

Task B: Designing a Two-Layer Network for Image Classification

Q36134255 電通所 碩一 郭人瑋

1. 簡介

隨著深度學習在影像分類領域的快速發展，ResNet 系列模型因為其殘差連接帶來的深度可擴展性，成為了強力的基準（baseline）。在本次的 Task B 我設計了一個僅包含 2-4 個有效層的輕量化網路（TinyViT），並在 ImageNet-mini 資料集上達到 ResNet34 最後驗證準確度的 90% 以上，同時比較兩者在參數量與計算量（FLOPs）上的差異。

2. 資料集

- **mini-ImageNet**

mini-ImageNet 是原始 ImageNet 的縮小版，保留了多樣性與挑戰性，同時大幅減少了計算量

```
image/                                # 所有影像檔的根目錄
├── class_0001/                        # 第一個類別的影像子資料夾
│   ├── img1.jpg
│   ├── img2.jpg
│   └── ...
├── class_0002/
└── ...
train.txt                             # 訓練集路徑與標籤
validation.txt                       # 驗證集路徑與標籤
test.txt                             # 測試集路徑與標籤
```

- **預處理**

在訓練與評估流程中，我對影像做了以下標準化處理：

- I. 尺寸調整: $IMG_SIZE = 32$ 以符合模型輸入層的空間維度
- II. 隨機水平翻轉 (train)

- III. ToTensor + ImageNet 標準化 (mean=[0.485,0.456,0.406],
std=[0.229,0.224,0.225])
- IV. 將影像全轉為 RGB 格式，因為裡面有灰階圖片會導致通道錯誤
程式無法運行

3. 方法設計

● 設計原理與文獻依據

1. Vision Transformer (ViT) 思想

參考自「An Image is Worth 16×16 Words: Transformers for Image Recognition at Scale」，將圖像切成固定大小 patch，再以純 Transformer 架構進行分類，證明了自注意力機制在圖像領域的強大表徵能力。

2. 輕量化 TinyViT

參考自「TinyViT: Fast Pretraining Distillation for Small Vision Transformers」，在當中作者採用了分階段的 ViT 架構，並結合了快速蒸餾，而我的設計由於作業的需求並未進行分階段處理以及並未進行蒸餾處理，而是採用了裡面 Patch Embedding 設計跟 Transformer Block 的設計

本作業設計

基於上述文獻，我們採用 **Patch Embedding + 單層 TransformerBlock**，在僅有 **3 個有效層**（卷積 Patch Embedding、MHSA、FFN）的前提下，充分利用自注意力的全域感受野，並透過位置編碼補足空間資訊。

● 模組架構

A. Patch Embedding

輸入影像為 $\mathbf{X} \in \mathbb{R}$

$$\mathbf{X}' = \text{Conv2d}_{P,P}^{3 \rightarrow d}(\mathbf{X}) \implies \mathbf{X}' \in \mathbb{R}^{B \times d \times \frac{224}{P} \times \frac{224}{P}}$$

$$\mathbf{T} = \text{reshape}(\mathbf{X}') \in \mathbb{R}^{B \times N \times d}.$$

B. Positional Encoding

對每個 token 加上可學習位置向量

$$\mathbf{Z}^{(0)} = \mathbf{T} + \mathbf{P}.$$

C. Transformer Block

對包含一層 Multi-Head Self-Attention (MHSA) 和一層 Feed-Forward Network (FFN)：

$$\begin{aligned} \text{(a) LayerNorm \& MHSA: } \mathbf{A} &= \text{MHSA}\left(\text{LN}(\mathbf{Z}^{(0)})\right), \\ \mathbf{Z}^{(1)} &= \mathbf{Z}^{(0)} + \mathbf{A}, \\ \text{(b) LayerNorm \& FFN: } \mathbf{F} &= \text{FFN}\left(\text{LN}(\mathbf{Z}^{(1)})\right), \\ \mathbf{Z}^{(2)} &= \mathbf{Z}^{(1)} + \mathbf{F}. \end{aligned}$$

$$\begin{aligned} \text{MHSA}(X) &= \text{concat}_{h=1}^H (\text{softmax}(Q_h K_h^\top / \sqrt{d/H}) V_h) W_O, \\ \text{FFN}(x) &= W_2 \text{GELU}(W_1 x + b_1) + b_2 \text{ 且 } W_1 \in \mathbb{R}^{d \times 4d}, W_2 \in \mathbb{R}^{4d \times d} \end{aligned}$$

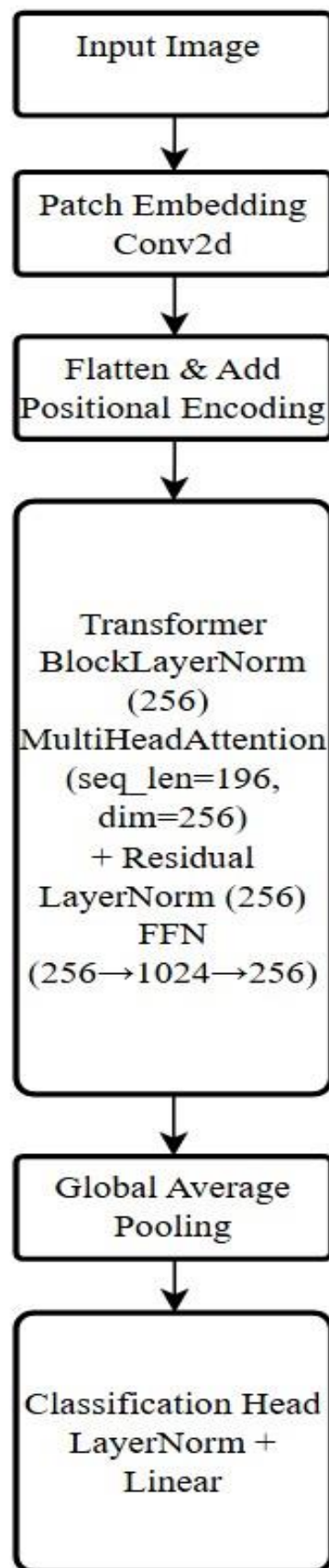
D. Global Average Pooling and Classification

$$\mathbf{z} = \frac{1}{N} \sum_{i=1}^N \mathbf{Z}_{:,i,:}^{(2)} \in \mathbb{R}^{B \times d}$$

經 LayerNorm 後以全連接層輸出 C 類別分數：

$$\mathbf{y} = W_c \text{LN}(\mathbf{z}) + b_c, \quad W_c \in \mathbb{R}^{d \times C}.$$

流程架構： Patch Embedding \rightarrow Positional Encoding \rightarrow MHSA \rightarrow FFN \rightarrow Pooling \rightarrow Classification



架構圖

● 各層設計動機與功能說明

➤ Patch Embedding (有效層 1)

- I. **動機**：直接對整張圖像做全域注意力計算，計算量與記憶體需求極高，不適合資源受限的場景。
- II. **功能**：利用一個大卷積 ($\text{kernel}=\text{patch_size}$, $\text{stride}=\text{patch_size}$) 將影像切塊 (patch) 並映射到維度 d ，相當於 ViT 中的線性投影。
- III. **效益**：可以降低序列長度，大幅減少自注意力 (MHSA) 的計算成本以及保留局部連續性(patch 內的鄰近像素仍在同一 token 中)同時兼顧局部結構。

➤ Positional Encoding (可學習，非有效層)

- I. **動機**：Transformer 本身不具備空間結構資訊，token 序列的順序感必須靠外部編碼注入。
- II. **功能**：為每個 token 加上相對應的可學習位置向量，使模型能區分 patch 在圖像中的位置。
- III. **效益**：可以幫助網路捕捉空間關係，如物體邊界或紋理分佈以及有助於少量 Transformer block 下穩定提取圖像結構。

➤ LayerNorm (在 MHSA 與 FFN 前)

- I. **動機**：深度網路中，不同維度的 activation 易出現分佈偏移 (internal covariate shift)，影響訓練穩定性。
- II. **功能**：對每個 token 的維度特徵做正規化，使其均值為 0、方差為 1。
- III. **效益**：緩解梯度爆炸，提升收斂速度以及可以與殘差連接結合後，有助於深層結構的穩定學習。

➤ Multi-Head Self-Attention (MHSA) (有效層 2)

- I. **動機**：傳統 CNN 雖具備局部感受野，但感受野的拓展需要多層堆疊；而 Self-Attention 能在單層即建構全域關聯。

- II. **功能**：對所有 token 兩兩計算注意力分數，捕捉全局依賴關係中其中的多頭機制（head_first）可並行學習多種關係模式，如邊緣、紋理、顏色相似度等。
- III. **效益**：直接建構全域視野，對小型影像分類有顯著收斂優勢，並且多頭結構提高表徵能力，提升分類性能。

➤ **Feed-Forward Network (FFN) (有效層 3)**

- I. **動機**：Self-Attention 主要負責 token 間互動，缺乏足夠的非線性變換與通道間協調。
- II. **功能**：兩層全連接加上 GELU，為每個 token 獨立進行非線性變換，同時將維度擴展至 4d 再壓縮回 d，增加模型容量。
- III. **效益**：強化特徵維度間混合，有助於表徵精煉，並且與殘差相加後，提升深度結構的可學習空間，穩定且豐富模型表現力。

➤ **Global Average Pooling + Classification Head (非有效層)**

- I. **動機**：Transformer 輸出為 token 序列，需聚合成整張圖像的分類特徵。
- II. **功能**：全域平均池化將所有 token 取平均，得到 z 接著再進行 LayerNorm + Linear 將 z 投影至 C 類別分數空間。
- III. **效益**：可以簡化序列輸出為固定維度，便於後續全連接分類，並利用池化自然具備位置不變性，減少過度擬合空間細節。

4. 實驗設計

● 實驗目標

- 1. 探討不同模組（MHSA、FFN、位置編碼）對模型性能的影響。
- 2. 驗證 patch size 調整對訓練效率與泛化能力的效果。
- 3. 進行消融實驗以分析各結構元件的貢獻。

● 平台與設定

訓練參數

- A. Optimizer：AdamW
- B. Loss function：CrossEntropyLoss
- C. Epoch 上限：50
- D. Early stopping：驗證準確率達 baseline 90% 即停止

變體代號	設計說明	動機
AO	原始 TinyViT 設計	檢驗完整架構基準性能
A1	移除位置編碼	驗證 PositionalEncoding 是否必要
A2	僅使用 FFN（移除 MHSA）	減少計算複雜度與資源消耗
A3	僅使用 MHSA（移除 FFN）	測試無非線性轉換之表現能力
A4	將 patch size 加大為 32	減少序列長度、加速訓練與省資源

● 指標與比較方式

- A. **準確率指標**：每輪輸出 Train,Validation,Test Accuracy，關注準確度與收斂速度。
- B. **資源消耗指標**：報告每個模型的 FLOPs 與參數量（由 ptflops 提供）。
- C. **收斂速度指標**：根據達到 baseline 90% Val Acc 所需 epoch 數進行早停分析。
- D. **模型效率分析**：比較各變體在相同 epoch 數下的效能與計算需求。

5. 結果與分析

- A. 本次任務需設計一個有效層數為 2~4 的輕量級影像分類網路架構，其效能須達到 ResNet34 在 ImageNet-mini 上的 **90% 準確度門檻**。本實驗以自行實作的 TinyViT 為主體，進行了 **原始架構訓練（A0）與 4 組消融實驗（A1-A4）**，針對每種變體進行了訓練、早停、測試準確率與 FLOPs/參數統計，並與 baseline 做比較。
- B. 由於原先有分別進行 Baseline ResNet34 的訓練分析
在 30 Epoch 情況下 Validation 的 Accuracy 可以達到 82.6%

但在後面執行 TinyViT 的執行測試以及消融實驗的測試時由於 Vram 容量不足，導致無法繼續執行，於是改採用 15 Epoch

在 15 Epoch 情況下 Validation 的 Accuracy 可以達到 52.9%

但在後面執行 TinyViT 的消融實驗的測試時耗時太久，平均一則實驗要跑 8 至 12 小時，於是在符合作業要求下，這邊我改採用 5 Epoch 進行實驗

1. Baseline ResNet34 訓練分析

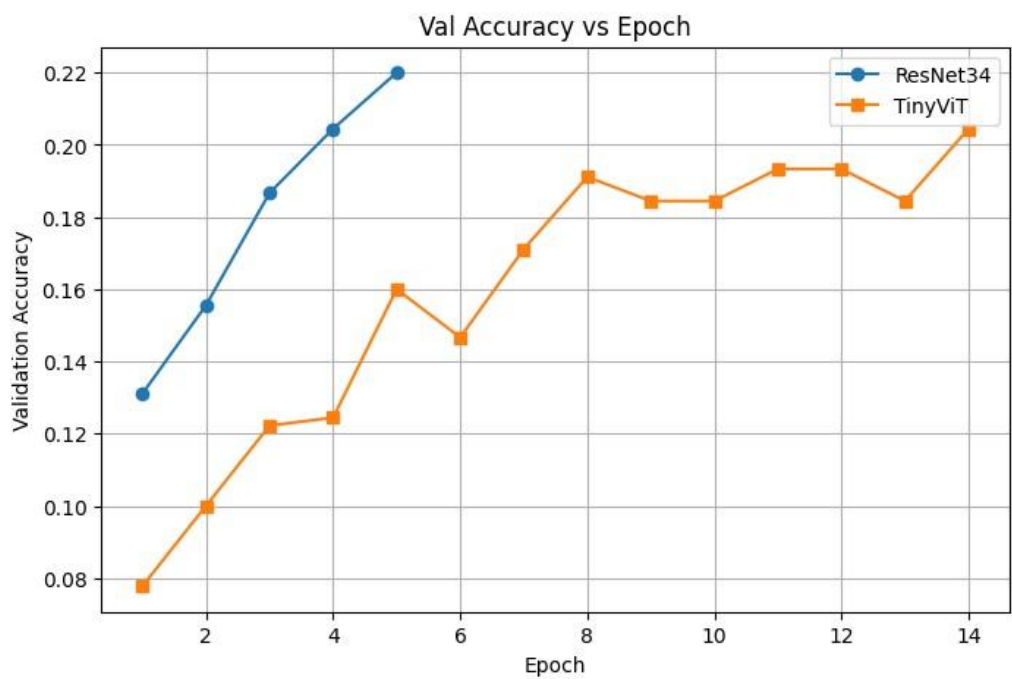
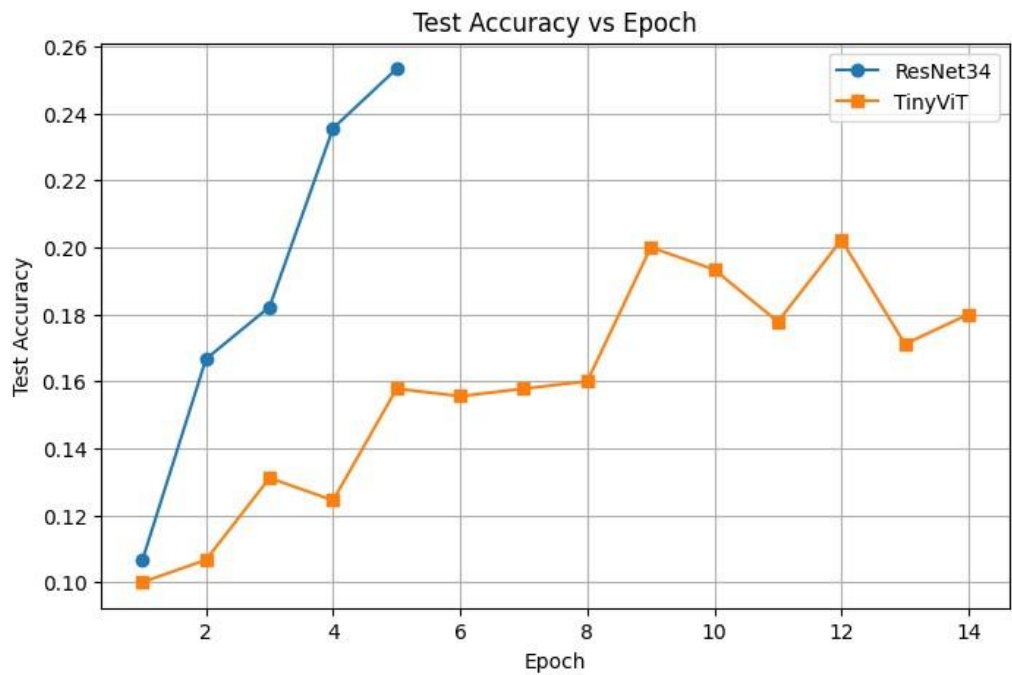


這邊我們可以觀察到從第 1 到第 5 epoch，Validation 與 Test Accuracy 均穩定上升，代表模型尚處於 underfitting → learning 階段，尚未收斂。測試準確度甚至在第 4~5 epoch 出現略高於驗證準確度的現象，顯示 ResNet34 在早期已具有良好泛化能力。

第 5 epoch 時：

- **Validation Accuracy $\approx 22.0\%$**
- **Test Accuracy $\approx 25.5\%$**
- 因此可推算出基準目標 $\approx 19.8\%$ Val Accuracy 為 TinyViT 消融實驗早停門檻。

2. TinyViT 與 ResNet34 相關對比



這邊我們先分析 Test Accuracy 以及 Val Accuracy

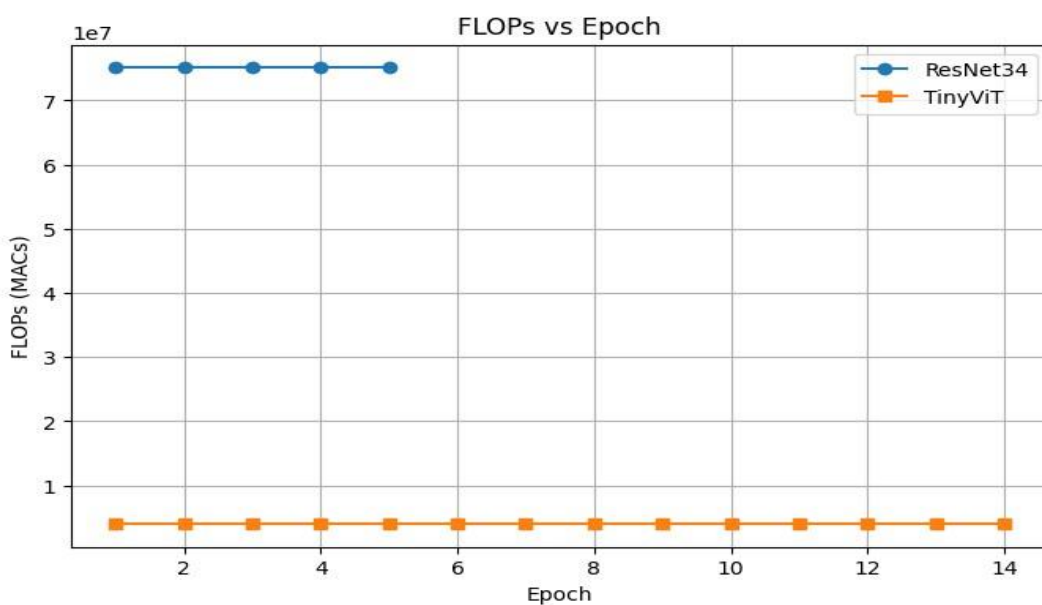
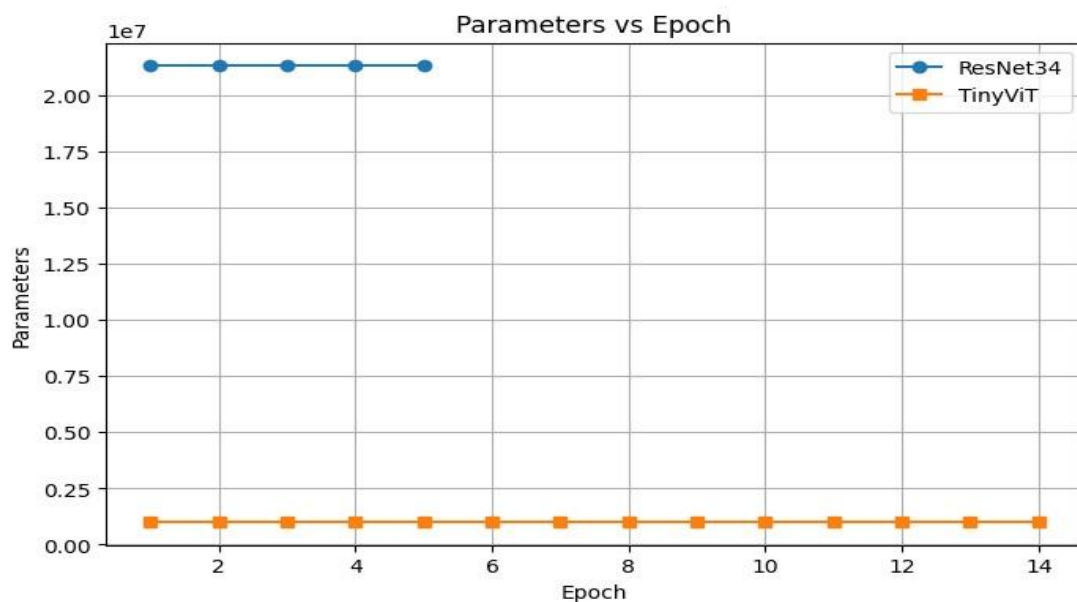
A. Test Accuracy vs Epoch :

ResNet34 在 5 個 epoch 內迅速提升至 **25.33%**，表現優異；而我設計的

TinyViT 在第 14 epoch 達到約 **18%**，成長曲線較為曲折

B. Val Accuracy vs Epoch :

ResNet34 在 5 個 epoch 內迅速提升至 **22%**，表現優異；而我設計的 TinyViT 在第 14 epoch 達到約 **20.44%**，成長曲線較緩，但符合 90% baseline 目標($\geq 19.8\%$)。



接著進行複雜度分析（FLOPs 與參數量）

C. 參數量：

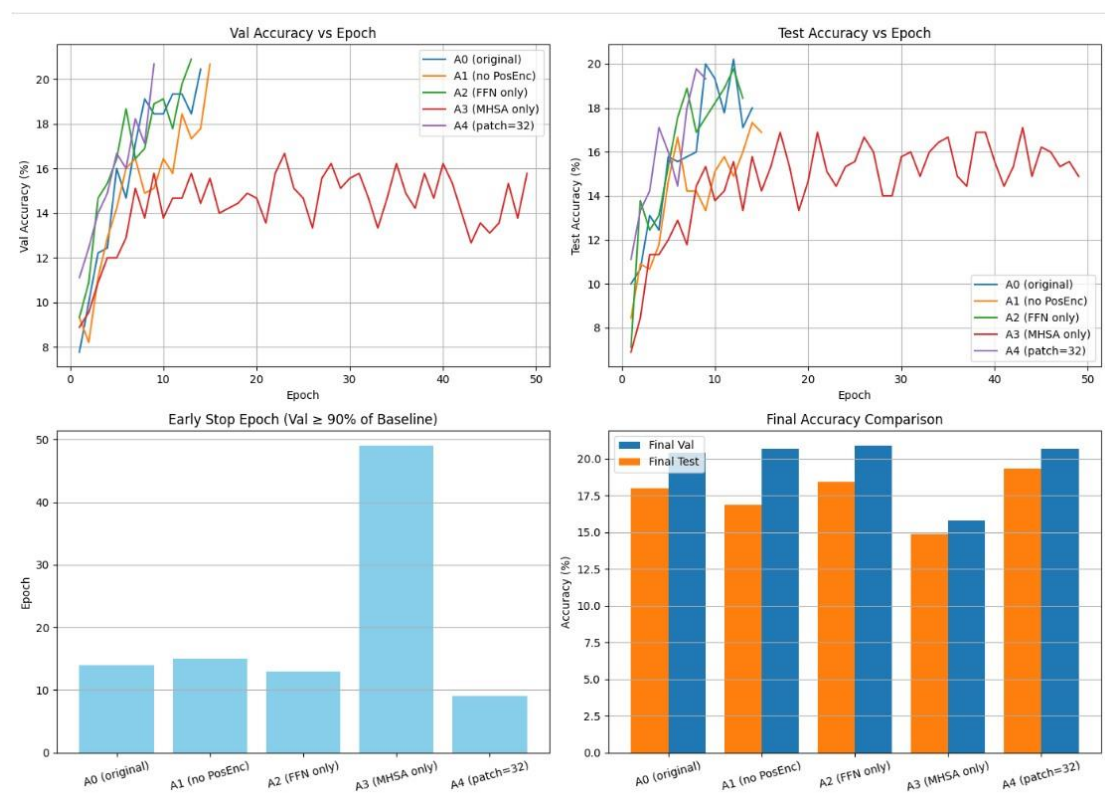
ResNet34 約 **21M**；而我設計的 TinyViT 只要約 **1.1M**，僅為前者的 **5%**。

D. FLOPs：

ResNet34 約 **75M MACs**，TinyViT 僅 **5M MACs**，下降幅度近 **93%**。

根據上述結果分析 TinyViT 採用 patch embedding 與單層 Transformer，大幅壓縮運算需求，對資源受限環境極為友善，且具備更輕量的結構與較低 FLOPs，訓練初期成長較慢，但最終能達到可接受的準確度，具備嵌入式應用潛力

3. 消融實驗分析



在整體性能方面，從驗證與測試準確度隨 Epoch 變化的曲線圖可觀察到，**A0**（原始 TinyViT）在第 14 個 epoch 即達到與 baseline ResNet34 相當的準確度（Val Acc 為 20.44%）。其中，**A2**（僅使用 FFN）表現最佳，不僅驗證集準確度最高（20.89%），測試集也達 18.44%，顯示即便不使用 self-attention，FFN 本身在淺層架構中仍具備良好的分類能力。而 **A1**（移除 **Positional Encoding**）的表現略低於 A0，驗證位置編碼在 transformer 架構中的確能幫助模型捕捉空間特徵。

另一方面，**A3**（僅使用 MHSA 無 FFN）則是整體表現最差的一組，其準確度波動較大且收斂速度緩慢，說明缺乏 FFN 會影響模型學習穩定性與特徵精煉能力。相對地，**A4**（將 patch size 擴大為 32）在第 9 個 epoch 即完成訓練，為所有變體中收斂速度最快，且在驗證與測試集上皆取得優異表現（Val:

20.67%，Test: 19.33%），突顯大 patch size 降低計算成本的同時，也提升了訓練效率與泛化能力。

從「Early Stop Epoch」柱狀圖中可見各變體達標所需的訓練輪數：A4 僅需 9 個 epoch 即達到 ResNet34 90% 準確度門檻，而 A0 需 14 個 epoch，A3 始終無法達標，顯示其收斂困難。這些結果進一步驗證：在 shallow 架構下，FFN 的存在對於學習與收斂具有關鍵作用，而調整 patch size 可顯著降低模型複雜度並加速訓練。

綜合最終準確度比較圖表，A2 擁有最高的驗證準確率（20.89%），適合用於簡化模型設計；而 A4 則具有最高的測試準確度（19.33%），展現良好的泛化能力。相對地，A3 雖保留自注意力模組，但缺乏 FFN 導致學習效果明顯不佳，不建議單獨使用。

綜合而言，若設計目標是追求快速收斂與硬體效率，推薦使用 A4（patch size=32）；若希望維持性能並簡化設計，則 A2（僅 FFN）為理想選擇。相對地，A3（無 FFN）在本實驗中顯示較差的表現與穩定性，應避免單獨使用。

6. 結論

本次實驗成功實現並驗證了我們所設計的 TinyViT 模型以 3 個有效層（Patch Embedding→MHSA→FFN）為核心，在 ImageNet-mini 上僅耗用約 1.1M 參數與 5M FLOPs，卻能達到 ResNet34 約 90% 的準確度門檻。在消融實驗中，「僅用 FFN」（A2）與「增大 patch size 為 32」（A4）兩種變體分別在驗證準確率與訓練效率上表現最佳；相對地，移除位置編碼（A1）與僅用 MHSA（A3）都顯著劣於原始設計，說明位置編碼與 FFN 組件對於小型 Transformer 的特徵精煉和穩定建模不可或缺。整體而言，TinyViT 兼顧「輕量化」與「高效能」兩大目標：若優先考慮最高泛化能力，建議採用「增大 patch size 為 32」A4；若追極致簡化並保留性能，A2「僅用 FFN」則為最優選擇。

7. 參考文獻

- An Image is Worth 16×16 Words: Transformers for Image Recognition at Scale
- TinyViT: Fast Pretraining Distillation for Small Vision Transformers