

1-20-2024

## Action Recognition Model of Directed Attention Based on Cosine Similarity

Chen Li

*Command & Control Engineering College, Army Engineering University of PLA, Nanjing 210007, China, 651220007@qq.com*

Ming He

*Command & Control Engineering College, Army Engineering University of PLA, Nanjing 210007, China, paper\_review@126.com*

Chen Dong

*Military Human Resource Support Center, Political Work Department of the Army, Beijing 100072, China*

Wei Li

*Command & Control Engineering College, Army Engineering University of PLA, Nanjing 210007, China*

Follow this and additional works at: <https://dc-china-simulation.researchcommons.org/journal>



Part of the Artificial Intelligence and Robotics Commons, Computer Engineering Commons, Numerical Analysis and Scientific Computing Commons, Operations Research, Systems Engineering and Industrial Engineering Commons, and the Systems Science Commons

---

This Paper is brought to you for free and open access by Journal of System Simulation. It has been accepted for inclusion in Journal of System Simulation by an authorized editor of Journal of System Simulation. For more information, please contact [xtfzxb@126.com](mailto:xtfzxb@126.com).

---

# Action Recognition Model of Directed Attention Based on Cosine Similarity

## Abstract

**Abstract:** Aiming at the lack of directionality of traditional dot product attention, this paper proposes a directed attention model (DAM) based on cosine similarity. To effectively represent the direction relationship between the spatial and temporal features of video frames, the paper defines the relationship function in the attention mechanism using the cosine similarity theory, which can remove the absolute value of the relationship between features. To reduce the computational burden of the attention mechanism, the operation is decomposed from two dimensions of time and space. The computational complexity is further optimized by combining linear attention operation. The experiment is divided into two stages : Four ablation experiments are carried out on each module of directed attention to show the best performance of DAM in accuracy and efficiency; the accuracy of the model is 7.3% higher than that of I3D-NL on the Sth-Sth V1(something something V1) dataset and 95.7% on the UCF101(101 human action classes from videos in the wild) dataset. The research results have a wide application prospect in safety monitoring, automatic driving, and so on.

## Keywords

action recognition, deep learning, attentional mechanism, cosine similarity, time-space decomposition

## Recommended Citation

Li Chen, He Ming, Dong Chen, et al. Action Recognition Model of Directed Attention Based on Cosine Similarity[J]. Journal of System Simulation, 2024, 36(1): 67-82.

# 基于余弦相似性的定向注意力行为识别模型

李晨<sup>1</sup>, 何明<sup>1\*</sup>, 董晨<sup>2</sup>, 李伟<sup>1</sup>

(1. 陆军工程大学 指挥控制工程学院, 江苏 南京 210007; 2. 陆军政治工作部 军事人力资源保障中心, 北京 100072)

**摘要:** 针对传统点积注意力缺乏方向性的问题, 建立了一种基于余弦相似性的定向注意力模型(directed attention model, DAM)。为有效表示视频帧时空特征间的方向关系, 运用余弦相似性理论, 定义了注意力机制中关系函数, 能够去除特征间关系绝对值; 为降低注意力机制计算量, 从时间和空间两个维度上对运算进行分解; 结合线性注意力运算, 进一步优化计算复杂度。实验分为两个阶段: 对定向注意力各模块开展了4个消融实验, 以表现DAM在精确度和效率方面的最佳性能; 该模型在 Sth-Sth V1(something something V1)数据集上的精确度较 I3D-NL(inflated 3D ConvNet non-local)高7.3%, 在UCF101(101 human action classes from videos in the wild)数据集上的识别精确率为95.7%。研究成果在安全监控、自动驾驶等方面应用前景广泛。

**关键词:** 行为识别; 深度学习; 注意力机制; 余弦相似性; 时空分解

中图分类号: TP391.4 文献标志码: A 文章编号: 1004-731X(2024)01-0067-16

DOI: 10.16182/j.issn1004731x.joss.22-0937

**引用格式:** 李晨, 何明, 董晨, 等. 基于余弦相似性的定向注意力行为识别模型[J]. 系统仿真学报, 2024, 36(1): 67-82.

**Reference format:** Li Chen, He Ming, Dong Chen, et al. Action Recognition Model of Directed Attention Based on Cosine Similarity[J]. Journal of System Simulation, 2024, 36(1): 67-82.

## Action Recognition Model of Directed Attention Based on Cosine Similarity

Li Chen<sup>1</sup>, He Ming<sup>1\*</sup>, Dong Chen<sup>2</sup>, Li Wei<sup>1</sup>

(1. Command & Control Engineering College, Army Engineering University of PLA, Nanjing 210007, China;

2. Military Human Resource Support Center, Political Work Department of the Army, Beijing 100072, China)

**Abstract:** Aiming at the lack of directionality of traditional dot product attention, this paper proposes a directed attention model (DAM) based on cosine similarity. To effectively represent the direction relationship between the spatial and temporal features of video frames, the paper defines the relationship function in the attention mechanism using the cosine similarity theory, which can remove the absolute value of the relationship between features. To reduce the computational burden of the attention mechanism, the operation is decomposed from two dimensions of time and space. The computational complexity is further optimized by combining linear attention operation. The experiment is divided into two stages: Four ablation experiments are carried out on each module of directed attention to show the best performance of DAM in accuracy and efficiency; the accuracy of the model is 7.3% higher than that of I3D-NL on the Sth-Sth V1(something something V1) dataset and 95.7% on the UCF101(101 human action classes from videos in the wild) dataset. The research results have a wide application prospect in safety monitoring, automatic driving, and so on.

**Keywords:** action recognition; deep learning; attentional mechanism; cosine similarity; time-space decomposition

收稿日期: 2022-08-09 修回日期: 2022-10-18

基金项目: 江苏省重点研发计划(BE20200729, BE2019762, BE2021729); 军内科研项目(LJ2021Z2010032); 军队重点课题(JYKYA2021029)

第一作者: 李晨(1994-), 男, 硕士生, 研究方向为计算机视觉。E-mail: 651220007@qq.com

通讯作者: 何明(1978-), 男, 教授, 博导, 博士, 研究方向为物联网、无人指控等。E-mail: paper\_review@126.com

## 0 引言

当前, AR(action recognition)逐渐成为计算机视觉领域的热门研究课题, 具备着广泛应用前景, 例如识别作业人员行为, 增强安全生产场所防护能力。AR 不仅要关注帧内空间特征, 还需注意帧间时间特征, 以有效提取视频中行为时空特征, 实现精确识别。虽然 NL(non-local)<sup>[1]</sup>利用自注意力机制提取全局特征, 但其通过点积乘法得到的特征相似关系缺乏方向性, 而 AR 任务较依赖视频帧的时空顺序, 例如同一组视频帧若时间排序不同则会呈现不同的行为。另外对于时空位置较多的特征图, NL 会耗费较大计算量, 这显然不利于即插即用于 CNN(convolutional neural network)。因此, 如何构造出能够提取方向关系特征的高效注意力行为识别模型, 是待需解决的问题。

本文基于余弦相似度<sup>[2]</sup>提出了 DAM(directed attention model)。DAM 的核心是 DA-NL(directed attention non-local), 其使用余弦相似性提取特征

间的方向性全局关系。另外通过时空分解方式将 DA-NL 运算进行分离以提高计算效率, 进而结合 Poly-NL(polynomials non-local)<sup>[3]</sup>和 NL<sup>[1]</sup>设计了 6 种具备定向注意力的多样性全局模块, 使其在低于原始计算代价的情况下, 实现时空交互, 并通过对各模块的消融实验进一步分析了各模块性能表征。最终基于各定向注意力模块构造的 DAM 与最近的 AR 模型在精确度、计算量、参数量方面进行了对比, DAM 较 NL 在总计算成本上具备更高的效率, 并在定向注意力的加持下较 NL 具备更高的识别精确度。

## 1 相关工作

AR 发展初期, 以改进密集轨迹(improved dense trajectories, iDT)<sup>[4]</sup>为代表的传统人工方法占据主导。但随着 CNN 发展, 深度学习方法逐渐成为主流, AR 深度学习方法可从时间特征提取、高效率优化和全局特征捕获 3 个角度概括, 如表 1 所示。

表 1 AR 方法总览概况  
Table 1 Overview of AR methods

角度	方法	代表模型	优、缺点
时间特征提取	双流网络	Two-stream Network <sup>[5-8]</sup>	可提取时间特征, 但计算量高、不稳定、特征局部
	3D CNN	C3D <sup>[9]</sup> 、I3D <sup>[10]</sup>	可提取时空特征, 但计算量高, 无法提取全局特征
	时态模块	TSM <sup>[11]</sup> 、TRN <sup>[12]</sup> 、NLNN <sup>[1]</sup>	可灵活提取时空特征, 但无法高效联系全局信息
高效率优化	输入数据优化	TSN <sup>[13]</sup>	可降低效率, 但影响识别准确性
	时空分解卷积	P3D <sup>[14]</sup> 、S3D <sup>[15]</sup>	可降低效率, 但不利于模型最优迭代
	深度分离卷积	CSN <sup>[16]</sup>	可降低效率, 但缺少跨通道信息
	混合卷积	ECO <sup>[17]</sup> 、X3D <sup>[18]</sup>	可降低效率, 但前期训练工作较困难
全局特征捕获	全局均匀采样	TSN <sup>[13]</sup>	可捕获时间全局特征, 但缺乏对空间全局特征
	LSTM	I3D-LSTM <sup>[19]</sup>	增强全局表征, 但训练效率较低
	自注意力机制	NLNN <sup>[1]</sup> 、DirectFormer <sup>[2]</sup>	可提取全局特征, 但缺乏特征方向性和运行效率

### (1) 时间特征提取

相较图像识别, AR 不仅要关注空间特征, 还要聚焦时间特征理解运动信息。具体包含以下 3 个方面: ①基于光流的双流网络<sup>[5-8]</sup>。两条 2D CNN 路径以视频帧和堆叠光流图为输入分别提取空间及时间特征, 光流表征图像变化携带运动信息, 可有效提取时间特征, 实现了与 iDT 比拟的

识别效果。但是易受光线变化影响, 对存储量和计算量要求较高, 且小位移特性不利于识别全局特征。②基于三维卷积的 3D CNN<sup>[9-10, 20]</sup>。其通过将二维卷积扩展为三维卷积, 直接提取行为时空特征。如膨胀三维卷积网络(inflated 3D ConvNet, I3D)<sup>[10]</sup>将经过图像数据集预训练的 2D CNN 中的二维卷积核参数沿时间轴复制, 得到初始化的 3D

CNN, 为AR使用图像识别成熟架构提供了便利。但3D CNN参数量和计算量较高, 并且三维卷积的局部性, 导致无法实现对全局特征的有效提取。③基于时态建模的时态模块<sup>[1, 11-12, 21-23]</sup>。其不需光流和3D卷积扩展时间维度, 可灵活嵌入至各2D CNN, 具备有效的时空特征提取能力。例如, 时间位移模块(temporal shift module, TSM)<sup>[11]</sup>通过移位操作高效提取帧间信息关系, NL<sup>[1]</sup>利用自注意力机制交互时空位置间的特征关系。但TSM只能实现局部帧间特征关系的提取, NL计算量受限于时空位置数量, 不利于模块灵活嵌入。在时间特征提取方面, 本文DAM聚焦于时态模块方法, 设计了灵活的DA-NL模块嵌于ResNet50-C2D(residual network 50-convolution 2D)实现行为识别。

## (2) 高效率优化

AR作为图像识别的时序扩展, 计算量极大, 双流网络中光流需额外存储和计算, 3D CNN计算量更为昂贵, 基于自注意力机制的时态模块同样需耗费较大计算量。具体包含以下3个方面: ①输入数据优化。时间段网络(temporal segment networks, TSN)<sup>[13]</sup>对视频帧均匀采样, 降低信息冗余, 以较低代价实现学习。隐式双流网络<sup>[24]</sup>将能够从视频帧中产生类光流的移动网络与时间流CNN连接, 缓解了光流计算开销。但从数据端进行效率优化, 会损失特性信息, 影响识别准确性。②时空分解卷积。其将3D卷积分解, 近似为空间2D卷积和时间1D卷积的外积。例如, 伪3D(pseudo 3D, P3D)<sup>[14]</sup>使用 $1 \times 3 \times 3$ 和 $3 \times 1 \times 1$ 卷积模拟 $3 \times 3 \times 3$ 的3D卷积, 参数量得到显著降低, 且可利用2D CNN初始化训练。分离3D(separable 3D, S3D)<sup>[15]</sup>还采用Top-heavy方式简化特征量优化效率。但是, 硬性的时空分解不利于模型最优迭代, 影响识别精细度。③深度分离卷积。如通道分离卷积网络(channel separated convolutional networks, CSN)<sup>[16]</sup>, 其将 $T \times H \times W \times C$ 的3D卷积分解为 $T \times H \times W \times 1$ 的逐通道卷积和 $1 \times 1 \times 1 \times C$ 的逐点卷积, 两个部分

可在Bottleneck结构基础上共同作用优化模型效率。但其逐通道卷积缺少跨通道信息, 缺乏空间关联, 影响时空特征提取。④混合卷积。例如, 高效卷积网络(efficient convolutional network for online, ECO)<sup>[17]</sup>通过2D CNN得到特征图后串联3D CNN实现分类、快慢网络(SlowFast)<sup>[25]</sup>并联配置2D和3D卷积的双分支提取空间和时间特征, 扩展3D(X3D)<sup>[18]</sup>在深度、空间、时间、通道、采样等指标上自动逐步扩展反馈评估以生成高效CNN。本文DAM聚焦于时空分解和混合思想, 将DA-NL的注意力机制进行时空分解, 并混合了线性注意力降低计算量, 缓解了自注意力机制受限于时空位置数量的问题, 更利于模块灵活嵌入。

## (3) 全局特征捕获

基于卷积操作的CNN缺乏对全局特征捕获的能力, 而AR涉及范围更广的时空特征。具体包含以下3个方面: ①全局均匀采样。如TSN<sup>[13]</sup>和TDN(temporal difference networks)<sup>[26]</sup>通过稀疏采样得到时序上的全局采样帧, 实现全局特征提捕获, 但此策略缺乏对空间特征的全局把握。②长短时记忆(long short term memory, LSTM)。例如, 时间分段LSTM(temporal segment LSTM, TS-LSTM)<sup>[27]</sup>和I3D-LSTM<sup>[19]</sup>首先使用2D CNN提取空间特征, 再输入LSTM进行融合实现时序特征提取, 增强全局表征能力。但LSTM时序先后顺序的严格迭代比较影响训练效率。③自注意力机制。其能直接关注到任意时空位置间的全局信息, 在运算上具备很强的并行性。基于点积注意力机制的NL<sup>[1]</sup>可灵活嵌于CNN捕获全局特征, 视频Transformer网络(video transformer network, VTN)<sup>[28]</sup>和联合Transformer(unified transformer, UniFormer)<sup>[29]</sup>将CNN与Transformer串联关注全局信息, 但点积注意力的特征权值会削减特征间方向性, 不利于行为的方向特征提取。视频视觉Transformer(video vision transformer, ViViT)等<sup>[2, 30-33]</sup>完全摒弃CNN, 使用纯Transformer进行AR任务, 具备极强的全局特征捕获能力。但局部归纳性能偏弱, 不具备

CNN的平移不变性和局部性，需要大量数据才能泛化到AR任务上。本文聚焦于自注意力机制，设计了基于余弦相似性的定向非局部模块(DA-NL)以关注特征间的方向关系，并且可与CNN灵活组合实现快速拟合。

综上所述，本文首先聚焦于时态模块方法，基于NL设计了灵活的DA-NL模块嵌于ResNet50-C2D实现行为识别；其次聚焦于时空分解和混合思想，将DA-NL的注意力机制进行时空分解，并混合了线性注意力降低计算量，缓解了自注意力机制受限于时空位置数量的问题，更利于模块灵活嵌入；最后聚焦于自注意力机制，设计了基于余弦相似性的定向非局部模块(DA-NL)以关注特征间的方向关系，可与CNN灵活组合实现快速拟合。

## 2 定向注意力设计

本节首先阐述了基于余弦相似性<sup>[2]</sup>的DA-NL，其次对注意力进行时空分解，进而结合NL、Poly-NL优化定向注意力效率。

### 2.1 定向注意力模块

NL虽然通过点积乘法得到了各时空位置间的相似性权重，并聚合各特征重新校准了每个时空位置特征，但其相似度权重包含的特征绝对值会对相似性功能造成影响。因为在绝对值较大的情况下，即使特征向量间差异较大(即向量间夹角较大)，也会出现相似度权重较大情况，从而导致特征向量较相似的假象。所以，绝对值造成的相似性作弊现象会减弱各位置特征间的方向性，并通过聚合最终影响输出的注意力特征矩阵。

因此，本文基于余弦相似度计算特征向量间的方向性相似关系，提出DA-NL提取定向注意力全局特征。因相似特征之间的角度应该更小，而余弦相似度可以作为一种正则化机制去除向量幅度，增加特征向量幅度的动机较小，学习的权重总体上也更小，所以其更利于衡量特征间的方向相似度，能更直接地最大化相似特征。

给定一个输入特征张量 $\mathcal{X} \in \mathbb{R}^{T \times H \times W \times C}$ ，将其折叠为 $\mathbf{X} \in \mathbb{R}^{N \times C}$ 作为DA-NL的输入特征矩阵，其中 $N = T \times H \times W$ 即时空维度合并到 $N$ 。如等式(1)所示， $\mathbf{X}$ 分别转换为查询特征矩阵 $\mathbf{Q} \in \mathbb{R}^{N \times C_q}$ 、关键特征矩阵 $\mathbf{K} \in \mathbb{R}^{N \times C_k}$ 、值特征矩阵 $\mathbf{V} \in \mathbb{R}^{N \times C_v}$ 。

$$\begin{aligned}\mathbf{Q} &= \mathbf{X}\mathbf{W}_\theta \\ \mathbf{K} &= \mathbf{X}\mathbf{W}_\phi \\ \mathbf{V} &= \mathbf{X}\mathbf{W}_g\end{aligned}\quad (1)$$

式中： $\mathbf{W}_\theta \in \mathbb{R}^{C \times C_q}$ ， $\mathbf{W}_\phi \in \mathbb{R}^{C \times C_k}$ ， $\mathbf{W}_g \in \mathbb{R}^{C \times C_v}$ 是3个可学习参数线性矩阵。 $\mathbf{Q}$ 和 $\mathbf{K}$ 具有相等特征维度即 $C_q = C_k$ ，但两者互为不同以避免相似性对称。 $\mathbf{V}$ 中每个时空位置以相似度矩阵为权值，对所有位置的值特征进行加权求和，得到输出特征。

基于余弦相似性<sup>[2]</sup>计算 $\mathbf{Q}$ 和 $\mathbf{K}$ 的相似性度矩阵，如式(2)和(3)所示：

$$\cos(\mathbf{Q}_i, \mathbf{K}_j) = \frac{\mathbf{Q}_i \mathbf{K}_j}{|\mathbf{Q}_i| \times |\mathbf{K}_j|} = \cos \theta_{ij} \quad (2)$$

$$\cos \theta_{ij} = \frac{\sum_{c=1}^C \mathbf{Q}_{i,c} \times \mathbf{K}_{j,c}}{\sqrt{\sum_{c=1}^C (\mathbf{Q}_{i,c})^2} \times \sqrt{\sum_{c=1}^C (\mathbf{K}_{j,c})^2}} \quad (3)$$

式中： $\mathbf{Q}_i$ 和 $\mathbf{K}_j$ 分别为 $i$ 和 $j$ 时空位置的特征向量； $\theta_{ij}$ 为 $\mathbf{Q}_i$ 和 $\mathbf{K}_j$ 的夹角； $\mathbf{Q}_{i,c}$ 和 $\mathbf{K}_{j,c}$ 分别为 $N=i$ 、 $C=c$ 和 $N=j$ 、 $C=c$ 的向量特征值， $C$ 表示通道数。

由式(2)和(3)可得，余弦相似性去除了特征绝对值 $|\mathbf{Q}_i| \times |\mathbf{K}_j|$ 的影响，由 $\cos \theta_{ij}$ 独立实现。进而，基于余弦相似性权重，对所有位置的值特征向量加权求和：

$$\mathbf{Y}_i = \sum_{j=1}^N \cos(\mathbf{Q}_i, \mathbf{K}_j) \times \mathbf{V}_j = \sum_{j=1}^N \cos \theta_{ij} \times \mathbf{V}_j \quad (4)$$

式中： $\mathbf{V}_j$ 为 $j$ 时空位置的值特征向量； $\mathbf{Y}_i$ 为 $i$ 时空位置的注意力特征向量。由式(4)可得， $\mathbf{Y}_i$ 是所有位置 $\mathbf{V}_j$ 的求和， $\mathbf{V}_j$ 的权重是 $\cos \theta_{ij}$ 值。 $\mathbf{Y}_i$ 聚合了全部位置的值特征向量，各位置值特征向量的权重是对应位置键值特征向量和 $i$ 位置查询特征向量的夹角余弦值 $\cos \theta_{ij}$ 。因此，最终输出的注意力特征矩阵中，每个 $i$ 位置上的特征向量都聚合了所有位置特征向量在 $i$ 位置特征向量方向上的投影。这

种方式可以避免绝对值影响各位置向量的相似性, 更有助于  $\cos$  值的发挥, 从而增强自注意力机制的方向性。

为了使定向注意力机制能够即插即用在 CNN 中, 构造了 2 种定向注意力模块: DA-NL-1(如图 1 所示, 符号 ‘C’ 表示余弦相似性运算, ‘×’ 表示点积运算, 灰色矩形表示核大小为 1 的卷积和归一化函数)和 DA-NL-2。

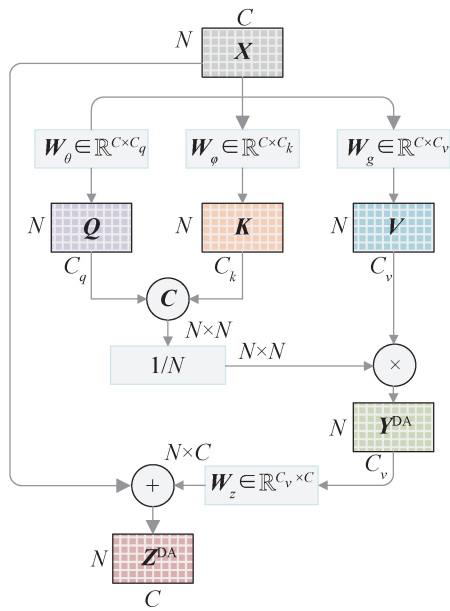


图1 DA-NL-1 定向注意力非局部块  
Fig. 1 Non-local module of directed attention DA-NL-1

(1) DA-NL-1: 使用余弦相似性函数  $\cos$ , 归一化函数  $U(x)$  如式(5)所示:

$$U(x) = \frac{x}{N} \quad (5)$$

$U(x)$  用来简化梯度计算,  $N$  是  $X$  的时空位置数。DA-NL-1 如式(6)和(7)所示:

$$Y^{DA-NL-1} = \frac{1}{N} \cos(Q, K) V =$$

$$\frac{1}{N} \cos(XW_\theta, XW_\phi) XW_g \quad (6)$$

$$Z^{DA-NL-1} = X + W_z Y^{DA-NL-1} \quad (7)$$

式中:  $W_z \in \mathbb{R}^{C_v \times C}$  将  $Y^{DA-NL-1}$  转换为与  $X$  相同的通道维度。 $Z^{DA-NL-1}$  中包含了原始  $X$  的特征信息和通过余弦相似性重新校准的各时空位置的特征。余弦相似性是每个时空位置贡献相似度权重, 权

重只包含双方位置间的  $\cos \theta_{ij}$  值, 学习的权重从整体上也较小, 能更直接地最大化相似特征。

(2) DA-NL-2: 与 DA-NL-1 区别在于归一化函数  $\text{softmax}$ , 如式 (8) 所示,  $\text{softmax}$  应用在  $A \odot \cos(Q, K)$  矩阵各行中。

$$\text{softmax}(A \odot \cos(Q, K)) = \frac{e^{A_{i,j} \cos(Q, K)}}{\sum_{n=1}^N e^{A_{i,n} \cos(Q, K_n)}} \quad (8)$$

DA-NL-2 如式(9)(10)所示:

$$Y^{DA-NL-2} = \text{softmax}(A \odot \cos(Q, K)) V =$$

$$\text{softmax}(A \odot \cos(XW_\theta, XW_\phi)) XW_g \quad (9)$$

$$Z^{DA-NL-2} = X + W_z Y^{DA-NL-2} \quad (10)$$

其中, 余弦相似度运算的结果范围为  $-1 \sim 1$ 。

仅使用余弦相似度项作为  $\text{softmax}$  函数的输入, 可能不足以对条件概率分布进行建模。因此, 为了增加每个条件概率项的可能概率值的范围, 添加了缩放超参数矩阵  $A$ ,  $\odot$  表示哈玛点积乘法。

## 2.2 时空分解定向注意力模块

DA-NL 计算复杂度为  $O(N^2)$ , 而 AR 任务中  $N = T \times H \times W$ , 并且包含去除绝对值操作的 DA-NL 计算量较 NL 更高。因此, 本节对 DA-NL 时空分解为 DA-DA(directed attention-directed attention)。

(1) 时空分解表示。前文将输入特征  $\mathcal{X} \in \mathbb{R}^{T \times H \times W \times C}$  折叠表示为  $X \in \mathbb{R}^{N \times C}$ , 本节将  $X$  时空分解表示为  $X \in \mathbb{R}^{(S,T) \times C}$ 。 $X_{s,t}$  表示  $X$  在时间位置  $t$  和空间位置  $s$  的  $C$  维向量, 其中  $1 \leq s \leq S (S = H \times W)$ ,  $1 \leq t \leq T$ 。

如图 2 所示, 符号 ‘C’ 和 ‘C’ 分别表示时间和空间的余弦相似性运算, ‘×’ 和 ‘×’ 表示时间和空间维度点积乘法运算。灰色矩形表示核大小为 1 的卷积和归一化函数, 以 DA-NL-1 为例, 基于  $X_{s,t}$  将其形式化为

$$\begin{cases} Q_{s,t} = \omega_\theta(X_{s,t}) \\ K_{s,t} = \omega_\phi(X_{s,t}) \\ V_{s,t} = \omega_g(X_{s,t}) \end{cases} \quad (11)$$

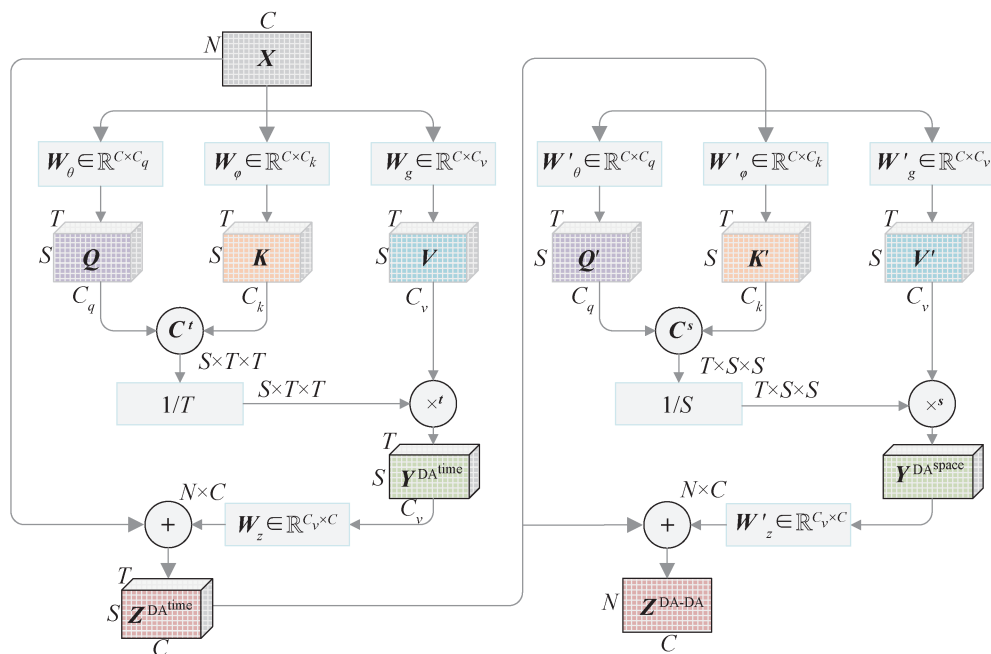


图 2 DA-DA 时空分解定向注意力非局部块

Fig. 2 Spatio-temporal decomposition of directed attention non-local module block DA-DA

$$Y_{s,t} = \sum_{s'=1}^S \sum_{t'=1}^T \frac{1}{N} \cos(Q_{s,t}, K_{s',t'}) V_{s',t'} \quad (12)$$

$$Z_{s,t} = X_{s,t} + \omega_z(Y_{s,t}) \quad (13)$$

式中： $\omega_\theta$ 、 $\omega_\phi$ 、 $\omega_\gamma$ 、 $\omega_z$ 是对应于 $W_\theta$ 、 $W_\phi$ 、 $W_\gamma$ 、 $W_z$ 的转换函数。 $Y_{s,t}$ 中需同时考虑时空维度，计算成本较大，因此将定向注意力在时间和空间维度上分解操作。

(2) DA<sup>time</sup>。将各空间位置的特征独立，DA-NL 限制在时间维度中，得到各空间位置的时间定向注意力特征，形式化为

$$Y_{s,t}^{\text{time}} = \sum_{t'=1}^T \frac{1}{T} \cos(Q_{s,t}, K_{s,t'}) V_{s,t'} \quad (14)$$

$$Z_{s,t}^{\text{time}} = X_{s,t} + \omega_z(Y_{s,t}^{\text{time}}) \quad (15)$$

(3) DA<sup>space</sup>。将各时间位置的特征独立，DA-NL 限制在空间维度中，得到各时间位置的空间定向注意力特征，形式化为

$$\begin{cases} Q'_{s,t} = \omega'_\theta(Z_{s,t}^{\text{time}}) \\ K'_{s,t} = \omega'_\phi(Z_{s,t}^{\text{time}}) \\ V'_{s,t} = \omega'_\gamma(Z_{s,t}^{\text{time}}) \end{cases} \quad (16)$$

$$Y_{s,t}^{\text{space}} = \sum_{s'=1}^S \frac{1}{S} \cos(Q'_{s,t}, K'_{s',t'}) V'_{s',t'} \quad (17)$$

$$Z_{s,t} = Z_{s,t}^{\text{time}} + \omega'_z(Y_{s,t}^{\text{space}}) \quad (18)$$

式中： $\omega'_\theta$ 、 $\omega'_\phi$ 、 $\omega'_\gamma$ 将 $Z_{s,t}^{\text{time}}$ 映射到新的查询特征 $Q'_{s,t}$ 、键值特征 $K'_{s,t}$ 和权值特征 $V'_{s,t}$ 。 $Y_{s,t}^{\text{space}}$ 通过与 $Z_{s,t}^{\text{time}}$ 残差连接形成最终的时空定向注意力特征图。

### 2.3 高效组合定向注意力模块

Poly-NL<sup>[3]</sup>和NL<sup>[1]</sup>的计算量较DA-NL更小。因此，如图3所示(图3(a)DA-Poly，对时空维度分别执行余弦相似性和哈玛点积运算点积乘法运算，‘C’表示时间维度余弦相似性运算，‘s’表示空间维度哈玛点积运算；图3(b)DA-DP，对时空维度分别执行余弦相似性和点积乘法运算，‘xs’表示空间维度点积乘法运算)，在DA-DA基础上，结合Poly-NL<sup>[3]</sup>、NL<sup>[1]</sup>共设计了6个定向注意力模块以改进效率，并在第4节进行了消融实验。

(1) DA-Poly: 首先，在时间维度进行DA-NL，如式(14)和(15)所示。

而后，在空间维度进行Poly-NL运算，如式(19)所示：



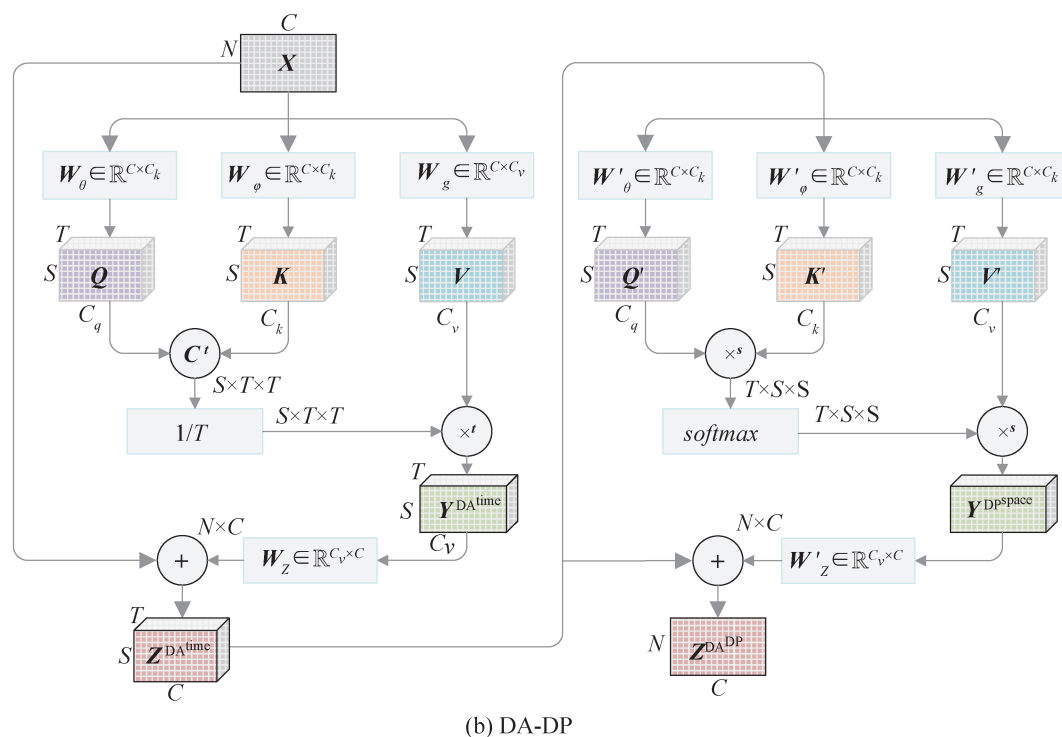
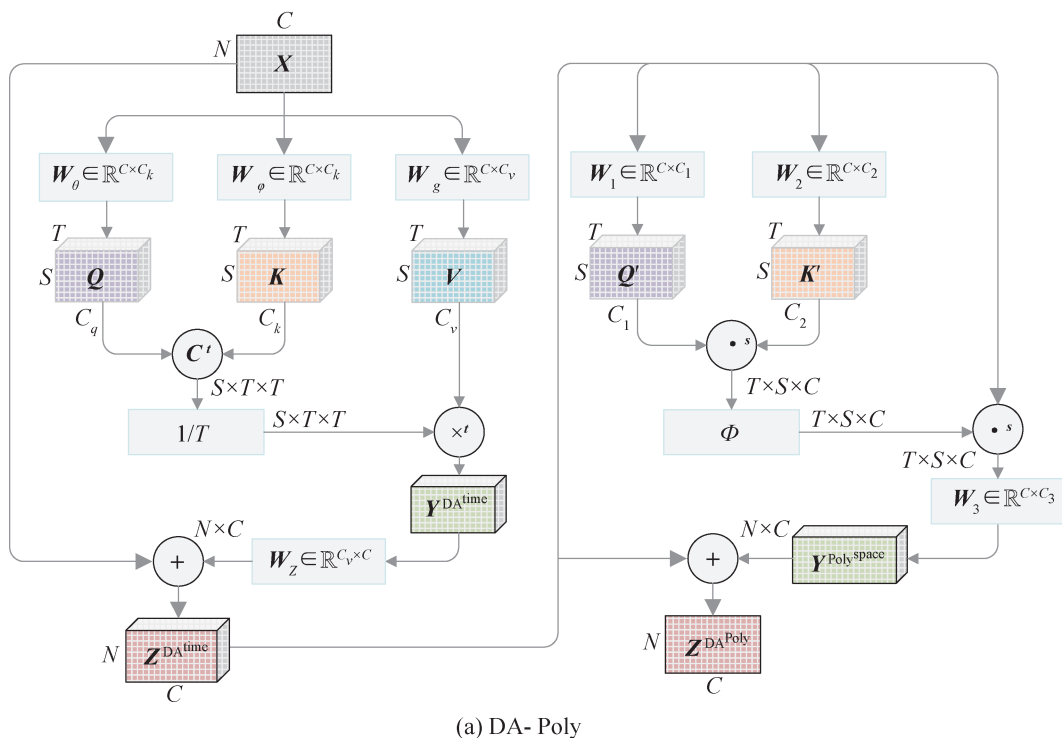


图 3 高效组合定向注意力模块 DA-Poly 和 DA-DP  
Fig. 3 Efficiently combined directed attention module DA-Poly and DA-DP

$$\begin{aligned}
 Y_{s,t}^{\text{space}} &= \omega_3 \left( \omega_e^{\text{space}} \left( \sum_{s'=1}^S \frac{1}{S} \times \right. \right. \\
 &\left. \left. \omega_1 \left( Z_{s',t}^{\text{time}} \right) \odot \omega_2 \left( Z_{s',t}^{\text{time}} \right) \odot Z_{s,t}^{\text{time}} \right) \right) \quad (19)
 \end{aligned}$$

式中:  $\omega_1$ 、 $\omega_2$ 、 $\omega_3$  为分别对应于  $W_1$ 、 $W_2$ 、 $W_3$  的转换函数 ( $W_1 \in \mathbb{R}^{C \times C_1}$ ,  $W_2 \in \mathbb{R}^{C \times C_2}$ ,  $W_3 \in \mathbb{R}^{C \times C_3}$  是 3 个可学习参数线性矩阵);  $\odot$  表示哈玛点积乘法;  $\omega_e^{\text{space}}$  是空间扩展功能将  $\mathbb{R}^{1 \times T \times C} \rightarrow \mathbb{R}^{S \times T \times C}$ 。  $\omega_3$  整合

最终的各位置特征,  $\mathbf{XW}_1$  与  $\mathbf{XW}_2$  具有相同的维度即  $C_1 = C_2$ 。

(2) DA-DP: 与 DA-Poly 类似, 同样在时间维度进行 DA-NL, 但在空间维度进行 NL, 如式(16)和(20)所示:

$$\mathbf{Y}_{s,t}^{\text{space}} = \sum_{s'=1}^S \text{softmax}(\mathbf{Q}'_{s,t}, \mathbf{K}'_{s',t}) \mathbf{V}'_{s',t} \quad (20)$$

(3) Poly-DA: 首先, 在时间维度进行 Poly-NL, 如式(21)所示:

$$\mathbf{Y}_{s,t}^{\text{time}} = \omega_3 \left( \omega_e^{\text{time}} \left( \sum_{t'=1}^T \frac{1}{T} \times \omega_1(\mathbf{X}_{s,t'}) \odot \omega_2(\mathbf{X}_{s,t'}) \right) \odot \mathbf{X}_{s,t} \right) \quad (21)$$

$$\mathbf{Z}_{s,t}^{\text{time}} = \mathbf{X}_{s,t} + \omega_z(\mathbf{Y}_{s,t}^{\text{time}})$$

式中:  $\omega_e^{\text{time}}$  是时间扩展功能将  $\mathbb{R}^{S \times 1 \times C} \rightarrow \mathbb{R}^{S \times T \times C}$ 。

而后, 在空间维度进行 DA-NL, 即式(16)和(17)。

(4) DP-DA: 首先, 在时间维度进行 NL, 如式(22)所示:

$$\mathbf{Y}_{s,t}^{\text{time}} = \sum_{t'=1}^T \text{softmax}(\mathbf{Q}_{s,t}, \mathbf{K}_{s,t'}) \mathbf{V}_{s,t'} \quad (22)$$

$$\mathbf{Z}_{s,t}^{\text{time}} = \mathbf{X}_{s,t} + \omega_z(\mathbf{Y}_{s,t}^{\text{time}})$$

而后, 在空间维度进行 DA-NL, 如式(16)和(17)所示。

### 3 定向注意力行为识别模型

本节使用了文献[1]中使用的 ResNet-C2D 和 ResNet-C3D 基准网络, 以便各版本定向注意力模块插入在 CNN 中形成 DAM。

#### 3.1 行为识别基准架构

(1) ResNet-C2D. 2D 卷积具有强大的空间特征表示能力, 并能通过池化层简单聚合时间特征。因此选择 2D ResNet<sup>[34]</sup> 为骨干网络, 并且对其调整以处理视频数据<sup>[1]</sup>。以 ResNet50 为例, 其适用于视频行为识别的 ResNet50-C2D 版本如表 2 所示(对每个块的第一个残差单元执行空间跨步<sup>[1, 34]</sup>)。

ResNet50-C2D 较原始 ResNet-50 的差异是, 所有  $k \times k$  的 2D 卷积都调整为  $1 \times k \times k$  的 3D 卷积形式, 并在 Conv<sub>1</sub> 之后添加时间池化 Pool<sub>1</sub> 将帧数减半, 全局平均池 Pool<sub>2</sub> 从 7×7 调整为 4×7×7 使得空间和时间特征同时聚合。类似的, 基于原始 ResNet101 改进了 ResNet101-C2D。

表 2 ResNet50-C2D 架构  
Table 2 Architecture of ResNet50-C2D

名称	输出尺寸	卷积核	步长
输入	32×224×224		
Conv <sub>1</sub>	16×112×112	1×7×7, 64	1, 2, 2
Pool <sub>1</sub>	8×56×56	3×3×3, max	1, 2, 2
Res <sub>1</sub>	8×56×56	$\begin{bmatrix} 1 \times 1 \times 1, & 64 \\ 1 \times 3 \times 3, & 64 \\ 1 \times 1 \times 1, & 256 \end{bmatrix} \times 3$	1, 1, 1
Pool <sub>2</sub>	4×56×56	3×1×1, max	2, 1, 1
Res <sub>2</sub>	4×28×28	$\begin{bmatrix} 1 \times 1 \times 1, & 128 \\ 1 \times 3 \times 3, & 128 \\ 1 \times 1 \times 1, & 512 \end{bmatrix} \times 4$	1, 2, 2
Res <sub>3</sub>	4×14×14	$\begin{bmatrix} 1 \times 1 \times 1, & 256 \\ 1 \times 3 \times 3, & 256 \\ 1 \times 1 \times 1, & 1024 \end{bmatrix} \times 6$	1, 2, 2
Res <sub>4</sub>	4×7×7	$\begin{bmatrix} 1 \times 1 \times 1, & 512 \\ 1 \times 3 \times 3, & 512 \\ 1 \times 1 \times 1, & 2048 \end{bmatrix} \times 3$	1, 2, 2
Pool <sub>3</sub>	1×1×1	4×7×7, average	1, 1, 1
Fc	1×1×1	2048×class	1, 1, 1

(2) ResNet-C3D. 通过引入额外的时间维度  $T$ , 将维度  $H \times W$  的 2D 卷积转换为维度  $T \times H \times W$  的 3D 卷积。传统 ResNet 中有两种主要类型的卷积核, 即 1×1 和 3×3。给定一个由 1×1 和 3×3 卷积组成的残差单元, 可将中间的 3×3 卷积核膨胀为 3×3×3(C3D<sub>3×3×3</sub>)<sup>[10]</sup>, 或将第一个 1×1 卷积核膨胀成 3×1×1(C3D<sub>3×1×1</sub>)<sup>[7]</sup>。文献[1]表明, C3D<sub>3×3×3</sub> 和 C3D<sub>3×1×1</sub> 实现了相当性能, 而后者参数较少, 计算效率更高。因此 ResNet-C3D 采用 C3D<sub>3×1×1</sub> 版本。

#### 3.2 定向注意力模型构建

基于 ResNet-C2D, 插入各定向注意力模块构建 DAM。如图 4 所示, DAM 是根据 4.3 节的消融实验设计, 基于 DA-NL、DA-DA 和 DA-Poly 模块

构建了多模块 DAM。预处理后的视频帧输入到 ResNet-C2D 基础架构中对每个帧单独处理, 最后对各帧的结果进行平均得到最终分类结果。在相同 2D CNN 框架下, DAM 较 C3D 具有更少的参数和计算成本。在卷积层的推理过程中, 特征图通过 2D CNN 进行独立运行, 不同之处是在架构的

不同位置插入了若干定向注意力模块, 较 3D 卷积而言能够提取到更强大的时空定向关系特征。DAM 的独特优势是它可轻松地将任何现成的 2D CNN 模型转换为可以同时处理空间和时间信息的伪 3D 定向全局模型, 与传统的 C3D 相比, 其在降低计算量的同时提取到了时空的定向全局关系。

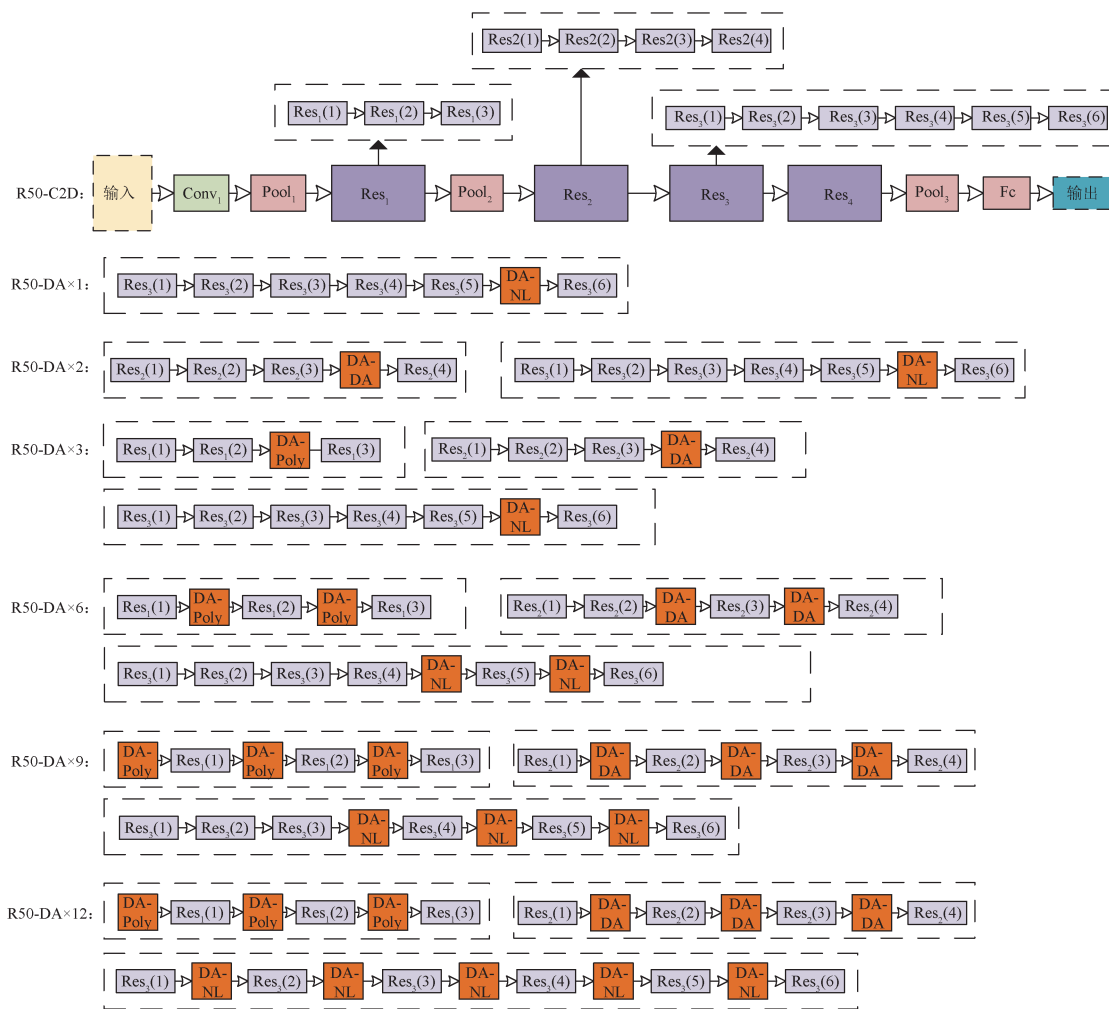


图4 基于不同数量定向模块的 DAM  
Fig. 4 DAM based on different number of orientation modules

## 4 实验评估

本节首先通过消融实验以在 ResNet 中合理使用各定向注意力模块, 从而构造 DAM; 进而将 DAM 与流行 AR 模型进行比较; 其次验证 DAM 的迁移学习能力。

### 4.1 数据集

Mini-Kinetics-400 数据集。考虑到数据量和硬件条件, 为提升实验速度基于 Kinetics-400<sup>[35]</sup> 轻量改造。包含 400 个动作类, 共有 23 156 个训练视频和 1 862 个验证视频, 共 12 G。

Something-Something V1<sup>[36]</sup> (Sth-Sth V1) 数据集

带有大量动作标签，更强调动作本身，包含人对日常对象执行的基本动作，动作类有174个，108 499个视频，持续时间2~6 s。

UCF101<sup>[37]</sup>共有13 320个剪辑视频，101个动作类，一个动作类的剪辑分为25组，每组包含4~7个剪辑，组中的剪辑具有一些共同的特征，例如背景或演员。动作类别包括人物交互、单人动作、人人交互、演奏乐器、运动。HMDB51<sup>[38]</sup>共包含6 849个剪辑视频，分为51个动作类别，每类别至少包含101个剪辑。类别包括面部动作、一般身体动作、物体交互动作、人类互动动作。本文使用 three train/test splits 策略<sup>[37]</sup>拆分UCF101和HMDB51数据集，以保持测试的一致性。此训练/测试拆分的设计方式是保持组分开，组内的剪辑是从单个长视频中获得的，因此在训练和测试中不共享来自同一组的剪辑。为了保证结果的客观准确，实验中显示了测试准确率平均值。

## 4.2 设置细节

DAM是在Ubuntu平台上使用Pytorch实现，具有一个NVIDIA RTX 3090显卡(24 GB RAM)。

(1) 训练。使用SGDM(SGD with Momentum)优化器，以最小化损失并更新包括权重、偏差等在内的参数。训练时为了获得更好的性能并加快处理速度，初始学习率设置为0.001，使用ExponentialLR按指数衰减调整学习率，dropout设置为0.5以减缓过拟合，动量设置为0.9，L2正则化权重衰减为 $5 \times 10^{-4}$ 。视频输入使用32个采样帧训练模型，这些采样帧是通过从原始视频中采样64帧，然后平均分成32个片段，从每个片段中挑选出一帧总共32帧。这些帧的较短边的大小固定为256，并利用边角裁剪和尺度抖动来增强数据，每个裁剪后的帧最终调整为 $224 \times 224$ 。用于训练模型的输入大小为 $N \times 32 \times 3 \times 224 \times 224$ ，其中 $N$ 是batch大小， $N$ 设置为16，所有模型都经过100个epoch的训练以进行微调，并选择最高的测试准确度进行报告。

(2) 测试。在消融实验中使用three-crop策略<sup>[26]</sup>进行测试，首先将每帧的较短边缩放为256，并从缩放帧中截取 $256 \times 256$ 的3个裁剪，从全长视频中随机采样10次，最终预测是所有剪辑的平均softmax分数。其他实验使用单采样和单中心 $224 \times 224$ 中心裁剪。

## 4.3 消融实验

本节对定向注意力模块在Mini-Kinetics-400数据集上消融实验。表3呈现了基于不同关系函数的单模块性能。表4和图5表现了各版本模块在不同位置上的准确性和效率。图6展示了基于各模块在不同数量情况下的各项性能，以构造DAM。

表3 基于不同关系函数的单模块性能对比  
Table 3 Single module performance comparison based on different relationship functions

模块	参数量/ MB	GFLOPS× View	Top1/ %	Top5/ %
基准模型	24.18	26.19×10×3	71.76	89.85
DP-NL <sup>[1]</sup>	25.56	26.47×10×3	72.68	90.47
Poly-NL <sup>[3]</sup>	25.33	26.20×10×3	72.65	90.46
DA-NL-1	25.53	27.02×10×3	73.21	90.73
DA-NL-2	25.57	27.04×10×3	73.18	90.71

如表3所示(插入位置是Res<sub>3</sub>中最后一个单元的前面，架构：Res50-C2D，帧数×view：32×256×256×10×3)，DA-NL在识别精度上较DP-NL<sup>[1]</sup>和Poly-NL<sup>[3]</sup>有明显提升，说明基于余弦相似性的定向注意力可更为有效表征各时空位置特征间的关系。但DA-NL较DP-NL和Poly-NL计算量稍大，不利于通过多帧输入和多模块途径，进一步提高识别精度。另外，DA-NL-1和DA-NL-2在精确度上未存在明显差异，这也许是因定向注意力去除了绝对值，不需softmax函数对向量幅度弱化。此后，DA-NL均为DA-NL-1。

表4(插入位置是各Res块中最后一个单元的前面，输入：32×256×256×10×3，架构：Res50-C2D)比较了不同Res层的各定向注意力时空分解模块：DA-NL、DA-DA、DA-Poly、DA-DP、Poly-DA、DP-DA。

表 4 各组合定向注意力模块在 Res 层的性能比较  
Table 4 Performance comparison of each combined directed attention module in Res layer

位置	模块	参数量/MB	GFLOPS× View	Top1/ %	Top5/ %
Res <sub>1</sub>	基准模型	24.18	26.19×10×3	71.76	89.85
	DA-NL	24.27	<b>237.92</b> ×10×3	<b>72.94</b>	<b>90.43</b>
	DA-DA	25.19	52.72×10×3	72.72	90.35
	DA-Poly	24.83	<b>26.26</b> ×10×3	<b>72.43</b>	90.27
	DA-DP	25.22	35.08×10×3	72.46	<b>90.26</b>
	Poly-DA	24.83	52.66×10×3	72.68	90.32
	DP-DA	25.22	52.68×10×3	72.65	90.31
Res <sub>2</sub>	DA-NL	24.51	<b>32.81</b> ×10×3	<b>73.07</b>	<b>90.58</b>
	DA-DA	26.23	27.85×10×3	72.94	90.41
	DA-Poly	26.08	<b>26.20</b> ×10×3	72.68	90.32
	DA-DP	26.25	26.75×10×3	<b>72.65</b>	<b>90.31</b>
	Poly-DA	26.08	27.84×10×3	72.90	90.38
	DP-DA	26.25	27.84×10×3	72.91	90.37
Res <sub>3</sub>	DA-NL	25.53	<b>27.02</b> ×10×3	<b>73.21</b>	<b>90.73</b>
	DA-DA	27.98	26.40×10×3	72.84	90.53
	DA-Poly	27.35	<b>26.19</b> ×10×3	72.53	90.39
	DA-DP	30.02	26.26×10×3	<b>72.49</b>	<b>90.37</b>
	Poly-DA	27.35	26.39×10×3	72.76	90.40
Res <sub>4</sub>	DA-NL	29.68	<b>26.29</b> ×10×3	<b>72.51</b>	<b>90.25</b>
	DA-DA	36.18	26.22×10×3	72.34	90.15
	DA-Poly	35.82	<b>26.19</b> ×10×3	<b>72.01</b>	<b>89.98</b>
	DA-DP	36.24	26.20×10×3	72.04	90.02
	Poly-DA	35.82	26.22×10×3	72.28	90.12
	DP-DA	36.24	26.22×10×3	72.31	90.10

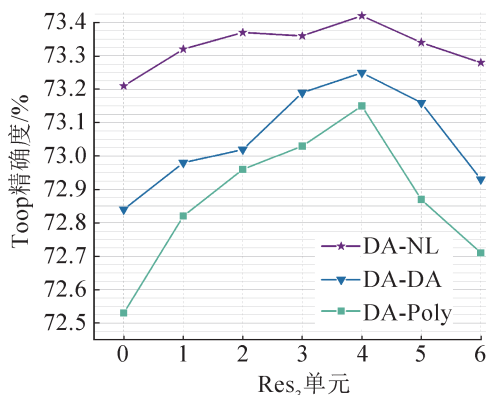


图 5 在 Res<sub>3</sub> 插入双模块的不同方案  
Fig. 5 Different options for inserting dual modules in Res<sub>3</sub>

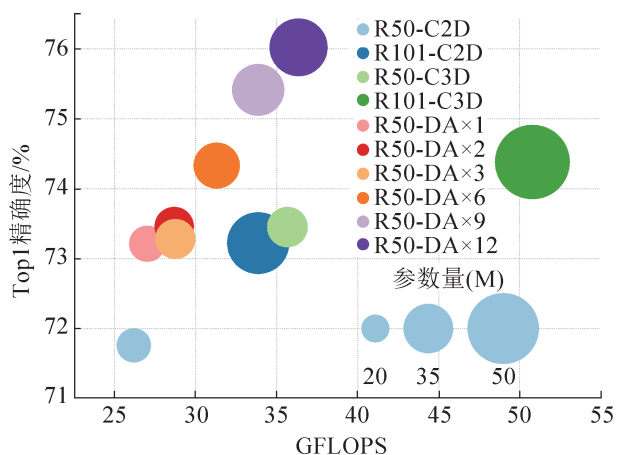


图 6 DAM 与基准架构的比较  
Fig. 6 Comparison of DAM with the benchmark architectures

根据 Res 层纵向分析, ResNet50-C2D 的性能指标有以下规律。参数量随模块插入位置的深度增加而逐渐提升,这可能是因深层特征通道数  $C$  较浅层特征变大,导致模块中变换矩阵参数量更大;计算量随插入层位置深度增加而逐渐降低,这可能是深层特征时空位置  $N$  逐渐减小的原因;精确度随插入位置深度增加逐渐提升,并在 Res<sub>4</sub> 层出现了回落,这可能是因 Res<sub>4</sub> 特征时空位置太少,导致特征间的关系特征提取不充分,不能充分发挥注意力机制效能。

另外,各 Res 层中出现共性情况。DA-NL 精确度最高,说明其提取特征间注意力关系最为充分;DA-DA 较 DA-NL 的计算量明显降低,精确度仅次于未分解的 DA-NL,证明时空分解操作的有效性;DA-Poly 的计算量最低且保证了精度,说明定向注意力和线性注意力结合的高效性;Poly-DA 和 DP-DA 较 DA-DA,并未在效率和精度上表现明显优势,这可能是因特征图的空间位置远大于时间位置,所以仅在时间维度进行 Poly-NL 不足以改进计算量。

其中:在 Res<sub>1</sub> 层,DA-NL 较基准架构精度提升了 1.22%,但 GFLOPS 高达 237.92 即 ResNet50-C2D 的近 10 倍计算量,这显然不能满足实际应用;DA-DA 将 GFLOPS 降低到了基准线的 2 倍即 52.72,并具备 1.02% 的精度提升;DA-Poly 的

GFLOPS与基线几乎相同即26.26,同时具有0.67%精度提升;在Res<sub>2</sub>层,DA-DA和DA-Poly较DA-NL在计算量上分别降低了4.96和6.61,在精确度上分别降低了0.13和0.42。Res<sub>3</sub>和Res<sub>4</sub>层,DA-DA和DA-Poly较DA-NL计算量降低并不明显。

图5(横坐标‘0’:模块1插入在第6个Res<sub>3</sub>单元前,横坐标‘1~6’:第2个模块插入在Res<sub>3</sub>单元前的位置,输入:32×256×256×30,架构:Res50-C2D)呈现了Res<sub>3</sub>中插入双模块的不同方案。无间隔插入(Res<sub>3</sub>单元为6)的精确度不如间隔插入(Res<sub>3</sub>单元为1~5)的精确度高,其中间隔1个res块(Res<sub>3</sub>单元为5)和间隔2个res块(Res<sub>3</sub>单元为4)的增益较好。这可能是因无间隔插入缺少了卷积的归纳偏置能力,导致双模块提取的关系特征较单模块增益不明显。

图6呈现了模块数量为1、3、6、9、12情况下的DAM效能。DAM的精确度随着模块的使用数量逐步提升,其中:DA×1的Top1精确度与R101-C2D相当,且参数量和GFLOPS更低,证明了使用单个定向注意力模块比加深CNN的深度更有利于提升识别效果,并且效率更高。DA×6的精确度超过R50-C3D,计算量是R50-C3D的87%,参数量是其113%,说明较于传统3D卷积,使用定向注意力模块能够取得更高的精确度以及计算效率。DA×9的精确度高于R101-C3D,且计算量是其67%,参数量是其70%。综上所述,多定向注意力模块设计的DAM,较传统的网络延伸和3D卷积方法,在精度和效率上都具备优势。

#### 4.4 与现有技术的比较

为了与其他方法公平比较,在Sth-Sth V1数据集<sup>[36]</sup>上使用单帧和单中心224×224裁剪的测试方案,预训练使用ImageNet数据集<sup>[39]</sup>和Kinetics-400数据集<sup>[35]</sup>。DAM采用了图6中12模块的版本,并测试了8、16和32帧的3种情况选择的对比模型算法,均是使用了ResNet50作为基准架构,具备相同的网络深度,计算量和参数量都在同一层级,便于比较各模型性能。另外,选择的对比模

型,涉及AR领域的多个设计方向,包括:纯2D卷积的C2D、具备光流的TSN;基于3D卷积的I3D、时空分解卷积的S3D和混合卷积的ECO;基于时态模块的TRN、TSM、TEA;基于自注意机制的TDN和NL。各对比模型算法的实验训练参数依照原配置,如表5所示(“\*”表示未发布源代码,“I”表示ImagNet,“K400”表示Kinetics-400,“S1M”表示Sports1M)。

如表6所示,传统2D CNN方法中,C2D<sup>[11]</sup>的Top1精确度仅为18.7%,TSN<sup>[13]</sup>虽利用稀疏采样策略提取全局特征,但精确度仅为19.7%,可知2D卷积时间建模较弱。而DAM基于定向注意力机制,在2D CNN架构下实现了48.2%~51.7%的精确度,较C2D和TSN实现了极大提升,说明定向注意力较2D卷积具备较强的时间特征提取能力。另外,3D CNN方法中,C3D<sup>[11]</sup>、ECO<sup>[17]</sup>、I3D<sup>[40]</sup>和S3D-G<sup>[15]</sup>虽然较传统2D CNN具备一定的时间特征提取能力,但其精确度较DAM仍有一定差距,其中只有S3D-G<sup>[15]</sup>在识别效果上与DAM相近,但其计算量远大于8帧输入的DAM,说明3D卷积只是实现了局部时空特征融合,而DAM通过定向注意力更好地实现了全局时空特征融合。最后在时态卷积方法中,在同样的32帧输入下,DAM的精确度较I3D+NL<sup>[40]</sup>高出7.3%,而计算量仅需其10%,说明定向注意力较传统注意力具备更好的全局特征提取能力,以及各模块的搭配设计使得效率得到了很好地优化。DAM较TRN<sup>[12]</sup>和TSM<sup>[11]</sup>识别效果更优异,较TEA<sup>[21]</sup>和TDN<sup>[26]</sup>在计算量上更为轻量,但精确度还有一定差距。

#### 4.5 UCF101的迁移学习

本节将4.4节预训练的DAM转移到更小的UCF101<sup>[37]</sup>和HMDB51<sup>[38]</sup>数据集,并在UCF101和HMDB51上进行了微调,以表明DAM学习的表示可推广到其他数据集。表7中包含的结果是three train/test splits平均类准确度<sup>[37]</sup>。很明显,DAM展示了非常强大的迁移学习能力,平均分类准确率高达95.7%和71.8%。

<http://www.china-simulation.com>

表 5 各对比模型在 Sth-Sth V1 数据集上的实验训练参数

Table 5 Experimental training parameters of each comparison model on Sth-Sth V1 dataset

模型	主干架构	预训练	初始学习率	Dropout	动量	权重衰减	帧数	分辨率	周期	Batch
C2D <sup>[1]</sup>	ResNet50	I+K400	0.001	0.5	0.9	$5 \times 10^{-4}$	32	224×224	100	16
TSN <sup>[13]</sup>	BNInception	I	0.020	0.8	0.9	$1 \times 10^{-4}$	8	224×224	50	8
TSN <sup>[13]</sup>	ResNet50	I	0.020	0.8	0.9	$1 \times 10^{-4}$	8	224×224	50	8
C3D <sup>[1]</sup>	ResNet50	I+K400	0.001	0.5	0.9	$5 \times 10^{-4}$	32	224×224	100	16
ECO <sup>[17]</sup>	BNIncep+R18	I+K400	0.001	0.5	0.9	$5 \times 10^{-4}$	8	224×224	100	32
I3D <sup>[40]</sup>	ResNet50	I+K400	0.001 25	0.3	0.9	$1 \times 10^{-4}$	32	224×224	100	8
I3D+NL <sup>[40]</sup>	ResNet50	I+K400	0.001 25	0.3	0.9	$1 \times 10^{-4}$	32	224×224	100	8
S3D-G <sup>[15]</sup>	BN-Inception	I	0.100	0.5	0.9	$1 \times 10^{-4}$	64	224×224	100	6
TRN <sup>[12]</sup>	BNInception	I	0.001	0.5	0.9	$1 \times 10^{-4}$	8	224×224	100	10
TRN <sup>[12]</sup>	ResNet50	I	0.001	0.5	0.9	$1 \times 10^{-4}$	8	224×224	100	10
TSM <sup>[11]</sup>	ResNet50	I+K400	0.010	0.5	0.9	$1 \times 10^{-4}$	8	224×224	50	64
TEA* <sup>[21]</sup>	ResNet50	I					16	224×224		
TDN <sup>[26]</sup>	ResNet50	I+S1M	0.020	0.5	0.9	$1 \times 10^{-4}$	8	224×224	60	128
DAM	ResNet50	I+K400	0.001	0.5	0.9	$5 \times 10^{-4}$	32	224×224	100	16

表 6 DAM 在 Sth-Sth V1 数据集上与其他方法的测试比较

Table 6 Comparison of DAM on Sth-Sth V1 dataset with tests of other methods

模型	主干架构	帧数×View	分辨率	参数量/MB	GFLOPS×View	Top1	Top5
C2D <sup>[1]</sup>	ResNet50	32×1×1	224×224	24.18	26.19×1×1	18.7	45.3
TSN <sup>[13]</sup>	BNInception	8×1×1	224×224	10.70	16×1×1	19.5	
TSN <sup>[13]</sup>	ResNet50	8×1×1	224×224	24.30	33×1×1	19.7	46.6
C3D <sup>[1]</sup>	ResNet50	32×1×1	224×224	28.51	35.67×1×1	32.8	60.3
ECO <sup>[17]</sup>	BNIncep+R18	8×1×1	224×224	47.50	32×1×1	39.6	
I3D <sup>[40]</sup>	ResNet50	32×2×1	224×224	28.00	153×2×1	41.6	72.2
I3D+NL <sup>[40]</sup>	ResNet50	32×2×1	224×224	35.30	168×2×1	44.4	76.0
S3D-G <sup>[15]</sup>	BN-Inception	64×1×1	224×224	11.56	71.38×1×1	48.2	78.7
TRN <sup>[12]</sup>	BNInception	8×1×1	224×224	18.30	16×1×1	34.4	
TRN <sup>[12]</sup>	ResNet50	8×1×1	224×224	31.80	33×1×1	38.9	68.1
TSM <sup>[11]</sup>	ResNet50	8×1×1	224×224	24.30	33×1×1	45.6	74.2
TEA <sup>[21]</sup>	ResNet50	16×1×1	224×224		70×1×1	<b>51.9</b>	<b>80.3</b>
TDN <sup>[26]</sup>	ResNet50	8×1×1	224×224		36×1×1	<b>52.3</b>	<b>80.6</b>
<b>DAM</b>	ResNet50	8×1×1	224×224	40.38	9.09×1×1	48.2	76.5
<b>DAM</b>	ResNet50	16×1×1	224×224	40.38	18.18×1×1	49.5	78.3
<b>DAM</b>	ResNet50	32×1×1	224×224	40.38	36.36×1×1	<b>51.7</b>	<b>80.1</b>

表 7 DAM 在 UCF101 和 HMDB51 上的平均类精度  
Table 7 Average class accuracy of DAM on UCF101 and HMDB51

模型	预训练	UCF101	HMDB51
C3D <sup>[1]</sup>	Kinetics400	84.6	54.3
I3D <sup>[40]</sup>	Kinetics400	91.8	66.4

为了更好地理解 DAM 如何学习用于行为识别的判别性时空描述符, 使用 CAM 方法<sup>[41]</sup>可视化模

型的特定类激活图。在实验中, 将 T 设置为 3 从 UCF101 测试集中的视频序列中均匀采样 3 个片段, 以获得特定类的激活图。作为比较还可可视化了 I3D 模型的激活图, 该模型利用局部时空信息, 而不是联合建模局部和全局时空信息。作为说明, 图 7(上: 原始视频帧, 中: I3D, 下: DAM)列出了模型的 3 个行为类别的激活图。DAM 关注人的

运动或者手部交互行为，如在打保龄球这个行为中，其重点关注了保龄球与抓球的手部。与 I3D 模型相比，DAM 能够更精确地关注行为中的与运

动位移相关的区域，而 I3D 缺乏对运动相关区域的精确关注，这依赖于 DAM 强大的时间建模能力。

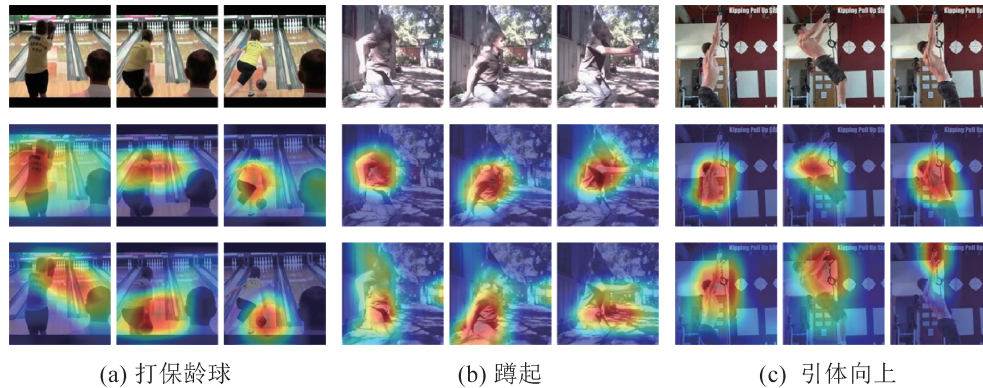


图 7 使用 CAM 可视化特定于行为的激活图  
Fig. 7 Activation diagrams of specific actions visualized by CAM

激活图表明，与无法联合建模视频中局部和全局时空特征的 I3D 相比，DAM 可以捕获视频帧中的特征关系，其更侧重于关注与行为密切相关的时空区域。例如，DAM 会更加关注附近有保龄球的人体，而只有人体的区域不会被激活。另外在“引体向上”示例中，DAM 通过单杠和手部激活，然而 I3D 被大面积人体区域激活。

## 5 结论

本文提出并实现了定向注意力模型(DAM)进行 AR 任务，其中基于余弦相似性的定向注意力模块可去除特征绝对值从而更有效表示特征间关系，且时空分解降低了注意力运算的计算量。DAM 较传统 2D、3D 和时态卷积，在不加大计算量的情况下表现出强大的时空特征表示能力。同时也存在一些不足：本文所提出的创新模块泛化性能和对于各模块的效率-精度平衡问题还有待进一步探索。在未来的工作中，将在更多的骨干网络上进行实验，验证所提出的模块的泛化能力；同时，将基于效率-精度平衡问题，探索如何自动生成各模块在骨干网络上的使用方案。

## 参考文献：

[1] Wang Xiaolong, Girshick R, Gupta A, et al. Non-local

Neural Networks[C]//2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. Piscataway, NJ, USA: IEEE, 2018: 7794-7803.

[2] Truong T D, Quoc-Huy Bui, Chi Nhan Duong, et al. DirecFormer: A Directed Attention in Transformer Approach to Robust Action Recognition[C]//2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2022: 19998-20008.

[3] Babiloni F, Marras I, Kokkinos F, et al. Poly-NL: Linear Complexity Non-local Layers with 3rd Order Polynomials[C]//2021 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2021: 10498-10508.

[4] Wang Heng, Cordelia Schmid. Action Recognition with Improved Trajectories[C]//2013 IEEE International Conference on Computer Vision. Piscataway, NJ, USA: IEEE, 2013: 3551-3558.

[5] Simonyan K, Zisserman A. Two-stream Convolutional Networks for Action Recognition in Videos[C]//Proceedings of the 27th International Conference on Neural Information Processing Systems. Cambridge, MA, USA: MIT Press, 2014: 568-576.

[6] Christoph Feichtenhofer, Axel Pinz, Zisserman A. Convolutional Two-stream Network Fusion for Video Action Recognition[C]//2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2016: 1933-1941.

[7] Christoph Feichtenhofer, Axel Pinz, Richard P Wildes. Spatiotemporal Residual Networks for Video Action Recognition[C]//Proceedings of the 30th International



- Conference on Neural Information Processing Systems. Red Hook, NY, USA: Curran Associates Inc., 2016: 3476-3484.
- [8] Christoph Feichtenhofer, Axel Pinz, Richard P Wildes. Spatiotemporal Multiplier Networks for Video Action Recognition[C]//2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2017: 7445-7454.
- [9] Tran D, Bourdev L, Fergus R, et al. Learning Spatiotemporal Features with 3D Convolutional Networks[C]//2015 IEEE International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2015: 4489-4497.
- [10] João Carreira, Zisserman A. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset[C]//2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2017: 4724-4733.
- [11] Lin Ji, Gan Chuang, Han Song. TSM: Temporal Shift Module for Efficient Video Understanding[C]//2019 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2019: 7082-7092.
- [12] Zhou Bolei, Andonian A, Oliva A, et al. Temporal Relational Reasoning in Videos[C]//Computer Vision-ECCV 2018. Cham: Springer International Publishing, 2018: 831-846.
- [13] Wang Limin, Xiong Yuanjun, Wang Zhe, et al. Temporal Segment Networks: Towards Good Practices for Deep Action Recognition[C]//Computer Vision-ECCV 2016. Cham: Springer International Publishing, 2016: 20-36.
- [14] Qiu Zhaofan, Yao Ting, Mei Tao. Learning Spatiotemporal Representation with Pseudo-3D Residual Networks[C]//2017 IEEE International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2017: 5534-5542.
- [15] Xie Saining, Sun Chen, Huang J, et al. Rethinking Spatiotemporal Feature Learning: Speed-accuracy Trade-offs in Video Classification[C]//Computer Vision-ECCV 2018. Cham: Springer International Publishing, 2018: 318-335.
- [16] Tran D, Wang Heng, Feiszli M, et al. Video Classification with Channel-separated Convolutional Networks[C]//2019 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2019: 5551-5560.
- [17] Mohammadreza Zolfaghari, Kamaljeet Singh, Thomas Brox. ECO: Efficient Convolutional Network for Online Video Understanding[C]//Computer Vision-ECCV 2018. Cham: Springer International Publishing, 2018: 713-730.
- [18] Feichtenhofer C. X3D: Expanding Architectures for Efficient Video Recognition[C]//2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2020: 200-210.
- [19] Wang Xianyuan, Miao Zhenjiang, Zhang Ruyi, et al. I3D-LSTM: A New Model for Human Action Recognition [C]//IOP Conference Series: Materials Science and Engineering. Bristol, United Kingdom: IOP Publishing, 2019: 032035.
- [20] Tran D, Ray J, Shou Z, et al. Convnet Architecture Search for Spatiotemporal Feature Learning[J]. Computing Research Repository, 2017, 16(8): 1-12.
- [21] Li Yan, Ji Bin, Shi Xintian, et al. TEA: Temporal Excitation and Aggregation for Action Recognition[C]//2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2020: 906-915.
- [22] 罗会兰, 陈翰. 时空卷积注意力网络用于动作识别[J]. 计算机工程与应用, 2023, 59(9): 150-158.
- Luo Huilan, Chen Han. Spatial-temporal Convolutional Attention Network for Action Recognition[J]. Computer Engineering and Applications, 2023, 59(9): 150-158.
- [23] 吴丽君, 李斌斌, 陈志聪, 等. 3D多重注意力机制下的行为识别[J]. 福州大学学报(自然科学版), 2022, 50(1): 47-53.
- Wu Lijun, Li Binbin, Chen Zhicong, et al. Action Recognition Under 3D Multiple Attention Mechanism[J]. Journal of Fuzhou University(Natural Science Edition), 2022, 50(1): 47-53.
- [24] Zhu Yi, Lan Zhenzhong, Newsam S, et al. Hidden Two-stream Convolutional Networks for Action Recognition [C]//Computer Vision-ACCV 2018. Cham: Springer International Publishing, 2019: 363-378.
- [25] Feichtenhofer C, Fan Haoqi, Malik J, et al. SlowFast Networks for Video Recognition[C]//2019 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2019: 6201-6210.
- [26] Wang Limin, Tong Zhan, Ji Bin, et al. TDN: Temporal Difference Networks for Efficient Action Recognition [C]//2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2021: 1895-1904.
- [27] Ma C Y, Chen M H, Kira Z, et al. TS-LSTM and Temporal-inception: Exploiting Spatiotemporal Dynamics for Activity Recognition[J]. Signal Processing: Image Communication, 2019, 71: 76-87.
- [28] Neimark D, Bar O, Zohar M, et al. Video Transformer Network[C]//2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW). Piscataway,

- NJ, USA: IEEE, 2021: 3156-3165.
- [29] Li Kunchang, Wang Yali, Gao Peng, et al. UniFormer: Unified Transformer for Efficient Spatiotemporal Representation Learning[EB/OL]. (2022-02-08) [2022-04-04]. <https://arxiv.org/abs/2201.04676>.
- [30] Arnab A, Dehghani M, Heigold G, et al. ViViT: A Video Vision Transformer[C]//2021 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2021: 6816-6826.
- [31] Fan Haoqi, Xiong Bo, Mangalam K, et al. Multiscale Vision Transformers[C]//2021 IEEE/CVF International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2021: 6804-6815.
- [32] Li Yanghao, Wu Chaoyuan, Fan Haoqi, et al. MViTv2: Improved Multiscale Vision Transformers for Classification and Detection[EB/OL]. (2022-03-30) [2022-04-04]. <https://arxiv.org/abs/2112.01526>.
- [33] Alfasly S, Lu Jian, Xu Chen, et al. Learnable Irrelevant Modality Dropout for Multimodal Action Recognition on Modality-specific Annotated Videos[EB/OL]. (2022-03-27) [2022-04-04]. <https://arxiv.org/abs/2203.03014>.
- [34] He Kaiming, Zhang Xiangyu, Ren Shaoqing, et al. Deep Residual Learning for Image Recognition[C]//2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2016: 770-778.
- [35] Kay W, João Carreira, Simonyan K, et al. The Kinetics Human Action Video Dataset[EB/OL]. (2017-05-19) [2022-04-04]. <https://arxiv.org/abs/1705.06950>.
- [36] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, et al. The "Something Something" Video Database for Learning and Evaluating Visual Common Sense[C]//2017 IEEE International Conference on Computer Vision (ICCV). Piscataway, NJ, USA: IEEE, 2017: 5843-5851.
- [37] Soomro K, Zamir A R, Shah M. UCF101: A Dataset of 101 Human Actions Classes from Videos in the Wild[EB/OL]. (2013-12-03) [2022-04-04]. <https://arxiv.org/abs/1212.0402>.
- [38] Kuehne H, Jhuang H, Garrote E, et al. HMDB: A Large Video Database for Human Motion Recognition[C]//2011 International Conference on Computer Vision. Piscataway, NJ, USA: IEEE, 2011: 2556-2563.
- [39] Deng Jia, Dong Wei, Socher R, et al. ImageNet: A Large-scale Hierarchical Image Database[C]//2009 IEEE Conference on Computer Vision and Pattern Recognition. Piscataway, NJ, USA: IEEE, 2009: 248-255.
- [40] Wang Xiaolong, Gupta A. Videos as Space-time Region Graphs[C]//Computer Vision-ECCV 2018: 15th European Conference. Heidelberg: Springer-Verlag, 2018: 413-431.
- [41] Zhou Bolei, Khosla A, Lapedriza A, et al. Learning Deep Features for Discriminative Localization[C]//2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Piscataway, NJ, USA: IEEE, 2016: 2921-2929.