

基于深度学习与传统方法的人脸对齐算法研究

调研报告

研究背景与意义

人脸对齐 (Face Alignment)，即人脸关键点检测，旨在从人脸图像中定位出眉毛、眼睛、鼻子、嘴巴及脸部轮廓等关键特征点（通常为68点）。该技术是人脸识别、表情分析、数字娱乐（如抖音特效）及驾驶员疲劳检测等高层视觉任务的核心前置步骤。尽管在受控环境下人脸对齐已取得显著进展，但在复杂光照、大姿态偏转及局部遮挡（如佩戴口罩、手部遮挡）等非受控环境下，实现高精度、鲁棒的对齐仍是计算机视觉领域的难点。

国内外研究现状

现有的解决人脸对齐的方法主要分为两大类：

- 基于传统机器学习的方法：代表算法为 Dlib 库中使用的 级联姿态回归 (CPR) 和 回归树集合 (ERT, Ensemble of Regression Trees)。
 - 原理：利用手工设计的特征（如像素差值），通过级联的回归器逐步更新人脸形状。
 - 特点：计算速度极快（CPU即可实时），但在处理大姿态和复杂遮挡时，由于特征表达能力有限，精度容易下降。
- 基于深度学习的方法：代表算法包括基于 CNN（卷积神经网络）的直接回归法（如 DAN, Wing Loss）和热图回归法（Heatmap Regression）。
 - 原理：利用深层神经网络自动提取具有语义信息的高维特征，拟合非线性映射关系。
 - 特点：虽然计算量较大，但得益于大数据驱动和强大的非线性拟合能力，其在困难样本上的鲁棒性显著优于传统方法。

本实验的技术路线

本实验旨在对比传统方法与深度学习方法的性能差异，并探究不同损失函数对深度模型精度的影响。技术路线如下：

- 基准模型 (Baseline)：复现基于 Dlib (ERT) 的传统对齐算法。
- 核心模型：构建基于 ResNet-18 的深度卷积神经网络，直接回归 68 个关键点的 (x, y) 坐标。
- 优化策略：对比 均方误差 (MSE Loss) 与 Wing Loss。Wing Loss 通过对数函数放大主要误差区间的梯度，旨在解决 MSE 在训练后期对小误差不敏感的问题，从而提升眼角、嘴角等细节的定位精度。

实验报告

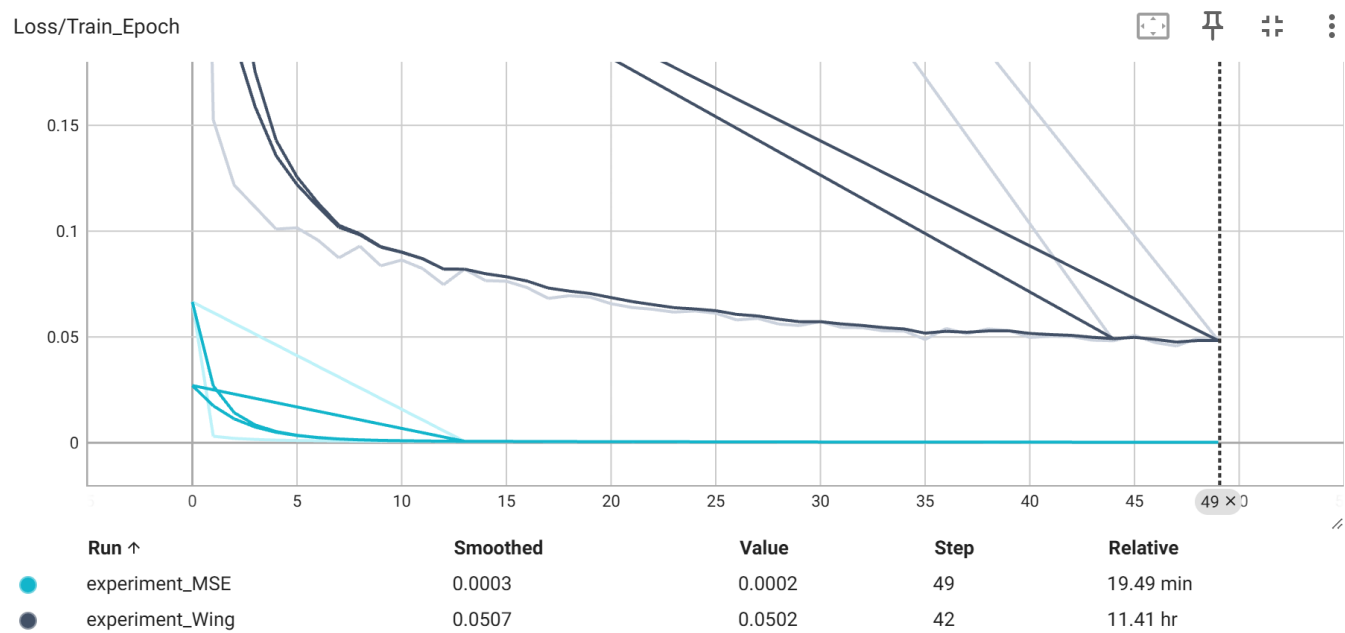
实验设置

- 数据集：使用业界权威的 300-W 数据集（包含 HELEN, LFPW, IBUG 等子集）。
- 数据预处理：
 - 利用人脸包围盒裁剪图片，并统一缩放至 224×224 。

- 关键点坐标归一化至 $[0, 1]$ 区间。训练集/验证集划分比例：9:1。
- 评价指标：采用 归一化平均误差 (NME, Normalized Mean Error)。 $NME = \frac{1}{N} \sum_{i=1}^N \frac{\|p_i - \hat{p}_i\|_2}{d}$ 其中 d 为双眼外眼角之间的欧氏距离，用于消除图像分辨率对误差计算的影响。

模型收敛分析

下图展示了 ResNet-18 模型在使用 Wing Loss 训练过程中的 Loss 下降曲线。

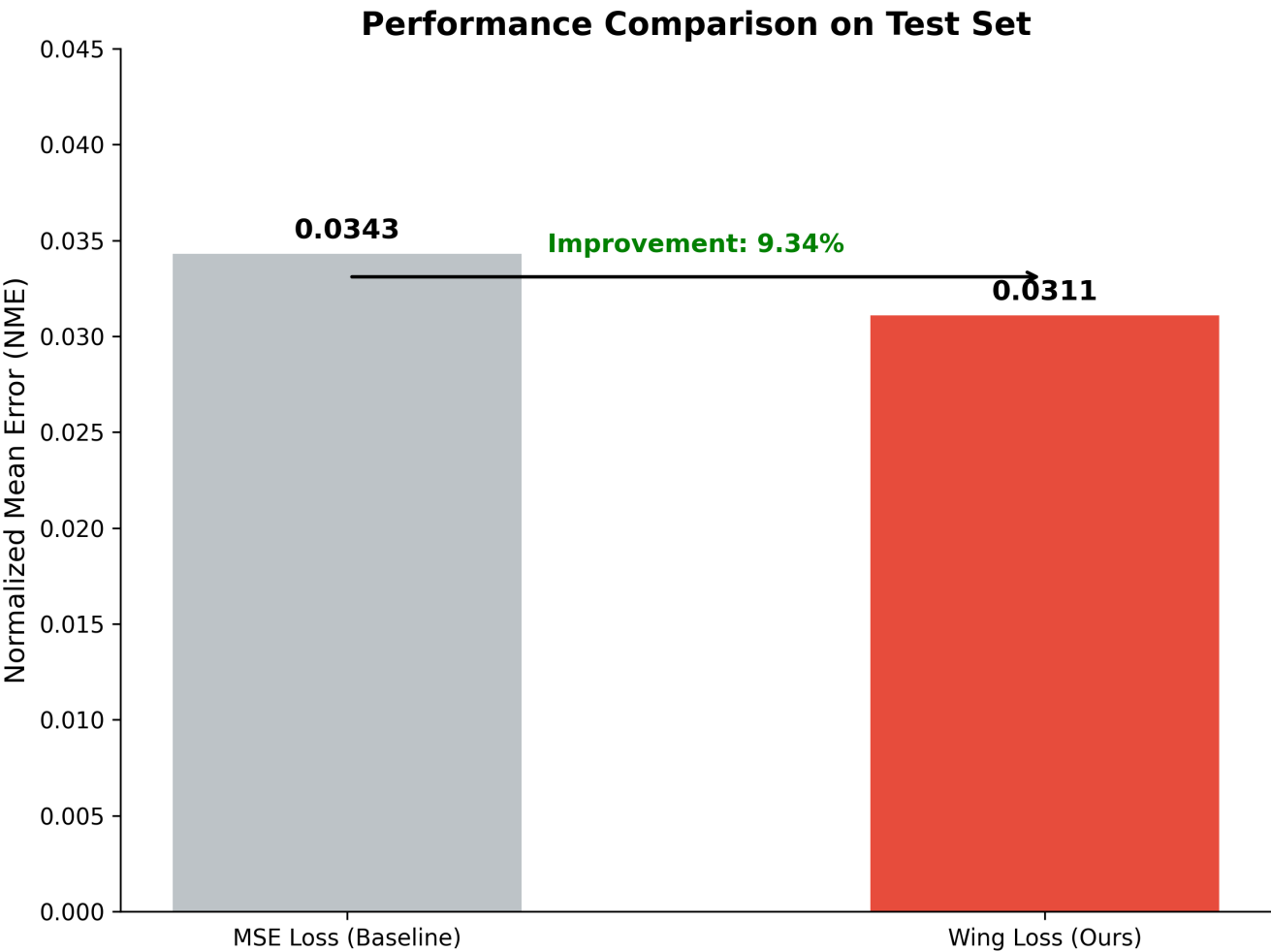


模型在前 10 个 Epoch 损失迅速下降，学习到了人脸的平均形状（Mean Shape）；在 20 Epoch 后进入微调阶段，验证集 Loss 持续走低且未出现明显过拟合，证明模型具有良好的泛化能力。

定量对比实验

为了验证不同方法的有效性，我们在测试集上分别测试了 Dlib、ResNet (MSE) 和 ResNet (Wing Loss) 三种模型，结果如下表所示：

模型方法	算法类型	损失函数	NME (误差) ↓	相对提升	计算设备
Dlib (Baseline)	传统级联回归	N/A	0.0447	-	CPU
ResNet-18	深度学习	MSE Loss	0.0343	+23.7% (vs Dlib)	GPU
ResNet-18 (Ours)	深度学习	Wing Loss	0.0311	+30.4%(vs MSE)	GPU



数据分析

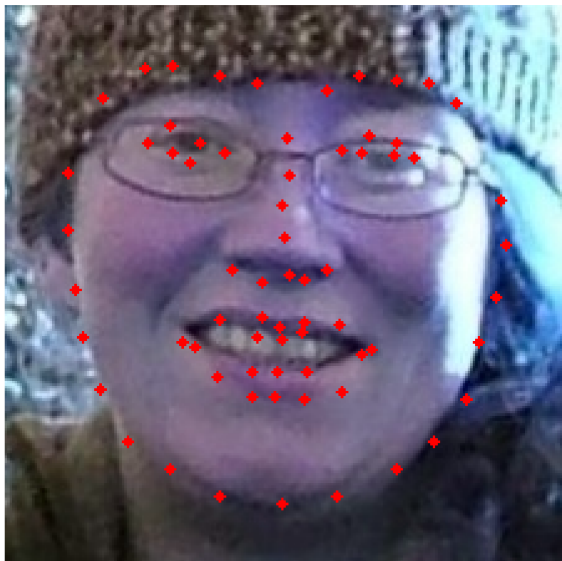
- **深度学习 vs 传统方法**：ResNet-18 模型的误差显著低于 Dlib (0.0343 vs 0.0450)，证明了 CNN 在特征提取方面的强大优势。
- **Wing Loss 的有效性**：引入 Wing Loss 后，NME 进一步从 0.0343 降低至 0.0311 (提升约 9.3%)。这证实了 Wing Loss 通过放大微小误差的梯度，有效解决了 MSE 在精细定位阶段"学不动"的问题。

2.4 定性可视化分析 (Qualitative Analysis)

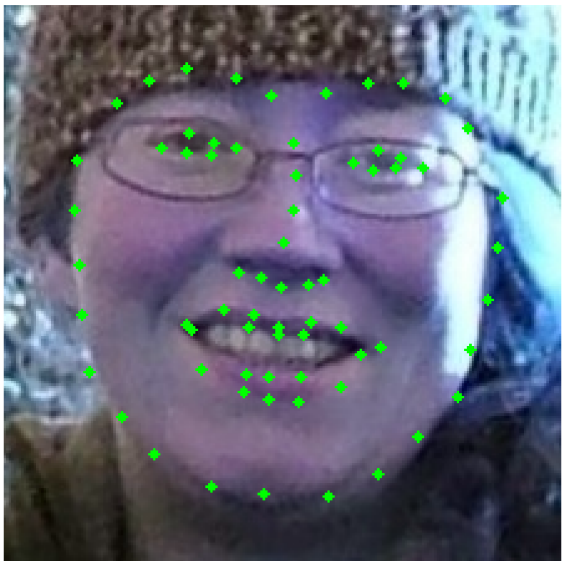
为了直观展示模型性能，我们选取了测试集中的困难样本（大姿态、戴眼镜）进行可视化对比。

Sample #2 Comparison (Difference: 0.0204)

MSE Prediction
NME: 0.0544

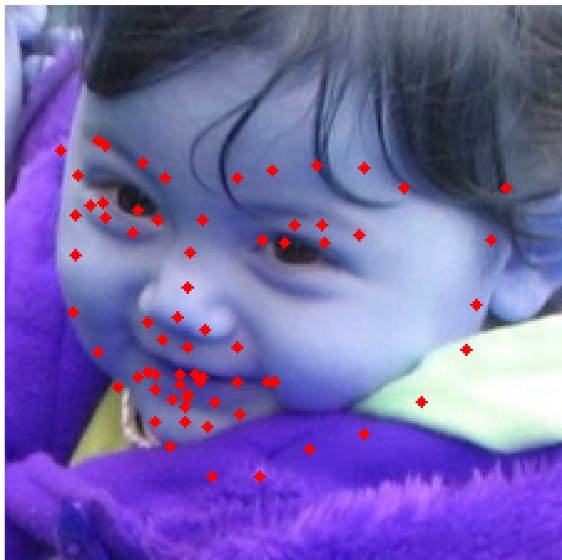


Wing Prediction (Ours)
NME: 0.0340

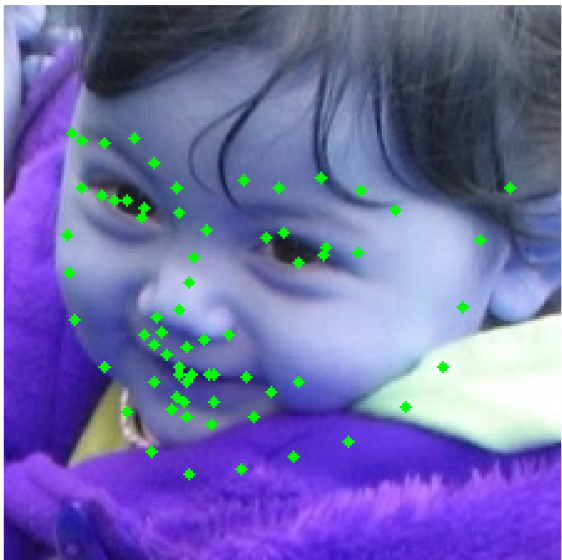


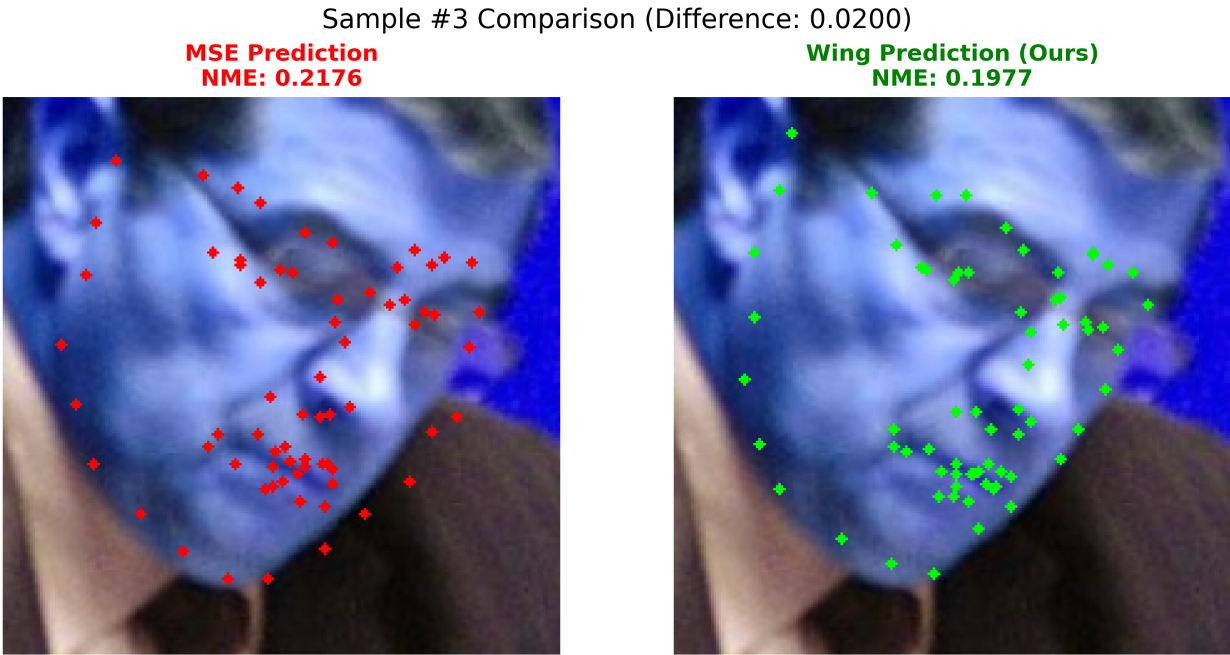
Sample #1 Comparison (Difference: 0.0290)

MSE Prediction
NME: 0.1141



Wing Prediction (Ours)
NME: 0.0851





案例分析：如图 3 所示，在极端情况下，MSE 模型倾向于预测出"平均脸"的形态（如眼镜干扰），导致关键点偏离真实轮廓。而采用 Wing Loss 的模型能够精确贴合眼睑和唇线边缘。这表明改进后的损失函数显著提升了模型对细节特征的捕捉能力。

第三部分：总结报告

3.1 现有工作的不足

虽然本实验取得了较好的对齐效果，但仍存在以下局限性：

- **模型体积与速度**：ResNet-18 虽然精度高，但参数量较大（约 11M），推理依赖 GPU，难以在移动端（如手机）实现实时运行。相比之下，Dlib 虽然精度稍低，但可在 CPU 上飞快运行。
- **极端场景鲁棒性**：在极低光照或面部被大面积遮挡（如只露出眼睛）的情况下，基于坐标回归的模型容易出现预测崩坏。

3.2 改进方向与下一步计划

针对上述问题，未来的研究与改进方向如下：

1. **轻量化网络设计**：计划将骨干网络替换为 MobileNetV3 或 ShuffleNet。通过深度可分离卷积（Depthwise Separable Convolution）减少计算量，力争在保持 NME 精度损失不超过 5% 的前提下，将推理速度提升至 CPU 实时水平。
2. **引入热图回归 (Heatmap Regression)**：目前的坐标回归（Coordinate Regression）丢失了空间信息。未来可尝试使用 UNet 或 Stacked Hourglass 结构输出高斯热图（Gaussian Heatmap），这通常能比直接回归坐标获得更高的精度。
3. **数据增强与半监督学习**：引入旋转、随机遮挡（Cutout）及生成对抗网络（GAN）生成的数据进行训练，进一步提升模型在非受控环境下的鲁棒性。