

DEEP ALTERNATIVE NEURAL NETWORK: EXPLORING CONTEXTS AS EARLY AS POSSIBLE FOR ACTION RECOGNITION



Jinzhuo Wang, Wenmin Wang, Xiongtao Chen, Ronggang Wang, Wen Gao
School of Electronics and Computer Engineering, Peking University

TASK AND MOTIVATIONS

We aim at automatically assigning high-level concepts to each untrimmed video, known as action recognition task. Our motivation comes from several limitations of existing approaches.

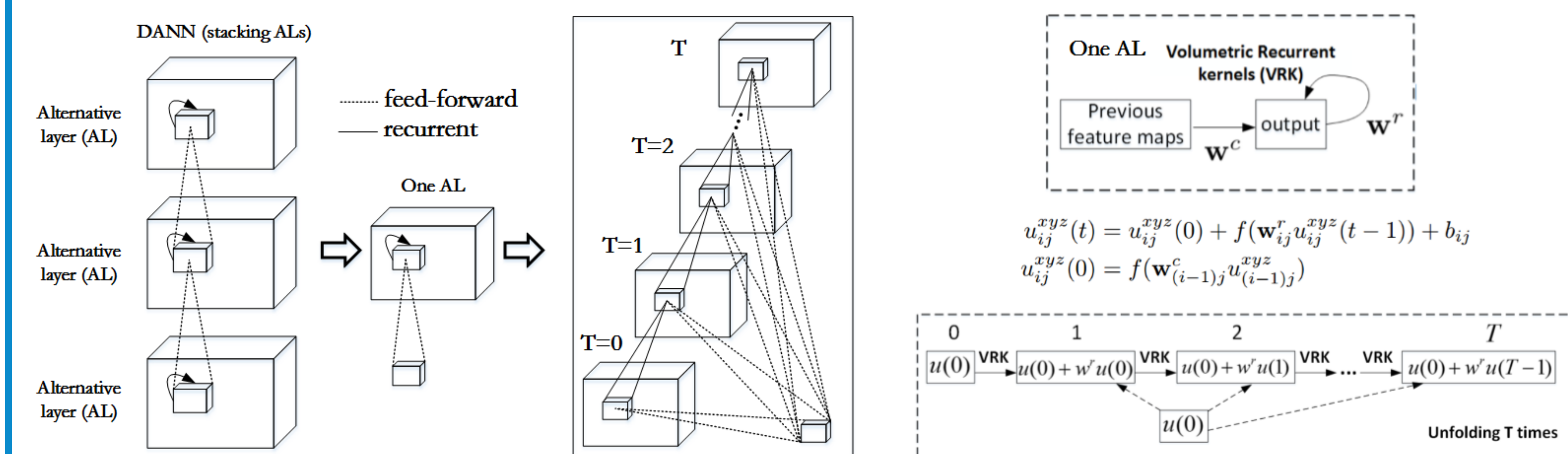
- Existing deep architectures often mine contexts in the latter layers, lacking contexts of low level features. In many computer vision tasks, especially action recognition, contexts of low-level features and its evolutions are crucial for distinguishing similar instances.
- For action recognition in videos, popular methods tend to extract features of entire frames which contain many irrelevant information. Such setting may burden deep network training.

CONTRIBUTIONS

Based on the left discussions about our motivation and existing limitations, we make several contributions as follows.

- A novel deep alternative neural network (DANN) for action recognition in videos, which mines contexts and its evolutions at the beginning of networks.
- An adaptive method to determine the temporal size of input video clip for deep network based on the density of optical flow energy, instead of manual choices used in relevant methods.
- A volumetric pyramid pooling layer to re-size the output of input video clips of arbitrary sizes to fixed-size before fc layers.

ARCHITECTURE

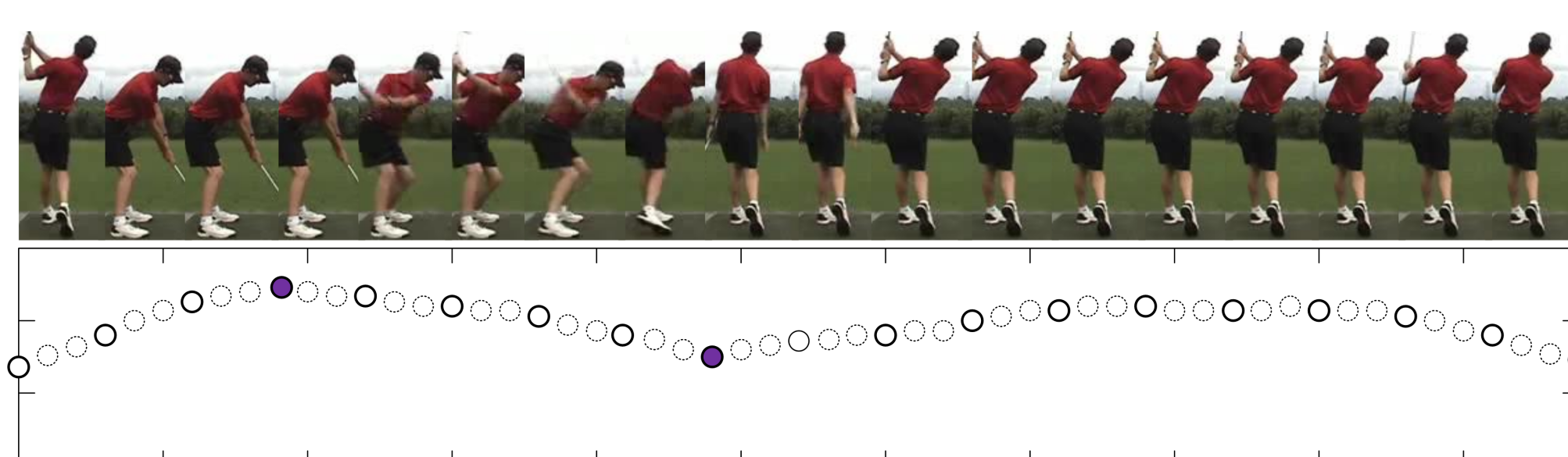


The key component of DANN is the alternative layer (AL), which consists of a standard volumetric convolutional layer followed by a designed recurrent layer. Volumetric convolution is performed to extract features from local spatiotemporal neighborhoods on feature maps in the previous layers. Then a recurrent layer is applied to the out-

put and iteratively proceeds for T times. This procedure makes each unit evolve over discrete time steps and aggregate larger RFs. The overall architecture of DANN has 6 alternative layers with 64, 128, 256, 256, 512 and 512 kernel response maps, followed by a VPPL and 3 fc layers of size 2048 each.

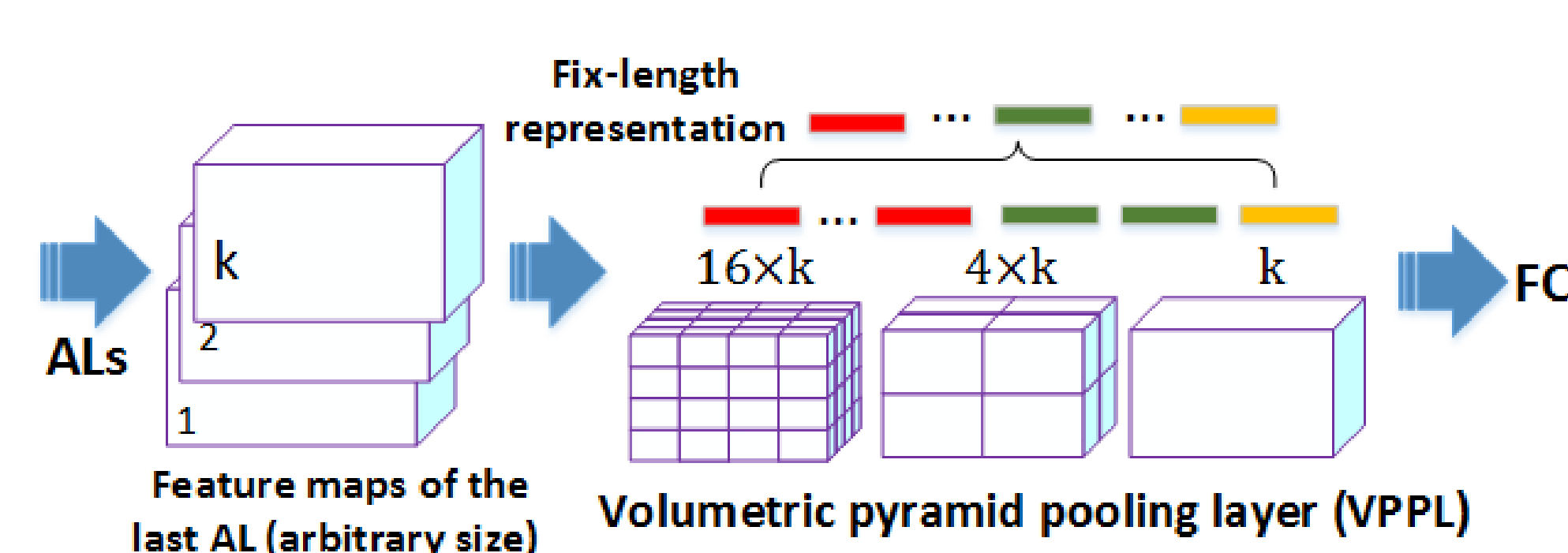
AL1	Pool1	AL2	Pool2	AL3	Pool3	AL4	Pool4	AL5	Pool5	AL6	VPPL	FC1	FC2	FC3	Softmax
64		128		256		256		512		512		2048	2048	2048	

ADAPTIVE INPUT



The input size of deep networks in temporal domain is often determined empirically since it is hard to evaluate all the choices. We instead introduce an adaptive method to automatically select the most discriminative video fragments using the density of optical flow energy. We attempt to preserve as much as motion information and appropriate range dependencies while not breaking their semantic structures in temporal domain. Many evidences show that motion energy intensity induced by human activity exhibits regular periodicity. We leverage this point for segmentation to pursue adaptive input.

VPP LAYER



The AL accepts input videos of arbitrary sizes and produces outputs of variable sizes. However, the fully connected layers require fixed-length vectors. To adopt DANN for input video clips of arbitrary sizes, we replace the last pooling layer with a volumetric pyramid pooling layer (VPPL) inspired by the success of spatial pyramid pooling layer (SPPL). The outputs of the VPPL are kM -dimensional vectors, where M is the number of bins and k is the number of kernels in the last alternative layer, which are then sent to the fc layers.

RESULTS

Table 4: Comparison with the state-of-the-art on HMDB51 and UCF101 (over three splits).

	Method	HMDB		UCF		
		51	51	101	101	
CNN	Slow fusion [13]	-	65.4	-	88.0	
	C3D [28]	-	85.2	-	88.6	
	Two-Stream(spatial) [22]	40.5	73.0	-	90.3	
	Two-Stream(temporal) [22]	54.6	83.7	-	91.5	
	LTC [29]	57.9	83.3	-	90.4	
	Very deep (temporal) [32]	-	87.0	-	91.4	
	Very deep (spatial) [32]	-	87.0	-	88.6	
Hand	IDT+FV [30]	57.2	85.9	-	89.2	
	IDT+HSV [19]	61.1	87.9	-	91.6	
	IDT+MIFS [16]	65.1	89.1	-	-	
	IDT+SFV [20]	66.8	-	-	-	
	Method	HMDB	UCF	Method	HMDB	UCF
	Two-stream [22]	59.4	88.0		59.4	88.0
	CNN+deep LSTM [35]	-	-		-	-
	TDD [31]	63.2	90.3		63.2	90.3
	TDD+IDT [31]	65.9	91.5		65.9	91.5
	C3D+IDT [28]	-	90.4		-	90.4
	Very deep (two-stream) [32]	-	91.4		-	91.4
	LTC+spatial	61.5	88.6		61.5	88.6
	DANN	63.3	89.2		63.3	89.2
	DANN+spatial	65.9	91.6		65.9	91.6

Table 4 reports the best DANN model and state-of-the-art approaches over three splits on UCF101 and HMDB51 datasets in terms of video-level accuracy. As can be seen, trajectory-based

features are still competitive in the area of deep learning, especially with the help of high-order encodings or deep architectures. Fusion strategies often outperform pure single deep networks. Note that all the other deep networks use a pre-defined temporal length to generate video clip as input such as 16-frame and 60-frame, while our DANN determines it in an adaptive manner. Combined with spatial stream, DANN achieves the accuracy of 65.9% and 91.6% on HMDB51 and UCF101, separately.

MAIN REFERENCES

- [1] K. Simonyan and A. Zisserman Two-stream Convolutional Networks for Action Recognition in Videos. *NIPS '14*
- [2] M. Liang, X. Hu and B. Zhang. Convolutional Neural Networks with Intra-layer Recurrent Connections for Scene Labeling. *NIPS '15*

FUTURE DIRECTION

Although adaptive temporal length is used in our model, the spatial size is still chosen in ad hoc manner. A natural direction as discussed in our paper is using a more compact video clip such as *tube* as the input of our DANN.

We expect the proposed architecture DANN be a general framework used for more CV tasks especially in which contexts are essential. Also, some related domains such as NLP have potential to enjoy the advantages of our DANN.

ACKNOWLEDGEMENT

This work was supported by Shenzhen Peacock Plan (20130408-183003656). The complete source code and pre-trained models are available at <http://github.com/wangjinzhuo/DANN>.