



Semi-Supervised image-to-Video Adaptation for Video Action Recognition

¹ Prof. Bhavana R Maale, ² Nikita Y

¹ Assistant Professor, Dept of Computer Science, VTU Center for PG Studies, Kalaburgi, Karnataka, INDIA

² Student, Dept of Computer Science, VTU Center for PG Studies, Kalaburgi, Karnataka, INDIA

Abstract: Human understanding of behavior has been well studied in computer vision applications. Many popular methods of action recognition have shown that knowledge of action can be gained from motion videos or still images effectively. Many of the current methods for identifying video action suffer from the problem of missing appropriate labeled videos for instruction. In such cases, over-fitting will be a potential issue, and action recognition efficiency is limited. We suggest an adaptation approach to improve video action identification by tailoring image information. The adapted information is used to learn the meanings of associated behavior by examining the components of both branded videos and images. We propose an IVA classifier to achieve good output of video action recognition, which can borrow the information adapted from images based on the common visual features.

INDEX : VIDEO, ADAPTION APPROACH, VA

I. NTRODUCTION

Due to its broad range of applications, such as automatic video tracking and video annotation, action recognition in personal videos created by users has become a significant research topic with the rapid advancement of the Internet and smart phone. Web-based consumer videos are posted by users and created by handheld cameras or mobile phones, which may involve considerable camera shake, occlusion, and background cluttered. Thus, these videos contain wide variations in intraclass within the same group of semance. Recognizing human behavior in such images is now a difficult task. The standard paradigm had been preceded by other methods of action identification. Next, it extracts from videos a large number of local motion features (e.g., space-time interest points (STIP), motion scale invariant feature transform (MoSIFT), etc.). Then, using bag-of-words (BoWs) representation, all the local features are quantized into a histogram map. Finally, the vector-based classifiers (e.g., help vector machine) are used in video testing to perform recognition. Such action-recognition methods have achieved promising results when the videos are clear. Noises and uncorrelated details may however be integrated into the BoW during the extraction and quantization of local characteristics. Therefore, when the videos involve significant camera shake, occlusion, cluttered backdrop and so on, these approaches are typically not reliable and could not be easily generalized. To improve accuracy of identification, relevant components of behavior, such as related objects, human presence, stance, and so on, should be used to create a simpler semanthetic understanding of human actions.

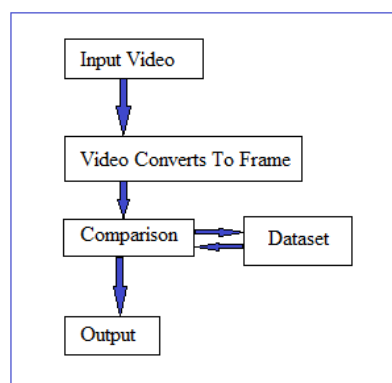


Fig. 1 Proposed System

II. RELATED WORK

[1] B. Ma, L. Huang, J. Shen, and L. Shao, "Discriminative tracking using tensor pooling," *IEEE Trans. Cybern.*, to be published, doi: 10.1109/TCYB.2015.2477879. Human activity acknowledgment has been all around investigated in utilizations of PC vision. Numerous fruitful activity recognition techniques have demonstrated that activity information can be successfully gained from movement recordings or still pictures. For a similar activity, the suitable activity information gained from various kinds of media, e.g., recordings or pictures, might be connected. Be that as it may, less exertion has been made to enhance the execution of activity acknowledgment in recordings by adjusting the activity learning conveyed from pictures to recordings. A large portion of the current video activity acknowledgment strategies experience the ill effects of the issue of lacking sufficient named preparing recordings. In such cases, we will have an issue of over fitting and the execution of activity acknowledgment is limited. Here in this paper, we propose an adjustment strategy to improve activity acknowledgment in recordings by adjusting information from pictures. The adjusted information is used to take in the cor-related activity semantics by investigating the basic parts of both named recordings and pictures. In the mean time, we stretch out the adjustment technique to a semi-directed system which can use both named and unlabeled recordings. In this way, the over-fitting can be eased and the execution of activity acknowledgment is made strides. Analyses on open benchmark datasets and certifiable datasets demonstrate that our technique outflanks a few other cutting edge activity acknowledgment strategies. [2] L. Liu, L. Shao, X. Li, and K. Lu, "Learning spatio-temporal representations for action recognition: A genetic programming approach," *IEEE Trans. Cybern.*, vol. 46, no. 1, pp. 158–170, Jan. 2016, Extracting discriminative and robust features from video sequences is the first and most critical step in human action recognition. In this paper, instead of using handcrafted features, we automatically learn spatio-temporal motion features for action recognition. This is achieved via an evolutionary method, i.e., genetic programming (GP), which evolves the motion feature descriptor on a population of primitive 3D operators (e.g., 3D-Gabor and wavelet). In this way, the scale and shift invariant features can be effectively extracted from both color and optical flow sequences. We intend to learn data adaptive descriptors for different datasets with multiple layers, which makes fully use of the knowledge to mimic the physical structure of the human visual cortex for action recognition and simultaneously reduce the GP searching space to effectively accelerate the convergence of optimal solutions. In our evolutionary architecture, the average cross-validation classification error, which is calculated by an support-vector-machine classifier on the training set, is adopted as the evaluation criterion for the GP fitness function. After the entire evolution procedure finishes, the best-so-far solution selected by GP is regarded as the (near-)optimal action descriptor obtained. The GP-evolving feature extraction method is evaluated on four popular action datasets, namely KTH, HMDB51, UCF YouTube, and Hollywood2. Experimental results show that our method significantly outperforms other types of features, either hand-designed or machine-learned. [3] A. Khan, D. Windridge, and J. Kittler, "Multilevel Chinese takeaway process and label-based processes for rule induction in the context of automated sports video annotation," *IEEE Trans. Cybern.*, vol. 44, no. 10, pp. 1910–1923, Oct. 2014. We propose four variants of a novel hierarchical hidden Markov models strategy for rule induction in the context of automated sports video annotation including a multilevel Chinese takeaway process (MLCTP) based on the Chinese restaurant process and a novel Cartesian product label-based hierarchical bottom-up clustering (CLHBC) method that employs prior information contained within label structures. Our results show significant improvement by comparison against the flat Markov model: optimal performance is obtained using a hybrid method, which combines the MLCTP generated hierarchical topological structures with CLHBC generated event labels. We also show that the methods proposed are generalizable to other rule-based environments including human driving behavior and human actions.

The proposed image-to-video adaptation process frame work s used to identify the action performed n the video that s taken from the user camera or personal mages. The dentification of movement n videos like that s very difficult. First, the user will upload the video as an nput and then the video will be given to the splitter of the fileThe frames that are obtained are then given to the process of comparison here the frames n which the object s present are obtained and that frames are compared with the data set present n the database. The dataset will be generated using SEMI-SUPERVISED methodology, and the functions will be stored for further analysis n the database. We use the SURF algorithm to remove and compare the frames with the dataset. RANSAC algorithm s used to calculate the frames distance, and the classification of the score base s based on the maximum classification value. f the obtained value s maximum then t will identify and show the operation.

III. PROPOSED SYSTEM

In this paper, we propose an adaptation method to enhance action recognition in videos by adapting knowledge from images. The adapted knowledge is utilized to learn the correlated action semantics by exploring the common components of both labeled videos and images. Meanwhile, we extend the adaptation method to a semi-supervised framework which can leverage both labeled and unlabeled videos. Thus, the over-fitting can be alleviated and the performance of action recognition is improved. Experiments on public benchmark datasets and real world datasets show that our method outperforms several other state-of-the-art action recognition methods

IV. METHODOLOGY

The proposed image-to-video adaptation process frame work is used to identify the action performed in the video that is taken from the user camera or personal images. The identification of movement in videos like that is very difficult. Next, the user uploads the video as an input and then the video is transferred to the file splitter. The frames that are obtained are then given to the process of comparison here the frames in which the object is present are obtained and that frames are compared with the data set present in the database. The dataset will be generated using SEMI-SUPERVISED methodology, and the functions will be stored for further analysis in the database. We use the SURF algorithm to remove and compare the frames with the dataset. RANSAC algorithm is used to calculate the frames distance, and the classification of the score base is based on the maximum classification value. If the obtained value is maximum then it will identify and show the operation.

V.EXPERIMENTALRESULTS

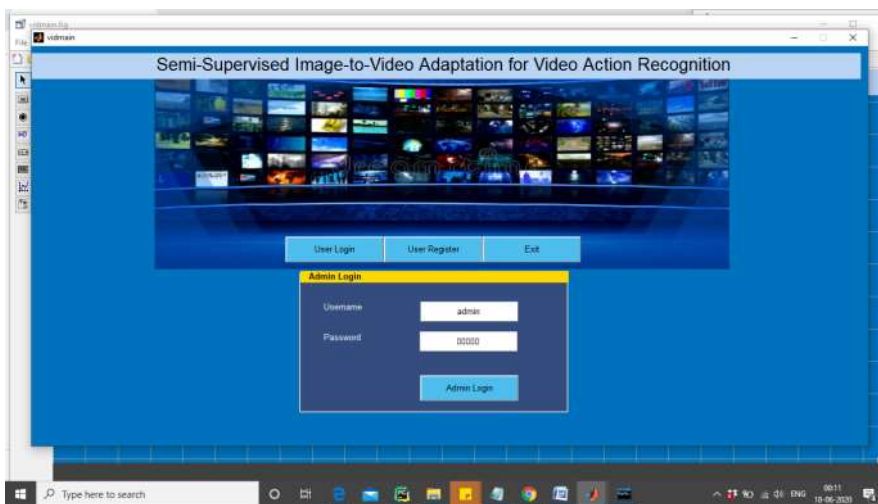


Fig2:HomeScreen

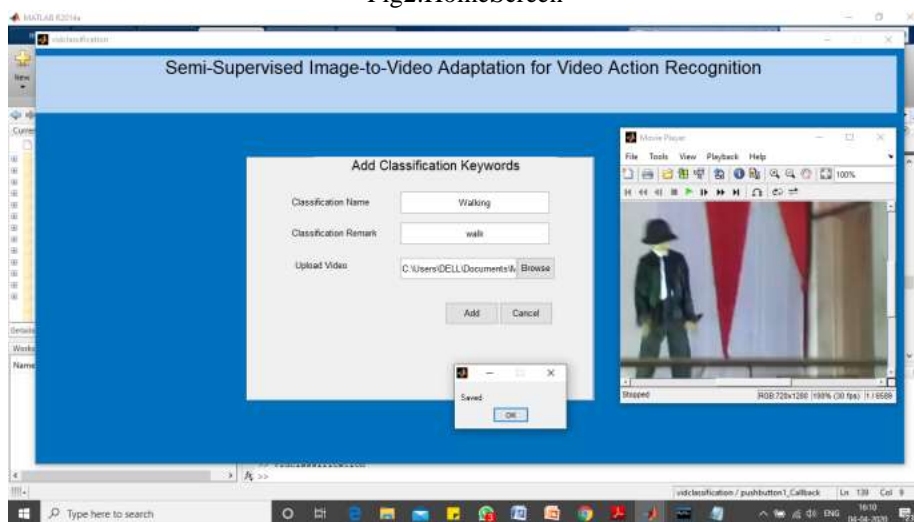


Fig3:AddKeywords&Upload Video

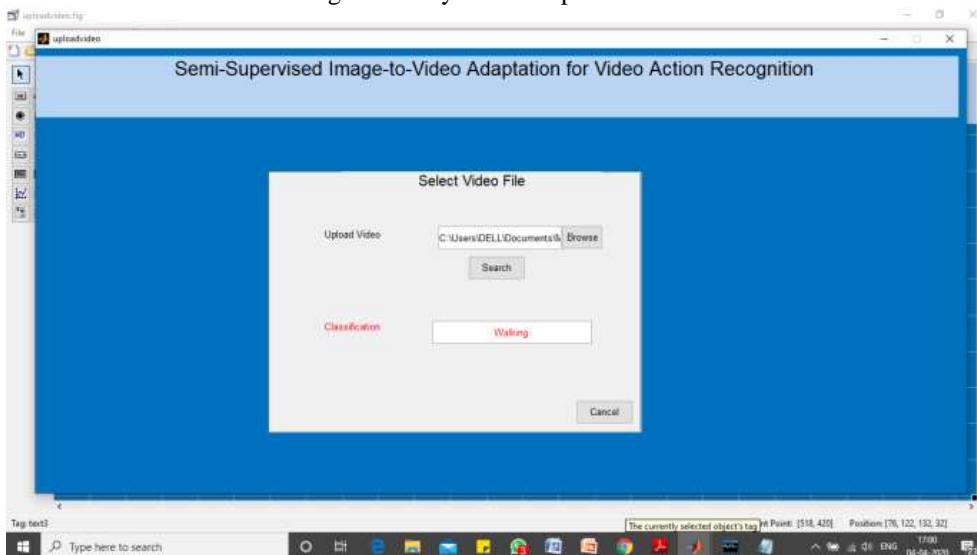


Fig4:Classification

VI.CONCLUSION

We propose an IVA classifier to achieve good output of video action recognition, which can borrow the information adapted from images based on the common visual features. Meanwhile, it can completely exploit the heterogeneous features of unlabeled videos to increase action recognition efficiency in videos. We affirm that information gained from images can affect the accuracy of video recognition and that different recognition results are obtained by the use of different visual indications. Compared with the state-of-the-art approaches, the proposed IVA has better performance in video action recognition. And IVA's success promises when there's only a few branded training videos available.

REFERENCE

- [1] B. Ma, L. Huang, J. Shen, and L. Shao, "Discriminative tracking using tensor pooling," IEEE Trans. Cybern., to be published, doi: 10.1109/TCYB.2015.2477879.
- [2] L. Liu, L. Shao, X. Li, and K. Lu, "Learning spatio-temporal representations for action recognition: A genetic programming approach," IEEE Trans. Cybern., vol. 46, no. 1, pp. 158–170, Jan. 2016.
- [3] A. Khan, D. Windridge, and J. Kittler, "Multilevel Chinese takeaway process and label-based processes for rule induction in the context of automated sports video annotation," IEEE Trans. Cybern., vol. 44, no. 10, pp. 1910–1923, Oct. 2014.
- [4] H. Wang, M. M. Ullah, A. Klaser, I. Laptev, and C. Schmid, "Evaluation of local spatio-temporal features for action recognition," in Proc. Brit. Mach. Vis. Conf., London, U.K., 2009, pp. 124.1–124.11.
- [5] L. Shao, X. Zhen, D. Tao, and X. Li, "Spatio-temporal Laplacian pyramid coding for action recognition," IEEE Trans. Cybern., vol. 44, no. 6, pp. 817–827, Jun. 2014. [6] M.-Y. Chen and A. Hauptmann, "MoSIFT: Recognizing human actions in surveillance videos," School Comput. Sci., Carnegie Mellon Univ., Pittsburgh, PA, USA, Tech. Rep. CMU-CS-09-161, 2009.
- [7] M. Yu, L. Liu, and L. Shao, "Structure-preserving binary representations for RGB-D action recognition," IEEE Trans. Pattern Anal. Mach. Intell., to be published, doi: 10.1109/TPAMI.2015.2491925. [8] L. Shao, L. Liu, and M. Yu, "Kernelized multiview projection for robust action recognition," Int. J. Comput. Vis., 2015, doi: 10.1007/s11263-015-0861-6. [9] C.-C. Chang and C.-J. Lin, "LIBSVM: A library for support vector machines," ACM Trans. Intell. Syst. Technol., vol. 2, no. 3, pp. 1–27, Apr. 2011.
- [10] Y. Han et al., "Semisupervised feature selection via spline regression for video semantic recognition," IEEE Trans. Neural Netw. Learn. Syst., vol. 26, no. 2, pp. 252–264, Feb. 2015.
- [11] C. Thurau and V. Hlaváč, "Pose primitive based human action recognition in videos or still images," in Proc. IEEE Comput. Vis. Pattern Recognit., Anchorage, AK, USA, 2008, pp. 1–8.
- [12] L. Liu, L. Shao, X. Zhen, and X. Li, "Learning discriminative key poses for action recognition," IEEE Trans. Cybern., vol. 43, no. 6, pp. 1860–1870, Dec. 2013.
- [13] A. Gupta, A. Kembhavi, and L. S. Davis, "Observing human-object interactions: Using spatial and functional compatibility for recognition," IEEE Trans. Pattern Anal. Mach. Intell., vol. 31, no. 10, pp. 1775–1789, Oct. 2009.
- [14] B. Yao et al., "Human action recognition by learning bases of action attributes and parts," in Proc. IEEE Int. Conf. Comput. Vis., Barcelona, Spain, 2011, pp. 1331–1338.
- [15] C. Liu, X. Wu, and Y. Jia, "Transfer latent SVM for joint recognition and localization of actions in videos," IEEE Trans. Cybern., to be published, doi: 10.1109/TCYB.2015.2482970.

