

Action Recognition in Videos – State of the Art

- Bag of Features video representations have become popular descriptor sampling is either:
- ▶ interest point based: e.g. SIFT, SURF; few stable descriptors only
- dense grid-based: provides robust recognition \rightarrow preferred
- In drawbacks of dense sampling: must process a large number of (often irrelevant) features \rightarrow huge computational load



Motivation – Biological Inspiration

- space-variant processing of the human visual system:
- first, potentially relevant (i.e. salient) regions are detected
- Further detailed processing happens only at salient locations
- **goal**: use saliency either to emphasize the most informative parts of the visual scene or to limit the processing to them

Saliency Masks

- **central mask**: filmmakers place the subject of interest in the center \rightarrow "distance of pixel to center" as simplest saliency measure
- analytical saliency mask: saliency model built on the structure tensor and its geometric invariants
- ▶ for a video f(x, y, t) the structure tensor J is:

$$\mathbf{J} = \int_{\Omega} \nabla f \otimes \nabla f \, \mathrm{d}\Omega = \int_{\Omega} \begin{bmatrix} f_x^2 & f_x f_y & f_x \\ f_x f_y & f_y^2 & f_y \\ f_x f_t & f_y f_t & f_y \end{bmatrix}$$

J's geometric invariants characterize typical video structures (uniform regions, edges, corners, transient corners):

$$H = 1/3 \operatorname{trace}(\mathbf{J})$$

 $S = M_{11} + M_{22} + M_{33}$
 $K = |\mathbf{J}|$

predict well eye movements on naturalistic videos (Vig et al. PAMI'12) **• empirical saliency mask**: a "ground truth" saliency map determined by measuring where humans actually look in the video \rightarrow fixation density map: superposition of Gaussians centered at each gaze sample



Space-Variant Descriptor Sampling for Action Recognition Based on Saliency and Eye Movements

Eleonora Vig¹, Michael Dorr², and David Cox¹

 1 The Rowland Institute at Harvard, Harvard University, 2 Schepens Eye Research Institute, Harvard Medical School

- peripheral mask: central mask inverted



Video Representations

- Motion Boundary Histograms

$$F(x;k,\lambda) = 1 - e^{-(x/\lambda)^k}$$

with k > 0 shape and $\lambda > 0$ scale parameters, to sample descriptors with certain probability

descriptor sets:

- > 5 hours Hollywood clips: 823 training and 884 test videos; 12 action classes
- binocular eye movements recorded at 1000 Hz from 3 subjects whose task was to identify the action(s) (gaze set is publicly available)







Results HOGHOF nter Ŭ 50 eriphery Ο S variant gaze 52 Conclusions

- many descriptors are unnecessary: discarding up to 70% has no effect enhanced with saliency-based descriptor sampling achieves best mAP (60%) on
- mimicking visual attention improves action recognition: Dense Trajectories Hollywood2 to date
- ► collected eye movements to probe the limits of saliency-based pruning (62%) ► feature combination is beneficial: separate representations for coarse scene gist and detailed foveal view of the scene
- strong center bias in man-made videos





x-axis: % descriptors kept, y-axis: mean Average Precision Legend: prune, concatenate, MKL, dashed line – baseline (no pruning) ► adopted the BoF processing pipeline of Wang et al., BMVC'09