

Table 3: Accuracy of action recognition techniques (Numbers are true recognition accuracy given in percentages. * Datasets in which the mean average precision is reported). The column *Type* indicates whether a method is purely Deep-net based(D), Representation Based(R) or Fused Solution(F).

| Reported Paper | Method | Type | Dataset | | | | | | |
|-------------------------------|---|------|---------|--------|-------|-------------|--------------|-----------------|----------|
| | | | HMDB51 | UCF101 | UCF50 | UCF-Sports* | Holly-wood2* | Olympic Sports* | Sports1M |
| Wang et al. (2011) | Dense Traj (Traj + HOG+HOF+MBH) | R | | | | 88.2 | 58.3 | | |
| Kliper-Gross et al. (2012) | Motion Interchange Patterns | R | 29.2 | | 68.5 | | | | |
| Sadanand and Corso (2012) | General | R | 26.9 | | | | | | |
| | Video Wise | | | | 76.4 | | | | |
| | Group Wise | | | | 57.9 | | | | |
| Oneata et al. (2013) | MBH + SIFT + Sqrt + L2 Normalization | R | 54.8 | | 90 | | 63.3 | 82.1 | |
| Wang and Schmid (2013) | Without Human Detector | R | 55.9 | | 90.5 | | 63 | 90.2 | |
| | With Human Detector | | 57.2 | | 91.2 | | 64.3 | 91.1 | |
| Jain et al. (2013) | Traj + HoG + HoF + MBH + DCS on w -flow | R | 52.1 | | | | 62.5 | | |
| Peng et al. (2014b) | Stacked FVs + FV | R | 66.8 | | | | | | |
| Peng et al. (2014a) | Hybrid-BoW | R | 61.1 | 87.9 | 92.3 | | | | |
| Kantorov and Laptev (2014) | MPEG-Flow : VLAD encodings of | R | 46.3 | | | | | | |
| Gaidon et al. (2014) | SDT tree ATEP | R | 41.3 | | | | 54.4 | 85.5 | |
| Simonyan and Zisserman (2014) | Two-stream (CNN-M-2048) | D | 59.4 | 88.0 | | | | | |
| Karpathy et al. (2014) | Transfer Learning on Sports 1M | D | | 65.4 | | | | | |
| | Clip Hit @ 1 - Slow Fusion | | | | | | | | 41.9 |
| | Video Hit @ 1 - Slow Fusion | | | | | | | | 60.9 |
| Sun et al. (2015) | Factorized Spatio Temporal Conv. Nets | D | 59.1 | 88.1 | | | | | |
| Wang et al. (2015b) | Two-Stream (ClarifaiNet) | D | | 88.0 | | | | | |
| | Two-Stream (GoogLeNet) | | | 89.3 | | | | | |
| | Two-Stream (VGG-16) | | | 91.4 | | | | | |
| Wang et al. (2015a) | TDD + Wang and Schmid (2013) | F | 65.9 | 91.5 | | | | | |
| | TDD (Only) | F | 63.2 | 90.3 | | | | | |
| Ng et al. (2015) | Conv Pooling Hit@1 (Best) | D | | | | | | | 72.4 |
| | LSTM Hit@1 (Best) | | | | | | | | 73.1 |
| | Conv Pooling (Image + Opt Flow) | | | 88.2 | | | | | |
| | LSTM (Image + Opt Flow) | | | 88.6 | | | | | |
| Fernando et al. (2015) | Rank Pooling | R | 63.7 | | | | 73.7 | | |
| Donahue et al. (2015) | LRCN- Weighted Avnerage of RGB + Flow | R | | 82.9 | | | | | |
| Wu et al. (2015) | Adaptive Multi-Stream Fusion | D | | 92.6 | | | | | |
| Jiang et al. (2015) | TrajShape+TrajMF | R | 48.4 | 78.5 | | | 55.2 | 80.6 | |
| | TrajShape+TrajMF+ Wang and Schmid (2013) | | 57.3 | 87.2 | | | 65.4 | 91 | |
| Lan et al. (2015) | Multi-Skip Feat. Stacking | R | 65.1 | 89.1 | 94.4 | | 68.0 | 91.4 | |
| Hoai and Zisserman (2015) | Proposed SSD + RCS | R | 62.2 | | | | 72.7 | | |
| Tran et al. (2015) | C3D on SVM | D | | 85.2 | | | | | |
| | C3D + Wang and Schmid (2013) on SVM | F | | 90.4 | | | | | |
| Misra et al. (2016) | ImageNet pretrain + tuple verification | D | 29.9 | | | | | | |
| | HMDB + UCF101 Labels Only | | 30.6 | | | | | | |
| Wang et al. (2016) | Proposed Only (RBG + Opt Flow Networks) | D | 62 | 92.4 | | | | | |
| Fernando and Gould (2016) | End to End Rank-pooling | D | | | | 87 | 40.6 | | |
| Fernando et al. (2016) | Hierarchical Rank-pooling (CNN Features) | D | 47.5 | 78.8 | | | 56.8 | | |
| | Hierarchical RP on CNN+ Fernando et al. (2015) | F | 65.0 | 90.7 | | | 74.1 | | |
| Li et al. (2016) | VLAD ³ | F | | 84.7 | | | | 90.8 | |
| | VLAD ³ + Wang and Schmid (2013) | F | | 92.2 | | | | 96.6 | |
| Varol et al. (2016) | LTC _{flow+RGB} | D | 64.8 | 91.7 | | | | | |
| | LTC _{flow+RGB} + Wang and Schmid (2013) | F | 67.2 | 92.7 | | | | | |
| Feichtenhofer et al. (2016) | Two Stream Fusion (VGG-16) | D | 65.4 | 92.5 | | | | | |
| | Two Stream Fusion (VGG-16) + Wang and Schmid (2013) | F | 69.2 | 93.5 | | | | | |
| de Souza et al. (2016) | Hybrid fusion of Wang and Schmid (2013) + Deep-nets | F | 70.4 | 92.5 | | | 72.6 | | |