

## 融合注意力機制與視覺特徵之視網膜病變自動分級研究

黃彥皓\*、陳漢成

### 摘要

本研究對視網膜病變影像進行多類別分類，並採用結合卷積神經網路 (CNN) 局部特徵提取能力與 Transformer 架構全域建模能力，提升模型對病灶區域的辨識能力。此外進一步生成病變區域的可視化熱力圖，提高深度學習模型的可解釋性。訓練結果顯示，模型於第 35 個訓練週期(epoch)達到最佳表現，驗證資料及準確率為 46.67%，Precision 為 37.22%，Recall 為 36.51%，F1-score 為 36.20%，顯示模型具備穩定的學習能力。在可視化結果中觀察到模型確實聚焦於影像中具病理特徵的區域。整體而言，本研究展示 CvT 模型在視網膜影像分析應用中的潛力，並透過可視化技術增強深度學習模型在醫學診斷中的解釋性。

**關鍵字：** Grad-CAM、Conformer、糖尿病視網膜病變

# Automated Retinal Disease Grading Based on the Integration of Attention Mechanisms and Visual Features

Yan-Hao Huang \*, Han-Cheng Chen

## Abstract

In this Research, a multi-class classification framework for retinal lesion images is proposed by leveraging Convolutional Neural Networks (CNNs) and Transformer architectures. The model integrates the local feature extraction capabilities of CNNs with the global context modeling power of Transformers. To further improve the interpretability of the deep learning model, class activation maps (heatmaps) were generated to visualize the regions of interest associated with pathological features.

Experimental results demonstrate that the model achieved its optimal performance at the 35th training epoch, attaining a validation accuracy of 46.67%, a precision of 37.22%, a recall of 36.51%, and an F1-score of 36.20%. These metrics indicate a consistent and stable learning process. Furthermore, the visualization results clearly show that the model accurately attends to clinically significant regions in the retinal images.

Overall, the findings highlight the potential of the Convolutional vision Transformer (CvT) in retinal image analysis, offering not only robust classification performance but also enhanced interpretability.

**Keywords:** Grad-CAM、Conformer、Diabetic Retinopathy

---

Yan-Hao Huang (Corresponding Author), Assistant Professor, Department of Green Energy and Information Technology National Taitung University. E-mail: yhhuang@nttu.edu.tw

Han-Cheng Chen, Student, Department of Green Energy and Information Technology National Taitung University. E-mail: 11302202@gm.nttu.edu.tw

## 壹、前言

隨著深度學習技術的迅速發展，人工智慧在醫學影像分析中的應用日益廣泛，特別是在眼科領域中，透過影像資料進行視網膜疾病的早期檢測與分類已成為輔助診斷的重要方向(Daanouni et al., 2021)。視網膜病變如糖尿病視網膜病變(Diabetic Retinopathy, DR)與黃斑部病變(Macular Degeneration, AMD)等疾病，若能及早發現並處理，可有效降低視力惡化甚至失明的風險。其中糖尿病視網膜病變是糖尿病患者最常見的視覺併發症之一，病灶可能包括微血管瘤、視網膜出血、硬性滲出物與新生血管增生等異常，初期症狀不明顯，卻可能隨著時間而導致永久視力損失，故及時檢測與分級診斷至關重要。

然而傳統診斷方法倚賴眼科醫師透過專業知識與臨床經驗逐張判斷眼底影像，不僅耗時費力，也容易因醫師間的主觀差異導致診斷不一致的問題。為此發展一套高效、準確且可解釋的自動化視網膜病變診斷系統，成為當前智慧醫療研究的重點目標之一。

在醫學影像領域中，機器學習(Machine Learning)技術的導入已為診斷流程帶來顯著突破。從傳統的支持向量機(SVM)、隨機森林(Random Forest)等經典演算法，到現今以卷積神經網路(CNN)為主體的深度學習架構，皆廣泛應用於影像分類、分割與偵測任務中(Daanouni et al., 2021; Gulshan et al., 2016)。這些模型能透過大量標註資料進行學習，自動提取影像中的病灶特徵，已成功應用於胸腔 X 光、腦部 MRI、乳房攝影與視網膜 OCT 等多種影像模態的診斷分析中。

結合卷積神經網路與自注意力機制(Self-Attention)的混合式架構模型 Convolutional vision Transformer (CvT) (Wu et al., 2021)，在多項電腦視覺任務中展現出優異表現。該模型保留了 CNN 在局部特徵擷取上的優勢，並融合 Transformer 的全域建模能力，使其在處理高解析度、病灶分佈不均的醫學影像時，能更全面掌握影像中的語意關聯與結構資訊，有效提升分類精度與穩健性。在醫學影像研究中，已有學者將類似架構應用於肺部 CT、乳房 X 光與病理切片等領域，並取得良好成效。

然而深度學習模型普遍存在「黑箱作業」的問題，即模型在做出分類決策時缺乏可解釋性。這種現象在醫學應用中特別敏感，因為臨床醫師往往需要明確理解模型的預測依據，才能將其納入實際診療流程。因此如何提升模型可解釋性(Interpretability)與信任度，已成為人工智慧醫療應用的重要研究方向之一(Tjoa & Guan, 2020)。

為了解決上述問題，本研究引入 Gradient-weighted Class Activation Mapping(Grad-CAM)技術(Selvaraju et al., 2017)，針對視網膜病變分類任務進行模型可視化分析。透過生成熱力圖，我們得以觀察模型在進行分類預測時所關注的區域，進一步判斷其是否與實際病灶特徵相符，增強模型在臨床上的可採用性與說服力。

本研究以公開資料集 RetinaMNIST 作為實驗基礎，建立基於 CvT 架構的 Conformer 模型，進行視網膜病變影像之多類別分類任務。為提升模型的透明性與臨床價值，本研究同時整合 Grad-CAM 可視化技術進行模型解釋，期望建構一套兼具準確性與可解釋性的 AI 輔助診斷系統，為未來智慧眼科與遠距診療提供技術基礎與應用參考。

## 貳、文獻探討

近年來深度學習技術在醫學影像分析領域的迅速發展，人工智慧系統已廣泛應用於視網膜疾病的自動分類與診斷。視網膜疾病如糖尿病視網膜病變 (Diabetic Retinopathy, DR) 和黃斑部病變 (Macular Degeneration, AMD) 等，若能及早偵測，對視力保存具有關鍵意義。傳統的卷積神經網路 (Convolutional Neural Networks, CNNs) 已被廣泛應用於眼底影像與 OCT (Optical Coherence Tomography) 影像中病灶的辨識，顯示出相當高的準確性 (Gulshan et al., 2016)。然而 CNN 雖擅長局部特徵擷取，但在全域結構建模上仍有所不足。

為克服此限制，近期研究引入 Transformer 架構處理醫學影像，並提出如 Vision Transformer (ViT) 與 Convolutional vision Transformer (CvT) 等模型，強化模型對長距離特徵關聯的理解能力 (Wu et al., 2021)。CvT 模型透過結合 CNN 的局部感受野與 Transformer 的全域注意力，提升模型對病灶結構的辨識能力與分類準確度。此外研究 (Daanouni et al., 2021) 提出將 Self-Attention 機制應用於糖尿病視網膜病變之診斷架構，結合 MobileNet CNN 與注意力模組，並透過 Grad-CAM 對模型進行可視化解釋，提升醫療影像的判讀可信度與臨床實用性。

除了用於分類與辨識任務外，Convolutional Vision Transformer (CvT) 也逐漸被應用於其他醫療影像分析任務，如影像配準。研究中提出一種結合 CvT 與 CNN 的新穎影像配準框架 CvTMorph (Chen et al., 2024)，專門用於建構 4D 醫學影像中的呼吸運動模型。此架構針對現有影像配準方法難以有效提取局部特徵的問題，設計了一個混合式深度學習模型，利用 CvT 的全域建模能力與 CNN 的局部特徵擷取能力，同時整合縮放與平方層以提升微分同胚 (diffeomorphic) 變換的穩定性。該研究於 4D-Lung 與 DIR-Lab 資料集上進行評估，並與多種傳統與深度學習影像配準方法進行比較，實驗結果顯示 CvTMorph 在準確度與結構保留性方面皆優於現有技術，展示 CvT 架構於醫學影像配準任務中的潛力。

這項工作為 CvT 在醫療影像領域的應用開拓了新方向，也間接驗證了該架構在其他任務 (如視網膜影像分類) 的泛化能力與靈活性，為本研究採用 CvT 架構處理視網膜病變分類任務提供有力支持。

另一方面，深度學習模型在醫學領域的應用仍面臨「黑箱問題」，即缺乏可解釋性，限制其臨床推廣。為了解釋模型判斷依據，Selvaraju 等人 (2017) 提出 Gradient-weighted Class Activation Mapping (Grad-CAM) 技術 (Selvaraju et al., 2017)，藉由計算輸出類別對特定卷積層特徵圖的梯度，生成關注區域的熱力圖。該方法已被證實可有效呈現模型判斷依據，增強模型透明度與臨床信任度 (Tjoa & Guan, 2020)。

綜合上述結合 CvT 模型與 Grad-CAM 可視化技術應用於視網膜疾病分類，不僅可提升模型對病灶區域的辨識能力，也可透過可視化結果輔助醫師理解模型決策依據，提升人工智慧在智慧醫療中的可用性與實用性。

### 參、研究方法

為提升模型對全域特徵的建模能力，Transformer 架構在許多研究中被引入，並衍生出多種視覺應用變體 (Dosovitskiy et al., 2020)。其中，Convolutional Vision Transformer (CvT) 模型結合 CNN 的局部感受野特性與 Transformer 的全域關聯建模能力，成為 Conformer 架構的代表之一，在圖像分類與醫學影像任務中展現出優異的性能(Wu et al., 2021)。此類混合架構不僅強化模型對複雜特徵的辨識能力，也提升在小型資料集上的泛化表現。

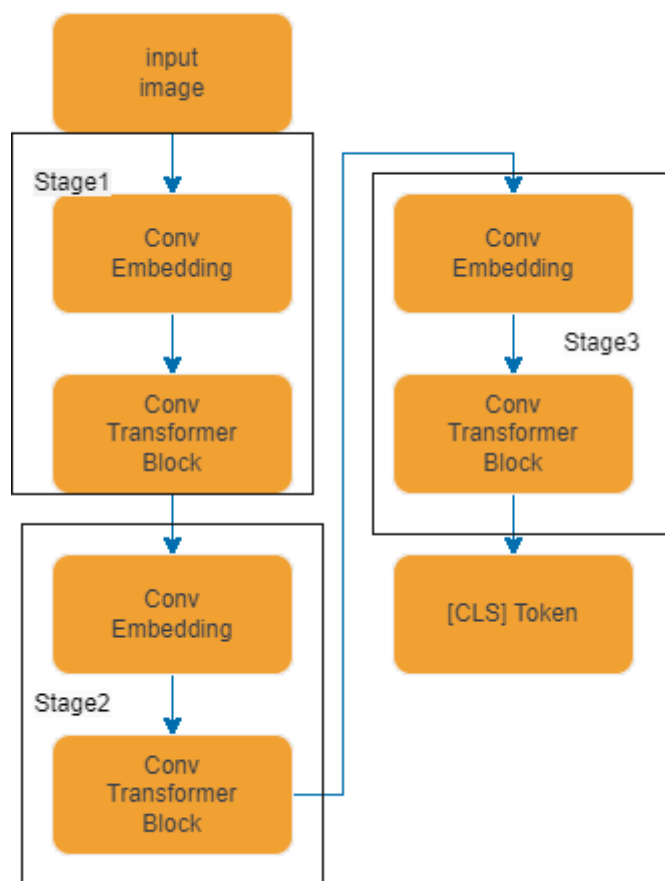


圖 1 CvT 流程圖

CvT 的整體架構可分為三個主要階段，Stage 1 是低層次特徵提取，提取基本圖像邊緣、紋理等局部特徵。Stage 2：中層次特徵表示，包含 Transformer Block 數量：2~6 個（根據模型規模），學習圖像中更複雜的區域關係和結構特徵。Stage 3：高層次語義建模，Transformer Block 數量：10~20 個（依 CvT 模型的大小），用來建構語義豐富

的全局特徵，並在此階段進行分類任務。

每個階段包含以下兩個部分：

1. 卷積式標記嵌入 (Convolutional Token Embedding)：
  - 將輸入影像或前一階段的特徵圖進行重塑為 2D 空間結構。
  - 使用具有重疊區域的卷積操作，提取局部空間特徵。
  - 應用層正規化 (Layer Normalization) 處理。
2. 卷積式變壓器區塊 (Convolutional Transformer Block)：
  - 在查詢 (Q)、鍵 (K)、值 (V) 嵌入中引入深度可分離卷積操作，稱為卷積式投影。
  - 取代傳統 ViT 中的位置線性投影，提升模型的局部感知能力。
  - 僅在最後一個階段加入分類標記 ([CLS] token)，進行最終分類。

這種設計使得 CvT 能夠在保留 Transformer 全局建模能力的同時，增強對局部特徵的捕捉，提升了模型在圖像分類等任務中的表現。

模型可解釋性分析：

為強化模型的臨床可解釋性，本研究採用 Grad-CAM (Gradient-weighted Class Activation Mapping) 作為可視化工具，有助於捕捉影像中具判別性的結構特徵。本研究進一步分析模型於糖尿病視網膜病變 (Diabetic Retinopathy) 分類任務中所關注的視覺區域，並將產生之熱力圖與實際病灶位置進行對比，以評估模型是否能正確聚焦於相關病理特徵。

Gradient-weighted Class Activation Mapping (Grad-CAM) 是一種用於深度卷積神經網路的可視化技術 (Selvaraju et al., 2017)，其主要目的是強化模型在圖像分類或辨識任務中的可解釋性。Grad-CAM 可產生一個「熱力圖」，用來表示模型在做出某一類別預測時，關注於輸入影像中的哪些區域。

Grad-CAM 的核心原理是透過反向傳播，取得模型對某一目標類別的預測分數對於特定卷積層特徵圖的梯度，進而量化各通道 (channel) 在該類別預測中的重要性，這些梯度反映了輸出對於每一個特徵圖的貢獻程度，因此可視為每個特徵圖對該預測類別的「權重」。

## 肆、實驗結果與討論

### 一、資料集

RetinaMNIST 是 MedMNIST 資料集中的一部分，主要包含來自眼底攝影 (Fundus Camera) 的視網膜影像，標註了不同程度的視網膜病變。該資料集常用於開發與驗證視網膜疾病的分類模型，如糖尿病視網膜病變的分級任務。由於視網膜病變常表現為局部細小病灶 (如微血管瘤、硬性滲出) 與整體視網膜色調改變 (如黃斑水腫或視神經盤變化)，因此需要一種能同時兼顧「局部細節」與「全域結構」理解的深度學習模型 (Yang et al., 2023)。RetinaMNIST 共包含 1600 張視網膜影像，分為五類：

No DR (沒有糖尿病視網膜病變的跡象)、Mild DR (輕度非增殖型糖尿病視網膜病變)、Moderate DR (中度非增殖型糖尿病視網膜病變)、Severe DR (嚴重非增殖型糖尿病視網膜病變) 及 PDR (增殖型糖尿病視網膜病變)。

## 二、評估指標

為了全面評估所提出模型的效能，本研究採用了多項常見的分類指標進行分析。首先，透過訓練與驗證準確率 (accuracy) 衡量模型整體的分類正確程度，以評估模型在不同資料階段的學習表現。另一方面為了更深入探討模型在各類別間的預測能力，亦納入精確率 (precision)、召回率 (recall) 與 F1 分數 (F1-score) 進行補充說明。精確率反映模型對正樣本預測的準確性，召回率則表示模型能成功找出所有正樣本的能力，而 F1 分數則為精確率與召回率的加權調和平均，可作為在類別分布不均時更穩健的整體指標。透過這些評估指標，本研究得以更完整掌握模型在視網膜病變分類任務中的效能與泛化能力。

## 三、Grad-CAM 可視化判讀結果

本研究以 RetinaMNIST 資料集為基礎，訓練 Conformer-CvT 模型進行糖尿病視網膜病變之多類別分類。實驗過程共訓練 80 個 epoch。

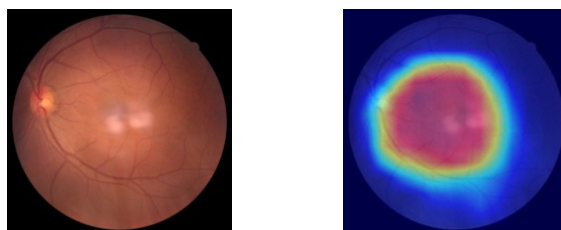


圖 2 分類 No DR 原圖與熱力圖

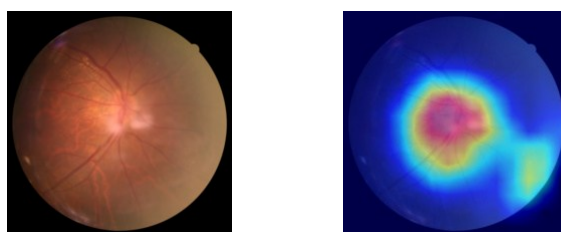


圖 3 分類 Mild DR 原圖與熱力圖

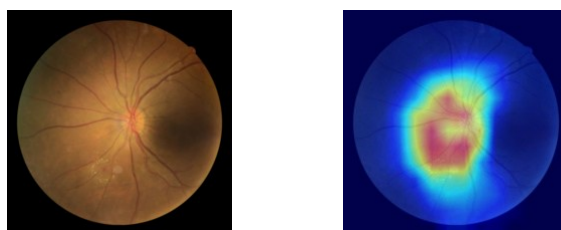


圖 4 分類 Moderate DR 原圖與熱力圖

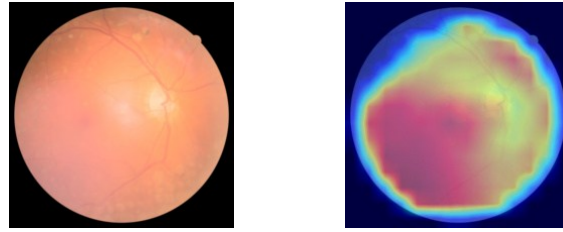


圖 5 分類 Severe DR 原圖與熱力圖

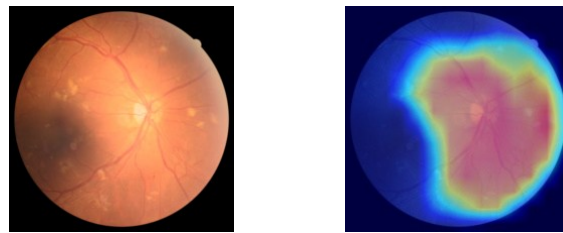


圖 6 分類 PDR 原圖與熱力圖

#### 四、討論

模型在訓練集上表現穩定，從第 21 個 epoch 起，訓練準確率逐步上升，最終於第 69 epoch 達到 100% 的訓練準確率，顯示模型已充分學習資料中的特徵模式。然而在驗證集上的表現相對較不穩定，訓練時整體驗證準確率約介於 42.50% 至 56.67% 之間，最高達到 56.67%，顯示模型可能存在過擬合 (overfitting) 的情形。

進一步觀察精確率、召回率與 F1-score，可以發現這些指標普遍落在 30.20% 至 38.17% 之間，表現與驗證準確率一致，未顯著提升。儘管模型在訓練階段學習良好，但對於驗證資料的泛化能力仍有待加強，可能受限於資料集規模、類別不平衡或部分類別特徵重疊等因素。

本研究採用 RetinaMNIST 資料集進行視網膜病變影像之多類別分類，資料集中共分為五個類別，分別標示為類別 0 至類別 4。如圖 2 所示，類別 0 為正常視網膜影像，未顯示明顯病理特徵；如圖 3 所示，類別 1 表示輕度視網膜病變，可能出現初步微血管瘤；如圖 4 所示，類別 2 為中度病變，影像中可觀察到局部滲出物與血斑；如圖 5 所示，類別 3 屬於嚴重病變，常伴隨明顯視網膜出血或大範圍微血管異常；如圖 6 所示，類別 4 為最嚴重的增殖性糖尿病視網膜病變 (Proliferative DR)，特徵包括新生血管生成與視網膜結構破壞。由於類別間在影像表現上呈現連續性變化，使得模型在辨別中度與嚴重病變時，容易發生類別混淆的情況。

在模型訓練完成後，本研究進一步應用 Grad-CAM (Gradient-weighted Class Activation Mapping) 技術進行特徵可視化分析，以探討模型在判斷過程中關注的區域

是否與實際病灶一致。可視化實作中，我們選定 Conformer 模型中屬於中階表徵層的 stage2.patch\_embed.proj 作為目標層，該層具有良好的空間感受能力，適合觀察中層語意特徵的活化情形。

由生成的 Grad-CAM 熱力圖結果可觀察到，在正常視網膜中(圖 2)，模型的注意力分布較為平均，並未集中於特定病灶區域。而在輕度與中度病變影像中(圖 3 與圖 4)，模型會逐漸聚焦於視網膜中央或靠近視神經盤的局部異常區域，如微血管瘤或細微出血點，顯示模型已學習到判斷病灶的重要位置。對於嚴重與增殖性病變(圖 5 與圖 6)，熱力圖則展現出更明顯且強烈的高亮活化區，聚焦於大面積滲出或新生血管密集分布區域，與臨床病理特徵具高度一致性。

然而在部分案例中，Grad-CAM 結果也顯示模型有時會錯誤聚焦於非病灶區域，或因病灶過小、邊界模糊導致模型辨識困難。尤其是在類別 2 與類別 3 的交界影像中，由於影像特徵的模糊過渡，可能導致模型誤判為相鄰等級，反映出模型對於中度與嚴重病變的邊界學習仍有進步空間。此外個別影像雜訊、亮度不均或資料前處理不一致也可能是導致分類不準確與注意力偏移的潛在因素。

整體而言，Grad-CAM 可視化結果不僅有助於理解模型的分類依據，也能作為臨床決策的輔助工具。透過視覺化呈現模型注意力集中區域，我們得以確認模型確實學習到了與醫學病理特徵高度相關的資訊，證實了 CvT 架構在視網膜影像分類任務中的可行性與應用潛力。同時這也提示後續研究可進一步針對難以區分的類別進行區塊重加權、資料擴增或多尺度特徵強化等機制，以提升整體模型的辨識準確度與穩定性。

## 伍、結論

本研究提出了一種基於 Convolutional Vision Transformer(CvT)的深度學習模型，用於糖尿病視網膜病變(Diabetic Retinopathy, DR)的自動化分類與分級。通過使用 RetinaMNIST 資料集，並應用了多種資料增強技術(如 Mixup、CutMix、AutoAugment)與正則化技術(如 Label Smoothing)，我們成功訓練了一個高效且具有良好泛化能力的分類模型。為了提高模型的可解釋性，我們採用了 Grad-CAM 技術來視覺化模型對不同病變級別的關注區域，並進一步應用了 aug\_smooth 與 eigen\_smooth 進行平滑處理，改善可視化結果的穩定性。

實驗結果顯示，該模型在糖尿病視網膜病變分類任務中表現出色，成功識別不同病變級別，並展示了較強的臨床可解釋性。透過可視化熱力圖的對比分析，模型能夠有效地聚焦於影像中的病灶區域，為臨床診斷提供支持。實驗表明 CvT 結構能夠有效率的分級視網膜病變並在提高診斷準確率，並且可視化熱力圖也說明了該模型能夠聚焦在視網膜的主要特徵點上。

本研究也存在一定的局限性。例如資料集規模相對較小，未來可以嘗試更大規模的資料集以進一步驗證模型的穩定性與泛化能力。此外未來可以探索將其他可解釋性技術與本研究中的 Grad-CAM 結合，進一步提升模型的臨床可解釋性和透明度。

## 引用文獻

- Chen, P., Zou, X., & Gou, Y. (2024). CvTMorph: Improving Local Feature Extraction in Medical Image Registration for Respiratory Motion Modeling with Convolutional Vision Transformer. *Current Medical Imaging*, 20(1), e15734056302592.
- Daanouni, O., Cherradi, B., & Tmiri, A. (2021). Self-attention mechanism for diabetic retinopathy detection. Paper presented at the Emerging Trends in ICT for Sustainable Development: The Proceedings of NICE2020 International Conference, 79–88.
- Gulshan, V., Peng, L., Coram, M., Stumpe, M. C., Wu, D., Narayanaswamy, A., Venugopalan, S., Widner, K., Madams, T., & Cuadros, J. (2016). Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama*, 316(22), 2402–2410.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. Paper presented at the Proceedings of the IEEE International Conference on Computer Vision, 618–626.
- Tjoa, E., & Guan, C. (2020). A survey on explainable artificial intelligence (xai): Toward medical xai. *IEEE Transactions on Neural Networks and Learning Systems*, 32(11), 4793–4813.
- Wu, H., Xiao, B., Codella, N., Liu, M., Dai, X., Yuan, L., & Zhang, L. (2021). Cvt: Introducing convolutions to vision transformers. Paper presented at the Proceedings of the IEEE/CVF International Conference on Computer Vision, 22–31.
- Yang, J., Shi, R., Wei, D., Liu, Z., Zhao, L., Ke, B., Pfister, H., & Ni, B. (2023). Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Scientific Data*, 10(1), 41.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., & Gelly, S. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv Preprint arXiv:2010.11929*,