

Recognizing fine-grained actions



Mallya and Lazebnik, ECCV'16 simplify it to only use full image as the context

Recent methods hard-code attention on the human

Our Model

Jointly predict an attention map to modulate the final convolutional feature before spatial average pool. The attention map can be unconstrained, or regularized using human pose keypoints.





OPTION OPTION 2

Improved action recognition performance! (MPII human-pose/action dataset)

Model (ResNet-101)	Val mAP
No attention	26.2%
Linear attention	30.3%
Pose-regularized attention	30.6%

Attentional Pooling for Action Recognition Rohit Girdhar and Deva Ramanan

Person Brush Lawn mower

Pose-regularized attention

Linear attention

What makes the model work? Rank-1 approximation of second order pooling

 \cong





- Bottom-up saliency: Certain image regions seem to pop-out
- Top-down task-guided attention: Similar to Ullman's Visual Routines (Cognition, 1984), extract desired information from base representation modulated by saliency

Why not only top-down attention?

already have top-down attention and

adding bottom-up saliency improves

Standard average pool models

performance significantly



Comparison with the state of the art

MPII pose dataset

Classify into one of 393 action classes

Method	Test mAP	Method	mAF
R*CNN (ICCV'15)	26.7%	R*CNN	28.5%
Mallya and Lazebnik (ECCV'16)	31.9%	Mallya & Lazebnik (w/o wtd loss)	33.8%
Ours (Linear attention)	36.0%	Ours (Linear attention)	35.0%
Ours (Pose regularized attention)	36.1%		
		Method (RGB stream only)	Accuracy
HMDB-51 dataset \rightarrow Classify short-video clips into 51 action classes.		TSN (ECCV'16)	51.0%
		ResNet-152 (NIPS'16)	46.7%
two-stream architectures) with attention pooling. At test time, we average prediction from uniformly sampled frames	attentional	TSN, ResNet-101 (ours)	47.1%
	predictions	Ours (Linear attention)	50.8%
		Ours (Pose reg. attention)	52.2%

Method	mAP
R*CNN	28.5%
Mallya & Lazebnik (w/o wtd loss)	33.8%
Ours (Linear attention)	35.0%
Method (RGB stream only)	Accuracy
TSN (ECCV'16)	51.0%
ResNet-152 (NIPS'16)	46.7%
TSN, ResNet-101 (ours)	47.1%
Ours (Linear attention)	50.8%
Ours (Pose reg. attention)	52.2%

Conclusion

- Consider replacing pooling with attentional pooling. Lightweight yet powerful!
- Connection to bilinear pooling suggests action recognition ~ fine-grained task!
- Self attention out-performs sequential attention. Our approach is an efficient implementation for self-attention.

HICO dataset

Detect 600 human-object interactions

Where does the model look?

Validation images in MPII that obtain strongest improvement in performance



The model learns to look for objects to recognize interactions

