

# 时空行为理解

章毓晋

清华大学电子工程系 100084 北京

## 总目录

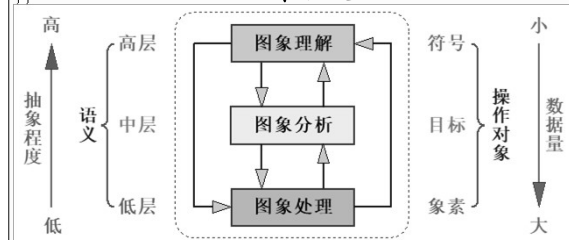
- ◆ 引言 (+ 课程历史)
- ◆ 正文
  - 时空行为理解层次
  - 从时空兴趣点到轨迹
  - 动作/活动识别 (示例)
  - 深度学习与行为识别 (示例)
- ◆ 结语

## 引言

- **图象:**  
用各种观测系统以不同形式和手段观测客观世界而获得的, 可以直接或间接作用于人眼并进而产生视觉的实体  
图象 (广义/抽象) ⊃ 图像 (狭义/具体)
- **图象和信息:**  
人类从外界 (客观世界) 获得的信息约有 75% 来自视觉系统

## 引言

- **图象工程:** 不同层次图象技术的有机结合及应用。 **三个层次**





## 引言

- **图象工程文献综述系列**
    - 每年《中国图象图形学报》5月那一期
    - 已历时 24 年, 涉及 15 种刊物, 15095 (62186) 篇论文
- 主要目的**
- (1) 概括我国图象工程发展现状
  - (2) 帮助读者查阅有关研究文献
  - (3) 对期刊编者和论文作者提供参考

## 引言

- **图象工程研究文献**
  - 图象处理 (图象 ⇒ 图象)
  - 图象分析 (图象 ⇒ 数据)
  - 图象理解 (图象 ⇒ 解释)
  - 技术应用 (图象处理、分析、理解技术的工程实现和应用)
  - 综述 (综合图象处理、分析、理解)

© Y.J.ZHANG.  



## 引 言

- 图象处理
  - 图象获取
  - 图象重建 (从投影重建)
  - 图象增强和恢复等
  - 图象 (视频) 压缩编码

从2000年起增加

- 图象数字水印和图象信息隐藏

20191216 7

© Y.J.ZHANG.  



## 引 言

- 2000 ~ 2004 (博士生学科前沿课)

### 国际标准MPEG-21 和 图象数字水印

- 2002年
- ✓ 图象数字水印 ⇒ 图象工程教材

20191216 8

© Y.J.ZHANG.  

## 引 言

### 不同时代标准化的驱动力

驱动力



机构组织驱动

技术驱动

市场驱动

70年代 | 80年代 | 90年代 | 21世纪

20191216 9

© Y.J.ZHANG.  

## 引 言

- 对2004年图象工程文献的统计



文献数量 (篇)

图象水印

A1 A2 A3 A4 A5 B1 B2 B3 B4 B5 C1 C2 C3 C4 D1 D2 D3 D4 D5 D6 E1

小类类别

20191216 10

© Y.J.ZHANG.  

## 引 言

- 对2007年图象工程文献的统计



文献数量 (篇)

图象水印

A1 A2 A3 A4 A5 A6 B1 B2 B3 B4 B5 C1 C2 C3 C4 C5 D1 D2 D3 D4 D5 D6 E1

小类类别

20191216 11

© Y.J.ZHANG.  

## 引 言

- 图象分析
  - 图象分割, 边缘及角点等基元的检测
  - 目标表达、描述、测量
  - 目标特性的提取分析
  - 目标检测和识别

从2000年起增加

- 人体生物特征提取和验证 (包括人体、人脸和器官的检测、定位与识别)

20191216 12

© Y.J.ZHANG.

## 引 言

- 2005~2008 (博士生学科前沿课)  
**人脸表情识别**
- 2007年
  - ✓ 人脸和表情识别 ⇒ 图象工程教材 (第2版, 中册)

20191216 13

© Y.J.ZHANG.

## 引 言

- 人脸表情分析  
人类情感交流中的地位

- 语言内容: 7%
- 语音语调: 38%
- 脸部表情: 55%

表情反映内心

- 测谎 (> 95%)

20191216 14

© Y.J.ZHANG.

## 引 言

- 对2009年图象工程文献的统计

20191216 15

© Y.J.ZHANG.

## 引 言

- 2009~2010 (博士生学科前沿课)  
**基于子空间的人脸识别**

程正东, 贾彗星  
李乐, 沈斌  
谭华春, 严严  
章毓晋, 朱云峰

20191216 16

© Y.J.ZHANG.

## 引 言

- 第1章 绪论      第2章 人脸检测
- 第3章 人脸跟踪    第4章 人脸描述
- 第5章 基本线性子空间方法
- 第6章 张量方法    第7章 核方法
- 第8章 非负矩阵(集)分解
- 第9章 分类器设计
- 第10章 评价指标与评测比较
- 附录A 张量      附录B 3-D人脸识别综述
- 附录C 相关识别概述    附录D 常用数据库

20191216 17



© Y.J.ZHANG.

## 引 言

- 2011~2013 (博士生学科前沿课)  
**人脸图象分析进展**  
—— 技术和应用



Yu-Jin ZHANG (Editor)  
33 experts from  
16 countries and regions

20191216 18

© Y.J.ZHANG.  **引言** 

- ◆ 背景介绍 Introduction and Background
- ◆ 特征提取 Facial Feature Extraction
- ◆ 特征降维 Feature Dimensionality Reduction
- ◆ 人脸识别 Face Recognition
- ◆ 表情分类 Facial Expression Classification
- ◆ 不变技术 Invariance Techniques

20191216 19



© Y.J.ZHANG.  **引言** 

- 图象理解
  - 图象匹配和融合
  - 场景恢复
  - 图象感知和解释
  - 基于内容的图象和视频检索

从2005年起增加



- 时空技术(高维运动分析、3-D姿态检测、跟踪、举止判断和行为理解)

20191216 20

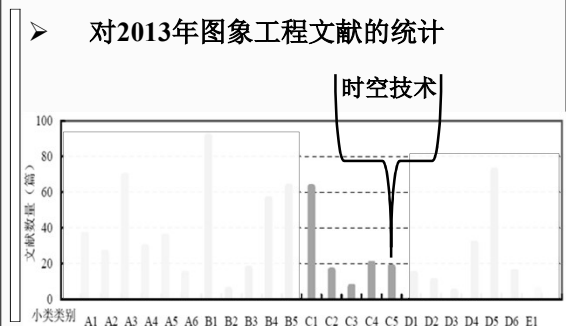
© Y.J.ZHANG.  **引言** 

- 2014~2018 (博士生学科前沿课)  
时空行为理解
- 2012年
  - ✓ 时空行为理解 ⇒ 图象工程教材 (第3版, 下册)



20191216 21

© Y.J.ZHANG.  **引言** 

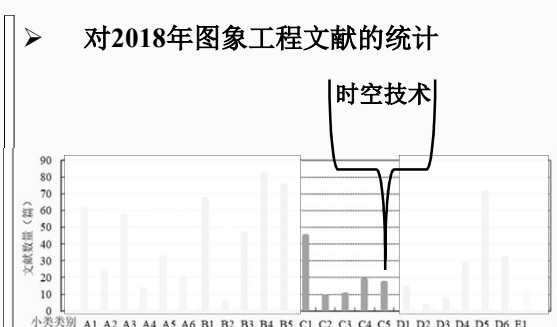
- 对2013年图象工程文献的统计





20191216 22

© Y.J.ZHANG.  **引言** 

- 对2018年图象工程文献的统计





20191216 23

© Y.J.ZHANG.  **引言** 

- 理解时空行为
  - ✓ 时空: 客观
    - 位置、轨迹、速度、外观、姿态、关系、...
    - 对群体目标活动中的聚合、消散、分化、合并等动态演变现象, ...
  - ✓ 行为: 主观
    - 举止、动向、态势、情感、...

20191216 24

© Y.J.ZHANG.  **引言** 



**计算机视觉前沿论坛 (2013~)**

徐光祐, 刘允才, 章毓晋. 计算机视觉——探索行为理解, 认知内心世界. 中国图象图形学报, 2013, 18(2): 131

章毓晋. 时空行为理解. 中国图象图形学报, 2013, 18(2): 141-151 (领跑者5000)

郑胤, 陈权崎, 章毓晋. 深度学习及其在目标和行为识别中的新进展. 中国图象图形学报, 2014, 19(2): 175-184 (领跑者5000)

20191216 25



© Y.J.ZHANG.  **时空行为理解层次** 

从动作到行为的五个层次

- (1) 动作基元 (action primitives)
- (2) 动作 (action)
- (3) 活动 (activity)
- (4) 事件 (events)
- (5) 行为 (behavior)

指用来构建动作的原子单元, 一般对应场景中局部短暂的运动信息

20191216 26



© Y.J.ZHANG.  **时空行为理解层次** 

从动作到行为的五个层次

- (1) 动作基元 (action primitives)
- (2) 动作 (action)
- (3) 活动 (activity)
- (4) 事件 (events)
- (5) 行为 (behavior)

由主体的一系列动作基元构成的有具体意义的集合体 (有序组合), 一般动作常代表由一个人所进行的简单运动模式, 且常仅持续秒的量级。人体动作的结果常导致人体姿态的改变

20191216 27



© Y.J.ZHANG.  **时空行为理解层次** 

从动作到行为的五个层次

- (1) 动作基元 (action primitives)
- (2) 动作 (action)
- (3) 活动 (activity)
- (4) 事件 (events)
- (5) 行为 (behavior)

为完成某个工作或实现某个目标而由主体执行的一系列动作的组合 (主要强调逻辑组合)。活动是相对大尺度的运动, 一般依赖于环境和交互人。活动常代表由多人参与的序列 (可能交互的) 复杂动作, 且常持续较长的时段

20191216 28



© Y.J.ZHANG.  **时空行为理解层次** 

从动作到行为的五个层次

- (1) 动作基元 (action primitives)
- (2) 动作 (action)
- (3) 活动 (activity)
- (4) 事件 (events)
- (5) 行为 (behavior)

指在特定时间段和特定空间位置发生的某种特定活动。通常其中的动作由多个主体/发起者执行 (群体活动)。对特定事件的检测常与异常活动有关

20191216 29

© Y.J.ZHANG.  **时空行为理解层次** 

从动作到行为的五个层次

- (1) 动作基元 (action primitives)
- (2) 动作 (action)
- (3) 活动 (activity)
- (4) 事件 (events)
- (5) 行为 (behavior)

主体/发起者主要指人或动物, 强调主体/发起者受思想支配而在特定环境/上下境中改变动作, 持续活动和描述事件等

20191216 30

© Y.J.ZHANG. 20191216 31

## 时空行为理解层次

五个层次 (示例)

动作基元

动作

活动

事件

行为

乒乓球比赛中的几个画面

© Y.J.ZHANG. 20191216 32

## 从时空兴趣点到轨迹

### 1. 空间兴趣点的检测

使用线性尺度空间表达对图象建模

$$L^{sp}(x, y; \sigma_i^2) = g^{sp}(x, y; \sigma_i^2) \otimes f^{sp}(x, y)$$

高斯核

$$g^{sp}(x, y; \sigma_i^2) = \frac{1}{2\pi\sigma_i^2} \exp[-(x^2 + y^2) / 2\sigma_i^2]$$

#### Harris兴趣点检测

$$\mu^{sp}(\cdot; \sigma_i^2, \sigma_j^2) = g^{sp}(\cdot; \sigma_i^2) \otimes \{[\nabla L(\cdot; \sigma_i^2)][\nabla L(\cdot; \sigma_j^2)]^T\}$$

图象点邻域中朝向分布矩阵

$$= g^{sp}(\cdot; \sigma_i^2) \otimes \begin{bmatrix} (L_x^{sp})^2 & L_x^{sp} L_y^{sp} \\ L_x^{sp} L_y^{sp} & (L_y^{sp})^2 \end{bmatrix}$$

© Y.J.ZHANG. 20191216 33

## 从时空兴趣点到轨迹

### 1. 空间兴趣点的检测

$\mu^{sp}$ 的本征值 $\lambda_1$ 和 $\lambda_2$  ( $\lambda_1 \leq \lambda_2$ )构成 $f^{sp}$ 沿两个图象坐标方向而变化的描述符

**检测角点函数的正极大值**

$$H^{sp} = \det(\mu^{sp}) - k \cdot \text{trace}^2(\mu^{sp}) = \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2$$

在感兴趣点, 本征值的比 $a = \lambda_2 / \lambda_1$ 应该很大

对 $H^{sp}$ 的正局部极值,  $a$ 应该满足 $k \leq a / (1+a)^2$

设 $k = 0.25$ ,  $H$ 的正极大值将对应理想的各向同性兴趣点 (此时 $a = 1$ , 即 $\lambda_1 = \lambda_2$ ),  $k$ 小 $\Rightarrow a$ 大

© Y.J.ZHANG. 20191216 34

## 从时空兴趣点到轨迹

### 2. 时空兴趣点的检测

检测在局部时空体中具有沿时和空都有图象值较大变化的位置

$$L(\cdot; \sigma_i^2, \tau_i^2) = g(\cdot; \sigma_i^2, \tau_i^2) \otimes f(\cdot)$$

空间方差 $\sigma_i^2$ 和时间方差 $\tau_i^2$

为了检测感兴趣点, 在 $f$ 中搜索具有 $\mu$ 的显著本征值 $\lambda_1, \lambda_2, \lambda_3$ 的区域

时空兴趣点示例

© Y.J.ZHANG. 20191216 35

## 从时空兴趣点到轨迹

通过对 (监控) 场景中各个运动目标行为的描述和刻画来提供对场景状态的把握

### 动态轨迹学习和分析:

先对感兴趣目标进行检测并跟踪, 接着用所获得的轨迹自动地构建场景模型, 最后用该模型描述监控的状况并提供对活动的标注

```

    graph LR
      Input[输入视频] --> Detection[目标检测]
      Detection --> Tracking[目标跟踪]
      Tracking --> Trajectory[轨迹]
      Trajectory --> Modeling[场景建模]
      Modeling --> Analysis[活动分析]
      Analysis --> Annotation[标注视频]
  
```

© Y.J.ZHANG. 20191216 36

## 从时空兴趣点到轨迹

### 3. 活动路径学习

活动定义在开始和结束的两个感兴趣点之间

3种基本结构的主要区别包括输入的种类, 运动矢量, 轨迹/视频片段, 以及运动抽象的方式

运动 隐含顺序 路径

单个轨迹 (a)

轨迹 路径

完整轨迹 (b)

视频片段 路径 运动单词

路径时序分解 (c)

轨迹和路径学习方案

© Y.J.ZHANG.

## 从时空兴趣点到轨迹

4. 活动路径建模  
 路径模型是对聚类的紧凑表达（图模型推理）  
 ① 考虑完整的路径，有平均的中心线，两边还有包络指示路径范围  
 ② 将路径分解为子路径（表示成子路径的树）

两种对路径建模的方式

20191216 37

© Y.J.ZHANG.

## 动作/活动识别（示例）

动作识别数据库

头顶击掌		挥手
侧向移动		挥手
弯腰		双手
行走		单脚前跳
跑		双脚前跳
		人原地跳

Weizmann 动作识别数据库中动作的示例图片

20191216 38

© Y.J.ZHANG.

## 动作/活动识别（示例）

Taking photo    Phoning    Riding horse    Using computer

Zheng Y, Zhang Y-J, Li X, Liu B D. "Action Recognition in Still Images Using a Combination of Human Pose and Context Information". Proc. ICIP, 785-788, 2012

20191216 39

© Y.J.ZHANG.

## 动作/活动识别（示例）

活动

姿态    Poselet based action classifiers    分类

上下文    Context based action classifiers    分类

静止图象     $\Sigma$     活动分类

Riding Horse  
Riding Bike  
Using Computer

20191216 40

© Y.J.ZHANG.

## 动作/活动识别（示例）

### 停车场的人车互动

特定目标、动作  
特殊环境、区域

20191216 41

© Y.J.ZHANG.

## 动作/活动识别（示例）

### 停车场的人车互动

人上车、人下车；车与人的相互位置关系

感兴趣的区域

20191216 42

© Y.J.ZHANG. 20191216

## 动作/活动识别 (示例)

停车场的人车互动

车的检测  
(位置, 状态, 环境, 相互关系, ……)

车外观相似?

车遮挡重叠?

20191216 43

© Y.J.ZHANG. 20191216

## 动作/活动识别 (示例)

停车场的人车互动 人的检测和定位

20191216 44

© Y.J.ZHANG. 20191216

## 动作/活动识别 (示例)

停车场的人车互动 类别和描述

**Table 1.** Generated events. In event signatures, “h” stands for a human ID, and “v” for a vehicle ID; “f” is the frame number at which the event occurs.

Event signature	Description
Appears(h,f)*	A human appears inside the ROI
Disappears(h,f)*	A human disappears inside the ROI
IsNear(h,v,f)**	The human is near the vehicle
IsStopped(v,f)**	The vehicle is currently stopped

\*: Not generated when the object intersects the ROI boundary.  
 \*\*: Generated only if “Appears” or “Disappears” event is generated at the same frame.

20191216 45

© Y.J.ZHANG. 20191216

## 动作/活动识别 (示例)

停车场的人车互动 规则和效果

- $Appears(h, f) \wedge IsStopped(v, f) \wedge IsNear(h, v, f) \Rightarrow GetsOut(h, v, f)$
- $Disappears(h, f) \wedge IsStopped(v, f) \wedge IsNear(h, v, f) \Rightarrow GetsIn(h, v, f)$
- $GetsIn(h, v1, f) \Leftrightarrow \neg GetsIn(h, v2, f)$
- $GetsOut(h, v1, f) \Leftrightarrow \neg GetsOut(h, v2, f)$

**Table 3.** Inference results on VIRAT dataset

Event	Precision	Recall	F-score
GetsIn	31,58%	85,71%	46,15%
GetsOut	30,00%	75,00%	42,86%

Zhang Y-J, Vincent T. “Spatial-Temporal Behaviour Understanding”. Proceedings of the 2nd International Conference on Electronics, Communications and Control, 2012

20191216 46

© Y.J.ZHANG. 20191216

## 动作/活动识别 (示例)

监控中的骑抢检测

监控摄像头遍布各处:  
公共汽车、超市、车载、地铁、酒店、小区、路口、商场、停车场、银行

20191216 47

© Y.J.ZHANG. 20191216

## 动作/活动识别 (示例)

监控中的骑抢检测

袁静, 章毓晋. 融合梯度差信息的稀疏去噪自编码网络在异常行为检测中的应用. 自动化学报, 2017, 43(4): 604-610

20191216 48



© Y.J.ZHANG. 动作/活动识别 (示例)

### 监控中的骑抢检测 视频序列

20191216 49

© Y.J.ZHANG. 动作/活动识别 (示例)

### 监控中的骑抢检测 定位和跟踪 (人、车、背景)

20191216 50

© Y.J.ZHANG. 动作/活动识别 (示例)

### 监控中的骑抢检测 流程和步骤

20191216 51

© Y.J.ZHANG. 动作/活动识别 (示例)

### 图像序列—预测序列 图像序列—被预测序列

20191216 52

© Y.J.ZHANG. 深度学习与行为识别

### 人脑分层次认知

20191216 53

© Y.J.ZHANG. 深度学习与行为识别

### 深度学习算法体系结构

20191216 54

© Y.J.ZHANG. 深度学习与行为识别

### 深度学习算法体系结构

1、自编码器的核心思想是将输入信号进行编码，使用编码之后的信号重建原始信号，目的是让重建信号与原始信号相比重建误差最小

特征：隐层单元  $h$

输入：可见层单元  $x$

20191216 55

© Y.J.ZHANG. 深度学习与行为识别

### 深度学习与行为识别

- The Classical Autoencoder
  - mapping visible input  $x$  to hidden representation  $y$   
 $y = f(x) = s(Wx + b)$
  - mapping hidden representation  $y$  back to input space  
 $z = g(y) = s(W'y + b')$

隐含  $y$

Reconstruction Error  $\min L(x,z)$  误差最小

输入  $x$   $z$  输入空间

20191216 56

© Y.J.ZHANG. 深度学习与行为识别

### 深度学习与行为识别

#### 深度学习算法体系结构

#### 去噪自编码器 (denoising autoencoder)

能发现输入中更多隐含的特征，能减小模型对输入数据的依赖性，能增强模型的鲁棒性

最小化重建误差  $L(x,z)$

随机破坏

20191216 57

© Y.J.ZHANG. 深度学习与行为识别

### 深度学习与行为识别

- Comparison between Classical Autoencoder and Denoising Autoencoder
- Four examples with different corruption rate
- Classical Autoencoder: the bases seem to be influenced by some spatial noise
- Denoising Autoencoder: smoother and clearer

Corruption Rate = 0% Corruption Rate = 20%

Corruption Rate = 40% Corruption Rate = 60%

You Q H Z, Zhang Y-J. "A New Training Principle for Stacked Denoising Autoencoders". Proceedings of the 7th International Conference on Image and Graphics, pp.384-389, 2013.

20191216 58

© Y.J.ZHANG. 深度学习与行为识别

### 深度学习与行为识别

#### 异常行为检测

融合梯度差信息的稀疏去噪自编码网络

模型学习过程-反向传播计算

训练

融合梯度差信息的稀疏去噪自编码网络

前向计算

20191216 59

© Y.J.ZHANG. 深度学习与行为识别

### 深度学习与行为识别

#### 检测

Patch\_1 Patch\_N

$N$ 个patch



检测

Loss=阈值-发生异常事件

Loss<阈值-正常事件

袁静, 章毓晋. 融合梯度差信息的稀疏去噪自编码网络在异常行为检测中的应用. 自动化学报, 2017, 43(4): 604-610.

20191216 60

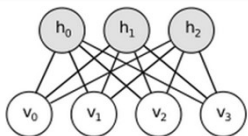
© Y.J.ZHANG.  **深度学习与行为识别** 

### 深度学习算法体系结构



2、限制/约束玻尔兹曼机是构成深度置信网络的基础单元，其本质是使得学习到的模型产生符合条件样本的概率最大

可以对全连通的玻尔兹曼机进行简化，使同层单元彼此独立（约束）

层内单元之间没有连接关系，层间单元是全连接关系



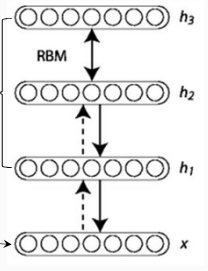
20191216 61

© Y.J.ZHANG.  **深度学习与行为识别** 

### 深度置信网络（Deep Belief Networks）

将约束玻尔兹曼机模型堆叠起来可构成一个深层置信网络，这样的几何模型可以用来提取训练数据中深层结构的特征

深层置信网络模型描述了可观测向量和所有隐藏层之间的联合概率分布



20191216 62

© Y.J.ZHANG.  **深度学习与行为识别** 

### 深度网络用于目标和行为识别

#### 谷歌的虚拟人脑



使用 1000 台电脑构造出了包含 10 亿个连接的深达9层的自编码器“神经网络”

具有一定的自动学习能力

可以模拟一些人脑的功能

在完全没有标签的情况下，该网络能够自动地从训练集中学习到某些高层概念（如去“认识”动物“猫”，将猫的图片与cat联系起来）

20191216 63

© Y.J.ZHANG.  **深度学习与行为识别** 

### 深度网络用于目标和行为识别

#### 大规模目标识别

在Image Net数据库（超过10,000,000张图象，超过1000个类别）上进行的大规模视觉识别比赛

早期最有效的算法

- 基于词袋模型（bag of words）
- 可变部件模型（deformable part model）

2012年基于深度学习的模型取得了最好的效果（超出了第二名至少10个百分点）

20191216 64

© Y.J.ZHANG.  **深度学习与行为识别** 



### 深度网络用于目标和行为识别

#### 图象的同时分类和标注

**图象分类：**指的是对图象内容做整体的描述，例如给定一幅图象确定它属于“海滩”、“厨房”、“卧室”等预先定义好的类别中的哪一类

**图象标注：**指的是对于图象中包含的内容作出判断，例如一幅图象中是否包含“天空”、“汽车”、“树木”等预先定义好的目标，“骑车”，“篮球比赛”等预先定义好的活动，……

20191216 65

© Y.J.ZHANG.  **深度学习与行为识别** 

### 深度网络用于目标和行为识别

#### 图象的同时分类和标注

Zheng Y, Zhang Y-J, Larochelle H. "Topic Modeling of Multimodal Data: An Autoregressive Approach". *Proceedings of Computer Vision and Pattern Recognition*, 1370-1377, 2014

Zheng Y, Zhang Y-J, Larochelle H. "A Deep and Autoregressive Approach for Topic Modeling of Multimodal Data". *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(6): 1056-1069, 2016

20191216 66

© Y.J.ZHANG.

## 深度学习与行为识别

**深度网络用于目标和行为识别**  
**图象的同时分类和标注**

使用基于神经自回归分布估计器 (NADE) 的  
 监督性神经自回归分布主题模型 (SupDocNADE)  
 来处理同时图象分类和标注问题

NADE (Neural Autoregressive Distribution Estimator)  
 DocNADE (Document)  
 SupDocNADE (Supervised)  
 SupDeepDocNADE (Deep)

20191216 67

© Y.J.ZHANG.

## 深度学习与行为识别

**SupDocNADE**  
 A single hidden layer

•Visual word  
 •Annotation  
 •Topic Feature  
 •Classification  
 •Annotation

20191216 68

© Y.J.ZHANG.

## 深度学习与行为识别

**SupDeepDocNADE** (Deep extension of SupDocNADE)

Retrieval Task  
 Topic Feature  
 Visual word  
 Annotation  
 Randomly Shuffle  
 $\omega(\rho)$

20191216 69

© Y.J.ZHANG.

## 深度学习与行为识别

Label Me  
 Open country  
 Mountain  
 UIUC-Sports  
 Bocce  
 Croquet

20191216 70

© Y.J.ZHANG.

## 深度学习与行为识别

**Flicker子集 (2.5万张图象, 38类) 上的结果**

Model	MAP
TF-IDF	0.384 ± 0.004
Multiple Kernel Learning SVMs [6]	0.623
TagProp [31]	0.640
Multimodal DBM [13]	0.651 ± 0.005
MDRNN [15]	0.686 ± 0.003
SupDeepDocNADE (1 hidden layer, 625 epochs pretraining)	0.654 ± 0.004
SupDeepDocNADE (2 hidden layers, 625 epochs pretraining)	0.671 ± 0.006
SupDeepDocNADE (3 hidden layers, 625 epochs pretraining)	0.670 ± 0.005
SupDeepDocNADE (2 hidden layers, 2325 epochs pretraining)	0.682 ± 0.005
SupDeepDocNADE (3 hidden layers, 2325 epochs pretraining)	0.686 ± 0.005
SupDeepDocNADE (2 hidden layers, 4125 epochs pretraining)	0.684 ± 0.005
SupDeepDocNADE (3 hidden layers, 4125 epochs pretraining)	0.691 ± 0.005

20191216 71

© Y.J.ZHANG.


## 深度学习与行为识别

**视频中的人体活动分类 问题和动机**

跳远  
 跳高

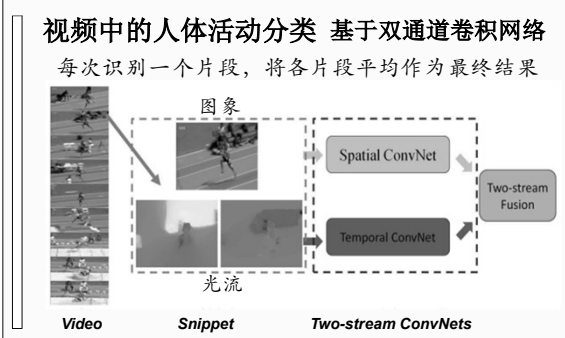
在助跑阶段很难区分, 需要考虑整个视频段的顺序结构  
 Chen Q Q, Zhang Y-J. "Sequential Segment Networks for Action Recognition". IEEE Signal Processing Letters, 24(5): 712-716, 2017

20191216 72

© Y.J.ZHANG. 


## 深度学习与行为识别

视频中的人体活动分类 基于双通道卷积网络  
每次识别一个片段，将各片段平均作为最终结果



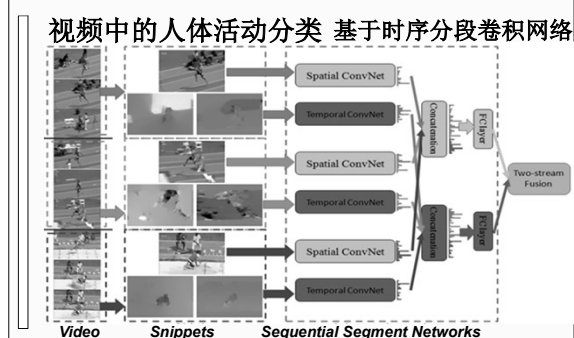
Video Snippet Two-stream ConvNets

20191216 73

© Y.J.ZHANG. 


## 深度学习与行为识别

视频中的人体活动分类 基于时序分段卷积网络



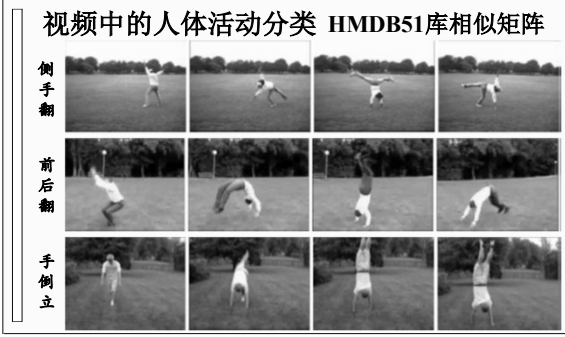
Video Snippets Sequential Segment Networks

20191216 74

© Y.J.ZHANG. 


## 深度学习与行为识别

视频中的人体活动分类 HMDB51库相似矩阵



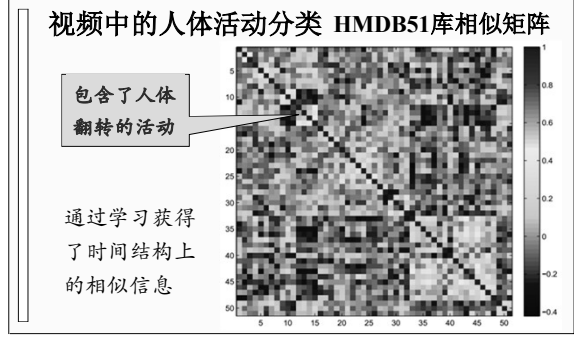
侧手翻  
前后翻  
手倒立

20191216 75

© Y.J.ZHANG. 

## 深度学习与行为识别


视频中的人体活动分类 HMDB51库相似矩阵



包含了人体翻转的活动

通过学习获得了时间结构上的相似信息

20191216 76

© Y.J.ZHANG. 

## 结语

- 引入一个有前途的研究方向
  - ✓ 对应高层的图象理解
- 介绍了一些相关研究和进展
  - ✓ 从客观到主观与客观结合
  - ✓ 从基元到行为 (五个层次)
  - ✓ 从点到线到面到体 (轨迹、路径、...)
  - ✓ 从数学物理到 (结合) 神经科学
- 未来发展趋势

20191216 77

© Y.J.ZHANG. 

## 联系信息

- ◆ 通信地址: 北京清华大学电子工程系 (100084)
- ◆ 办公地址: 清华大学, 罗姆楼, 6层305室
- ◆ 办公电话: (010)62798540
- ◆ 传真号码: (010)62770317
- ◆ 电子邮箱: [zhang-yj@tsinghua.edu.cn](mailto:zhang-yj@tsinghua.edu.cn)
- ◆ 个人主页: [oa.ee.tsinghua.edu.cn/~zhangyujin/](http://oa.ee.tsinghua.edu.cn/~zhangyujin/)

20191216 78