

Journal of Image  
and Graphics

# 中国图象图形学报



ISSN1006-8961  
CN11-3758/TB

2012 **7**  
Vol.17 No.

中国科学院遥感应用研究所  
中国图象图形学学会主办  
北京应用物理与计算数学研究所

# 中国图象图形学报

Zhongguo Tuxiang Tuxing Xuebao

2012年7月 第17卷 第7期(总第195期)

## 目次

### 综述

中国多媒体技术研究: 2011 ..... 史元春, 徐光祐, 高原(741)

### 图像处理和编码

图像引导滤波的局部多尺度 Retinex 算法 ..... 方帅, 杨静荣, 曹洋, 武鹏飞, 饶瑞中(748)

基于第一主成分方向稳定性的图像零水印 ..... 付剑晶, 王珂(756)

小波变换估计非线性扩散最优停止时间 ..... 蒋平, 张建州(770)

### 图像分析和识别

改进的 Beamlet 与 Canny 相结合提取复杂图像线特征 ..... 曾接贤, 周沥沥, 符祥(775)

动静态信息融合及动态贝叶斯网络的步态识别 ..... 杨旗, 薛定宇(783)

融合音频单词与视觉特征的成人视频检测 ..... 刘毅志, 唐胜, 王向东, 林守勋, 张勇东(791)

基于平面区域跟踪的目标位姿参数自动测量 ..... 回丙伟, 文贡坚, 赵竹新, 钟金荣(798)

3 维图像中边界曲面的分类追踪及抽取 ..... 丁德福, 程柳航, 王利生(806)

局部时空域模型的核密度估计目标检测方法 ..... 王兴宝, 刘纯平, 费兰英, 王朝晖, 季怡(813)

基于相位谱和调谐幅度谱的显著性检测方法 ..... 李崇飞, 高颖慧, 卢凯, 曲智国(821)

### 图像理解和计算机视觉

“目标-场景”语境关联的生成图模型分析 ..... 谢昭, 李姍琦, 高隼(828)

融合上下文信息的场景结构恢复 ..... 武晖, 于昕, 隋尧, 张利(839)

带有局部控制因子的图割光流估计 ..... 路子赟, 唐土生, 高隼, 沈琳, 刘伟(846)

改进的 Harris 亚像素角点快速定位 ..... 何海清, 黄声享(853)

### 计算机图形学

应用最小生成树实现点云分割 ..... 孙金虎, 周来水, 安鲁陵(858)

### 医学图像处理

病毒进化的离散差分进化超声图像特征选择 ..... 张巧荣, 朱长明, 倪军, 刘海波(866)

分水岭优化的 Snake 模型肝脏图像分割 ..... 兰红, 张璐(873)

### 遥感图像处理

Normalized Cut 与分水岭变换在高光谱影像混合像元端元提取中的应用 ..... 许菡, 李小娟(880)

改进非局部均值滤波的 SAR 图像降噪方法 ..... 郑永恒, 程建, 曹宗杰(886)

~~~~~  
第四届国际遥感考古会议 ..... (892)

## 中国图象图形学报

刊名题字: 宋 健

月刊(1996 年创刊)

第 17 卷 第 7 期

2012 年 7 月 16 日出版

**主管单位** 中国科学院

**主 办** 中国科学院遥感应用研究所  
中国图象图形学学会

北京应用物理与计算数学研究所

**主 编** 李小文

**编辑出版** 《中国图象图形学报》编辑出版委员会

北京 9718 信箱 邮编 100101

电子信箱:jig@irsa.ac.cn

电话:010-64807995 010-82614429

网 址:www.cjig.cn

**印刷装订** 北京北林印刷厂

**广告经营许可证** 京朝工商广字第 0346 号

**总 发 行** 北京报刊发行局

**订 购** 全国各地邮局

**国外发行** 中国国际图书贸易总公司

(中国国际书店)

(北京 399 信箱 邮编 100044)

**Superintended by** Chinese Academy of Sciences

**Sponsored by** Institute of Remote Sensing Application,  
CAS China Society of Image and Graphics  
Institute of Applied Physics and Computational  
Mathematics

**Chief editor** LI Xiaowen

**Editor, Publisher** Editorial and Publishing Board  
of Journal of Image and Graphics  
(P. O. Box 9718, Beijing 100101, China)  
E-mail:jig@irsa.ac.cn

**Distributed by** Beijing Bureau for Distribution of Newspapers  
and Journals

**Domestic** All Local Post Offices in China

**Foreign** China International Book Trading Corporation  
(P. O. Box 399, Beijing 100044, China)

**Printed by** Beijing Beilin Printing House

ISSN 1006-8961 CN11-3758/TB CODE ZTTFXZ 国内邮发代号: 82-831 国外发行代号: M1406 国内定价: 45.00 元

# Journal of Image and Graphics

( Monthly , Started in 1996 )

Vol. 17 No. 7 July 2012

## Contents

### Review

Researches on multimedia technology in China, 2011 ..... Shi Yuanchun, Xu Guangyou, Gao Yuan (741)

### Image Processing and Coding

Local multi-scale Retinex algorithm based on guided image filtering  
..... Fang Shuai, Yang Jingrong, Cao Yang, Wu Pengfei, Rao Ruizhong (748)

Image zero-watermark based on direction stability of first principal component vector ..... Fu Jianjing, Wang Ke (756)

Stopping-time estimation for anisotropic diffusion using discrete wavelet transform ..... Jiang Ping, Zhang Jianzhou (770)

### Image Analysis and Recognition

Complex image line feature extraction based on improved Beamlet transform and the Canny operator  
..... Zeng Jiexian, Zhou Lili, Fu Xiang (775)

Gait recognition based on dynamic & static information fusion and dynamic bayesian network ..... Yang Qi, Xue Dingyu (783)

Fusing audio-words with visual features for adult video detection  
..... Liu Yizhi, Tang Sheng, Wang Xiangdong, Lin Shouxun, Zhang Yongdong (791)

Automatic measurement for an object's position and attitude via tracking planar regions  
..... Hui Bingwei, Wen Gongjian, Zhao Zhuxin, Zhong Jinrong (798)

Detection and extraction of boundary surface patches within 3D images ..... Ding Defu, Cheng Liuhang, Wang Lisheng (806)

Foreground object detection method using kernel density estimation of a local spatio-temporal model  
..... Wang Xingbao, Liu Chunping, Fei Lanying, Wang Zhaohui, Ji Yi (813)

Saliency detection method based on phase spectrum and amplitude spectrum tuning  
..... Li Chongfei, Gao Yinghui, Lu Kai, Qu Zhiguo (821)

### Image Understanding and Computer Vision

"Object-Scene" contextual associated generative graph model analysis ..... Xie Zhao, Li Shanqi, Gao Jun (828)

Structure recovery algorithm using contextual information ..... Wu Hui, Yu Xin, Sui Yao, Zhang Li (839)

Graph cut optical flow estimation with a local control factor ..... Lu Ziyun, Tang Tusheng, Gao Jun, Shen Lin, Liu Wei (846)

Improved algorithm for Harris rapid sub-pixel corners detection ..... He Haiqing, Huang Shengxiang (853)

### Computer Graphics

Research on point cloud segmentation using a minimum spanning tree ..... Sun Jinhu, Zhou Laishui, An Luling (858)

### Medical Image Processing

Virus-evolutionary discrete differential evolution algorithm for feature selection of cervical lymph nodes in ultrasound images  
..... Zhang Qiaorong, Zhu Changming, Ni Jun, Liu Haibo (866)

Liver image segmentation algorithm based on the Snake model And optimized by watershed transformation  
..... Lan Hong, Zhang Lu (873)

### Remote Sensing Image Processing

Endmember extraction for hyperspectral image based on normalized cut and watershed transformation  
..... Xu Han, Li Xiaojuan (880)

SAR image denoising via improved non-local means filter ..... Zheng Yongheng, Cheng Jian, Cao Zongjie (886)

中图法分类号: TP391.41 文献标志码: A 文章编号: 1006-8961(2012)07-0791-07

论文引用格式: 刘毅志,唐胜,王向东,林守勋,张勇东. 融合音频单词与视觉特征的成人视频检测[J]. 中国图象图形学报,2012,17(7): 791-797

## 融合音频单词与视觉特征的成人视频检测

刘毅志<sup>1,2,3</sup>, 唐胜<sup>2</sup>, 王向东<sup>2</sup>, 林守勋<sup>2</sup>, 张勇东<sup>2</sup>

1. 湖南科技大学计算机科学与工程学院, 湘潭 411201; 2. 中国科学院计算技术研究所, 北京 100190;
3. 湖南省知识处理与网络化制造重点实验室, 湘潭 411201

**摘要:** 基于多模态的检测方法是过滤成人视频的有效手段,然而现有方法中缺乏准确的音频语义表示方法。因此本文提出融合音频单词与视觉特征的成人视频检测方法。先提出基于周期性的能量包络单元(简称EE)分割算法,将音频流准确地分割为EE的序列;再提出基于EE和BoW(Bag-of-Words)的音频语义表示方法,将EE的特征描述为音频单词的出现概率;采用复合加权方法融合音频单词与视觉特征检测结果;还提出基于周期性的成人视频判别算法,与基于周期性的EE分割算法前后配合,以充分利用周期性进行检测。实验结果表明,与基于视觉特征的方法相比,本文方法显著提高了检测性能。当误检率为9.76%时,检出率可达94.44%。

**关键词:** 成人视频检测;多模态融合;音频单词;视觉特征;能量包络单元

## Fusing audio-words with visual features for adult video detection

Liu Yizhi<sup>1,2,3</sup>, Tang Sheng<sup>2</sup>, Wang Xiangdong<sup>2</sup>, Lin Shouxun<sup>2</sup>, Zhang Yongdong<sup>2</sup>

1. Institute of Computer Science and Engineering, Hunan University of Science and Technology, Xiangtan 411201, China;
2. Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China;
3. Key Laboratory of Knowledge Processing and Networked Manufacturing, College of Hunan Province, Xiangtan 411201, China

**Abstract:** Multi-modality based adult video detection is an effective approach for filtering pornographic information. However, existing methods lack accurate representation methods of audio semantics. Therefore, a novel method is presented in this paper to fuse audio-words with visual features for adult video detection. First, we propose a periodicity-based segmentation algorithm of units of energy envelope (EE). Audio streams are divided into sequences of EE. Second, audio semantics representation method based on EE and BoW (Bag-of-Words) is presented to describe the features of the EE as the occurrence probabilities of audio-words. Integrated weighting methods are used to fuse the detection results of audio-words and visual features. Furthermore, we propose a periodicity-based decision algorithm to judge adult videos to cooperate with the preceding periodicity-based segmentation algorithm. Therefore, we make full use of the periodicity. Our experiments show that our approach remarkably improves the detection performance compared with the method based on visual features. The true positive rate achieves 94.44% while the false positive rate is 9.76%.

**Key words:** adult video detection; multi-modality fusion; audio-words; visual features; units of energy envelope

## 0 引言

互联网上的成人信息已经泛滥成灾,其中,成人

视频表达直观、内容丰富,对于青少年的危害不容忽视。目前,成人视频的过滤方法大多是基于视频关键帧(KF)的视觉特征。即:先提取视频关键帧及其视觉特征,再用分类器分析关键帧是否为成人图像,

收稿日期:2011-03-03;修回日期:2011-08-29

基金项目:国家重点基础研究发展计划(973)项目(2007CB311105);国家自然科学基金项目(60873165);北京市科技新星计划项目(2007B071);北京市教育委员会共建项目

第一作者简介:刘毅志(1973—),男,副教授,2011年于中国科学院计算技术研究所获计算机应用技术专业博士学位,主要研究方向为多媒体内容分析与检索、服务计算。E-mail:liuyizhi928@gmail.com

再用阈值判决算法获得结果。阈值判别算法是指当视频中成人关键帧的数目超过设定阈值时,就判定此视频为成人视频。

基于视觉特征的检测方法实质上是一种静态图像检测方法。静态图像检测的代表性工作有 Forsyth 等人提出的基于人体轮廓的检测方法<sup>[1]</sup>、Jones 等人提出的基于肤色模型的检测方法<sup>[2]</sup>、以 WIPE 系统<sup>[3]</sup>和 Google 系统为代表的基于多特征融合的检测方法<sup>[4]</sup>。基于视觉单词的检测方法<sup>[5-6]</sup>已被用于降低肤色检测或全局特征带来的误检率。视觉单词具有一定的语义,但特征提取速度较慢。

随着互联网上低质视频的快速增长,单纯依靠视觉特征难以准确地检测成人视频。因此,融合视频中运动和音频等时域信息的检测方法逐渐成为人们关注的焦点。Endeshaw 等人<sup>[7]</sup>和 Jansohn 等人<sup>[8]</sup>较为深入地研究了利用成人视频中运动信息的周期性进行检测的方法,获得了较好效果。但是,文献[8]采用周期性检测来描述运动特征,其效果不如运动直方图。可见,运动信息的周期性没有充分发挥作用。另一方面,成人视频中运动信息的区分性不强。例如,鼓掌、俯卧撑和仰卧起坐等都难以与色情行为区分开来。

众所周知,成人视频中往往伴随着周期性的喘息、呻吟或尖叫。Rea 等人<sup>[9]</sup>通过音频流能量的自相关性分析表明,成人视频中音频流能量具有周期

性。Zuo 等人<sup>[10]</sup>分别利用音频特征与视觉特征检测成人视频,分别通过最近中心分类器与高斯混合模型进行分类,再基于 Bayes 理论进行融合。其中,音频特征是采用音频帧的低层特征。

综上所述,现有的成人视频检测方法中,缺乏准确的音频语义表示方法。具体表现在 3 个方面:1)没有充分利用音频流能量的周期性;2)音频帧的尺度小,区分性较差;3)低层特征与高层概念之间存在语义鸿沟。

针对上述问题,提出融合音频单词(Audio-words)与视觉特征的成人视频检测方法。其中包括:1)提出用于音频流分割和成人视频判别的两种基于周期性的算法,充分利用音频流能量的周期性,将音频内容表示为能量包络单元(EE)的序列;2)提出基于 EE 和 BoW (Bag-of-words) 的音频语义表示方法,将 EE 的特征描述为音频单词的出现概率;3)由于 EE 是一种不定长的、较大尺度的音频片段,并且音频单词是音频信息的一种中层表达,所以,音频单词直方图不仅区分性较强,而且能够有效缩小语义鸿沟。文献[11]的研究表明,大尺度音频片段的分类正确率要明显高于小尺度音频片段的分类正确率,并且这个趋势与分类器选择无关。

图 1 给出了本文方法的框架,包括基于音频单词的检测方法、基于视觉特征的检测方法和两者的融合方法 3 个部分。

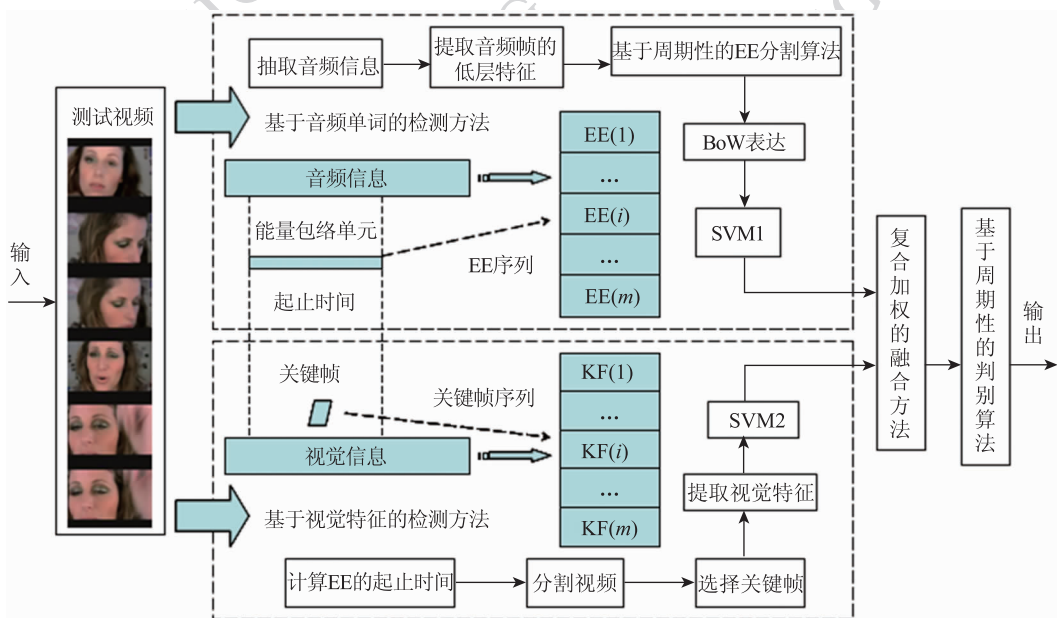


图 1 融合音频单词与视觉特征的成人视频检测方法的框架

Fig. 1 The framework of fusing audio words and visual features for adult video detection

## 1 基于音频单词的成人视频检测方法

本文首次将音频单词用于成人视频检测。因为音频单词在音频内容分析<sup>[12-13]</sup>、语义事件检测<sup>[14]</sup>和语义概念分类与标注<sup>[15-16]</sup>等领域表现出优良的性能。图 1 上半部分显示了基于音频单词的检测方法的主要过程:从视频中抽取音频信息,提取音频帧的低层特征,根据基于周期性的 EE 分割算法将音频流切分为 EE 的序列,采用 BoW 模型将 EE 的特征表示为音频单词的出现概率,再用支持向量机(SVM)分类。SVM 是当前性能最好的分类器之一<sup>[8]</sup>。下面详细叙述音频特征提取、音频单词生成、基于周期性的 EE 分割算法、以及基于 EE 和

BoW 的音频语义表示方法等。

### 1.1 音频特征提取和音频单词生成

将视频文件转换为音频文件后,无重叠地每隔 20 ms 提取一个音频帧,每帧包含 36 维音频低层特征。这些特征包括 13 维的美尔频率倒谱系数(MFCC)及其 13 维的 MFCC 差分系数、1 维过零率、1 维短时能量、4 维分带短时能量和 4 维分带短时能量比。

图 2 说明了音频单词的生成过程。在训练集的成人视频中,先粗略地剪切部分子视频,提取它们的音频信号;再从中剪辑典型的成人音频,提取其中音频帧的特征;聚类音频特征,获得音频词表  $V_{\text{words}}$ ,其中的  $M$  个向量称为音频单词,即

$$V_{\text{words}} = \{W_1, \dots, W_M\} \quad (1)$$

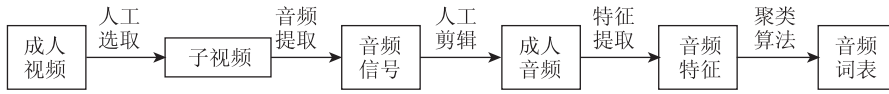


图 2 音频单词的生成过程

Fig. 2 The procedure of creating audio-words

### 1.2 基于周期性的 EE 分割算法

音频流分割有定长分割与变长分割两种方法。长度是指分割得到的音频片段中包含的音频帧数量。它在很大程度上影响着检测性能,所以,定长分割方法需要进行长度调节才能获得较好的性能。EE 分割是一种变长分割方法。现有方法根据短时能量将音频流分割为一个完整的声学信号单元,速度快且较鲁棒<sup>[17]</sup>。

但是,在分割成人视频时,经常出现长度较大的 EE。因此,提出基于周期性的 EE 分割算法。具体如下:

输入:音频特征文件,音频片段的最大长度  $L_{\text{max}}=100$ , 阈值  $T_1$  和  $T_2$ ;

输出:切分点的位置(即对应的音频帧号)。

1) 采用中值滤波的方法对能量包络波形图进行平滑处理;

2) 根据式(2)的检测函数,计算每一个音频帧  $i$  对应的检测函数值

$$d_i = \max_{j=1, \dots, J} (E_{i+j}/E_i) \quad (2)$$

式中,  $E_i$  表示第  $i$  帧对应的短时能量值,  $J$  为求检测函数时选取的后续帧数(经验值为 10);

3) 两个切分点之间的能量包络称为 EE。当  $d_i \leq T_1$  时,当前点不是切分点;  $d_i \geq T_2$  时,当前点一

定是切分点;  $T_1 < d_i < T_2$  时,  $P(i)$  为切分点,即

$$P(i) = \begin{cases} 1 & d_i \geq T_2 \\ \frac{d_i - T_1}{T_2 - T_1} & T_1 < d_i < T_2 \\ 0 & d_i \leq T_1 \end{cases} \quad (3)$$

4) 如果两个切分点之间的音频片段长度(即音频帧数目)大于  $L_{\text{max}}$ , 那么通过统计方法预测 EE 内部的周期  $T$ , 并按照  $T$  将其分割为长度较小的音频片段。周期  $T$  是指已经切分的 EE 的平均长度。

该算法利用成人视频中音频流能量的周期性,更准确地将音频流分割为 EE 的序列,为准确表达音频内容奠定了基础。

### 1.3 基于 EE 和 BoW 的音频语义表示方法

该方法将 EE 特征描述为音频单词出现概率。其目的有两个方面:其一,音频单词具有一定的语义,能够有效缩小语义鸿沟,其二,将不同长度的 EE 转化为相同特征维数的音频单词直方图,以便于 SVM 进行分类。

将训练集  $V_{\text{train}}$  分为  $a$  个正样本视频  $\{V_{p1}, \dots, V_{Pa}\}$  与  $b$  个负样本视频  $\{V_{N1}, \dots, V_{Nb}\}$ 。提取它们的音频特征,切分为无重叠的 EE。设其中任意一个视频  $V$ , 包含  $p$  个 EE(用  $U$  表示), 则

$$\mathbf{V} = \{U_1, \dots, U_p\} \quad (4)$$

设其中某个  $U$  包含  $L$  个音频帧, 则  $U$  的长度为  $L$ 。

$$U = \{F_1, \dots, F_L\} \quad (5)$$

将  $U$  的每一个音频帧和音频词表  $V_{\text{words}}$  的每一个音频单词相匹配, 得到结果  $S$ 。即

$$S = \text{Match}(U, V_{\text{words}}) = \{\text{Match}(F_1, V_{\text{words}}), \dots, \text{Match}(F_L, V_{\text{words}})\} \quad (6)$$

对于任意一个音频帧  $F_j$ , 根据式(1)有

$$\text{Match}(F_j, V_{\text{words}}) = \text{Match}(F_j, \{W_1, \dots, W_M\}) = \{D_{j,1}, \dots, D_{j,M}\} \quad (7)$$

如果音频帧  $F_j$  与音频单词的欧氏距离中, 第  $k$  个 ( $k = \{1, \dots, 10\}$ ) 最小, 表明音频帧  $F_j$  与第  $k$  个音频单词最相似, 则该音频帧的  $M$  维特征中, 第  $k$  维记作 1, 其他维记作 0。即

$$D_{j,k} = \begin{cases} 1 & k = \text{argmin}_k E(F_j, \{W_1, \dots, W_k, \dots, W_M\}) \\ 0 & \text{其他} \end{cases} \quad (8)$$

式中,  $E(A, B)$  表示  $A$  与  $B$  的欧氏距离。这样, 音频帧  $F_j$  就可以表示为一个  $M$  维的特征序列, 每一维由 0 或 1 组成, 表示对应的音频单词是否出现。即

$$D_j = \{D_{j,1}, \dots, D_{j,k-1}, D_{j,k}, D_{j,k+1}, D_{j,M}\} = \{0, \dots, 0, 1, 0, \dots, 0\} \quad (9)$$

为了避免长度的影响, 将不定长的 EE 的特征或定长的音频片段的特征都表示为各个音频单词在其中的出现概率。即将  $U$  中  $L$  个音频帧的特征按列叠加, 其和除以长度  $L$ 。

$$U = \left\{ \frac{\sum_{j=1}^L D_{j,1}}{L}, \dots, \frac{\sum_{j=1}^L D_{j,M}}{L} \right\} \quad (10)$$

训练集  $V_{\text{train}}$  中, 正样本视频的所有 EE 特征前标注为 1, 负样本视频的所有 EE 特征前标注为 -1, 将它们的顺序打乱后进行 SVM 训练以获取 EE 的 SVM 分类模型。

## 2 音频单词与视觉特征的融合方法

图 1 下半部分说明了基于视觉特征的检测方法。其主要过程是: 根据 EE 对应的起止时间, 选择离其中心最近的关键帧 (KF), 提取视觉特征后, 用 SVM 分类。视觉特征是指图像的各种特征, 包括全局特征和视觉单词直方图等。为了不失一般性, 选择颜色矩与边缘纹理的前融合作为视觉特征。图 1

右部采用复合加权的融合方法, 还提出基于周期性的成人视频判别算法。

### 2.1 复合加权的融合方法

在实验中发现: 基于音频单词的方法查准率较高, 查全率偏低; 基于视觉特征的方法查全率较高, 查准率偏低。因此, 先将两种检测方法的 SVM 结果进行加权平均, 再采用最大值融合方法。本文称之为复合加权的融合方法。

如果  $\text{Method}_{\text{EE}}$  表示基于音频单词的方法,  $\text{Method}_{\text{KF}}$  表示基于视觉特征的方法,  $R(V)$  表示视频  $V$  的检测结果,  $W$  表示权值, 那么

$$R_{\text{Method}_{\text{EE}}}(V) = (v_1^{(1)}, \dots, v_m^{(1)}) \quad (11)$$

$$R_{\text{Method}_{\text{KF}}}(V) = (v_1^{(2)}, \dots, v_m^{(2)}) \quad (12)$$

式中,  $(v_1^{(1)}, \dots, v_m^{(1)})$  分别为 SVM1 分类器的  $m$  个预测值,  $(v_1^{(2)}, \dots, v_m^{(2)})$  分别为 SVM2 分类器的  $m$  个预测值。加权融合过程为

$$R_W(V) = W \times R_{\text{Method}_{\text{EE}}}(V) + (1 - W) \times R_{\text{Method}_{\text{KF}}}(V) = W \times (v_1^{(1)}, \dots, v_m^{(1)}) + (1 - W) \times (v_1^{(2)}, \dots, v_m^{(2)}) \quad (13)$$

最终融合结果  $R(V)$  是式(11)(13)的最大值。为了获得更好的性能, 结合下面的成人视频判别算法调整权值。

$$R(V) = \text{MAX}(R_{\text{Method}_{\text{EE}}}(V), R_W(V)) \quad (14)$$

### 2.2 基于周期性的成人视频判别算法

为了配合基于周期性的 EE 分割算法, 更充分地利用音频流能量的周期性获取更好的性能, 提出基于周期性的成人视频判别算法。因为含有成人内容的音频流能量在一段时间内总是连续出现, 所以, 当含有成人内容的 EE 连续出现  $N$  ( $1 \leq N \leq 10$ ) 次时, 则认为它们所在的视频就是成人视频; 否则, 为非成人视频。调节连续次数  $N$ , 可以减少个别音频片段误判带来的干扰。

## 3 实验与分析

从互联网上收集视频建立训练集和测试集。训练集有 48 个成人视频和 300 个非成人视频, 用于训练 SVM 分类模型。测试集有 50 个成人视频和 150 个非成人视频。实验使用的机器配置是 1.86 GHz 的双核 CPU 和 2 G 内存, 在 visual studio 2003 环境下实现。性能评价指标选用 ROC(R) 曲线, 其横坐标是误检率 (FPR), 纵坐

标是检出率(TPR)。

分类器都采用 RBF 核的 SVM。

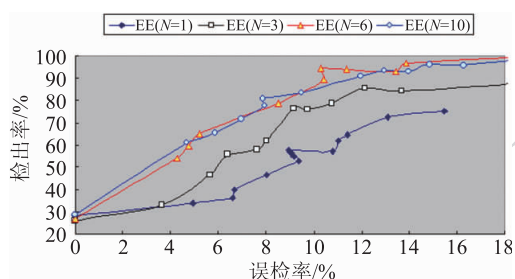
### 3.1 基于音频单词的检测方法

将视频转换为音频,所有音频数据均为采样率为 44 100/s、量化位数为 16 的混和单声道 WAV 格式。如表 1 所示,采用基于周期性的 EE 分割算法将训练集与测试集分别切分为 60 941 个和 69 211 个 EE。

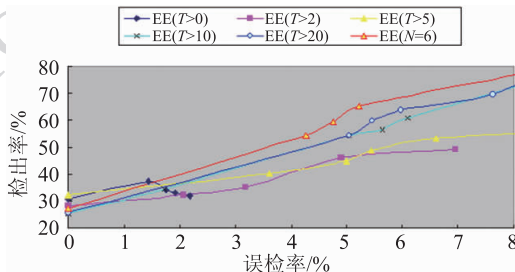
表 1 本文数据库的相关统计  
Table 1 The statistics of our dataset

|     | 成人视频<br>数量/个 | 非成人视频<br>数量/个 | 作用          | EE<br>数量/个 |
|-----|--------------|---------------|-------------|------------|
| 训练集 | 20           | 0             | 构建音频单词      | 0          |
|     | 28           | 300           | 训练 SVM 分类模型 | 61 941     |
| 测试集 | 50           | 150           | 测试检测性能      | 69 211     |

首先,评估基于音频单词的检测方法,其 ROC 曲线如图 3 所示。图 3(a)给出了基于周期性的判别算法中调节连续次数  $N(1 \leq N \leq 10)$  获得的性能,而图 3(b)显示了传统的阈值判别算法在阈值  $T$  从 0 调整到 30 的性能。其中,EE ( $N = 6$ ) 效果较好。不难看出,与传统的阈值判别算法相比,基于周期性的判别算法能够较大地提高检测性能。



(a) 采用 EE 和基于周期性的判别算法

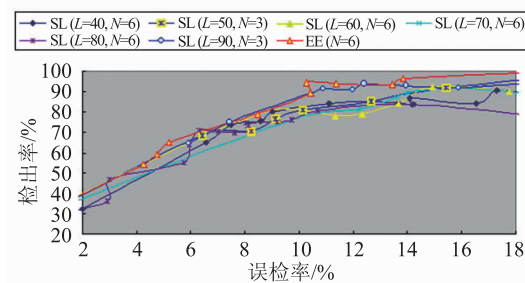


(b) 采用 EE 和阈值判别算法

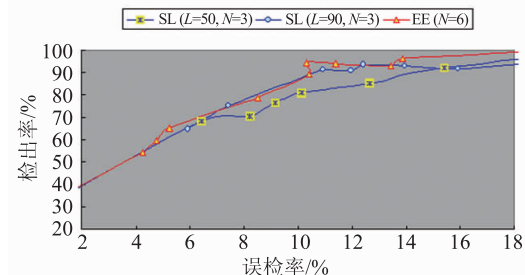
图 3 基于音频单词的检测方法中两种成人视频判别算法的比较

Fig. 3 Comparisons of decision algorithms in adult video detection based on audio-words

其次,评估两种音频片段的分割算法。如图 4(a)所示,定长的音频片段分割算法中,分别选取切分长度  $L = 10, 20, \dots, 100$  和连续次数  $N = 1, 3, 6, 10$  进行实验,其 ROC 曲线表示为  $SL(L, N)$ 。图 4(b)显示了两种分割算法的性能对比,基于周期性的 EE 分割算法优于定长的音频片段分割算法,无需调节切分长度即可获得更好性能。



(a) 采用定长音频片段和基于周期性的判别算法



(b) 两种音频片段分割算法的比较

图 4 基于音频单词的检测方法中两种分割算法的性能比较

Fig. 4 Comparisons of segmentation algorithms in adult video detection based on audio-words

但是,由于每个音频帧要和音频词表内的每个音频单词相匹配,这会增加计算时间。

### 3.2 基于视觉特征的检测方法

从训练集视频中抽取关键帧,经人工标注后得到 8 083 幅成人图像与 31 616 幅非成人图像,将它们的颜色矩与边缘纹理前融合,再用 SVM 训练分类模型。

如图 5(a)所示,在基于视觉特征的检测方法中采用基于周期性的判别算法,分别选取连续次数  $N = 1, 3, 6, 10$  进行实验,其中, KF ( $N = 10$ ) 效果较好。图 5(b)采用阈值判别算法,阈值  $T$  由 0 到 20 变化,随着阈值的增加,开始时性能提高较快,阈值达到某个值时性能不再提高,并有可能下降。这里, KF ( $T > 5$ ) 效果较好。

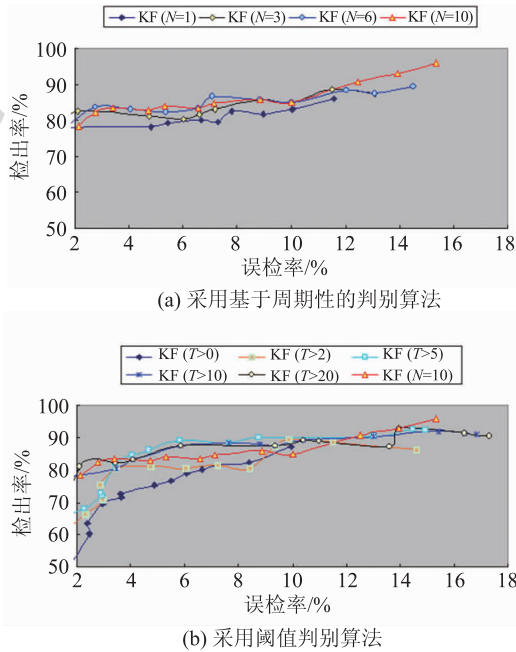


图 5 基于视觉特征的成人视频检测方法  
Fig. 5 Adult video detection based on global features of keyframes

### 3.3 融合音频单词与视觉特征的检测方法

利用 2.1 节方法融合音频单词和视觉特征的 SVM 分类结果,再用 2.2 节方法判别成人视频。其 ROC 曲线表示为 EE + KF。

由图 6(a)可知,当权值  $W = 0.39$  且连续次数  $N = 10$  时,EE + KF ( $N = 10$ )性能优良。

从图 6(b)看出,基于音频单词的检测方法能够获得较高的检出率,但随着误检率的变化波动很大;基于视觉特征的检测方法的检出率偏低,但性能较稳定。

融合两种方法之后,性能显著提高,当误检率为 9.76% 时,检出率达到了 94.44%。

## 4 结 论

由于现有的成人视频检测方法缺乏准确的音频语义表示方法,造成检测性能仍然偏低。因此,本文提出融合音频单词与视觉特征的成人视频检测方法,主要创新点有基于 EE 和 BoW 的音频语义表示方法、基于周期性的 EE 分割算法与成人视频判别算法。

利用 EE 特征中音频单词的出现概率表示音频内容,再和视觉特征进行复合加权融合。实验结果表明,与基于视觉特征的检测方法相比,该方法显著

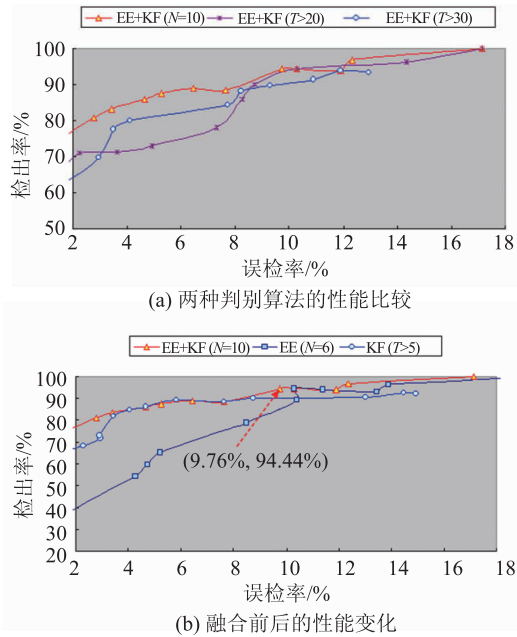


图 6 融合音频单词和视觉特征的成人视频检测方法  
Fig. 6 Fusing audio-words and global features for adult video detection

提高了成人视频的检测性能。当误检率为 9.76% 时,检出率可达 94.44%。

下一步工作是,充分利用视频中的运动信息,将它与音频、视觉等模态信息进行融合,进一步提高检测性能。

### 参考文献 (References)

- [ 1 ] Forsyth D A, Fleck M M. Automatic detection of human nudes [J]. International Journal of Computer Vision, 1999, 32 (1): 63-77.
- [ 2 ] Jones M J, Rehg J M. Statistical color models with application to skin detection [J]. International Journal of Computer Vision, 2002, 46 (1): 81-96.
- [ 3 ] Wang J Z, Li J, Wiederhold G, et al. System for screening objectionable images[J]. Computation Communication, 1998, 21: 1355-1360.
- [ 4 ] Rowley H A, Yushi J, Baluja S. Large scale image-based adult-content filtering[C] // Proceedings of the 1st International Conference on Computer Vision Theory and Applications. New York City: Springer, 2006: 290-296.
- [ 5 ] Wang Y S, Li Y N, Gao W. Detecting pornographic images with visual words[J]. Transactions of Beijing Institute of Technology, 2008, 28 (5): 410-413. [王宇石, 李运宁, 高文. 基于局部视觉单词分布的敏感图像检测[J]. 北京理工大学学报, 2008, 28 (5): 410-413.]

- [ 6 ] Deselaers T, Pimenidis L, Ney H. Bag-of-visual-words models for adult image classification and filtering[ C ] // Proceedings of the 19th International Conference on Pattern Recognition. Tampa, USA: University of South Florida, 2008: 1-4.
- [ 7 ] Endeshaw T, Garcia J, Jakobsson A. Classification of indecent video by low complexity repetitive motion detection[ C ] // Proceedings of 37th Applied Imagery Pattern Recognition Workshop. Washington DC, USA: IEEE Computer Society Press, 2008.
- [ 8 ] Jansohn C, Ulges A, Breuel T M. Detecting pornographic video contents by combining image features with motion information [ C ] // Proceedings of the 17th ACM International Conference on Multimedia. New York City, USA: ACM Press, 2009: 601-604.
- [ 9 ] Rea N, Lacey G, Lambe C, et al. Multimodal periodicity analysis for illicit content detection in videos[ C ] // Proceedings of 3rd European Conference on Visual Media Production. Washington D. C. , USA: IEEE Computer Society Press, 2006: 106-114.
- [ 10 ] Zuo H Q, Wu O, Hu W M, et al. Recognition of blue movies by fusion of audio and video[ C ] // Proceedings of IEEE International Conference on Multimedia & Expo. Washington DC, USA: IEEE Computer Society Press, 2008: 37-40.
- [ 11 ] Zhang Y B, Zhou J, Bian Z Q, et al. A two-stage content-based audio segmentation algorithm[ J ]. Chinese Journal of Computers, 2006, 29(3): 457-465. [张一彬, 周杰, 边肇祺, 等. 一种基于内容的音频流二级分割方法[ J ]. 计算机学报, 2006, 29(3): 457-465.]
- [ 12 ] Lu L, Hanjalic A. Towards optimal audio “keywords” detection for audio content analysis and discovery[ C ] // Proceedings of the 14th ACM International Conference on Multimedia. New York City, USA: ACM Press, 2006: 825-834.
- [ 13 ] Lu L, Hanjalic A. Audio keywords discovery for text-like audio content analysis and retrieval[ J ]. IEEE Transactions on Multimedia, 2008, 10 ( 1 ): 74-85.
- [ 14 ] Xu M, Xu C, Duan L Y, et al. Audio keywords generation for sports video analysis[ J ]. ACM Transactions on Multimedia Computing, Communications and Applications, 2008, 4 ( 2 ): Article 11.
- [ 15 ] Peng Y X, Lu Z W, Xiao J G. Semantic concept annotation based on audio PLSA model[ C ] // Proceedings of the 17th ACM International Conference on Multimedia. New York City, USA: ACM Press, 2009: 841-844.
- [ 16 ] Zeng Z, Zhang S W, Li H P, et al. A novel approach to musical genre classification using probabilistic latent semantic analysis model [ C ] // Proceedings of IEEE International Conference on Multimedia & Expo. Washington DC, USA: IEEE Computer Society Press, 2009: 486-489.
- [ 17 ] Zhao D, Wang X D, Qian Y L, et al. Fast commercial detection based on audio retrieval [ C ] // Proceedings of IEEE International Conference on Multimedia & Expo. Washington DC, USA: IEEE Computer Society Press, 2008: 1185-1188.