

國立陽明交通大學百川學士學位學程

112 學年度專題探索二

實作 Detr-物體辨識模型

學生：孫奇霆

指導老師：帥宏翰

113.6.25

目錄

壹、 摘要.....	1
貳、 研究動機.....	1
參、 文獻探討.....	2
(一) 傳統的物件辨識方法.....	2
(二) DETR 的創新思路.....	2
肆、 研究目標.....	3
伍、 研究工具.....	4
陸、 實作細節.....	4
(一) 實現 DETR 模型.....	4
(二) 資料集與前處理(dataloader).....	5
(三) Loss Function.....	5
(四) 優化器(optimizer).....	6
柒、 實驗結果與分析.....	6
(一) 改變超參數.....	6
(二) backbone 可視化.....	7
(三) 訓練過程視覺化.....	8
捌、 結論與討論.....	9
(一) DETR 模型的優勢與局限性.....	9
(二) 性能優化策略.....	9
(三) 未來改進方向.....	10
參考文獻.....	11
附錄.....	12

壹、摘要

本專題實作了基於 Transformer 架構的 DETR (Detection Transformer) 3D 物體辨識模型，並利用 COCO 數據集進行訓練與測試。研究首先使用 PyTorch 實現了 DETR 的基本架構，包括 ResNet backbone、Transformer 編碼器和解碼器、以及雙邊匹配損失函數。隨後，通過調整模型的超參數，如改變 backbone 從 ResNet50 到 ResNet18、減少 Transformer 層數等，比較了 DETR 在不同配置下的性能表現。實驗結果顯示，在有限的訓練週期內，較小的模型反而能獲得更好的效果。

本專題還探討了位置編碼對模型的影響，發現即使移除位置編碼，模型仍能保持不錯的表現。此外，通過分析不同尺寸物體的檢測效果，發現 DETR 在大物體檢測上表現優異，而對小物體的檢測相對較弱。為了深入理解模型的工作原理，將 backbone 提取的局部特徵和 Transformer 的注意力機制進行了可視化，揭示了模型如何逐步改進預測結果。

貳、研究動機

物體檢測是計算機視覺領域的核心任務之一，長期以來 R-CNN 系列為代表的兩階段檢測器是主流。然而，這類方法存在明顯的局限性：首先，它們依賴於複雜的後處理步驟，如非極大值抑制（NMS），增加了模型的複雜度；其次，由於採用順序處理方式，並行性較差，影響了訓練和推理的效率。

2020 年，Facebook AI 研究團隊提出的 DETR (Detection Transformer) 模型為物體檢測任務帶來了革命性的變化。DETR 巧妙地將 Transformer 架構應用於物體檢測，利用其強大的自注意力機制來捕捉圖像中不同位置間的相互關係。這種方法不僅簡化了檢測流程，還大大提高了模型的並行性能。

DETR 的出現引發了一個重要問題：它是否真的如論文所述，能夠在不需要大量超參數調整的情況下，達到或超越傳統 R-CNN 模型的性能？這個問題直接關係到 DETR 的實用性和適應能力。

因此，本研究的主要動機在於深入探討 DETR 模型的性能特點和適應性。我們希望通過調整超參數，如 backbone 結構、Transformer 層數等，來驗證 DETR 是否確實具備良好的適應能力，還是僅在特定參數配置下才能展現其優越性。

此外，我們也對 DETR 的內部工作機制深感興趣。通過可視化分析，我們希望能夠揭示 DETR 如何利用 Transformer 架構來理解和處理圖像信息，特別是在不同尺寸物體的檢測上的表現差異。這些洞察不僅有助於我們更好地理解和優化 DETR 模型，還可能為未來物體檢測算法的設計提供新的思路和方向。

參、文獻探討

(一)傳統的物件辨識方法

最簡單的物體檢測流程是：(1)輸入圖像，用一個滑動的窗口擷取一小塊圖片；(2)將各區域傳遞給卷積神經網絡[4](CNN)，並將區域分類；(3)一旦將每個區域劃分為相對應的類別後，就可以組合這些區域，來檢測原始圖像。(4)重複 1~3 直到每個位置、每種大小的窗口都嘗試完畢。

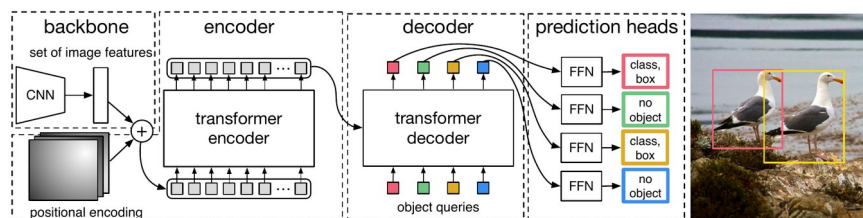
這種方法的優點是直接使用分類模型就可以做到物件辨識；但這種方法有兩大缺點，首先這花費了大多數的時間在不必要的圖片上，因此如何將圖片有效率切分就是一項挑戰。再者，一個物體可能會被不同的窗口多次檢測到，如何融合這些數據又是一項挑戰。

為了解決窗口切分的效率問題，R-CNN[3](RNN+CNN)引入了區域提議 (Region Proposals)，先找出有可能的候選位置，然後對候選區域進行分類。後來 Fast-RCNN[2]與 Faster RCNN 也是類似的作法，只是速度上加速非常多。但無論是哪一種方法，都沒有解決時間浪費在不必要之圖片的問題，且 RCNN 的並行性極差，導致訓練和推論效率低下。

隨著 Transformer[1]的出現，由於良好的並行性，很大程度上取代了 RNN 這種遞迴機制，針對物體檢測的 DETR[6](Detection Transformer)也應運而生。

(二)DETR 的創新思路

DETR 使用了 CNN 作為 backbone 作為特徵提取器，嵌入 positional encoding 後交給 transformer encoder 處理。為了生成邊界框，DETR 使用 100 個 object queries，每一個 query 會生成一組類別與邊界框預測。其中有些預測為「無類別」，代表該邊界框中沒有物體。只要畫出有類別的邊界框就得到了成功辨識的模型。



DETR 架構圖(來自[6])

DETR 融合了 CNN 與 Transformer 的優勢，前者善於從像素點中分離出不同的局部特徵，而後者善於從大量的局部訊息中找到全局特徵。DETR 不需要區域

提議，物體也不會同時被不同的窗口多次檢測到，一次解決傳統 RCNN 的兩大痛點。以下是 DETR[6]模型的一些主要特點和組件：

1. End-to-End：DETR 是一個端到端訓練的模型，它不依賴於傳統的區域提議網絡（RPN），因此和非極大值抑制（NMS）等手段，省下了調整許多參數的麻煩，因而簡化了物體檢測流程。
2. Transformer 架構：該架構最初被設計用於自然語言處理（NLP）任務，但其注意力機制在圖片辨識上同樣取得巨大成功。Transformer 幫助模型理解圖像中的上下文信息，並處理物體之間的關係。
3. Bipartite Matching Loss：DETR 引入了一種名為雙邊匹配損失（Bipartite Matching Loss）的新型損失函數，用於在訓練過程中匹配預測邊界框和真實邊界框。這種方法從訓練的角度解決了同一個物體重複檢測的問題。
4. 物體和背景的區分：由於 decoder 一次的輸出為固定長度，因此 DETR 使用一種稱為「無物體類別」（"no object" class）來區分圖像中的物體和背景，這有助於模型更準確地進行物體檢測。
5. 全局理解：通過自注意力機制，DETR 能夠一次全局地理解圖像，不像 CNN 需要一層一層傳遞，這使得模型較 CNN 能夠處理一些複雜的場景，尤其在圖片像素量很多時尤為明顯。

DETR 提供了一種簡單且統一的檢測架構。由於其較不需要依賴人類經驗且擁有優異的並行性能，DETR 已成為計算機視覺領域的重要研究對象。

肆、研究目標

本研究具體目標如下：

1. 利用 PyTorch 框架實現 DETR 模型的完整架構，包括 ResNet backbone、Transformer 編碼器和解碼器、以及雙邊匹配損失函數等核心組件。藉由這個過程，深入理解 DETR 的工作原理和實現細節。
2. 系統地調整和優化 DETR 模型的超參數，包括但不限於更改 backbone 網路（如從 ResNet50 到 ResNet18）、調整 Transformer 的層數和隱藏單元數量、探究位置編碼的影響等。透過這些實驗，評估 DETR 模型的適應性和性能穩定性。
3. 使用 COCO 資料集進行訓練和測試，對比不同配置下 DETR 模型的性能表現。特別關注模型在不同尺寸物體（小、中、大）檢測上的表現差異。
4. 開發和應用視覺化工具分析 DETR 模型的內部工作機制。具體包括：

5. 視覺化 ResNet backbone 提取的局部特徵，了解模型在特徵提取階段的表現。
6. 透過視覺化訓練過程中的預測結果，觀察模型如何逐步改進其檢測和分類能力。
7. 基於實驗結果和分析，提出 DETR 模型的優化建議和潛在的改進方向，為未來的研究和應用提供參考。

伍、研究工具

研究工具和資源的組合，將使我們能夠全面實現研究目標，深入探討 DETR 模型的性能和特性，並為未來的改進提供可靠的實驗基礎。本研究主要採用以下工具和資源來實現目標：

- COCO 資料集：作為本研究的核心資料來源，COCO (Common Objects in Context) 是一個大型開源圖片資料集，廣泛用於物體檢測、分割和關鍵點檢測等電腦視覺任務。COCO 資料集包含超過 33 萬張影像，其中 20 萬張已經過標註，涵蓋 80 個物體類別和 150 萬個物體實例。這個資料集能夠為我們的 DETR 模型提供全面的訓練和測試基礎。
- PyTorch 框架：使用 PyTorch 作為主要的深度學習框架，用於實現 DETR 模型。PyTorch 豐富的神經網絡模組，使得複雜模型的構建和實驗變得更加靈活和高效。
- COCO API：利用 COCO 官方提供的 API 來加載、解析和視覺化 COCO 資料集中的標註資訊。這個 API 支援對物體實例、關鍵點和圖像標題等多種標註類型的處理，大大簡化了資料預處理的工作。
- ResNet 預訓練模型：使用在 ImageNet 上預訓練的 ResNet 模型作為 DETR 的 backbone，以提取圖像的初始特徵。我們將探索 ResNet50 和 ResNet18 等不同版本對模型性能的影響。
- 評估指標：採用 COCO 評估指標，如平均精確度 (AP) 和不同物體尺寸的 AP (AP_S 、 AP_M 、 AP_L) 等，來評估模型性能。

陸、實作細節

(一)實現 DETR 模型

借用 resnet50 作為特徵提取的主幹(backbone)，resnet 最後的池化層與線性投射層是作為結果輸出之用途，因此將其排除。接著使用了 1x1 的卷積核從 resnet50 的 2048 維投射到 512 維。backbone 的結果再嵌入位置資訊(position

encoding)後，作為 encoder 的輸入。長度為 100、可學習的 query position 將作為 decoder 輸入並產生對應的 100 個獨立的預測。這些預測再經過 linear_class 與 linear_bbox 得到 91 個類別的置信度及邊界框的位置(cxcywh)，並將邊界框用 sigmoid 轉換到(0, 1)區間，以避免超過圖像大小。程式碼請參閱 *DETR 簡易實作*。

(二)資料集與前處理(dataloader)

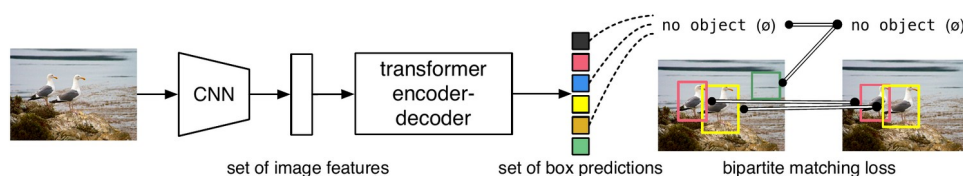
dataloader 會至少回傳以下三種資料（以 batch size=4 為例）：

1. image：統一縮放成 512x512。shape(4, 512, 512)。
2. id：遵守 COCO 原始 ID，但有少許 ID 沒有出現過。shape(4, 100)。
3. bbox：以 cxcywh 格式儲存，範圍為[0, 1]。shape(4, 100, 4)。

為了程式的並行性，需要將一個 batch 的資料儲存在 tensor 裡，但 tensor 不允許每筆資料的物體數量不一，因此必須填充(padding)「無類別」至 100。規定無類別的 ID=0，cxcywh=(1, 1, 1, 1)。

(三)Loss Function

在 DETR 的實作中，Loss Function 是非常重要的部分。它負責衡量模型預測結果和真實標註之間的差異，以此來指導模型的學習。由於 DETR 會固定輸出 100 個預測，然而其順序不一定符合 dataloader 預設的順序，因此需要找出最佳的匹配方式，使預測和 ground truth 形成一一對應。最佳匹配如下圖所示：



匹配示意圖(來自[6])

實作 DETR 的 Loss Function 主要由三部分組成：

1. 計算 Cost Matrix：類別 Cost 和邊界框 Cost 分別通過 _cat_cost 和 _bbox_cost 方法計算，相加後得到 100*100 的 Cost Matrix。

$$-\mathbb{1}_{\{c_i \neq \emptyset\}} \hat{p}_{\sigma(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\sigma(i)}).$$

_cat_cost 與 cross_entropy 過程類似，但缺少了對數運算，原因是這裡想要計算的是「距離」，而非損失。

`_bbox_cost` 則沒有這種區別，距離與損失皆為同一種函數如下：

$$\lambda_{\text{iou}} \mathcal{L}_{\text{iou}}(b_i, \hat{b}_{\sigma(i)}) + \lambda_{\text{L1}} \|b_i - \hat{b}_{\sigma(i)}\|_1$$

需要注意的是，即使是「無類別」也會有 `bbox`，在計算 `bbox cost` 時必須排除目標類別為空類別的匹配。程式碼請參閱 *cost 計算實作*。

2. 建立配對 (Matching)：

使用匈牙利算法來進行預測和真實標籤之間的最佳匹配。這部分在 `_assign` 方法中實現，並使用了 `linear_sum_assignment` 函數，為了最佳化性能，這部分用多線程處理每一個 `sample`。

3. 計算匹配成功的 Loss：

在完成配對之後，計算匹配成功的預測和真實標籤之間的最終損失。這部分包括類別損失和邊界框損失，分別通過 `_cat_loss` 和 `_bbox_loss` 方法計算。與 `bbox cost` 相同，必須排除目標類別為空類別的匹配。

(四)優化器(optimizer)

根據論文，`resnet` 作為 `backbone` 最好要避免變動，因此設定 `backbone` 的學習率為 `1e-5`，其餘部分為 `1e-4`，使用 `AdamW` 作為優化器，前 `100epoch` 不需太多調整即可有不錯的梯度下降效果（左圖），另外設定 `scheduler`，使學習率在 `200epoch` 後減少 10 倍，在後期有顯著的梯度下降效果（右圖）。

柒、實驗結果與分析

(一)改變超參數

根據原始論文，`DETR` 應該擁有良好的適應性，不容易出現完全不可用的參數，加大模型也不容易出現性能飽和，因此改變模型的參數，觀察訓練出來的模型有什麼變化？這裡以 `Average Precision(AP)` 作為判斷根據。調整的方向有以下三點：

1. 改變 `backbone` 從 `resnet50` 至 `resnet18`。
2. 減半原始論文中 `hidden_dim` `nheads` `num_encoder_layers` `num_decoder_layers` 的值。
3. 拿掉 `positional encoding`

經過以上調整的模型在 200 個 `epoch` 後的表現結果如下表：

調整代號	Params	AP	AP50	AP _S	AP _M	AP _L
原始	31316512	0.09	0.21	0.01	0.05	0.18
原始(400epoch)	31316512	0.19	0.36	0.04	0.15	0.39
(1)	12962656	0.26	0.48	0.07	0.23	0.53
(2)	25392480	0.25	0.48	0.07	0.22	0.51
(1)+(2)	12962656	0.3	0.53	0.08	0.28	0.6
(1)+(2)+(3)	12959456	0.2	0.41	0.04	0.15	0.44

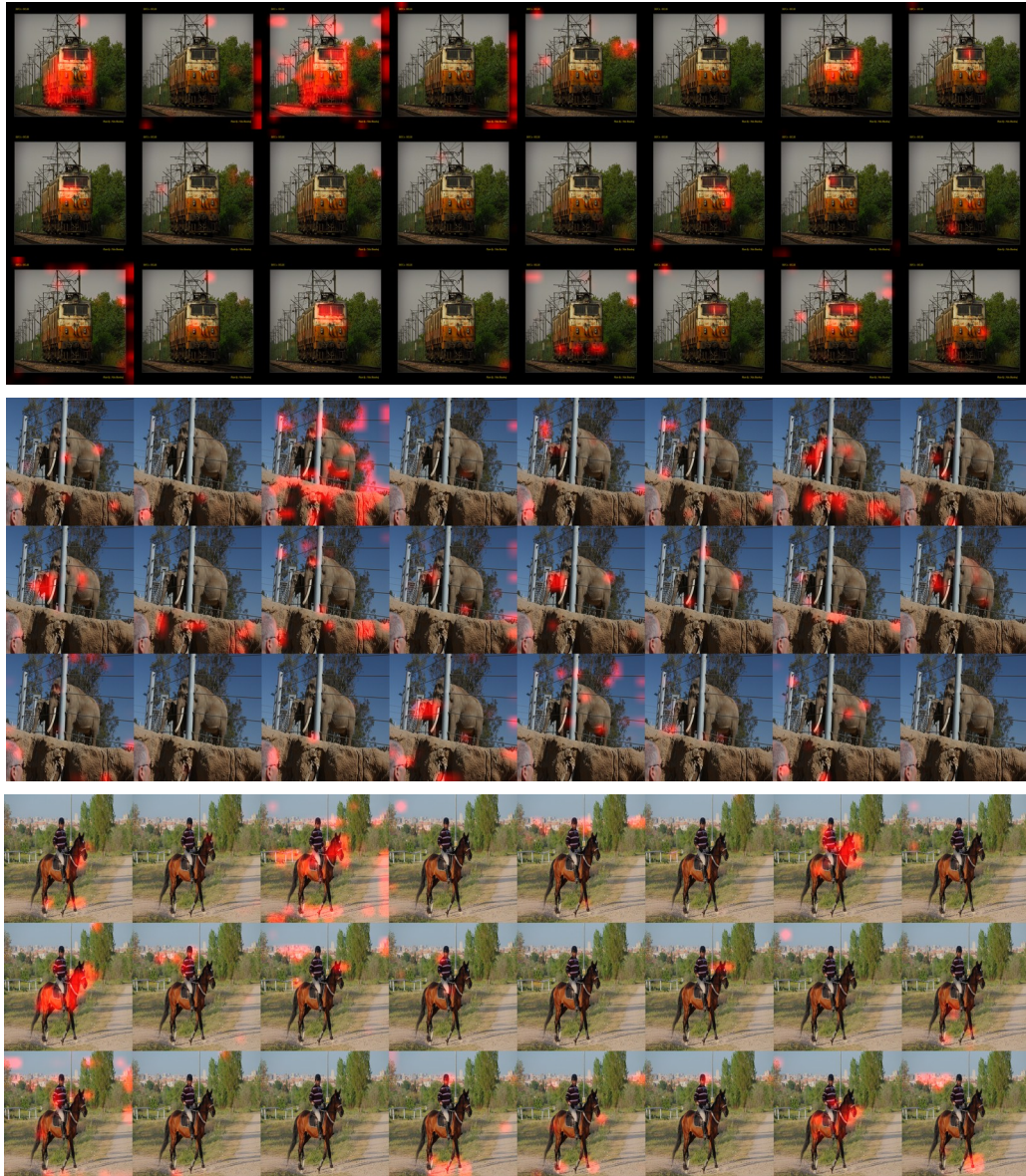
可以看出模型越大效果不一定越好，推測是因為模型參數較多，需要更多迭代次數與訓練資源才能收斂到足夠好的結果，從上表經過 400 個 epoch 訓練後的結果可以發現，模型確實有很大的進步空間。在只有 200 個 epoch 的限制下，使用 resnet18，並設定 hidden_dim=64, nhead=4, num_encoder_layers=3, num_decoder_layers=3 可以得到最好的結果，若訓練資源充足，增加模型大小並不會導致瓶頸。

值得注意的是，即使拿掉 positional encoding 仍然可以得到不差的結果，可能的原因有兩種，首先 dataloader 已經將每一張圖片轉成 512*512，拉直後的順序一致，因此經過一段時間的訓練，transformer 可以演化出類似 positional encoding 的模式。另一種推測是，將圖片拉直後，會出現週期性的數據，transformer 可以根據上下文的週期性推測出這張圖片原始的長寬，進而發展出類似 positional encoding 的模式。

從表中可以發現數據按照 AP_S/AP_M/AP_L 的順序遞增，這意味著越大的物體精準度越高，這與原始論文[6]的結果非常相似，推測是因為相較於以往的 R-CNN 模型，transformer 藉助自注意力機制，更善於發現大範圍的模式，反而對於只出現幾個像素的小物體不太敏感。

(二)backbone 可視化

512*512 的三通道彩色圖片經過 backbone(resnet)輸出為 32*32 的 512 通道熱點圖，這些熱點圖將是 transformer 用來生成類別預測及邊界框的根據。取前 24 通道進行可視化，將 32*32 重新放大為 512*512 並用紅色表示一個通道的熱點，形成下圖。可以發現有些通道用來探測圓弧，有些善於發現對稱等等。這個時候的 resnet 已經具備區分不同物體的功能。



resnet 對三張圖片的熱點圖

(三)訓練過程視覺化

下圖是模型的預測結果，與正確答案匹配後，隱藏 98 個未匹配的結果，其中一個框辨識出物體，但是類別錯誤，另一個框尚未辨識出物體，但預測框的位置最接近正確答案，因此被成功匹配，loss function 會設法讓這兩個預測框更接近他們所分配到的正確答案，同時讓剩下 98 個預測靠近「無類別」。



訓練中的模型預測及匹配結果

捌、結論與討論

(一) DETR 模型的優勢與局限性

DETR 模型作為物體檢測領域的創新方法，展現出了多方面的優勢。首先，其端到端的訓練方式顯著簡化了物體檢測流程，無需複雜的後處理步驟，大大提高了模型的易用性。其次，基於 Transformer 的架構使得 DETR 具有高度的並行性能，這不僅提升了訓練效率，也加快了推理速度。此外，DETR 的自注意力機制能夠有效捕捉整個圖像的上下文資訊，使得模型在理解複雜場景時表現出色。最後，DETR 的靈活性使其易於擴展到其他相關任務，如全景分割，展現出廣闊的應用前景。

然而，DETR 模型也存在一些局限性。最顯著的問題是其較長的訓練時間，相比傳統方法，DETR 通常需要更多的訓練週期才能達到收斂。同時，自注意力機制的計算複雜度隨輸入大小呈平方增長，這在處理高解析度圖像時可能造成效率瓶頸。另一個值得注意的問題是 DETR 在小物體檢測上的表現不如大物體，這可能限制其在某些特定應用場景的效果。最後，DETR 模型通常需要大量的訓練資料才能充分發揮其潛力，這在資料有限的情況下可能會成為一個挑戰。

(二) 性能優化策略

針對 DETR 模型的優化，可以從多個角度進行改進。首先，可以考慮改進 backbone 網路，使用更強大或更適合的特徵提取方法，如 Swin Transformer，以提供更優質的初始特徵。其次，增強位置編碼的效果也是一個重要方向，例如採

用可變形 DETR 中的可變形注意力機制，可以更好地捕捉物體的空間資訊。為了提高對小物體的檢測能力，可以考慮多尺度特徵融合的策略，結合不同層級的特徵圖來增強模型的表達能力。

在訓練過程中，優化損失函數設計可以加速收斂並提高準確率。同時，採用更多樣化的資料增強技術能夠有效提升模型的泛化能力。對於計算資源有限的情況，可以考慮使用知識蒸餾技術，將大型模型的知識轉移到小型模型中，在保持性能的同時減少計算開銷。這些優化策略的綜合應用，有望顯著提升 DETR 模型的整體性能。

(三)未來改進方向

展望 DETR 模型的未來發展，有多個潛在的改進方向值得探索。首要任務是開發計算效率更高的 DETR 變體，以減少訓練時間和推理延遲，使其更適合實際應用場景。同時，改進 DETR 在小樣本學習場景下的表現也很重要，這可以減少對大規模標註資料的依賴，擴大模型的應用範圍。

另一個有前景的方向是跨模態整合，結合圖像和文本等多模態資訊，可能會大幅提高檢測的準確性和靈活性。探索動態調整模型結構的方法，使 DETR 能夠根據不同複雜度的場景自適應調整，也是一個有趣的研究方向。此外，深入研究 DETR 的決策過程，提高模型的可解釋性，對於增進我們對模型行為的理解至關重要。

参考文献

1. VASWANI, Ashish, et al. Attention is all you need. Advances in neural information processing systems, 2017, 30.
2. GIRSHICK, Ross. Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. 2015. p. 1440-1448.
3. GIRSHICK, Ross, et al. Rich feature hierarchies for accurate object detection and semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. 2014. p. 580-587.
4. LECUN, Yann, et al. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998, 86.11: 2278-2324.
5. COCO Dataset. [online]. 2017. COCO Consortium. Available from: <http://cocodataset.org> [Accessed 12 April 2024].
6. CARION, Nicolas, et al. End-to-end object detection with transformers. In: European conference on computer vision. Cham: Springer International Publishing, 2020. p. 213-229.

附錄

1. DETR 簡易實作

```
class DETR(nn.Module):
    def __init__(self, num_classes, hidden_dim, nheads,
                  num_encoder_layers, num_decoder_layers):
        super().__init__()
        self.backbone = nn.Sequential(*list(resnet50(weights=ResNet50_Weights.DEFAULT).children())[:-2])
        self.conv = nn.Conv2d(2048, hidden_dim, 1)
        self.transformer = nn.Transformer(hidden_dim, nheads,
                                           num_encoder_layers, num_decoder_layers)
        self.linear_class = nn.Linear(hidden_dim, num_classes + 1)
        self.linear_bbox = nn.Linear(hidden_dim, 4)
        self.query_pos = nn.Parameter(torch.rand(100, hidden_dim))
        self.row_embed = nn.Parameter(torch.rand(50, hidden_dim // 2))
        self.col_embed = nn.Parameter(torch.rand(50, hidden_dim // 2))

    def forward(self, inputs):
        x = self.backbone(inputs)
        h = self.conv(x)
        H, W = h.shape[-2:]
        pos = torch.cat([
            self.col_embed[:W].unsqueeze(0).repeat(H, 1, 1),
            self.row_embed[:H].unsqueeze(1).repeat(1, W, 1),
        ], dim=-1).flatten(0, 1).unsqueeze(1)
        h = self.transformer(pos + h.flatten(2).permute(2, 0, 1),
                             self.query_pos.unsqueeze(1))
        return self.linear_class(h), self.linear_bbox(h).sigmoid()
```

2. cost 計算實作

```
def _cat_cost(self, targ_cat: torch.Tensor, pred_cat: torch.Tensor):
    pred_cat = pred_cat.softmax(dim=-1)
    cat_cost = -torch.stack([pred_cat[i, :, targ_cat[i]] for i in range(pred_cat.shape[0])]).mT
    return cat_cost

def _bbox_cost(self, targ_bbox, pred_bbox, targ_cat):
    pred_bbox = pred_bbox.unsqueeze(1).expand(-1, 100, -1, -1)
```

```
targ_bbox = targ_bbox.unsqueeze(2).expand(-1, -1, 100, -1)
cost_bbox = bboxLoss(pred_bbox, targ_bbox)
mask = (targ_cat != 0).unsqueeze(-1).expand(-1, -1, 100)
cost_bbox *= mask
return cost_bbox
```