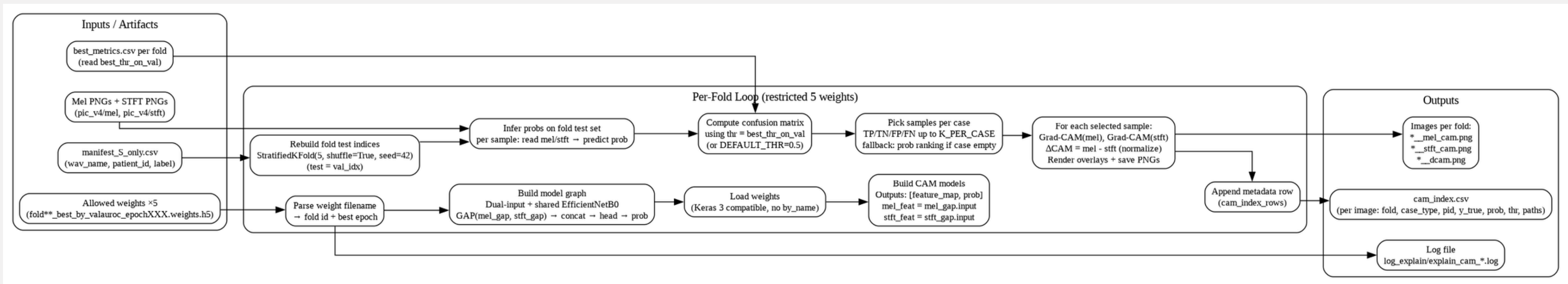
07 模型推論與可解釋性分析架構

- 使用已訓練完成之模型權重進行可解釋性分析
- 同一個模型推論結果，同時支援模型層級與特徵層級兩種可解釋性分析



07 模型推論與可解釋性分析架構

- (英文版本)
- 本流程圖描述模型推論與可解釋性分析；模型訓練流程不包含於此圖中

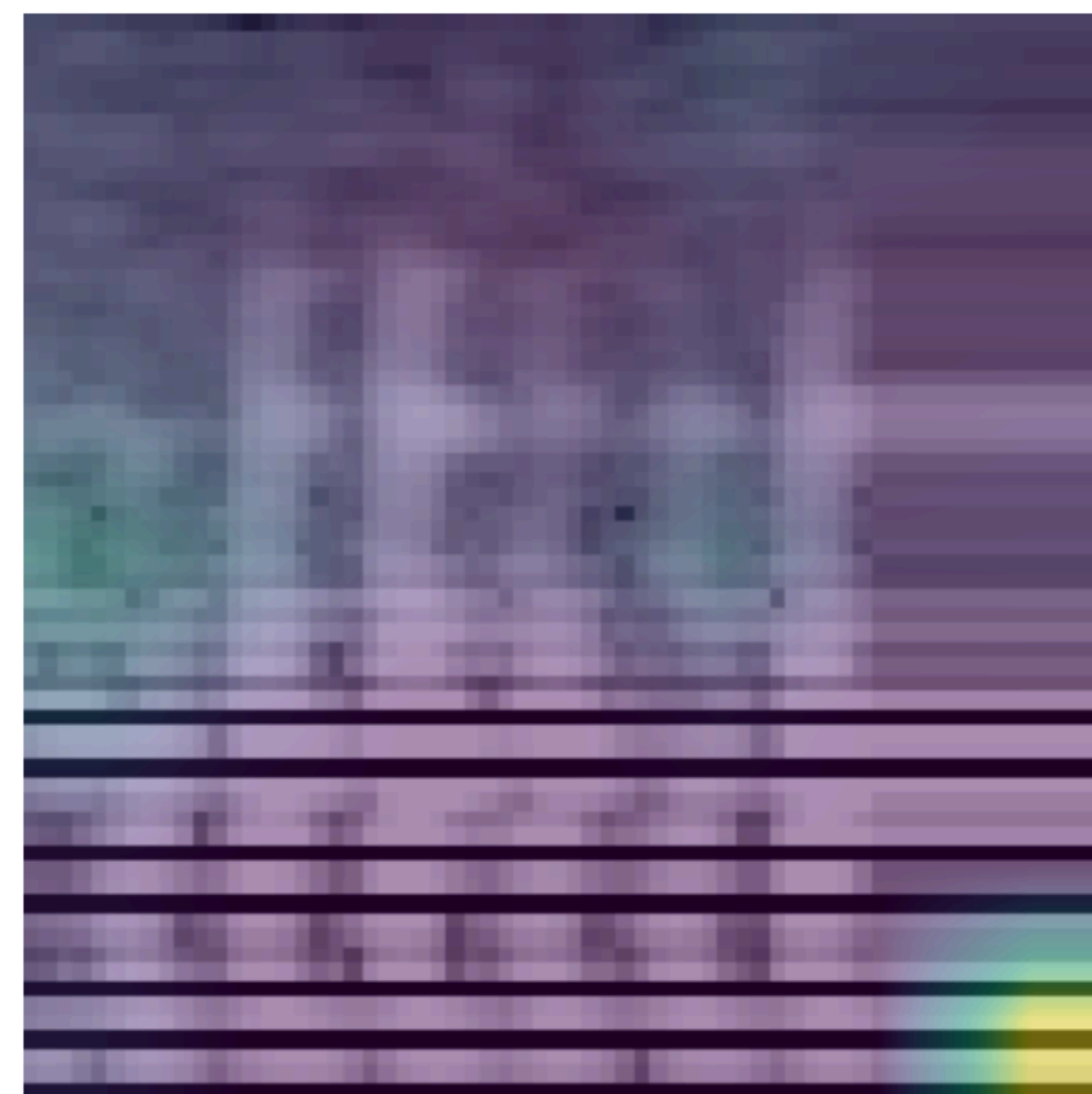


07 Grad-CAM (1/4) 以fold01_FN_pid14_14_S1_10_為例

梅爾時頻圖的 Grad-CAM

- 模型主要關注：低到中頻區段、非連續的時間窗
- 模型在梅爾時頻圖中，捕捉到部分疑似病理相關的頻帶特徵，但這些特徵強度不足或持續時間不夠長

縱軸
||
頻率



橫軸＝時間

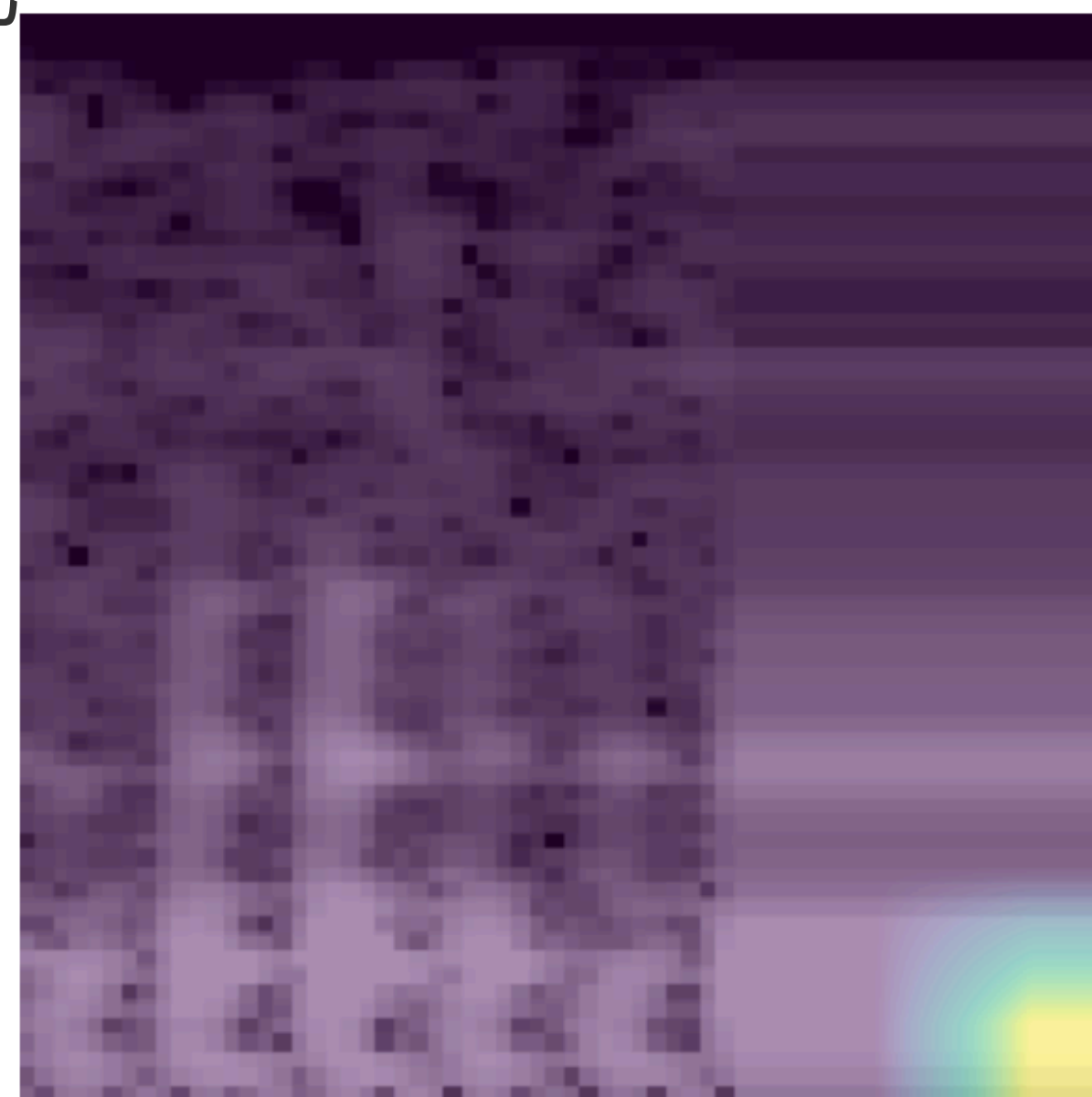
- 顏色越亮 → 對預測為正類的貢獻越大

07 Grad-CAM (2/4) 以fold01_FN_pid14_14_S1_10_為例

STFT 的 Grad-CAM

- 顯示模型在短時間頻譜中關注的關鍵時間 × 頻率區域
- 模型利用 STFT 捕捉 瞬時異常或突發變化

縱軸
||
頻率



橫軸＝時間

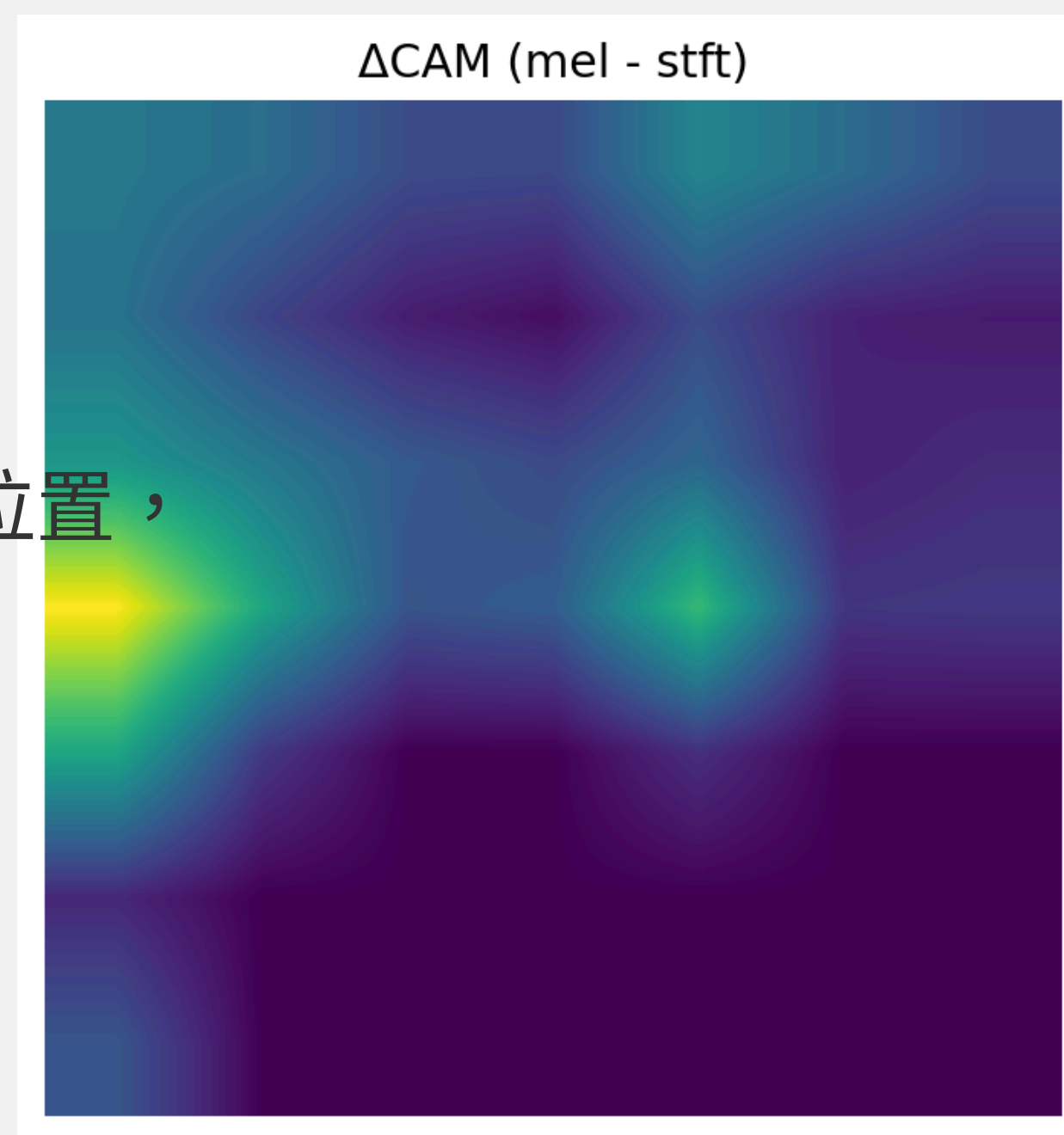
- 顏色越亮對模型判斷越重要；暗色是貢獻低

07 Grad-CAM (3/4) 以fold01_FN_pid14_14_S1_10_為例

- 模型在梅爾時頻圖比STFT更依賴的位置與較不依賴的位置
- 在不同時間與頻段，注意力明顯偏向某一種圖
- 梅爾時頻圖在低頻、局部時間窗出現較強主導，頻譜圖整體影響較弱

亮色 / 偏黃區域 ($\Delta > 0$)

→ 模型在這些時間 × 頻率位置，
更依賴梅爾時頻圖



暗色 / 偏紫區域 ($\Delta < 0$)

→ 這些區域STFT的相對貢獻較高

07 Grad-CAM (4/4) 以fold01_FN_pid14_14_S1_10_為例

- Mel CAM

- 模型在 低到中頻區段 有注意力，聲學結構已偏離健康樣態
- 但注意力零散，沒有長時間維持
- 模型推斷是 **早期或間歇性病理特徵**

- STFT CAM

- 病人可能偶發不穩定聲音，但 **頻率與持續度未構成明確風險訊號**

- Δ CAM

- 有些區域偏向 mel、有些偏向 stft：模型不確定要選結構性特徵還是瞬時特徵

這位病人已經開始出現一些異常聲音的跡象

模型在不同的聲學表示（mel 與 STFT）中都有看到異常狀況

但這些異常出現得不夠久、不夠集中

所以模型判斷風險偏高但仍落在臨界值附近，沒有直接判成陽性

07 遮罩式特徵可解釋分析

- mel_branch_mean：在該 fold 中，把 mel 分支遮掉時，模型預測機率平均下降多少
- stft_branch_mean：在該 fold 中，把 stft 分支遮掉時，模型預測機率平均下降多少
- Δprob 越大（正值）＝該分支越重要
- 負值＝遮掉反而讓預測更高，代表在整體上不是主要決策依據

fold	mel_branch_mean	stft_branch_mean	解讀
1	+0.078	−0.111	強烈依賴 mel，stft 整體偏干擾
2	+0.102	−0.151	mel 主導非常明顯
3	+0.150	−0.011	mel 為主，stft 幾乎中性
4	+0.095	+0.017	唯一一折 stft 有輕微正貢獻
5	+0.104	−0.045	mel 主導，stft 次要

跨 5 折交叉驗證，模型決策主要基於梅爾時頻圖
符合人類對呼吸或聲音異常的感知方式

07 遮罩式特徵可解釋分析

- 頻帶層級 (Frequency-wise)
 - 把 Mel / STFT 圖沿著 **頻率方向切成 10 段**，一次遮掉一段看模型信心掉多少
 - 模型對 **不同頻率區段** 的依賴 **不平均**，主要集中在 Mel 頻譜上
 - 頻帶重要性主要來自 Mel，STFT 的頻帶遮罩影響小、不穩定
- 時間層級 (Time-wise)
 - 把頻譜圖沿著 **時間軸切成 10 段**，遮掉某一段時間觀察模型預測機率變化
 - 結果發現 **模型依賴一段時間內的聲學型態**

**跨 5 折交叉驗證，模型決策主要基於梅爾時頻圖
符合人類對呼吸或聲音異常的感知方式**

08 結果與討論 (1/2)

- 深度學習模型測試集準確度皆近八成 → 具一定預測能力
- 雙輸入架構中，模型主要依賴 Mel 頻譜進行判斷，STFT 提供輔助資訊
- 我認為 STFT 更像是正則化的存在，單用 STFT 的效果雖不好，但結合梅爾時頻圖一起有助於模型準確判斷
- 此分析僅採 a 母音音檔，納入其他母音有助於辨識更準確完整
- 再修改輸入輸出可以預測出需要使用多少的玻尿酸劑量，對臨床上的資源利用率會有正面效益

08 結果與討論 (2/2)

限制與心得

- 資料量很少，每段音檔長度已經用到極致，待資料量更充足，希望能用常見的 ResNet, LeNet 等 CNN 測試看看效果
- 從曲線圖可看出模型尚未完全收斂，也因為資料量小，epoch 設很小，這次期末報告試用兩種圖來訓練模型，在閱讀文獻就有看到梅爾時頻圖的效果會比較好，如果可以應該還可以再找別種可以從音檔的擷取的特徵，建構多模態模型，讓 AI 在臨床上更具信度

114-1

期末報告

Thanks for your attention

Q&A

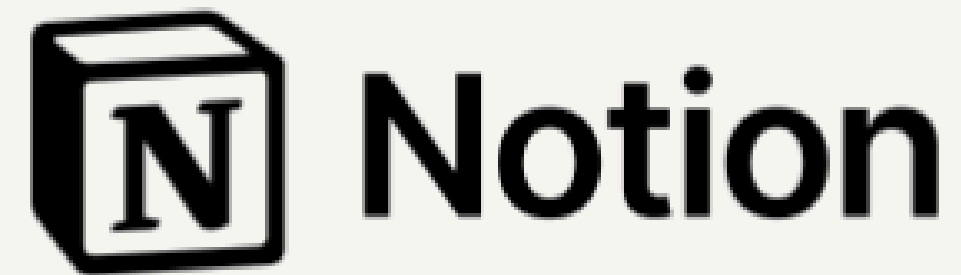
報告
合作醫師：
指

醫師

報告日期：2025. 01. 12

01 附錄

更多補充文字的連結



114-1 最佳化期末報告 | Notion

資料介紹

 statuesque-arrow-b28 on Notion

https://statuesque-arrow-b28.notion.site/114-1-2e436e7f278c801995cfd6439d14df89?source=copy_link

02 附錄

資料的連結（持續上傳）

hiimsharon/ai-laryngeal-injection-outcome-...



應用人工智慧方法預測注射式喉成型術後狀況

1 Contributor 0 Issues 0 Stars 0 Forks

hiimsharon/ai-laryngeal-injection-outcome-prediction: 應用人工智慧方法預測注射式喉成型術後狀況

應用人工智慧方法預測注射式喉成型術後狀況. Contribute to hiimsharon/ai-laryngeal-injection-outcome-prediction development by creating an account on GitHub.

GitHub