

整合全域幾何先驗：在關鍵點估計網路中應用單應性損失之可行性分析與實施框架

第一部分：幾何特徵偵測的基礎典範

在深入探討將單應性(Homography)約束整合至損失函數此一高階研究問題之前，必須先為任務本身建立一個堅實的理論與實踐基礎。此基礎的核心在於正確地建構問題，並選擇最適合的深度學習典範。本部分將深入論證，為何對於如羽球場交點這類抽象幾何特徵，關鍵點估計(Keypoint Estimation)是遠比標準物件偵測(Object Detection)更為優越的方法，並詳細解析我們將以此為基礎進行修改的YOLOv8-Pose框架。

第一節：抽象幾何特徵偵測的獨特挑戰

1.1 標準物件偵測的不足之處

標準的物件偵測框架，如YOLO系列，其設計初衷與成功基石在於處理具有豐富視覺特徵的實體物件，例如球員、球拍或車輛¹。這些模型的深度卷積神經網路(CNN)擅長從影像中學習階層式的特徵，包括顏色(如球衣顏色)、紋理(如布料或皮膚)以及複雜的形狀與部件組合(如人體輪廓)¹。然而，當偵測目標從具體物件轉變為抽象的幾何圖形，如羽球場地線的交點時，此類方法的根本性設計瓶頸便暴露無遺。

場地線交點(如L型、T型、十字型)在視覺上是極度稀疏的特徵。它們缺乏可供模型學習的內在視覺屬性，其可辨識性幾乎完全由周圍白色或黃色線條的空間排列所定義¹。此一特性直接衝擊了標準物件偵測方法的核心優化目標——交並比(Intersection over Union, IoU)。

IoU透過計算預測邊界框(Bounding Box)與真實邊界框之間的面積重疊度來衡量定位的準確性²。對於一個具有一定體積的常規物件，這是一個合理且有效的指標。然而，當應用於羽球場交點這類幾乎是點狀的特徵時，IoU的物理意義便開始瓦解。一份針對此問題的深入分析報告指出，一個理想的交點真實邊界框可能僅有幾個像素的大小¹。在這種極端情況下，預測框即便只有微不足道的幾個像素偏移，也可能導致IoU值從一個接近完美的值(例如0.9)災難性地驟降至0¹。這種現象揭示了一個根本性的問題：IoU-based損失函數(如CIoU、GIoU)為此任務創造了一個極其不利於學習的損失地貌(Loss Landscape)。這個地貌可以被想像成一片廣闊的平原(代表IoU為0的區域)，中間僅存在一個極其微小且陡峭的深坑(代表高IoU的區域)。在訓練過程中，梯度

訊號會變得極其「陡峭」和嘈雜，模型收到的回饋是一種不穩定的「全有或全無」訊號，使其預測可能在「幾乎完美」和「完全錯誤」兩個極端狀態之間劇烈震盪¹。這不僅阻礙了模型進行平滑的、亞像素級的微調，更從根本上使學習過程變得低效且難以收斂到最佳的定位點。

1.2 典範轉移：以關鍵點估計作為幾何原生的解決方案

為了解決上述困境，我們必須進行一次典範轉移：將問題的表述從「框選一個物件」回歸到其幾何本質——「定位一個點」。關鍵點估計 (Keypoint Estimation) 正是實現此轉移的理想典範¹。此方法不再試圖用一個「框」去包圍一個「點」，而是直接將任務建構為對關鍵點 \$(x,y)\$ 座標的直接回歸。YOLOv8框架透過其姿態估計 (Pose Estimation) 模式，原生支援此類任務¹。這種方法的優越性體現在其損失函數的設計上。關鍵點估計的損失函數不再基於面積重疊，而是基於距離，例如簡單的L1/L2距離，或是更為先進的「物件關鍵點相似度」(Object Keypoint Similarity, OKS)¹。採用這種基於距離的損失函數，從根本上解決了IoU所面臨的所有問題。相較於IoU的陡峭深坑，距離損失的損失地貌可以被想像成一個平滑的碗狀盆地。無論預測點與真實點相距多遠，損失函數總能提供一個平滑、連續且有意義的梯度訊號¹。一個偏離5像素的預測會比偏離2像素的預測受到更大的懲罰，而模型也能清晰地知道需要朝哪個方向進行修正以減小損失。在這樣的地貌中，優化器可以輕鬆地沿著梯度方向穩定下降，最終收斂到盆地的最低點，即亞像素級別的精確座標。這極大地提升了任務的「可學習性」(learnability)，並完美地統一了任務的最終目標(最小化座標誤差)和模型的優化過程。

下表總結了兩種典範在處理幾何特徵偵測任務時的根本差異，並闡明了為何本報告選擇關鍵點估計作為後續研究的基礎。

表1: 幾何特徵偵測的典範比較：物件偵測 vs. 關鍵點估計

特徵	物件偵測 (IoU-based)	關鍵點估計 (Distance-based)	對於場地交點偵測的意涵
核心任務	定位並分類一個有面積的「物件」	直接回歸一個或多個「點」的座標	將交點視為「點」更符合其幾何本質 ¹
主要度量	面積重疊度 (IoU)	歐幾里得距離 (如L1, OKS)	距離度量更能直接反映定位的精確度
損失地貌	不連續、陡峭的「深坑」	平滑、連續的「碗狀」盆地	關鍵點估計的優化過程更穩定、更易收斂 ¹
梯度訊號品質	嘈雜、不穩定、「全有或全無」	平滑、有意義、具方向性	模型能從梯度中獲得清晰的修正方向，有利於學習
亞像素精度適用性	極差，微小偏移導致IoU驟降至零 ¹	極佳，為微小偏移提供平滑的梯度訊號	關鍵點估計是實現高精度定位的必要條件

第二節：YOLOv8-Pose框架：高精度定位的機制

在確立了關鍵點估計作為正確的解決典範後，我們選擇了Ultralytics開發的YOLOv8框架中的姿態估計模式(YOLOv8-Pose)作為實現此典範的技術基礎。其高度模組化的設計和卓越的性能，使其成為我們後續進行損失函數客製化的理想平台¹。

2.1 架構與預測流程

YOLOv8-Pose模型並非一個全新的架構，而是在成熟的YOLOv8物件偵測器基礎上進行的擴展。其核心架構由三個部分組成：

1. 骨幹網路 (**Backbone**)：通常採用基於CSPDarknet的設計，負責從輸入影像中提取一系列由淺到深、由粗到細的階層式特徵圖(feature maps)¹。
2. 頸部網路 (**Neck**)：採用如特徵金字塔網路(FPN)和路徑聚合網路(PANet)相結合的結構，其作用是融合來自骨幹網路不同層級的特徵圖，將高層次的語義資訊(有助於分類)與低層次的空間細節資訊(有助於定位)結合起來¹。
3. 預測頭 (**Head**)：這是YOLOv8-Pose與標準YOLOv8最關鍵的區別所在。YOLOv8-Pose採用了一個「雙頭」並行設計。在頸部網路輸出的每一個尺度特徵圖上，都並行連接了兩個預測頭：
 - 偵測頭 (**Detect Head**)：與標準YOLOv8的偵測頭完全相同，負責預測物件的類別(class)和邊界框(bounding box)。
 - 姿態頭 (**Pose Head**)：專為姿態估計新增的預測頭，負責預測與每個偵測到的物件相關聯的關鍵點座標(keypoints)¹。

這種並行設計意味著模型在一次前向傳播中，能夠同時完成「這是一個什麼類型的交點？」、「它大概在影像的哪個區域？」以及「它精確的幾何中心座標在哪裡？」這三個子任務。最終，模型的輸出是一個結構化的列表，其中每個元素都代表一個被偵測到的交點，並包含其類別標籤、一個邊界框，以及最關鍵的一個高精度的\$(x,y)\$關鍵點座標¹。

2.2 標準複合損失函數的解構

為了後續能在此基礎上進行擴展，我們必須深入理解YOLOv8-Pose的標準損失函數。這個總損失是一個由多個部分加權構成的複合損失，它協同工作，共同引導模型向著期望的目標優化。其形式可表示為¹：

$$L_{total} = w_{cls} \cdot L_{cls} + w_{bbox} \cdot L_{bbox} + w_{kpts} \cdot L_{kpts}$$

其中 w 是平衡各部分損失的權重係數。

- 分類損失 (**Lcls**)：負責懲罰模型對交點類別(L型、T型、十字型)的錯誤分類。YOLOv8通常採用帶有Logits的二元交叉熵損失(BCEWithLogitsLoss)，它將多類別分類問題視為多個獨立的二元分類問題，具有良好的靈活性¹。
- 邊界框損失 (**Lbbox**)：負責優化邊界框的定位。YOLOv8採用了如完整交並比損失(Complete IoU, CIoU)或分佈焦點損失(Distribution Focal Loss, DFL)的組合¹。
- 關鍵點損失 (**Lkpts**)：這是我們任務的核心，專門負責懲罰關鍵點座標的預測誤差。YOLOv8-Pose的關鍵點損失是基於物件關鍵點相似度(Object Keypoint Similarity, OKS)來

計算的¹。

2.3 邊界框在關鍵點估計中的共生作用

一個常見的困惑是：「既然基於IoU的損失函數不適用於點狀特徵，為何我們仍需費力標註並讓模型預測邊界框？」這個問題的答案揭示了YOLOv8-Pose架構設計的精妙之處。邊界框在此任務中扮演著兩個不可或缺的共生角色。

首先，也是最關鍵的作用，是作為尺度歸一化器(**Scale Normalizer**)。關鍵點定位的黃金標準OKS，其核心公式為 $OKS = \exp(-d^2 / 2s^2k^2)$ ¹。在這個公式中，

d 是預測點與真實點之間的歐幾里得距離。如果沒有歸一化，一個5像素的預測誤差對於遠景中只有 10×10 像素的交點來說是致命的，但對於近景特寫中 100×100 像素的交點來說則可能是個微不足道的誤差。變數 s (物體尺度)正是為了解決這個問題而引入的。在實踐中，由於分割掩碼(segmentation mask)標註成本高昂，通常使用邊界框的面積(即 $width * height$)作為 s^2 的一個高效且合理的近似¹。因此，通過在OKS計算中用邊界框面積來歸一化距離

d ，損失函數變得具有尺度不變性。模型受到的懲罰不再是絕對的像素誤差，而是相對於交點自身大小的相對誤差。邊界框的存在，是實現這種高級、自適應優化機制的基礎。

其次，邊界框扮演著**隱式注意力機制(Implicit Attention Mechanism)**的角色。模型並非直接在整張圖上暴力搜索所有可能的關鍵點，而是採用了一種更為高效的、分而治之的策略。偵測頭首先透過預測邊界框來回答「影像的哪些區域可能包含一個交點？」，將注意力從無關的背景轉移到包含場地線的候選區域。然後，姿態頭就在這個高度相關的局部區域內，集中精力解決「在這個小區域內，精確的交點中心在哪裡？」的問題¹。這種「先找框，再找點」的流程，實際上是將一個複雜的全局搜索問題分解為多個簡單的局部搜索問題，極大地提高了模型的學習效率和最終的定位精度。

下表詳細解構了標準YOLOv8-Pose損失函數的各個組成部分及其在本任務中的具體作用。

表2: 標準YOLOv8-Pose損失函數解構

損失組件	公式/概念	在場地交點偵測中的目的	關鍵參數
Lcls (分類損失)	二元交叉熵 (BCE)	準確地區分交點類型 (L-junction, T-junction, cross-junction)。	類別權重
Lbox (邊界框損失)	CIoU / DFL	1. 學習預測一個能準確反映交點尺度的邊界框。2. 作為隱式注意力機制，縮小關鍵點搜索範圍。	IoU閾值, DFL分佈範圍
Lkpts (關鍵點損失)	物件關鍵點相似度 (OKS)	實現亞像素級的精確交點定位。OKS損失利用 $L_{\{box\}}$ 提供的尺度資訊進行歸一化，實現尺度不變的優化。	物體尺度 s (來自邊界框), 關鍵點常數 k

第二部分：單應性作為全域幾何約束

在建立了以關鍵點估計為核心的技術基礎後，我們現在轉向使用者查詢的核心問題：將單應性(Homography)整合到損失計算中。本部分將首先闡述單應性的理論基礎，然後基於此理論，提出一個具體的、可計算的單應性一致性損失項 $L_{homography}$ 。

第三節：二維單應性的理論基礎

3.1 定義與性質

單應性是一種描述兩個平面在不同視角下投影關係的幾何變換³。在電腦視覺中，它通常由一個 3×3 的矩陣

H 來表示，這個矩陣可以將一個影像中的點對應到另一個影像中的對應點⁴。其數學表達式為：
 $p' = Hp$

其中 p 和 p' 分別是兩個影像中對應點的齊次座標(homogeneous coordinates)， \sim 符號表示等式在一個非零尺度因子下成立⁶。

單應性變換具有以下關鍵性質：

- 射影變換 (Projective Transformation): 它是一種射影變換，這意味著它能保持共線性，即原來在同一直線上的點，經過單應性變換後依然在同一直線上⁵。
- 8個自由度 (8 Degrees of Freedom): 儘管 H 矩陣有9個元素，但由於它是定義在一個尺度因子之下的，所以它只有8個獨立的自由度⁵。這意味著，要唯一地確定一個單應性矩陣 H ，至少需要4對不共線的點對應關係。

在我們的任務中，羽球場地本身是一個平面。因此，無論攝影機從哪個角度拍攝，所得到的影像與一個標準的、鳥瞰的2D場地圖之間，都存在一個單應性關係。這個關係捕捉了所有由透視造成的幾何畸變。

3.2 透過直接線性變換(DLT)進行估計

直接線性變換(Direct Linear Transform, DLT)是從點對應關係中計算單應性矩陣 H 的經典演算法⁹。DLT的巧妙之處在於，它將一個看似非線性的問題(由於齊次座標中的投影片除法)轉化為一個可以透過標準線性代數方法求解的齊次線性系統⁶。

其推導過程如下：

1. 對於每一對點對應關係 (p, p') ，其中 $p = [x, y, 1]^T$, $p' = [x', y', 1]^T$ ，它們的單應性

關係 $p' \sim H p$ 可以寫成 $p' \times (H p) = 0$, 其中 \times 是向量叉積⁸。

2. 將這個叉積展開, 並進行代數重排, 可以為 H 矩陣的 9 個未知元素 h_i 導出兩條線性方程式⁶。這兩條方程式可以寫成矩陣形式:

$$A h = 0$$

其中 A_i 是一個 2×9 的矩陣, 其元素僅由點座標 (x, y, x', y') 構成, h 是由 H 矩陣元素拉直而成的 9×1 向量。

3. 當我們有 $N \geq 4$ 對點對應時, 可以將所有 A_i 矩陣垂直堆疊起來, 形成一個 $2N \times 9$ 的大矩陣 A 。整個問題就轉化為求解齊次線性系統:

$$A h = 0$$

4. 這個系統的解 h 存在於矩陣 A 的零空間(null space)中。為了找到在存在測量噪聲時的最佳解(即最小化 $\|Ah\|^2$), 我們使用奇異值分解(Singular Value Decomposition, SVD)⁷。對矩陣

A 進行SVD分解 $A = U \Sigma V^T$, 所求的解 h 正是與最小奇異值對應的右奇異向量, 即 V 矩陣的最後一列⁶。

DLT演算法的優雅之處在於它提供了一個非迭代的、封閉形式的解, 將一個幾何問題轉化為一個經典的線性代數問題。這為我們將其整合到深度學習的損失函數中提供了堅實的數學基礎。

第四節: 單應性一致性損失(Lhomography)的構建

基於單應性的理論, 我們可以設計一個新的損失項 $L_{homography}$, 其核心思想是:懲罰模型預測出的、不符合單一全域平面投影幾何約束的關鍵點集合。如果模型預測出的一組關鍵點是幾何上正確的, 那麼它們必須能夠透過一個唯一的單應性矩陣 H 映射到一個標準的、鳥瞰的場地圖上。

4.1 概念框架

我們的損失函數將衡量一種「自我一致性」。它不直接將模型預測的單應性矩陣 H_{pred} 與一個真實的 H_{gt} 進行比較(因為獲取每張影像的 H_{gt} 是不切實際的)。相反, 它提出一個問題:「模型自己預測出的這組關鍵點, 在多大程度上符合一個嚴格的平面射影幾何約束?」一個低的一致性損失意味著所有預測點的空間排列高度結構化, 彷彿它們確實是從一個平面上投影而來。一個高的一致性損失則意味著這些點的排列是隨機的、雜亂的, 在幾何上是不可能同時出現的。

4.2 數學構建

以下是計算 $L_{homography}$ 的具體步驟:

1. 定義標準點集(Canonical Points): 我們首先需要一個參考基準。根據羽球世界聯合會(

BWF)的官方規則，我們可以創建一個標準化的2D鳥瞰場地圖，並記錄下所有場地交點的精確座標¹。這個點集

$P_{\text{canonical}} = \{p_{c,1}, p_{c,2}, \dots, p_{c,N}\}$ 在整個訓練過程中是固定不變的。

2. 預測影像點集 (**Predicted Points**): 在每一次訓練迭代中，YOLOv8-Pose模型會對輸入的影像進行前向傳播，輸出一組預測的關鍵點座標。我們收集所有被高信賴度偵測到的 N 個交點的座標，形成預測點集 $P_{\text{pred}} = \{p_{p,1}, p_{p,2}, \dots, p_{p,N}\}$ 。
3. 估計單應性矩陣: 利用這 N 對點對應關係 $(P_{\text{canonical}}, P_{\text{pred}})$ ，我們使用前述的 DLT演算法來估計一個單應性矩陣 H_{pred} ，該矩陣能將標準點集映射到預測點集。
4. 計算重投影誤差 (**Reprojection Error**): 得到 H_{pred} 後，我們用它將標準點集 $P_{\text{canonical}}$ 重新投影回影像空間，得到一組新的重投影點集 $P_{\text{reproj}} = H_{\text{pred}} \cdot P_{\text{canonical}}$ 。理想情況下，如果 P_{pred} 本身是完全幾何一致的，那麼 P_{reproj} 應該與 P_{pred} 完全重合。它們之間的差異，即重投影誤差，就構成了我們的損失值。這個概念在自監督學習和幾何驗證中被廣泛使用¹³。

最終的損失函數可以被定義為所有點的預測位置與其重投影位置之間的平均L1距離：

$$L_{\text{homography}} = \frac{1}{N} \sum_i \|p_{\text{pred},i} - p_{\text{reproj},i}\|_1$$

這個損失函數的設計極為巧妙。它利用模型自身的輸出作為監督訊號的一部分，來約束其輸出的內部結構。如果模型預測的點在幾何上不一致（例如，某個點的偏移破壞了整體的透視關係），DLT演算法會找到一個「最佳擬合」的 H_{pred} 。在這種情況下，那個不一致的點經過重投影後的位置 $p_{\text{reproj},i}$ 將會偏離其原始預測位置 $p_{\text{pred},i}$ ，從而產生一個非零的損失值。這個損失值會產生梯度，指導模型調整其預測，使其更符合全域的幾何約束。

第三部分：實現端到端訓練：可微分幾何層

將前述的 $L_{\text{homography}}$ 整合到端到端的訓練流程中，面臨著一個巨大的技術障礙：傳統的幾何估計算法，特別是其核心的SVD步驟，本質上是不可微分的，或在特定條件下梯度不穩定。本部分將深入分析這一障礙，並探討三種可行的解決路徑，以構建一個可微分的幾何層。

第五節：古典幾何估計中的可微分性障礙

5.1 RANSAC與SVD的問題

在傳統的電腦視覺流程中，為了提高對噪聲和離群點的穩健性，通常會在DLT之前使用隨機抽樣一致性(RANSAC)演算法來篩選出內點(inliers)³。然而，RANSAC是一個基於隨機採樣和迭代投票的演算法，其過程是不連續的，因此不可微分，無法直接嵌入到基於梯度下降的深度學習訓練環節中。

即便我們暫不考慮RANSAC，只看DLT演算法本身，其核心的求解步驟 $Ah=0$ 依然存在問題。如

前所述，這個齊次線性系統是透過SVD求解的。這正是可微分性的主要障礙所在。

PyTorch的官方文件 `torch.linalg.svd` 中明確警告：只有當輸入矩陣沒有重複的奇異值時，梯度才是有限的。如果任意兩個奇異值之間的距離接近於零，梯度將會變得數值不穩定¹⁵。這種不穩定性源於梯度計算公式中包含了類似

$\frac{1}{\sigma_i^2 - \sigma_j^2}$ 這樣的項，當奇異值 σ_i 和 σ_j 趨於相等時，分母趨近於零，導致梯度爆炸¹⁵。

這引出了一個深刻的悖論：**SVD反向傳播的不穩定性**，恰恰在模型表現優異時最為嚴重。回顧DLT演算法，當我們有 $N \geq 4$ 個完全精確、無噪聲的共面點對應時，構造出的 $2N \times 9$ 矩陣 A 的零空間維度恰好為1，其最小的奇異值 σ_9 將精確地等於0⁶。在訓練過程中，隨著模型學習的進行，其預測的關鍵點

P_{pred} 會越來越準確，越來越符合共面約束。這將導致矩陣 A 越來越接近秩虧（rank-deficient），其最小的幾個奇異值會不斷趨近於零，並可能彼此非常接近。這意味著，一個成功的、產生精確預測的前向傳播，反而會導致一個失敗的、梯度不穩定的反向傳播。DLT演算法所追求的理想條件（一個清晰的零空間），恰恰是標準反向傳播演算法崩潰的條件。這個根本性的衝突，是將古典幾何演算法融入深度學習時必須克服的核心挑戰。

第六節：通往可微分單應性估計的路徑

為了解決上述的可微分性障礙，學術界和工業界已經探索出多條路徑。我們可以將其歸納為三種主要策略：直接迴歸、使用可微分函式庫，以及實現先進的可微分SVD。

6.1 路徑A：直接迴歸（黑箱方法）

此方法完全繞過了DLT/SVD的求解過程。最具代表性的工作是DeTone等人提出的HomographyNet¹⁸。他們設計了一個類似VGG的標準CNN架構，其輸入是堆疊在一起的兩張影像塊，輸出則直接是一個8維的向量。這個8維向量代表了單應性變換的一種參數化形式，即源影像四個角點相對於目標影像四個角點的位移（4-point parameterization）。

在訓練時，網路的輸出（預測的8維向量）與真實的位移向量之間計算一個簡單的L2損失。在推斷或損失計算的最後一步，這個8維向量可以透過一個固定的、不可學習的DLT轉換層，計算出最終的 3×3 單應性矩陣 H ¹⁸。

分析：這種方法的優點是其端到端的完全可微分性，因為它避免了任何有問題的SVD操作。然而，它的缺點也同樣明顯：它用一個「黑箱」的迴歸網路取代了一個有著堅實幾何基礎、數學上清晰的DLT演算法。這種黑箱方法可能難以學習到幾何問題的內在結構，其泛化能力和最終能達到的精度，可能不如那些明確將幾何約束編碼到模型結構中的方法。

6.2 路徑B：可微分函式庫（實用的工程解決方案）

這條路徑是目前最為實用和直接的解決方案。它依賴於專為深度學習設計的電腦視覺函式庫，這些函式庫提供了常用視覺演算法的可微分實現。Kornia是最著名的一個，它為PyTorch提供

了大量可微分的幾何與影像處理工具²⁰。

Kornia函式庫中包含了一個關鍵函式：

`kornia.geometry.homography.find_homography_dlt(points1, points2)`²⁵。這個函式接收兩組對應點的張量(Tensor)作為輸入，並輸出一個或一批單應性矩陣的張量。

分析：這個函式可以被視為一個標準的、可微分的PyTorch層。在模型的前向傳播過程中，它內部執行了DLT演算法來計算單應性矩陣。更重要的是，Kornia的開發者為這個操作實現了一個客製化的、數值穩定的反向傳播函式(backward function)。這意味著梯度可以安全地流過這個DLT計算層，而不會遇到前面討論的梯度爆炸問題。對於我們的任務而言，這是實現所提出的`$L_{homography}`最直接、最可靠的方法。它完美地結合了DLT演算法的幾何嚴謹性和深度學習模型的可訓練性，讓我們可以站在巨人的肩膀上，而無需從頭解決複雜的SVD反向傳播問題。

6.3 路徑C：先進的可微分SVD(前沿研究領域)

這條路徑深入到問題的數學核心，旨在從根本上使SVD操作本身變得穩健可微。這是目前學術界的一個前沿研究方向。

近期的研究指出，SVD的不可微性本質上源於在存在重複奇異值時，求解梯度所需的線性方程組變成了欠定(underdetermined)系統²⁶。一個創新的解決方案是，利用摩爾-彭若斯偽逆(Moore-Penrose pseudoinverse)來求解這個欠定系統，從而得到一個唯一的、最小範數的最小二乘解。這個解在數學上是良定義的，從而使得整個SVD操作在理論上變得可微分²⁶。

其他的先進方法還包括使用不同的技術來近似SVD的梯度，例如：

- 幕迭代法 (Power Iteration): 利用幕迭代法來近似梯度計算，避免直接求解不穩定的線性系統²⁷。
- 帕德近似 (Padé Approximants): 使用有理函式(帕德近似)來逼近梯度函式，以獲得更平滑和穩定的梯度估計²⁹。

分析：這些方法代表了將幾何視覺完全融入深度學習框架的最前沿探索。它們提供了最為數學嚴謹的解決方案，但從零開始實現的複雜度非常高。在本報告的範疇內，它們的主要價值在於驗證了「可微分DLT」這個概念在理論上是成立且可行的，從而為我們選擇使用Kornia這一工程解決方案提供了堅實的理論背書。

第四部分：分析、綜合與建議

綜合前述所有部分的分析，本部分將提出一個最終的、整合了單應性約束的複合損失函數，對其與其他幾何損失進行深入的比較分析，並為使用者提供一個清晰的戰略性實施路線圖和結論。

第七節：用於全域一致性姿態估計的新型複合損失函數

7.1 最終構建

基於上述討論，我們提出一個全新的複合損失函數 L_{final} ，它將標準的YOLOv8-Pose損失與我們設計的單應性一致性損失相結合：

$$L_{final} = w_{pose} \cdot L_{yolo_pose} + w_{homography} \cdot L_{homography}$$

- **L_{yolo_pose}:** 這是標準的YOLOv8-Pose複合損失，即 $L_{cls} + L_{box} + L_{kpts}$ ，其詳細定義見第二節。這個損失項是模型的基礎，確保模型能夠學習到基本的分類能力和準確的、具備尺度感知的關鍵點定位能力。
- **L_{homography}:** 這是我們在第四節中定義的、基於重投影誤差的單應性一致性損失。在實際實現中，其核心的單應性矩陣估計步驟將透過第六節討論的Kornia可微分DLT層來完成。這個損失項作為一個全域正則化項，強制模型預測出的所有關鍵點在幾何上必須符合單一平面投影的約束。

7.2 損失權重的考量與策略

在上述公式中，超參數 w_{pose} 和 $w_{homography}$ 的設定至關重要。它們平衡了模型的兩個學習目標：一是局部定位的精確性（由 L_{yolo_pose} 驅動），二是全域結構的幾何一致性（由 $L_{homography}$ 驅動）。

一個關鍵的考量是， $L_{homography}$ 只有在模型的關鍵點預測 P_{pred} 已經達到一定準確度時，才能提供有意義的梯度訊號。在訓練初期，當模型的預測還處於隨機和嘈雜的階段時，用這些噪聲點去估計單應性矩陣 H_{pred} 將會得到一個同樣無意義的結果。基於這個無意義的 H_{pred} 計算出的重投影誤差，將會是混亂的噪音，反而可能干擾模型的基礎學習過程。因此，一個合理的訓練策略是採用損失退火（Loss Annealing）。具體而言，我們可以在訓練開始時將 $w_{homography}$ 設為0，讓模型完全專注於透過 L_{yolo_pose} 學習基礎的、粗略的定位能力。隨著訓練的進行，當驗證集上的 L_{yolo_pose} 損失開始下降並趨於穩定時，我們再逐步地、平滑地將 $w_{homography}$ 的值從0增加到其最終的目標值（例如1.0）。這種課程學習（Curriculum Learning）式的策略，允許模型先學會「走」（粗略定位），再學會「跑得優雅」（精確且幾何一致的定位），從而確保了訓練過程的穩定性和最終的收斂效果。

第八節：比較分析：全域 vs. 局部幾何先驗

使用者提出的 $L_{homography}$ 是一種**全域（Global）幾何約束。與此同時，研究文獻中也提出了其他形式的幾何約束，例如在¹中提到的，基於交點角度的局部（Local）**幾何損失 L_{geom} ¹。

L_{geom} 的思想是，對於每個預測的L型或T型交點，計算其內部線條構成的角度，並懲罰其與90度的偏差。為了全面評估 $L_{homography}$ 的效用，有必要將這兩種方法進行深入的比較分析。

$L_{homography}$ 無疑是比 L_{geom} 更為強大和全面的約束。一個滿足全域單應性約束的

點集，其內部的所有局部角度、長度比例、共線性等幾何屬性必然都是正確的。反之則不然，即使所有局部角度都正確，也不能保證整個點集符合一個全域的透視變換。

然而，這兩種約束並非完全的競爭關係，而更應被視為互補關係。它們在計算成本、對噪聲的穩健性以及提供訊號的尺度上各有優劣。

- 計算成本與穩健性： L_{geom} 的計算非常廉價，它只需要在一個小的局部影像塊 (patch) 內偵測幾條線段並計算夾角。它對離群點的穩健性也更強。即使影像中只有一個交點被準確預測，它也能提供有效的梯度訊號。相比之下， $L_{homography}$ 的計算成本更高 (涉及SVD)，並且它要求至少有4個分佈良好的、大致準確的點預測才能估計出一個有意義的單應性矩陣。如果4個點中有一個是嚴重的離群點，估計出的 H_{pred} 將會被完全扭曲。
- 訓練階段的適用性：鑑於上述特性， L_{geom} 更適合在訓練早期使用，因為它能在模型預測還很嘈雜時提供穩定的局部幾何修正訊號。而 $L_{homography}$ 則更適合在訓練中後期，當模型的預測已經比較穩定後，再引入作為一個強大的全域「精修」工具。

一個真正先進的系統，甚至可以考慮將兩者結合，在訓練的不同階段賦予它們不同的權重，以充分利用各自的優勢。

下表對這兩種幾何一致性損失進行了多維度的比較。

表3: 幾何一致性損失的比較分析

屬性	L_{geom} (局部角度約束)	$L_{homography}$ (全域平面約束)
約束類型	局部 (Local)	全域 (Global)
計算成本	低 (局部線段偵測)	中等 (矩陣構建與SVD)
對噪聲/離群點的穩健性	較高，僅受局部預測影響	較低，易受單個離群點影響
強制執行的資訊	局部角度正確性 (例如，90度角)	所有幾何屬性 (角度、長度比、共線性等)
實現複雜度	中等 (需實現局部角度估計器)	高 (需處理可微分SVD，推薦使用Kornia)
資料需求	每個交點的局部影像塊	至少4個跨影像的對應點

第九節：結論與戰略性建議

9.1 核心問題的回答

綜合本報告的全面分析，對於使用者提出的核心問題：「在YOLO-Pose的損失計算中加入 homography 會有用嗎？」，我們的結論是肯定的，並且極具研究與應用價值。

將單應性一致性作為損失項，不僅僅是「有用」，它代表了一種從局部特徵匹配向全域結構理解的範式提升。它為模型訓練引入了一個強大的、基於物理世界幾何規律的先驗知識，能夠強制模型學習到遠比單純的點位迴歸更為豐富和結構化的場景表徵。相較於局部角度約束，單應性約束更為全面和嚴格，有望在模型的定位精度和幾何一致性上帶來質的飛躍。

9.2 推薦實施路徑

我們為使用者規劃了一條清晰、分步的戰略性實施路線圖：

1. 建立基準模型 (**Establish Baseline**): 首先，使用精心標註的羽球場交點數據集，訓練一個標準的YOLOv8-Pose模型。此模型僅使用其內建的標準複合損失 L_{yolo_pose} ¹。在獨立的測試集上，使用正確的評估指標——基於OKS的mAP(特別是 $mAP@[.50:.05:.95]$)和PCK——來嚴格評估其性能，建立一個強有力的性能基準線。
2. 實現可微分單應性層 (**Implement Differentiable Homography Layer**): 整合Kornia函式庫到您的PyTorch專案中。創建一個新的、可微分的PyTorch模組，該模組接收模型預測的關鍵點集 P_{pred} 和預定義的標準點集 $P_{canonical}$ 作為輸入，內部調用 `kornia.geometry.homography.find_homography_dlt` 來計算 H_{pred} ，並最終計算並返回基於重投影誤差的 $L_{homography}$ 。
3. 訓練複合損失模型 (**Train with Composite Loss**): 將新的單應性損失項整合到訓練流程中，形成最終的複合損失 L_{final} 。在訓練過程中，採用第七節中提出的損失退火策略，即逐步增加 $w_{homography}$ 的權重，以確保訓練的穩定性。
4. 進行嚴格評估 (**Evaluate Rigorously**): 使用與第一步完全相同的測試集和評估指標，對新訓練出的模型進行評估。將其性能與基準模型進行量化比較，以科學地驗證引入單應性損失所帶來的性能增益。

9.3 未來展望：從損失函數到實際應用

本研究的意義遠不止於提升模型的評估指標。我們為損失函數設計的機制，本身就為最終的實際應用鋪平了道路。

在第七節中，我們設計了一個可微分的模組，它能從模型預測的關鍵點中即時計算出單應性矩陣 H 。這個模組不僅在訓練時用於計算損失，在模型訓練完成後，它就變成了一個強大的推斷工具。在對比賽視訊進行分析時，我們可以將這個模組直接應用於模型在每一幀上的輸出。模型預測出一組關鍵點，這個模組便能立刻計算出將當前視角「拉平」到標準鳥瞰圖所需的單應性矩陣 H ¹。

利用這個即時計算出的 H ，我們可以：

- 實現全自動的視角校正：將任意角度的比賽畫面實時轉換為統一的2D鳥瞰圖，極大地簡化後續所有基於位置的分析¹。
- 賦能高階戰術分析：在統一的座標系下，精確地量化分析球員的跑動覆蓋、回球落點分佈、反應速度和戰術模式。
- 提升轉播觀看體驗：在電視轉播中實時疊加AR視覺效果，如球員跑動軌跡、擊球速度和落點預測。

這種訓練機制與應用工具之間的優雅協同，是本研究框架最引人注目的特點之一。透過將全域幾何約束直接編碼到學習目標中，我們不僅訓練出一個更精確的模型，更在訓練過程中直接構建了實現下一代智慧體育分析系統所需的關鍵技術模組。

引用的著作

1. 羽球場交點YOLO偵測與分類_.pdf
2. A Review of Homography Estimation: Advances and Challenges - MDPI, 檢索日期 : 7月 31, 2025, <https://www.mdpi.com/2079-9292/12/24/4977>
3. What is Homography? How to estimate homography between two images? - GeeksforGeeks, 檢索日期 : 7月 31, 2025, <https://www.geeksforgeeks.org/computer-vision/what-is-homography-how-to-estimate-homography-between-two-images/>
4. Homography examples using OpenCV (Python / C ++) | - LearnOpenCV, 檢索日期 : 7月 31, 2025, <https://learnopencv.com/homography-examples-using-opencv-python-c/>
5. The Ultimate Guide to Homography - Number Analytics, 檢索日期 : 7月 31, 2025, <https://www.numberanalytics.com/blog/ultimate-guide-to-homography>
6. Homography Estimation*, 檢索日期 : 7月 31, 2025, https://cseweb.ucsd.edu/classes/wi07/cse252a/homography_estimation/homography_estimation.pdf
7. 41 Homographies - Foundations of Computer Vision, 檢索日期 : 7月 31, 2025, <https://visionbook.mit.edu/homography.html>
8. How do you find the homography matrix given 4 points in both images? - AI Stack Exchange, 檢索日期 : 7月 31, 2025, <https://ai.stackexchange.com/questions/21042/how-do-you-find-the-homography-matrix-given-4-points-in-both-images>
9. Direct linear transformation - Wikipedia, 檢索日期 : 7月 31, 2025, https://en.wikipedia.org/wiki/Direct_linear_transformation
10. Estimating the Homography Matrix with the Direct Linear Transform (DLT) - Medium, 檢索日期 : 7月 31, 2025, <https://medium.com/@insight-in-plain-sight/estimating-the-homography-matrix-with-the-direct-linear-transform-dlt-ec6bbb82ee2b>
11. Direct Linear Transform - Carnegie Mellon University, 檢索日期 : 7月 31, 2025, https://www.cs.cmu.edu/~16385/s17/Slides/10.2_2D_Alignment_DLT.pdf
12. Image homographies - Introduction and course overview, 檢索日期 : 7月 31, 2025, <https://www.cs.cmu.edu/~16385/lectures/lecture9.pdf>
13. Self-Supervised Camera Pose Estimation With Geometric Consistency 1 Introduction - Anirudh Chakravarthy, 檢索日期 : 7月 31, 2025, <https://anirudh-chakravarthy.github.io/files/Monodepth+RelPose.pdf>
14. Deep Homography Estimation for Visual Place Recognition - arXiv, 檢索日期 : 7月 31, 2025, <https://arxiv.org/html/2402.16086v1>
15. torch.linalg.svd — PyTorch 2.7 documentation, 檢索日期 : 7月 31, 2025, <https://pytorch.org/docs/stable/generated/torch.linalg.svd.html>
16. torch.svd — PyTorch 2.7 documentation, 檢索日期 : 7月 31, 2025, <https://pytorch.org/docs/stable/generated/torch.svd.html>
17. Matrix Backpropagation for Deep Networks With Structured Layers - CVF Open Access, 檢索日期 : 7月 31, 2025, https://openaccess.thecvf.com/content_iccv_2015/papers/Ionescu_Matrix_Backpr

[opagation_for_ICCV_2015_paper.pdf](#)

18. (PDF) Deep Image Homography Estimation - ResearchGate, 檢索日期:7月 31, 2025,
https://www.researchgate.net/publication/305881252_Deep_Image_Homography_Estimation
19. [1606.03798] Deep Image Homography Estimation - arXiv, 檢索日期:7月 31, 2025 , <https://arxiv.org/abs/1606.03798>
20. [1910.02190] Kornia: an Open Source Differentiable Computer Vision Library for PyTorch, 檢索日期:7月 31, 2025, <https://arxiv.labs.arxiv.org/html/1910.02190>
21. What is Kornia library, 檢索日期:7月 31, 2025,
<https://kornia.readthedocs.io/en/latest/get-started/introduction.html>
22. kornia/kornia: Geometric Computer Vision Library for Spatial AI - GitHub, 檢索日期:7月 31, 2025, <https://github.com/kornia/kornia>
23. raw.githubusercontent.com, 檢索日期:7月 31, 2025,
<https://raw.githubusercontent.com/kornia/kornia/master/README.md>
24. Kornia: an Open Source Differentiable Computer Vision Library for PyTorch, 檢索日期:7月 31, 2025,
https://openaccess.thecvf.com/content_WACV_2020/papers/Riba_Kornia_an_Open_Source_Differentiable_Computer_Vision_Library_for_PyTorch_WACV_2020_paper.pdf
25. kornia.geometry.homography - Read the Docs, 檢索日期:7月 31, 2025,
<https://kornia.readthedocs.io/en/latest/geometry.homography.html>
26. Differentiable SVD based on Moore-Penrose Pseudoinverse for Inverse Imaging Problems - arXiv, 檢索日期:7月 31, 2025, <https://arxiv.org/pdf/2411.14141>
27. cvlab-epfl/Power-Iteration-SVD: Backpropagation-Friendly-Eigendecomposition - GitHub, 檢索日期:7月 31, 2025,
<https://github.com/cvlab-epfl/Power-Iteration-SVD>
28. Backpropagation-Friendly Eigendecomposition - NIPS, 檢索日期:7月 31, 2025,
<http://papers.neurips.cc/paper/8579-backpropagation-friendly-eigendecomposition.pdf>
29. KingJamesSong/DifferentiableSVD: A collection of differentiable SVD methods and ICCV21 "Why Approximate Matrix Square Root Outperforms Accurate SVD in Global Covariance Pooling?" - GitHub, 檢索日期:7月 31, 2025,
<https://github.com/KingJamesSong/DifferentiableSVD>