

面向开放场景的多模态

步态识别鲁棒性增强策略研究

Research on Robustness Enhancement Strategies for  
Multi-Modal Gait Recognition in Open Scenarios

# 简介

- 步态识别作为一种非侵入式生物特征识别技术，在智慧安防、人机交互等领域具有广阔的应用前景。然而，在开放监控场景下，复杂多变的协变量(如衣着、光照、遮挡及视角变化)对传统单模态识别模型构成了严峻挑战，导致其鲁棒性和泛化能力显著下降。为此，融合多模态信息(如耦合人体的外观表征与运动学结构信息)被视为突破识别瓶颈的关键路径。多模态数据能提供更具判别力的互补信息，从而增强模型对环境干扰的抵抗能力。本课题的核心在于探索何种鲁棒性增强策略能最有效地利用这些异构信息。本研究将聚焦于多模态融合策略本身，将其作为提升鲁棒性的核心手段。研究将系统性地比较和评估不同的融合范式，包括但不限于从经典的特征级联、决策仲裁到先进的基于注意力机制的动态对齐与加权策略。本研究将在模拟开放场景的基准数据集上进行实证分析，旨在为构建高鲁棒性的多模态步态识别系统提供算法选型和策略参考。

# 研究

- (1) 搜索步态识别(Gait Recognition)的基础知识，包括其定义、生物特征识别原理、及其在安防和人机交互中的具体应用。
- (2) 深入研究单模态步态识别(Single-Modal Gait Recognition)的主流方法(例如基于剪影的模型 GaitSet)，并重点查明“开放场景”下的关键挑战，即协变量(Covariates)如衣着、光照、遮挡和视角变化如何导致模型鲁棒性下降。
- (3) 学习多模态步态识别(Multi-Modal Gait Recognition)的概念，特别是研究所涉及的异构信息源，包括“外观表征”(如RGB视频、轮廓剪影)和“运动学结构信息”(如人体骨骼关节点数据)的提取技术。
- (4) 系统性地学习多模态融合(Multi-Modal Fusion)策略，这是本课题的核心。首先，研究简介中提到的经典融合范式：
  - (a) 特征级联(Feature-level fusion)
  - (b) 决策仲裁(Decision-level fusion)

- (5) 深入研究先进的多模态融合技术，特别是简介中强调的“基于注意力机制的动态对齐与加权策略”(Attention-based dynamic alignment and weighting)，搜索跨模态注意力(Cross-Modal Attention)等机制如何用于融合步态数据。
- (6) 调研和学习用于评估开放场景步态识别的基准数据集(Benchmark Datasets)，例如CASIA-B、OU-MVLP或其它包含复杂协变量的数据集，并了解其评价标准和性能指标(如Rank-1准确率)。
- (7) 搜索近期的学术论文和综述，内容聚焦于“多模态步态识别”和“鲁棒性增强”，分析不同研究如何比较和评估多种融合策略的有效性，为课题的“实证分析”部分收集参考。
- (8) 综合以上信息，总结归纳完成此毕业设计所需掌握的核心知识点：即理解单模态的局限性、掌握多模态数据的处理方法，并能深入比较(从经典到先进)不同多模态融合策略在提升开放场景鲁棒性方面的原理和表现。

# 1. 引言：开放场景下单模态步态识别的脆弱性

- 1.1 核心背景
- 步态作为一种独特的生物特征步态，即个体行走时的姿态变化，被视为一种关键的生物特征。其核心优势在于，它是唯一可在远距离、非受控状态下被获取的生物信息，且具有非接触性和难以模仿的特点。随着视频监控设备在机场、商场等公共场所的普及，步态识别在智慧安防、视频监控和法律援助等领域的应用价值日益凸显。
- 然而，步态识别的独特优势与其面临的核心挑战紧密相连。使其变得有价值的“远距离”特性，恰恰意味着传感器(摄像头)获取的数据质量必然受限——目标通常较小、分辨率较低、光照和视角均处于非受控状态。因此，使步态识别变得有用的条件，也正是使其脆弱的原因。

- 1.2 “开放场景”的定义与问题阐述
- 本课题所指的“开放场景”(Open Scenarios)，在学术界通常被称为“野外”或“真实世界”场景("in-the-wild")。这指的是在超市、城市街道等部署了数百个摄像头的完全不受控的真实环境。
- 在这些场景中，会引入大量复杂的“协变量”(Covariates)。这些因素，如衣着变化、携带物品、光照、遮挡及视角变化，会“显著影响识别系统的准确性”并导致“性能下降”。正如课题简介所述，这导致传统单模态模型的鲁棒性和泛化能力显著下降。
- 这种脆弱性源于传统识别流程的“串行故障点”。经典的步态识别流程通常是串行的：行人检测→行人分割→行人追踪→行人识别。在开放场景中，流程的早期阶段(如依赖背景扣除的轮廓分割)极易受到光照、阴影或衣着的影响而失败。一旦轮廓信息损坏，无论后续的识别模型多么先进，整个识别链条都会断裂。这就为多模态方法提供了理论基础：通过并行的、互补的数据通道(如轮廓和骨骼)，一个通道的失败可由另一个通道补偿，从而实现系统级的鲁棒性。

- 1.3 多模态论点：以融合作为鲁棒性路径
- 面对单模态的瓶颈，本课题的核心论点是：融合多模态信息是提升鲁棒性的“关键路径”。具体而言，研究聚焦于耦合“外观表征”(如轮廓)与“运动学结构信息”(如骨骼)。
- 其理论基础在于，异构数据能提供“互补增强”和“多特征的相互验证”。目标是构建一个系统，使其能够利用一种模态的信息来纠正或补充另一种模态在特定干扰下的缺失。本课题的核心并非发明新的模态，而是系统性地探索如何最有效地融合这些已有模态，以抵抗环境干扰。

## 2. 协变量挑战：真实世界干扰因素的解构

- 2.1 协变量的定义
- 协变量是导致步态识别系统效率降低或准确性显著下降的外部因素。它们是阻碍步态识别技术从实验室走向实际应用的核心障碍。

## ● 2.2 关键协变量分析

- **视角(View Angle)**: 视角变化被认为是影响步态识别性能的“最主要因素”之一，也是计算机视觉领域“最本质和最困难的任务”之一。原因在于，视角的变化会导致提取的人体轮廓(外观)产生“巨大差异”。
- **衣着(Clothing)**: 衣着或穿着条件(如穿大衣、裙子)是公认的主要挑战。它直接改变了人体的视觉外观和轮廓形状，使得基于外观的模型难以区分“衣着的变化”和“身份的变化”。
- **遮挡(Occlusion)**: 遮挡是实际应用中的一个“关键限制”，在真实世界中“不可避免”。遮挡源可能来自：1) **个体因素**，如衣物或携带物品(例如背包)；2) **环境因素**，如行人走过柱子或被其他行人遮挡；3) **摄像机部分视角**。过去，相关研究受限于“缺乏带有可量化遮挡指标的公共步态数据库”，而OccGait和Gait3D等新数据集的出现正是在回应这一挑战。
- **其他协变量**: 还包括携带条件(如背包)、行走速度和光照条件等

- 在分析这些协变量时，必须认识到它们对不同模态的影响是非对称的。例如，衣着变化对基于轮廓的外观模态是灾难性的，因为它直接改变了轮廓特征；但对于运动学模态来说，只要姿态估计算法(如OpenPose)能够“看透”衣物定位到关节点，衣着的影响就相对较小。相反，严重的肢体遮挡可能会使姿态估计算法完全失效(例如无法检测到腿部关节点)，导致运动学模态崩溃；但此时，外观模态可能仍能从可见的上半身轮廓中提取到有效的识别信息。
- 这种“模态特定的脆弱性”是理解多模态融合为何至关重要的核心。没有任何一种数据模态是永远可靠的。一个真正的鲁棒系统必须具备动态识别和权衡不同模态信号质量的能力，在特定协变量(如衣着)干扰外观时，更“信任”运动学数据；而在另一些协变量(如遮挡)干扰运动学时，转而更“信任”外观数据。这直接指向了基于注意力机制的动态融合策略，而非简单的静态特征相加。
- 同时，学术界对协变量的评估标准已经发生了根本性转变。早期的数据库(如CASIA-B)在受控的实验室环境中模拟协变量(例如，让受试者“穿上大衣”或“携带背包”）。然而，研究已经证实这种模拟是不充分的。当前的SOTA(State-of-the-Art)研究必须在“in-the-wild”基准(如Gait3D和GREW)上进行评估。这些数据集的挑战不在于单独的“衣着”或“视角”，而在于同时面对“穿着大衣”、“被部分遮挡”、“视角刁钻”和“光线昏暗”的复合型真实挑战。

### 3. 多模态范式：利用异构数据增强鲁棒性

- 3.1 融合的理论优势
- 多模态融合的目的是实现“互补增强和有效融合”。它通过整合来自不同传感源或不同表征的信息，提供了“更丰富的数据支持”。其核心价值在于，通过“多种特征的相互验证”，可以“减少”环境变化等因素对识别稳定性的影响。

- 3.2 两种数据族：外观(“是什么”) vs. 运动学(“如何动”)
- 当前的步态识别研究，特别是多模态方法，已广泛收敛于融合两大类异构数据：
  - **基于外观的表征(Appearance-based)**：这类方法“利用信息丰富的视觉特征”。它包括原始的RGB图像、二值化的人体轮廓(Silhouettes)以及更精细的人体解析(Human Parsing)图。这族数据回答了“这个人看是什么样子”的问题。
  - **基于模型的表征(Model-based)**：这类方法“利用人体的潜在结构”。它包括2D或3D的人体骨骼(Skeletons)和3D人体网格(Meshes)。这族数据专注于运动学(Kinematics)，回答了“这个人是如何运动的”的问题。

- 3.3 实践中的互补性：优势与权衡
- 精确总结这两种模态间的核心权衡：
  - **外观模态：**优势在于它直接来自包含“更丰富信息”的原始RGB帧。劣势在于它对衣着、速度等协变量非常敏感。它信息丰富，但信噪比低(易受干扰)。
  - **运动学模态：**优势在于它对某些协变量(特别是“衣着”)具有“鲁棒性”。劣势在于它依赖于姿态估计算法，而这些算法本身需要“相对较高的图像分辨率”才能可靠运行，并且带来了“更高的计算成本”。它信息更纯粹(只有结构和运动)，但可能信息量较少(丢失了体型、轮廓等信息)。

- 多模态融合的目标，就是设计一个架构，能够同时利用外观模态的“丰富性”和运动学模态的“鲁棒性”，实现“鱼与熊掌兼得”。
- 这种融合策略在概念上是一种“解耦”(Disentanglement)。在单模态的外观模型中，“步态特征”和“外观特征”(如衣着)是高度纠缠(Entangled)在一起的。提到有工作试图“自动地将姿态/步态特征从外观特征中分离出来”。而多模态架构，通过设立一个独立的运动学分支，显式地强迫模型从骨骼数据中学习“运动”；同时，外观分支则可以更专注于学习“体型”。先进的融合策略(如C2Fusion)不仅是简单地相加特征，而是在一个解耦的潜在空间中重组这些信息。
- 此外，选择模态也是一种在预处理和后处理之间分配计算复杂度的架构权衡。外观模态的预处理可能很简单(如背景减除)，但需要一个极其复杂的识别模型来处理后续的噪声。相反，运动学模态的预处理极其复杂(例如，运行OpenPose 本身就是一个庞大且计算昂贵的神经网络)，但在预处理完成后，其输出(例如17个关节点的坐标)是低维、干净的数据，这使得后续的识别模型可以设计得相对简单。

## 4. 外观表征模态分析(即“是什么”)

- 4.1 轮廓(Silhouettes) 和步态能量图(GEI)
- 从历史上看，步态轮廓是“步态识别中大多采用”的特征。其核心是提取“完整封闭的运动人体轮廓”，这通常通过高斯背景模型或深度学习分割网络来实现。为了处理序列，这些轮廓帧常被聚合成一个单一的模板，如“步态能量图”(Gait Energy Image, GEI)，或“分段帧差能量图”(SFDEI)。轮廓模态的脆弱性(如对噪声和阴影的敏感 )是推动多模态研究的根本原因。

- 4.2 人体解析(Human Parsing)
- 人体解析是近期SOTA研究中的一个重要趋势，它被视为一种“细粒度的分割”。与二值化的轮廓不同，Parsing模型将人体分割为不同的语义部分(如头、躯干、左臂、右腿等)。这种模态被用于MultiGait++等前沿工作中。其理论优势在于：如果一个协变量(如背包)在二值轮廓中会与“躯干”融合，导致特征污染；而一个语义级的Parsing模型则有潜力将“背包”(外部物体)与“躯干”(身体部分)区分开。

- 4.3 光流(Optical Flow)
- 光流(即帧间运动矢量)是GREW等大型基准数据集提供的模态之一，并被用于SOTA模型的比较中。光流的理论优势是它对静态外观(如衣服的颜色和纹理)完全不敏感，因为它只捕捉运动。然而，它对摄像机自身的运动(抖动)非常敏感。
- 外观表征的发展展现出一条清晰的“演化路径”：每一个新模态的出现，都是为了解决前一个模态的缺陷。
  - 轮廓：因噪声和衣着问题而失败。
  - GEI：试图通过平均化来“去噪”，但牺牲了时序动态信息。
  - 人体解析：通过引入语义知识来解决衣着/背包问题。
  - 光流：通过只关注运动来解决静态外观(如衣物颜色)的干扰。
- 这表明，没有任何一种外观表征是完美的。当前的SOTA模型(如MultiGait++)往往会融合多种外观模态(如轮廓+解析+光流)，这暗示了单一外观表征的局限性。

# 5. 运动学与结构模态分析(即“如何动”)

- 5.1 2D/3D骨骼(Keypoints)
- 这是最主流的“基于模型”的方法。它涉及使用开源算法(如OpenPose或AlphaPose)从视频帧中提取人体关节点(Keypoints)，或直接从深度相机获取。这些稀疏的坐标点序列随后被编码为高级特征，如“关节点角度轨迹”或肢体长度。这些特征因其对衣着变化的鲁棒性而备受青睐。

- 5.2 3D人体网格(3D Meshes)
- 这是一种更稠密、更完整的3D表征，例如SMPL模型。Gait3D数据集 在这方面具有开创性，它提供了从视频中恢复的3D SMPL模型，这些模型能捕捉到“关于体型、视角和动态的稠密3D信息”。这可以说是目前最完善的表征，因为它隐式地同时包含了运动学信息(骨架)和一种“理想化”的体型外观信息。

- 5.3 从坐标到图像：“骨骼图”(Skeleton Maps)
- 直接处理坐标序列(一种时序数据)与处理轮廓图像(一种视觉数据)需要截然不同的网络架构(如RNNvs.CNN)。这为后续的特征融合带来了挑战。因此，近期的SOTA工作(如SkeletonGait )采用了一种巧妙的技巧：将骨骼坐标(稀疏的点列表)转换(渲染)为2D图像表征，即“骨骼图”(Skeleton Maps)。这一转换的意义重大：它统一了数据格式。通过将骨骼转换为“图像”，研究者现在可以使用相同的强大CNN架构(例如ResNet变体 )来并行处理轮廓图分支和骨骼图分支。这极大地简化了多模态融合的网络设计。
- 然而，运动学模态的鲁棒性建立在一个“隐藏的依赖”之上：即上游的姿态估计算法(如OpenPose)的准确性。但姿态估计算法本身也受协变量(特别是遮挡和运动模糊)的严重影响。因此，运动学分支并非“银弹”，它在某种程度上是将“鲁棒的步态识别”问题转移到了“鲁棒的姿态估计”问题。如果OpenPose因遮挡而出错，运动学分支将接收到无意义的“噪声”数据。这再次凸显了动态融合机制的必要性——该机制必须能够评估姿态数据的实时质量，并在姿态估计失败时动态降低其权重。
- 推动这一领域发展的另一个关键催化剂是数据集的演进。在早期，研究者需要自己在大规模数据集(如GREW)上运行昂贵的姿态估计算法。而Gait3D和GREW等现代基准数据集，在发布时直接提供了预先计算好的多模态数据流(轮廓、2D/3D骨骼、光流等)。这极大地降低了研究门槛，使得研究者(如本课题)能够将精力聚焦于融合策略本身，而不是数据预处理。

# 6. 核心研究轴：多模态融合策略的分类与SOTA分析

- 6.1 经典融合范式：早期、晚期与混合
- 传统的融合策略根据融合发生的阶段进行划分：
  - 早期融合(Early Fusion/**特征级**)：在数据输入层或浅层特征层进行融合。优势是允许模型学习模态间底层的复杂关联。劣势是要求数据严格对齐。
  - 晚期融合(Late Fusion/**决策级**)：各个模态独立进行预测，最后融合决策结果(如对Softmax概率求和或投票)。优势是灵活性高，各分支可独立训练。劣势是完全丢失了特征交互信息。
  - 对于深度学习，这种二元划分过于粗糙。当前的研究重点是中层融合(Middle Fusion)，即在网络的哪个深度阶段(输入层、中层或高层)进行特征交互。

- 6.2 深度学习基线：静态融合策略
- 在比较先进策略之前，必须建立强有力的基线(Baselines)。在SOTA研究中，两种最常见的静态特征级融合方法是：
  - **逐元素相加(Add Fusion)**：将来自不同分支(如轮廓和骨骼图)的特征图(Feature Maps)在逐个元素上相加。这假设两个特征空间是语义对齐的。
  - **通道拼接(Concat Fusion)**：沿着通道维度将特征图“堆叠”在一起。这种方法更通用，它通常会立即跟随一个 $1 \times 1$ 卷积层，目的是在通道间“混合”信息并降低维度。这是一个计算廉价且效果惊人强大的基线。

- 6.3 高级融合：基于注意力机制的动态策略
- 这是当前SOTA的研究前沿。与Add或Concat这样的静态融合(即融合权重是固定的)相反，基于注意力的融合是动态的——融合权重是根据输入数据实时计算的。
- 其核心是跨模态注意力(Cross-Modal Attention, CMA) 机制。CMA专为“跨模态特征交互”和“高效联合建模”而设计。它允许一个模态的特征(作为Query)去“查询”另一个模态的特征(作为Key和Value)。这使得模型能够“专注于每种模态中最具信息量的方面”。
- 如第2节所述，这是对“模态特定脆弱性”问题的直接回应。如果当前帧的轮廓因遮挡而损坏，CMA机制可以通过学习，在生成最终特征时动态地降低对轮廓分支的注意力权重，转而提高对骨骼分支的权重。这就是通过动态融合实现的“鲁棒性”。

- 6.4 SOTA融合架构案例研究
- 本课题旨在比较不同的融合策略。通过分析近期的SOTA模型，我们发现“基于注意力的策略”本身并非铁板一块，而是存在多种不同的实现路径。

- 案例 1：SkeletonGait++及其“Attention Fusion”
- 融合模态：轮廓(Silhouette)+骨骼图(Skeleton Maps)。
- 融合架构：多个文献共同清晰地定义了其“Attention Fusion”机制。关键点：它不是基于Transformer的跨模态注意力。它是一种卷积注意力模块(Convolutional Attention Module)：
  - 首先，在通道上Concatenate轮廓和骨骼分支的特征图。
  - 将其送入一个“小型网络”(small network)。
  - 这个“小型网络”被定义为：一个 $1 \times 1$ 的压缩卷积→一个 $3 \times 3$ 的普通卷积→一个 $1 \times 1$ 的扩张卷积。
  - 最后通过一个Softmax层生成逐元素的注意力权重(Mask)，再乘回原特征。
- 分析：这是一种轻量级、高效的空间/通道注意力实现。它学习的是融合后的特征图中“哪些区域或通道更重要”。这是本课题极佳的比较对象之一。

- 案例 2：MultiGait++及其"C2Fusion"
- 融合模态：轮廓(Silhouette)+人体解析(Parsing)+光流(Optical Flow)。
- 融合架构：提出了一种名为"C2Fusion"的新策略。
- 核心思想：C2Fusion的目标不仅是混合特征，而是要同时“提取跨模态的共享特征(Commonalities)并“鼓励每个模态强调其独特属性(Unique Attributes)”。
- 分析：这在概念上是更高级的融合。它主动地试图对特征进行“解耦”(见3.3节的分析)。它试图分离出“所有模态都同意的信息”(例如，这个人正在向前走)和“只有特定模态知道的信息”(例如，Parsing模态知道“这是一个背包，不是躯干”)。实验数据表明，这种方法在多个基准上优于SkeletonGait++。

- 案例 3：GaitCSF及其"Channel Shuffle Module"
- 融合模态：轮廓(Contour)+姿态热图(Heatmap)。
- 融合架构：使用“基于通道洗牌的特征选择性调节模块(CFS)”。
- 核心思想：它将通道分为数组，在组内并行处理通道统计和空间统计，然后使用(源自ShuffleNet的)“通道洗牌”(Channel Shuffling) 操作在不同组之间交换信息。
- 分析：这是第三条路径，其设计重点是效率。它旨在通过“轻量级的参数设计”实现“自适应的特征增强”。

- 本课题的核心任务(比较融合策略)必须认识到“注意力”是一个广谱的术语。通过上述案例，我们至少可以识别出三个不同的“注意力”家族：
- **卷积注意力(Conv-Attention)**: (如SkeletonGait++)，使用小型CNN生成注意力图。
- **解耦注意力(Disentanglement-Attention)**: (如MultiGait++)，使用特定目标(共享/独特)来指导融合。
- **高效注意力(Efficient-Attention)**: (如GaitCSF)，使用分组卷积和通道洗牌等高效操作实现自适应。
- 本课题的一个强大贡献将是对这些不同“注意力”家族进行直接的实证比较。

- 下表总结了这些SOTA融合策略的关键差异，这将是本课题进行算法选型和策略参考的核心依据。
- 表 1：SOTA融合架构的对比分析

模型/架构	融合的模态	核心融合机制	机制原理	假设优势
SkeletonGait++	轮廓+骨骼图	“Attention Fusion” (卷积注意力)	Concat→小型网络 ( $1 \times 1 \rightarrow 3 \times 3 \rightarrow 1 \times 1$ Conv)→Softmax生成Mask。	轻量级的空间/通道注意力，学习融合特征的局部重要性。
MultiGait++	轮廓+人体解析 +光流	“C2Fusion” (解耦注意力)	提取跨模态的“共享特征” “(Commonalities)，同时保留各模态的“独特属性” “(Uniqueness)。	更优的特征解耦， 获得信息更丰富的表征，SOTA性能。
GaitCSF	轮廓+姿态热图	CFS模块 (高效注意力)	通道分组→并行(通道/空间)统计→“通道洗牌”(Channel Shuffle)。	“轻量级参数设计”下的“自适应特征增强”。
(基线)	轮廓+骨骼图	Concat Fusion	Concatenate→ $1 \times 1$ Conv。	简单、高效、强大的基线。

- 在所有这些研究中，OpenGait框架扮演了至关重要的角色。它不仅是一个开源项目，更是一个旨在提供“公平比较”(Fair Conditions)的基准测试框架。它在同一套代码库中实现了SkeletonGait++和MultiGait++等模型，并允许研究者在不同数据集上对比不同的融合机制(如Add, Concat, Attention)。因此，OpenGait是执行本课题研究(即“比较和评估不同的融合范式”)的理想工具。

# 7. 鲁棒性评估：基准数据集与实验协议

- 7.1 "In-the-Wild"基准的必要性
- 学术界长期以来受到“在受控环境中捕获的”数据库的限制。经典的CASIA-B数据集 擅长隔离单一协变量(如11个固定视角 )，但它无法代表真实的“开放场景”。为了在“不受控的场景”中训练和评估识别器，新的“in-the-wild”基准是必需的。

- 7.2 SOTA "In-the-Wild"多模态数据集分析
- **Gait3D:**
  - 定位：一个“in-the-wild”数据集，包含4,000名受试者，在一个不受控的室内场景中捕获。
  - 多模态数据：其核心贡献是同时提供了多种模态：2D轮廓、Keypoints(骨骼)和3D SMPL网格模型。
  - 分析：Gait3D完美契合本课题的需求，因为它原生支持外观模态(轮廓)和运动学模态(骨骼、SMPL)的融合研究。

- **GREW(Gait REcognition in the Wild):**

- 定位：第一个大规模"in-the-wild"步态数据集。其规模空前：包含26,345名受试者和128,671个序列，数据来自数百个摄像头。它还包含一个庞大的“干扰项集”(Distractor set)，这使其极度接近真实世界的安防应用。
- 多模态数据：同样提供了丰富的模态，包括轮廓、GEIs、光流以及2D/3D姿态。
- 分析：由于其巨大的规模和“多样化且实际的视角变化”以及“自然的挑战性因素”，GREW是评估鲁棒性的终极基准。
- SOTA模型(第6节)和SOTA数据集(第7节)的发展是“共同进化”的。没有像Gait3D和GREW这样提供多模态数据的基准，研究界就无法系统地开发和验证像SkeletonGait++这样的融合模型。反过来，这些数据集的出现也推动了研究领域转向融合策略的研发。

- 7.3 实用实验协议：**OpenGait**框架
- 如前所述，OpenGait是一个用于步态识别的软件基准框架。它在统一的训练和评估流程下，复现了SkeletonGait++和MultiGait++等模型。它被设计用于在“公平的条件”下进行比较。对于本课题而言，必须使用OpenGait框架。它是唯一能确保对不同融合策略(Add, Concat, Attention)进行科学有效和公平比较的工具。

- 7.4 评估指标：超越平均准确率
- 标准的评估指标是Rank-k准确率(例如 Rank-1)。然而，要证明“鲁棒性”，仅报告一个平均的 Rank-1准确率是不够的。本课题的评估必须利用GREW等数据集提供的“丰富的属性”。一个完整的评估应包括：
  - **全局准确率：**在整个测试集上的平均Rank-1准确率。
  - **协变量特定准确率：**这才是衡量鲁棒性的关键。应在数据集的特定子集上单独报告Rank-1准确率。例如，比较不同融合策略在“正常衣着(NM)”vs.“穿大衣(CL)”；“未携带(NM)”vs.“带包(BG)”；以及不同视角变化范围下的性能。

# 8. 总结与实施建议

- 8.1 文献综合结论
- 问题已明确：单模态步态识别在开放场景中(面对协变量时)不具备鲁棒性。
- 方向已收敛：多模态融合，特别是融合“外观表征”(轮廓、解析)和“运动学结构”(骨骼、姿态)，是当前学术界和工业界公认的、实现鲁棒性的SOTA路径。
- 前沿已确定：当前的研究前沿不是“是否要融合”，而是“如何融合”。本课题的定位精准地处于这一前沿。
- 趋势已显现：简单的静态融合(Add/Concat)作为基线是有效的，但动态的、基于注意力的策略(如CMA，C2Fusion，Conv-Attention)在性能上更具优势。
- 目标已清晰：存在多种高级融合策略(如表1所示)，本课题的核心任务就是对它们进行系统性的、公平的实证比较。

- 8.2 毕业设计推荐实验方案
- 为实现本课题“为构建高鲁棒性系统提供算法选型和策略参考”的目标，推荐以下四步实验方案：
- **第 1 步：工具(Framework)**
  - 选用OpenGait框架。
  - 理由：它提供了“公平比较”所需的统一平台，并已内置了SkeletonGait++和MultiGait++等SOTA 实现，极大降低了复现难度。
- **第 2 步：数据集(Datasets)**
  - 选用GREW或Gait3D。
  - 理由：它们是唯二满足本课题全部要求的基准：1) 包含真实的“in-the-wild”协变量；2) 提供(预先计算好的)多模态数据流(轮廓、骨骼、解析等)。

### • 第3步：模型(Variables to Test)

- 目标：隔离融合策略作为唯一变量。在OpenGait中，使用相同的骨干网络(Backbone)和相同的输入模态(例如，轮廓+骨骼图)来测试以下融合策略：
  - 基线 1: **Concat Fusion**(通道拼接+ $1 \times 1$ Conv)。
  - 基线 2: **Add Fusion**(逐元素相加)。
  - SOTA 1: **SkeletonGait++** "Attention Fusion"(即卷积注意力模块)。
  - SOTA 2: **MultiGait++** "C2Fusion"(实现其“共享+独特”的融合逻辑)。
  - (可选 SOTA 3): 一个标准的跨模态 Transformer融合模块(基于Q/K/V)。

## ● 第4步：评估(Result Measurement)

- 指标1(全局)：报告在整个测试集上的平均Rank-1准确率。
- 指标2(鲁棒性/核心)：这是本课题最重要的成果。利用GREW的属性标签，报告按协变量细分的Rank-1准确率。必须创建表格，清晰对比上述(第3步)所有策略在以下条件下的性能：
  - 衣着变化：正常(NM)vs.穿大衣(CL)。
  - 携带条件：正常(NM)vs.背包(BG)。
  - 视角变化：小角度vs.大角度(跨视角)。
  - (若数据可得)：不同遮挡水平下的性能。

- 通过执行这一方案，本课题将能产出经验性的、有数据支撑的结论。例如，不仅能得出“C2Fusion的平均准确率最高”，还能得出更深入的见解，如：“尽管C2Fusion平均最高，但SkeletonGait++的卷积注意策略在(CL)衣着变化下的性能下降最少，表明它在冬季安防场景下具有更强的鲁棒性”。这正是本课题标题所承诺的——提供“算法选型和策略参考”。