

deconvolutions, forming the decoding part of the network. The KDN is organized sim-ilar to [1]. The first five convolutional subnetworks (conv1 to conv5) each consist of multiple temporal convolutional layers, where each temporal convolution is followed by batch normalization and ReLU activation. Each convolutional subnetwork is fol-lowed by a temporal max-pooling. The subnetworks conv6 and conv7, each consists of a temporal convolution followed by ReLU and dropout. The subnetwork conv8 consists of 1×1 convolution with batch normalization. We apply deconvolution along the time axis on the output of conv8. We also apply a 1×1 convolution and batch normaliza-tion to the output of pool4 and add (element-wise) the result with deconv1 features. This skip connection, used in semantic segmentation to produce better visual features,



Fig. 2: Visualization of the self-supervised learned attention model on two videos from the task 'assemble clarinet' (left) and 'perform CPR' (right) from ProceL. Notice that our method successfully learns to focus on important region of each frame. For example, for clarinet, it focused on cork, ligature, screws, lower and upper joints in the associated key-steps.

is also useful in key-step discovery, as it helps to recover temporal information for key-step and video classification. Finally, we apply a temporal deconvolution and obtain the final predictions, which are T_{ℓ} outputs each of dimension M.

More Results from Visual Attention. Figure 2 shows more results for the visualization of our self-supervised learned attention model on two videos from the task 'assemble clarinet' (left) and 'perform CPR' (right) from ProceL. Notice that our method success-fully learns to focus on important region of each frame. For example, for clarinet, it focuses on cork, ligature, screws, lower and upper joints in the associated key-steps. This is particularly vital for localization of key-steps and recognition of the task, as also demonstrated by our quantitative results.

References

069	1. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation.
070	IEEE Conference on Computer Vision and Pattern Recognition (2015)