

Article Info

Volume 83

Page Number: 5041 - 5047

Article Received: 24 July 2019

Revised: 12 September 2019 Accepted: 15 February 2020 Publication: 27 March 2020

Publication Issue:

March - April 2020

Article History

Contextual Action Recognition in Videos using Tube Convolutional Neural Network

S. Venkata Kiran¹, Dr. S.Venkatnarayanan² ¹Research Scholar (SSSEC1527), ECE Dept., SSSUTMS, Bhopal, (India) ²Professor, ECE Dept., SSSUTMS, Bhopal, (India)

Abstract

In an image classification and object detection Deep learning has been exhibited to accomplish great results.But deep learning on video analysis has been limited due to complexity of video data and lack of annotations. In this paper,we propose Tube Convolutional Technique (T-CT) for action detection in videos. The proposed architecture is a unified deep network that is able to identify and localize action based on 3D convolution features. A video is first divided into equal length eight frame clips and next for each clip a set of tube proposals are generated based on 3D TCT features. Finally, the tube proposals of differents are coupled along using network flow and spatio-temporal action detection is performed victimisation these linked video proposals.

Index Terms; action recognition; T-CT

I. INTRODUCTION

Human action recognition is a fundamental yet challenging task with considerable efforts having been investigated for decades. Motivated by the notable-success of convolutional neural networks (CNNs) for visual recognition in still images, many recent works take advantage of deep models to train end-to-end networks for recognizing actions in videos.

Human actions in video sequences are three dimensional (3D) spatio-temporal signals. Jointly modeling spatiotemporal information via a TCT in an end-to-end deep network provides a natural and efficient approach for action recognition. The objective of activity location is to identify each event of a given activity inside a long video, and to restrict every identification both in space and time. Deep learning getting the hang of learning based methodologies have essentially enhanced video activity acknowledgment execution. Contrasted with activity acknowledgment, activity location is an all the more difficult assignment because of adaptable volume shape and huge spatio-fleeting inquiry space.

In addition, so as to top true both spatial and fleeting data of an activity, two-stream arranges (a spatial CNN and a movement CNN) are utilized. In this way, the spatial and movement data are prepared independently. Region Convolution Neural Network (R-CNN) for protest discovery in pictures was proposed by Girshick et al. [4]. It was trailed by a quick R-CNN proposed in [3], which incorporates the classifier too. Afterward, speedier R-CNN [20] was created by presenting an area proposition organize. It has been widely used to create great outcomes for question recognition in pictures. A characteristic speculation of the RCNN from 2D pictures to 3D spatio-fleeting volumes is to contemplate their adequacy for the issue of activity discovery in recordings. A direct spatio-



fleeting speculation of the R-CNN approach is treat activity recognition in recordings as an arrangement of 2D picture location utilizing quicker RCNN. Nonetheless, tragically, this approach does not take the worldly data into account and isn't adequately expressive to activities. Motivated recognize bv the spearheading work of speedier R-CNN, we propose Tube Convolutional Technique (T-CT) for activity discovery. To better catch the spatioworldly data of video, we misuse 3D ConvNet for activity recognition, since it can catch movement attributes in recordings and shows promising outcome on video activity acknowledgment.We propose a novel system by utilizing the clear energy of 3D ConvNet. In our approach, an info video is partitioned into approach length cuts first. At that point, the cuts are nourished into Tube Proposal Network (TPN) and an arrangement of tube proposition are gotten. Next, tube proposition from each video cut are connected by their actionness scores what's more, cover between nearby proposition to shape a total tube proposition for spatio-transient activity restriction in the video. At long last, the Tube-of-Interest (ToI) pooling is connected to the connected activity tube proposition to produce a settled length include vector for activity mark expectation.

II. LITERATURE SURVEY/RELATED WORK

2D CNN based. To explore the spatio-temporal information in human actions, the two-stream architecture is first proposed in where two 2D CNNs are applied to the appearance

(RGB frames) and motion (stacked optical flow)

domains, respectively. Based on this architecture, several mechanisms are presented to fuse the two networks over the appearance and motion .extend the architecture via the multi-granular structure .

A key volume mining deep framework is designed by Zhu et al. to identify key video clips and perform classification simultaneously Temporal segment networks is proposed which adopts a sparse temporal sampling strategy to enable longrange temporal observations . network for modeling spatio-temporal informationin action recognition . LSTM networks are employed to combine the frame-level features of 2D CNNs to explicitly model spatio-temporal relationships .make use of LSTMs in an encoderdecoder framework for unsupervised video representation. Attention models are also presented based on the recurrent networks to weight the important frames or highly relevant spatio-temporal locations as well.



Figure 1: Overview of the proposed Tube Convolutional Technique(T-CT).

III. R-CNN FROM 2D TO 3D

To better catch the spatio-fleeting data in video,



we misuse 3D CNN for activity proposition age also, activity acknowledgment. One favorable position of 3D CNN over 2D CNN is that it catches movement data by applying convolution in both time and space. Since 3D convolution

what's more, 3D max pooling are used in our approach, not just in the spatial measurement yet in addition in the fleeting measurement, the span of video cut is decreased while recognizable data is concentrated. To deliver a settled length include vector, we propose another pooling layer - Tubeof- Intrigue (ToI) pooling layer. The ToI pooling layer is a 3D speculation of Region-of-Interest (RoI) pooling layer of R-CNN. The exemplary max pooling layer characterizes the part size, walk and cushioning which decides the state of the vield. Interestingly, for RoI pooling layer, the yield shape is settled initially, at that point the portion size and walk are resolved in like manner. Contrasted with RoI pooling which takes 2D highlight guide and 2D districts as information, ToI pooling manages highlight 3D shape and 3D tubes. Signify the span of a component 3D square as $d \times h \times w$, where d, h and w individually speak to profundity, stature and width of the element solid shape. Back-propagation of ToI pooling layer routes the derivatives from output back to the input. Assume xi is the i-th activation to the ToI pooling layer, and yi is the j-th output. Then the partial derivative of the loss unction (L) with respect to each input variable xi can be expressed as:

$$\frac{\partial L}{\partial x_i} = \frac{\partial y_j}{\partial x_i} \frac{\partial L}{\partial y_j} \to \frac{\partial L}{\partial x_i} = \sum_j [i = f(j)] \frac{\partial L}{\partial y_j}.$$
 (1)

IV. T-CT PIPELINE

Our work makes the following contributions:

• We present a Tube Proposal Network, which influences skip pooling in worldly space to safeguard transient data for activity confinement in 3D volumes.

• We propose another pooling layer – Tube-of-Interest (ToI) pooling layer in T-CNN. The ToI pooling layer is a 3D speculation of Region-of-Interest (RoI) pooling layer of R-CNN. It successfully eases the issue with variable spatial and transient sizes of tube proposition. We demonstrate that ToI pooling can enormously enhance the acknowledgment comes about.

A graphical illustration of ToI pooling is presented in Figure 1.



Figure 2: Tube of interest pooling.

As shown in Figure 20ur T-CNN is an end-to-end deep learning framework that takes video clips as input. The core component is the Tube Proposal Network (TPN) to produce tube proposals for each clip. Linked tube proposal sequence represents spatio-temporal action detection in the video and is also used for action recognition.

A. Tube Proposal Network

Our 3D ConvNet comprises of seven 3D convolution layers and four 3D max-pooling layers. We signify the piece state of 3D convolution/pooling by $d \times h \times w$, where d, h,w are profundity, stature and width, individually. In all



convolution layers, the bit sizes are $3\times3\times3$, cushioning and walk stay as 1. The quantities of channels are 64, 128 and 256 individually in the initial 3 convolution layers also, 512 in the rest of the convolution layers. The part measure is set to $1 \times 2 \times 2$ for the initial 3D max-pooling layer, also, $2\times2\times2$ for the rest of the 3D max-pooling layers.



Figure 3: Tube proposal network

B. Linking Tube Proposals

We obtain a set of tube proposals for each video clip after the TPN. We then link these tube proposals to form a proposal sequence for spatiotemporal action localization of the entire video. Each tube proposal from different clips can be linked in a tube proposal sequence for action detection.

C. Action Detection

In the wake of connecting tube recommendations, we get an arrangement of connected tube which speak proposition arrangements, to potential activity cases. The following stage is to proposition group these connected tube arrangements. The tube proposition in the connected groupings may have distinctive sizes. Keeping in mind the end goal to extricate a settled length include vector from every one of the connected proposition grouping, our proposed ToI pooling is used. At that point the ToI pooling layer is trailed by two completely associated layers and a dropout layer. The measurement of the last completely associated layer is N + 1 (N activity classes and 1 foundation class).

V. EXPERIMENTS

To verify the effectiveness of the proposed T-CNN for action detection, we evaluate T-CNN on three trimmed video datasets including UCF-Sports [21], J-HMDB [8], UCF-101 [11] and one un-trimmed video dataset – THUMOS' 14 [12].

We implement our method based on the Caffetoolbox [10]. The TPN and recognition network share weights in their common layers.





Figure 5: Action detection results by T-CT



Figure 6: The ROC and AUC curves for UCF-Sports Dataset [21] are shown in (a) and (b), respectively. The results are shown for Jain et al. [6] (green), Tian et al. [26] (purple), Soomro et al. [23] (blue), Wang t al. [28] (yellow), Gkioxari et al. [5] (cyan) and Proposed Method (red). (c) shows the mean ROC curves for four actions of THUMOS'14. The results are shown for Sultani et al. [24] (green), proposed method (red) and proposed method without negative mining (blue).

VI.CONCLUSION

In this paper we propose a conclusion to-end Tube



Convolutional Technique (T-CT) for activity recordings. recognition in It abuses 3D convolutional system to extricate viable spatiohighlights worldly and perform activity confinement and acknowledgment in a bound together structure. Coarse proposition boxes are thickly tested in view of the 3D convolutional include 3D square and connected for activity acknowledgment and limitation. Broad examinations on a few benchmark datasets show the quality of T-CNN for spatiotemporal confining activities, even in untrimmed recordings.

REFERENCES

- 1. P. Felzenszwalb, D. McAllester, and D. discriminatively Ramanan.A trained. multiscale, deformable part model. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 1-8, 2008.
- 2. A. Gaidon. Z. Harchaoui. and C Schmid.Temporal localization of actions with actoms. IEEE Transactions on Pattern Analysis and Machine Intelligence, 35(11):2782-2795,2013.
- 3. R. Girshick. Fast r-cnn.In IEEE International Conference on Computer Vision (ICCV), December 2015.
- 4. R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2014.
- 5. G. Gkioxari and J. Malik.Finding action tubes. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 759-768, 2015.
- 6. M. Jain, J. Van Gemert, H. J'egou, P. Bouthemy, and C. G.Snoek.Action localization with tubelets from motion. In IEEE Conference on Computer Vision and

Pattern Recognition (CVPR), pages 740-747, 2014.

- 7. M. Jain, J. C. van Gemert, and C. G. Snoek. What do 15,000 object categories tell us about classifying and localizing actions? In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 46-55, 2015.
- 8. H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black.Towards understanding action recognition. IEEE International In Conference on Computer Vision (ICCV), pages 3192-3199, 2013.
- 9. S. Ji, W. Xu, M. Yang, and K. Yu.3d convolutional neural networks for human action recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence, 35(1):221-231, 2013.
- 10. Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. and T. Darrell. Guadarrama. Caffe: Convolutional architecture for fast feature embedding. arXiv preprint arXiv:1408.5093, 2014.
- 11. Y.-G. Jiang, J. Liu, A. RoshanZamir, I. Laptev, M. Piccardi, M. Shah, and R. Sukthankar. THUMOS challenge: Action recognition with a large number of classes. /ICCV13-Action- Workshop/, 2013.
- 12. Y.-G. Jiang, J. Liu, A. RoshanZamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. THUMOS challenge:Action recognition with a large number of classes.http://crcv.ucf.edu/THUMOS14/, 2014.
- 13. R. Joseph and F. Ali. Yolo9000: Better, faster, stronger. InIEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.
- 14. A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with 5045



convolutionalneural networks. In IEEE conference on Computer Vision and Pattern Recognition (CVPR), pages 1725– 1732,2014.

- 15. Y. Ke, R. Sukthankar, and M. Hebert.Event detection incrowded videos. In IEEE International Conference on ComputerVision (ICCV), pages 1–8, 2007.
- 16. T. Lan, Y. Wang, and G. Mori. Discriminative figure-centric models for joint action localization and recognition.In IEEE International Conference on Computer Vision (ICCV), pages 2003– 2010, 2011.
- 17. Y. LeCun, Y. Bengio, and G. Hinton.Deep learning. Nature,521(7553):436–444, 2015.
- F. Negin and F. Bremond. Human action recognition invideos: A survey. INRIA Technical Report, 2016.
- X. Peng and C. Schmid.Multi-region twostream r-cnnforaction detection. In European Conference on Computer Vision(ECCV), pages 744–759, 2016.
- 20. S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towardsreal-time object detection with region proposal networks. In Advances in Neural Information Processing Systems (NIPS),pages 91–99. 2015.
- 21. M. Rodriguez, A. Javed, and M. Shah. Action mach: a spatio-temporal maximum average correlation height filter for action recognition. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 1–8, 2008.
- 22. K. Simonyan and A. Zisserman.Twostream convolutional networks for action recognition in videos. In Advances in Neural Information Processing Systems (NIPS), pages 568–576, 2014.
- 23. K. Soomro, H. Idrees, and M. Shah.Action localization in videos through context

walk.In IEEE International Conference on Computer Vision (CVPR), pages 3280– 3288, 2015.

- 24. W. Sultani and M. Shah. What if we do not have multiple videos of the same action? – video action localization using web images. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2016.
- 25. L. Sun, K. Jia, D.-Y. Yeung, and B. E. Shi. Human action recognition using factorized spatio-temporal convolutional networks.In IEEE International Conference on Computer Vision (ICCV), pages 4597– 4605, 2015.
- 26. Y. Tian, R. Sukthankar, and M. Shah. Spatiotemporal deformable part models for action detection. In IEEE Conference on Computer Vision and Pattern Recognition, pages 2642–2649, 2013.
- 27. D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri.Learning spatiotemporal features with 3d convolutional networks. In IEEE International Conference on Computer Vision (ICCV), pages 4489–4497, 2015.
- 28. L. Wang, Y. Qiao, and X. Tang.Video action detection with relational dynamicposelets. In European Conference on Computer Vision (ECCV), pages 565–580, 2014.
- 29. L. Wang, Y. Qiao, X. Tang, and L. V. Gool.Actionness estimation using hybrid fully convolutional networks. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 2708– 2717, June 2016.
- 30. P. Weinzaepfel, Z. Harchaoui, and C. Schmid.Learning to track for spatiotemporal action localization.In IEEE International Conference on Computer Vision (ICCV), pages 3164–3172, 2015.



31. J. Yue-Hei Ng, M. Hausknecht, S. Vijayanarasimhan,O. Vinyals, R. Monga, and G. Toderici. Beyond short snippets: Deep networks for video classification. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 4694–4702, 2015.