



Sequence to Sequence – Video to Text (S2VT)

Subhashini Venugopalan¹, Marcus Rohrbach², Jeff Donahue²,

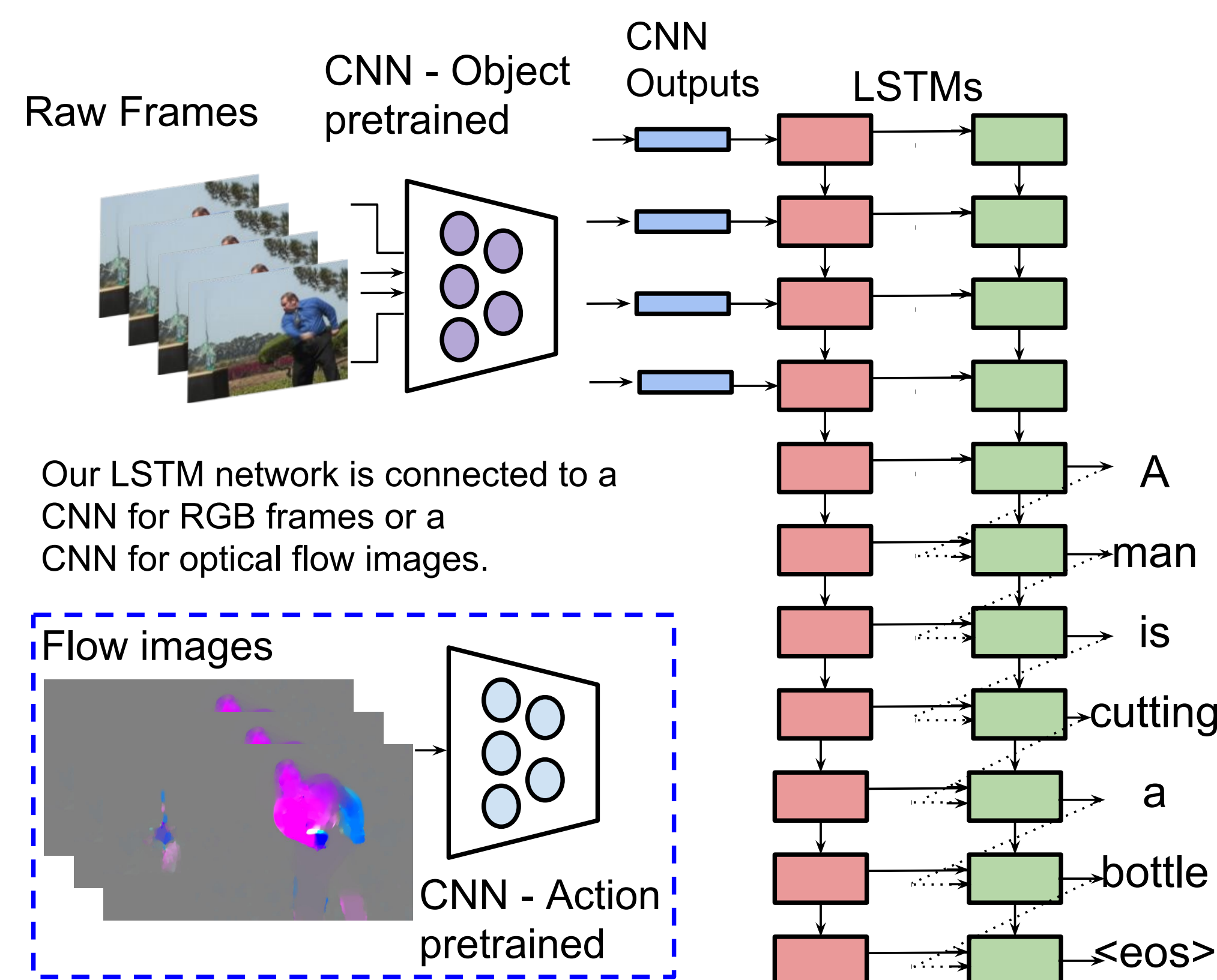
Raymond Mooney¹, Trevor Darrell², Kate Saenko³

¹ UT-Austin ² UC-Berkeley ³ UMass-Lowell



GOALS

Given a video clip, output a natural language sentence that describes the event depicted in the video.



Our model uses a CNN-RNN based encoder-decoder approach and learns to associate a sequence of video frames to a sequence of words in order to generate a description of the event in the video clip. The model is trained on raw RGB frames as well as optical flow images.

DATASETS

We demonstrate our approach on large, realistic collections of YouTube videos and Hollywood movie clips.



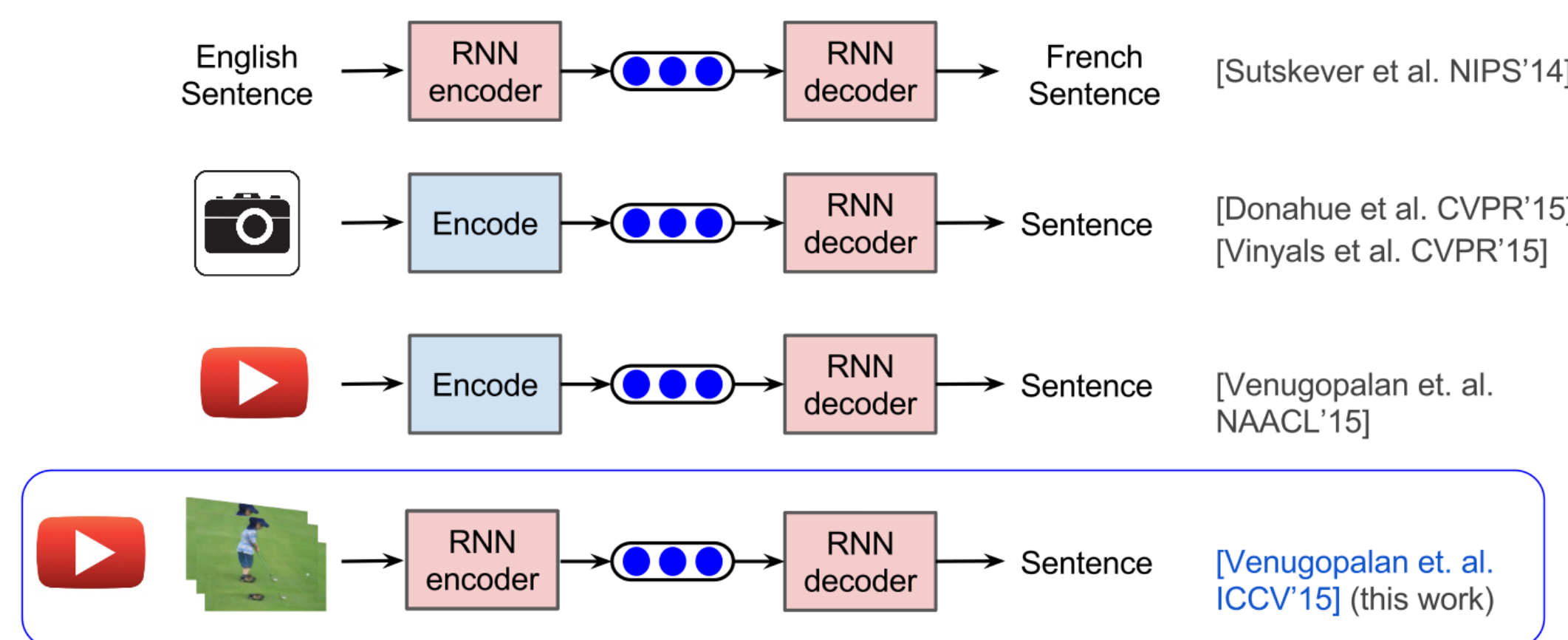
(a) YouTube Video corpus



(b) MPII Movie Description Dataset

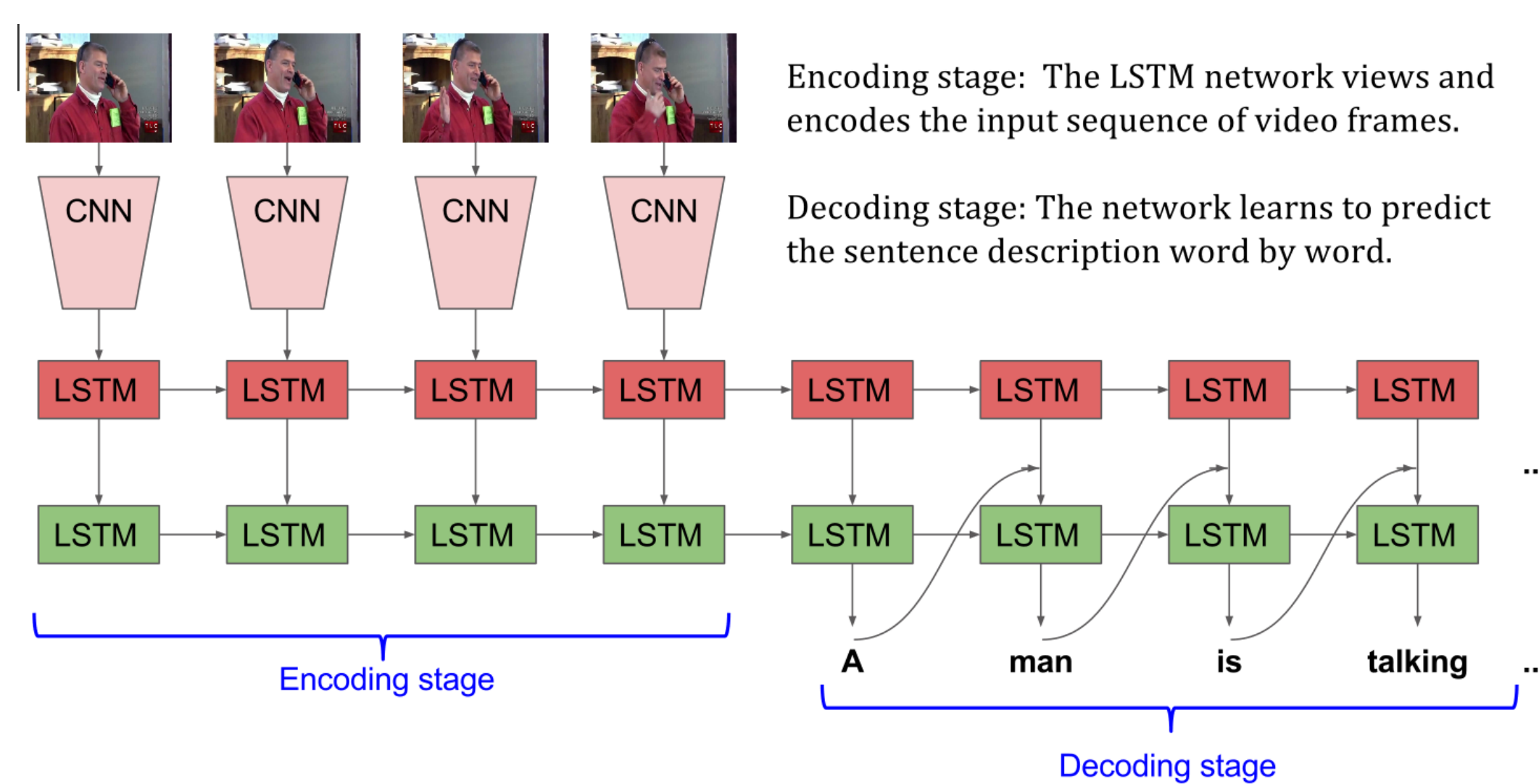
The YouTube dataset, collected by Chen and Dolan (ACL 2011) consists of 1970 videos, where each video is accompanied by about 41 human descriptions (sentences), see (a) above. We also show results on two large movie description corpora - the Montreal (M-VAD) and MPII movie description datasets consisting of nearly 200 Hollywood movies with DVS sentences, see (b) above.

INSIGHT



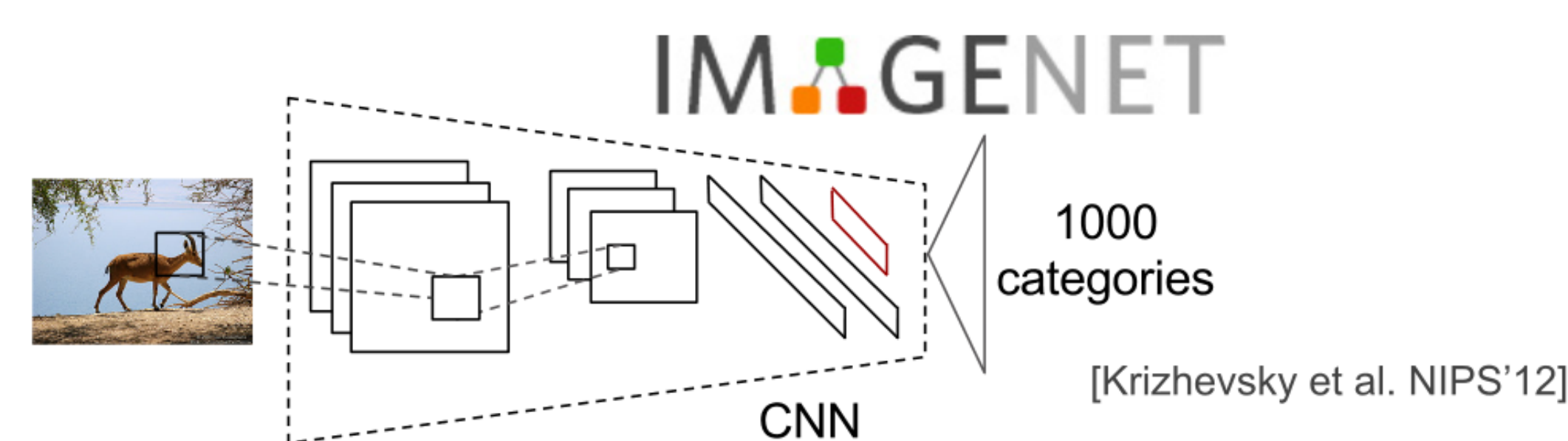
The broad idea of our approach is to encode a video frame sequence and decode it to a sequence of english words (sentence) using LSTMs.

OVERVIEW



The image is forward propagated through a CNN. The activations of the fully connected layer just before classification forms the input to the LSTM network.

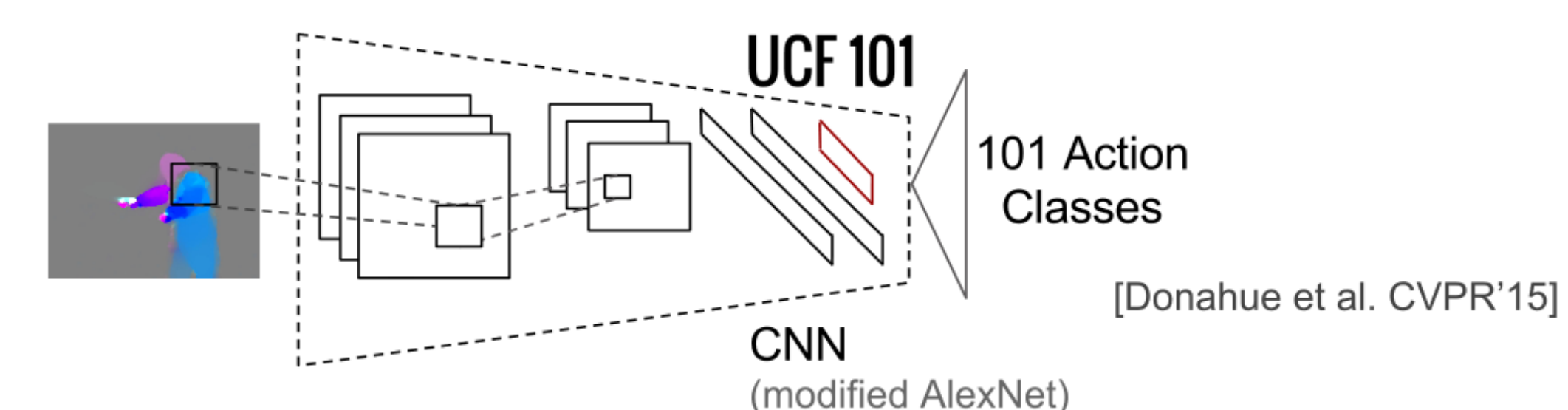
Raw RGB Frames: CNN is initialized with weights from a model trained on the ImageNet classification task.



Optical Flow Frames: Additionally, our model incorporates activity information using optical flow.

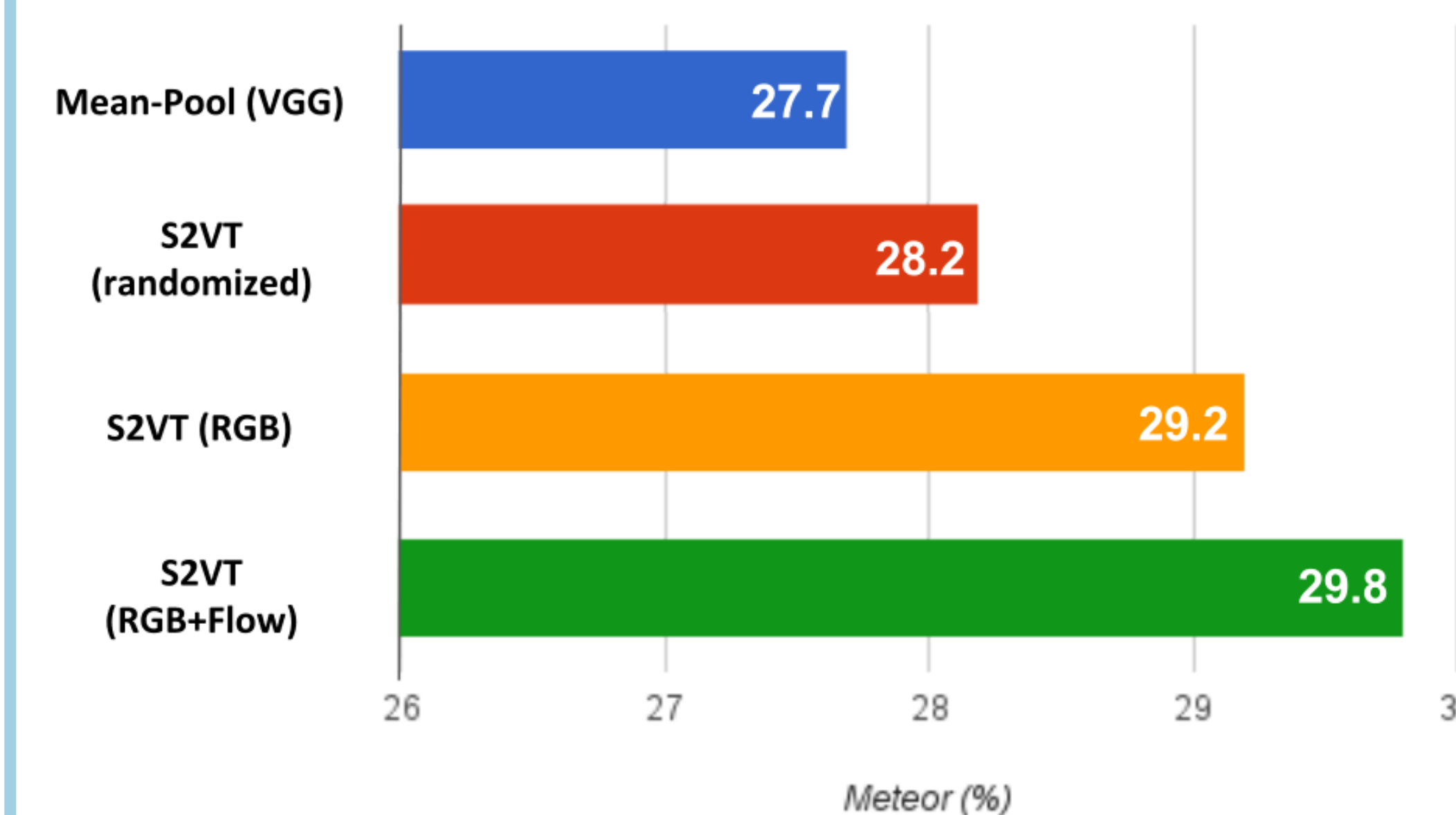


Optical flow is extracted by considering two consecutive frames in the video (above). The flow CNN in the S2VT network is initialized with weights from a model trained on the UCF101 activity recognition task (below).



RESULTS - YOUTUBE DATASET

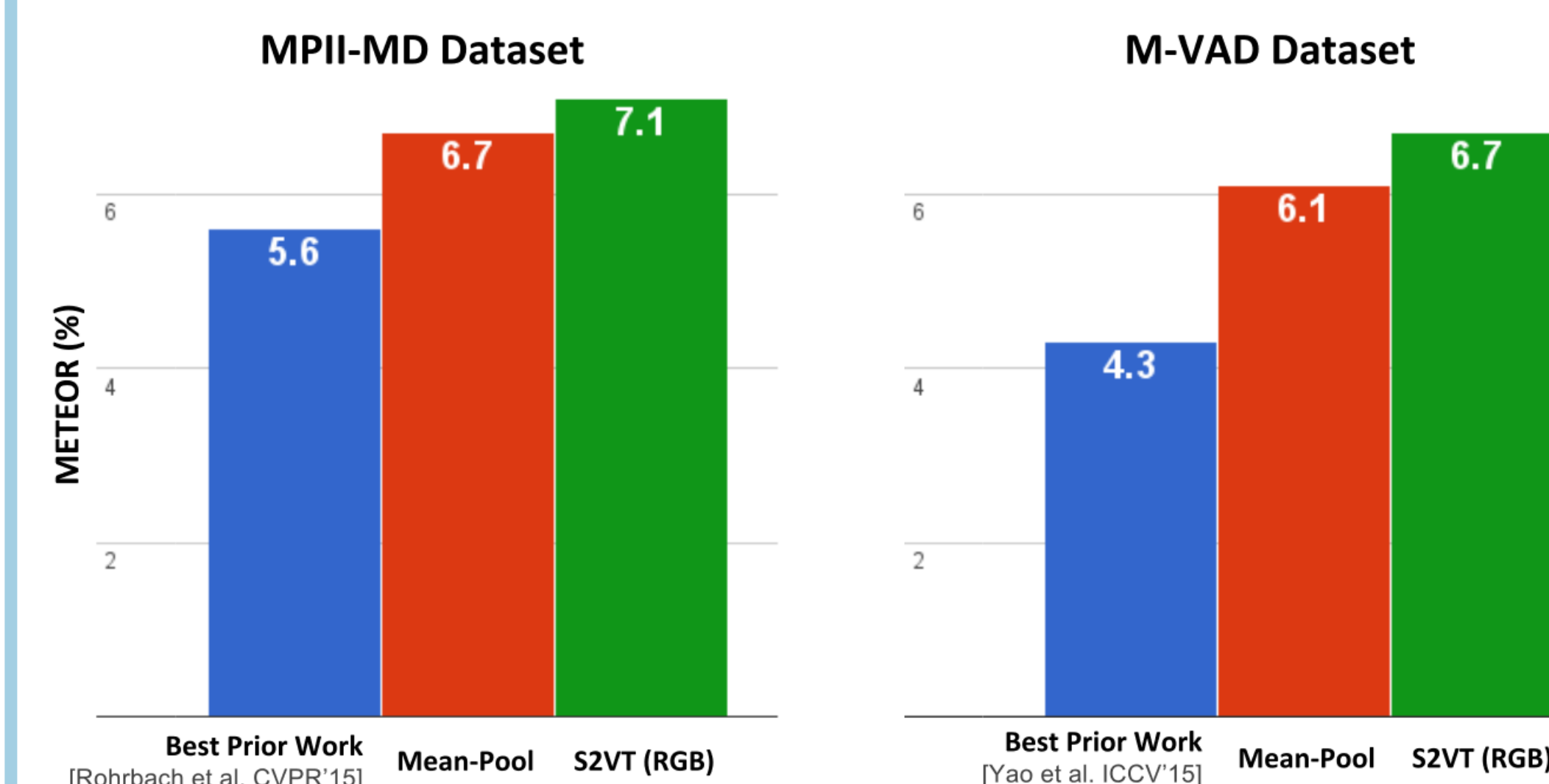
Results on the MSVD corpus of Youtube video clips.



We use the machine translation metric METEOR to compare the quality of the generated description against the multiple ground truth reference sentences.

RESULTS - MOVIE DESCRIPTION

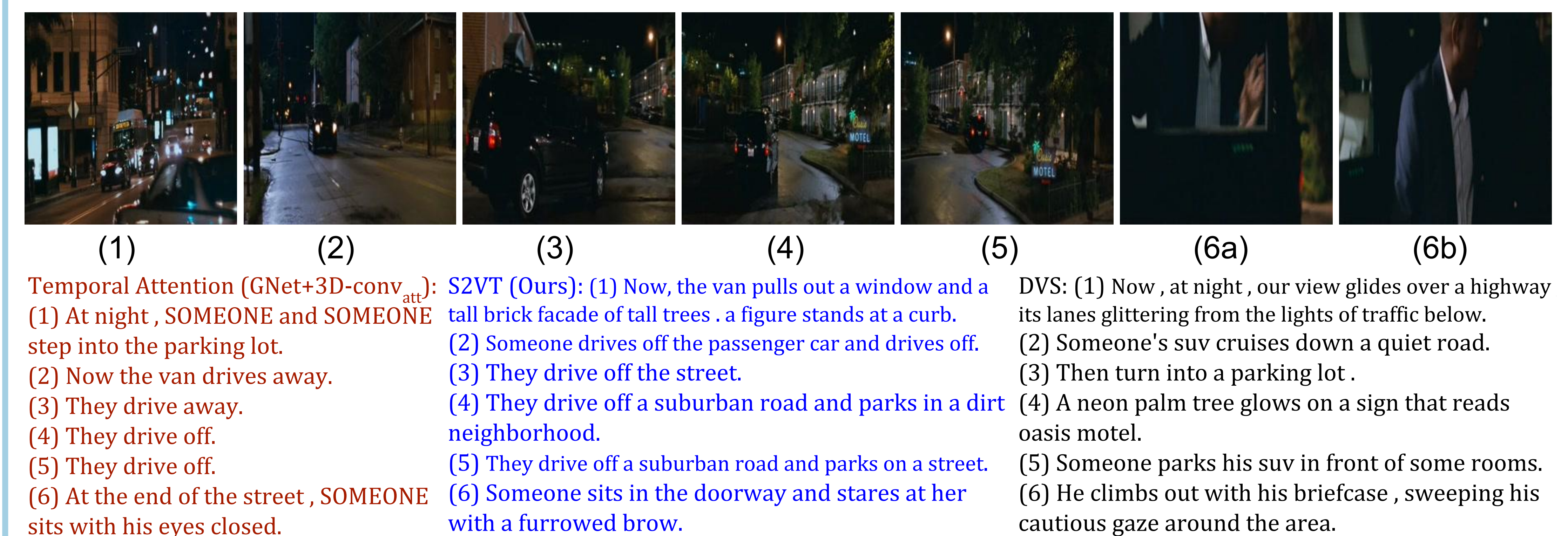
Results on the MPII-MD and M-VAD movie corpora.



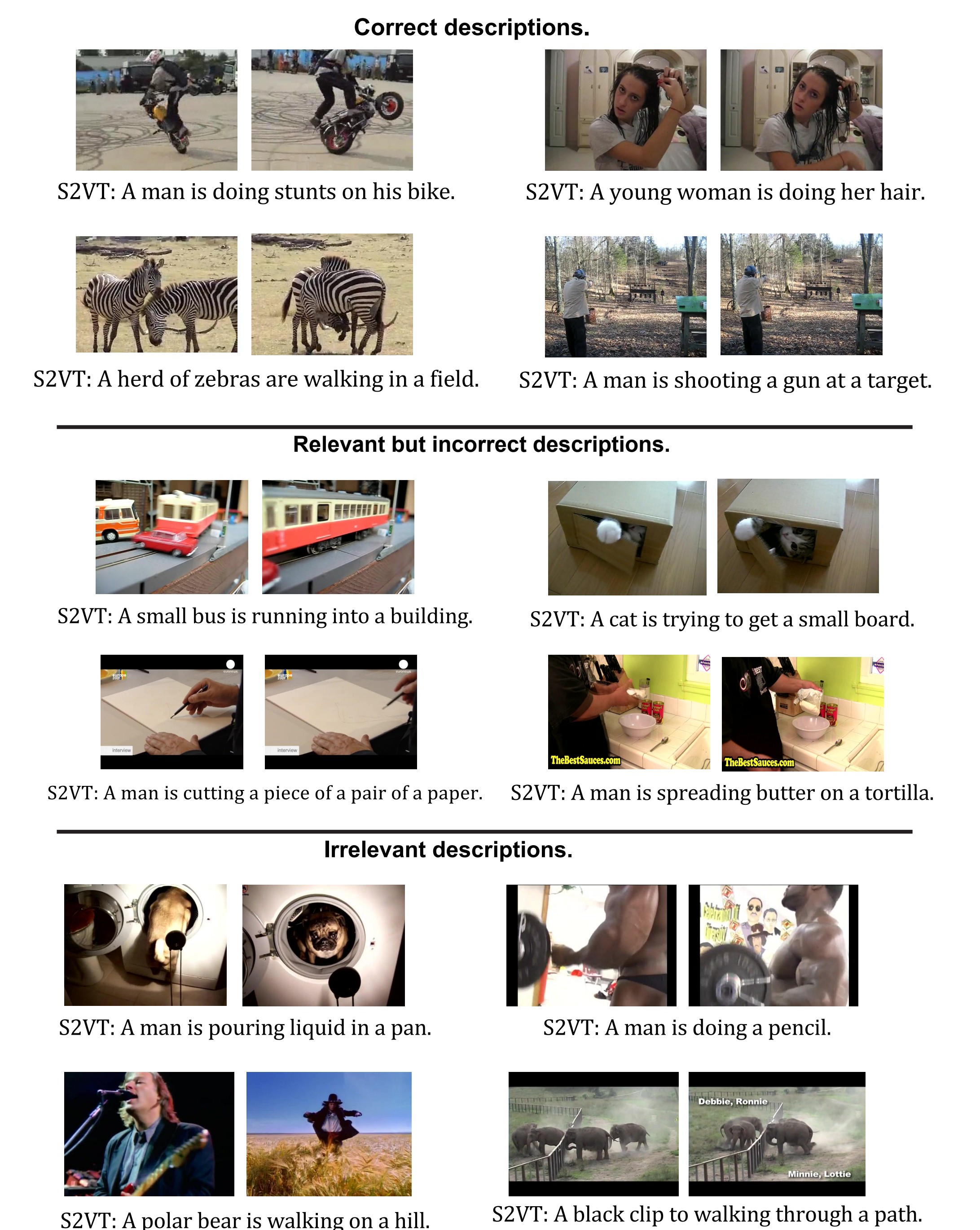
We use METEOR to compare the predicted sentence against the ground truth DVS description of the clip.

QUALITATIVE RESULTS - MOVIE DESCRIPTION

Representative frames of 6 contiguous clips from the movie "Big Mommas: Like Father, Like Son". Below are descriptions generated by prior art (Temporal Attention), and our model (S2VT), as well as the groundtruth (DVS) sentences.



QUALITATIVE RESULTS - YOUTUBE



LINKS

Project Page:
<http://vsubhashini.github.io/s2vt.html>
 Code: <https://github.com/vsubhashini/caffe/tree/recurrent/examples/s2vt>