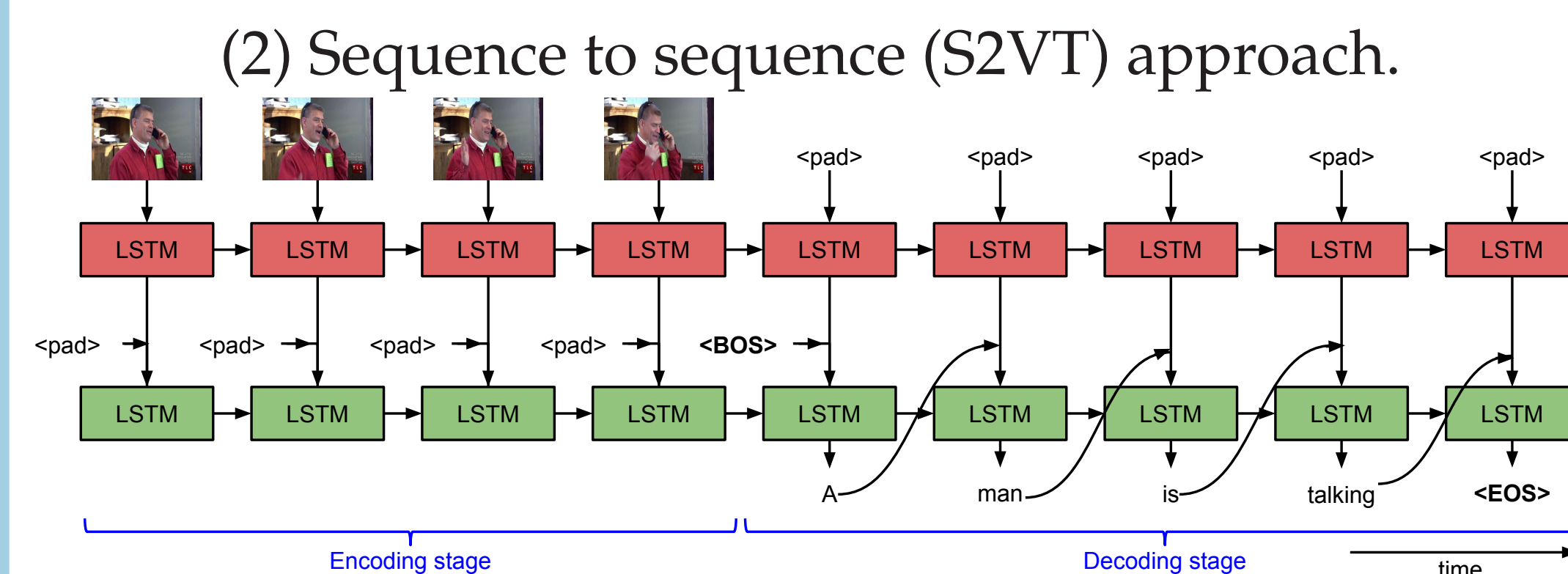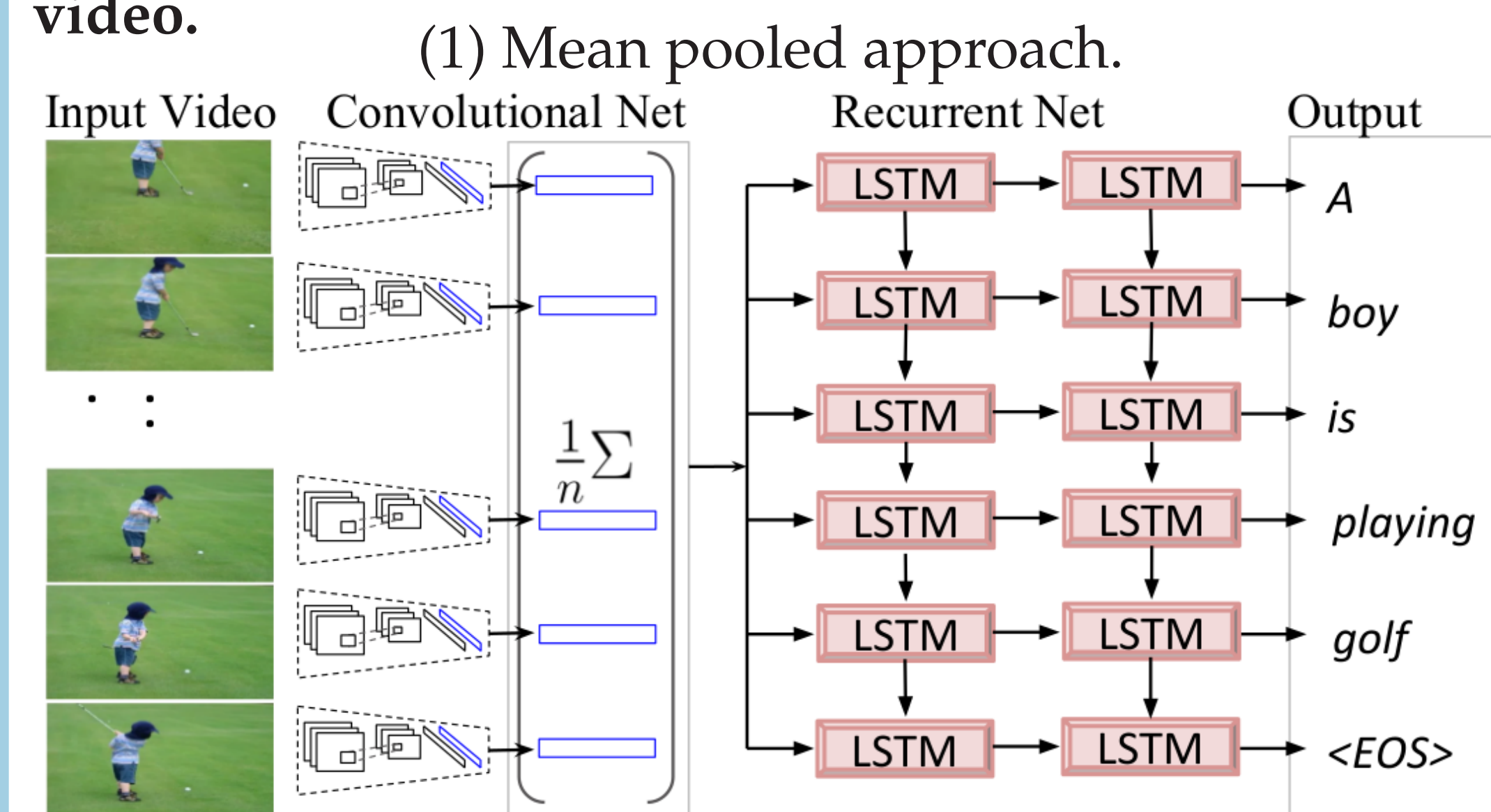# Translating Videos to Natural Language using Deep Recurrent Neural Networks

Subhashini Venugopalan[1], Huijuan Xu[3], Jeff Donahue[2],
Marcus Rohrbach[2], Raymond Mooney[1], Trevor Darrell[2], Kate Saenko[3]

[1] `UT-Austin` [2] `UC-Berkeley` [3] `UMass-Lowell`

## GOALS

**Given a short YouTube video, output a natural language sentence that describes the event depicted in the video.**

#### (1) Mean pooled approach.



Input Video   Convolutional Net   Recurrent Net   Output

#### (2) Sequence to sequence (S2VT) approach.



Encoding stage   Decoding stage   time

We present methods to generate descriptions for events depicted in videos using Convolutional Neural Networks (CNNs) and Long Short Term Memory (LSTM) networks.

## DATASETS

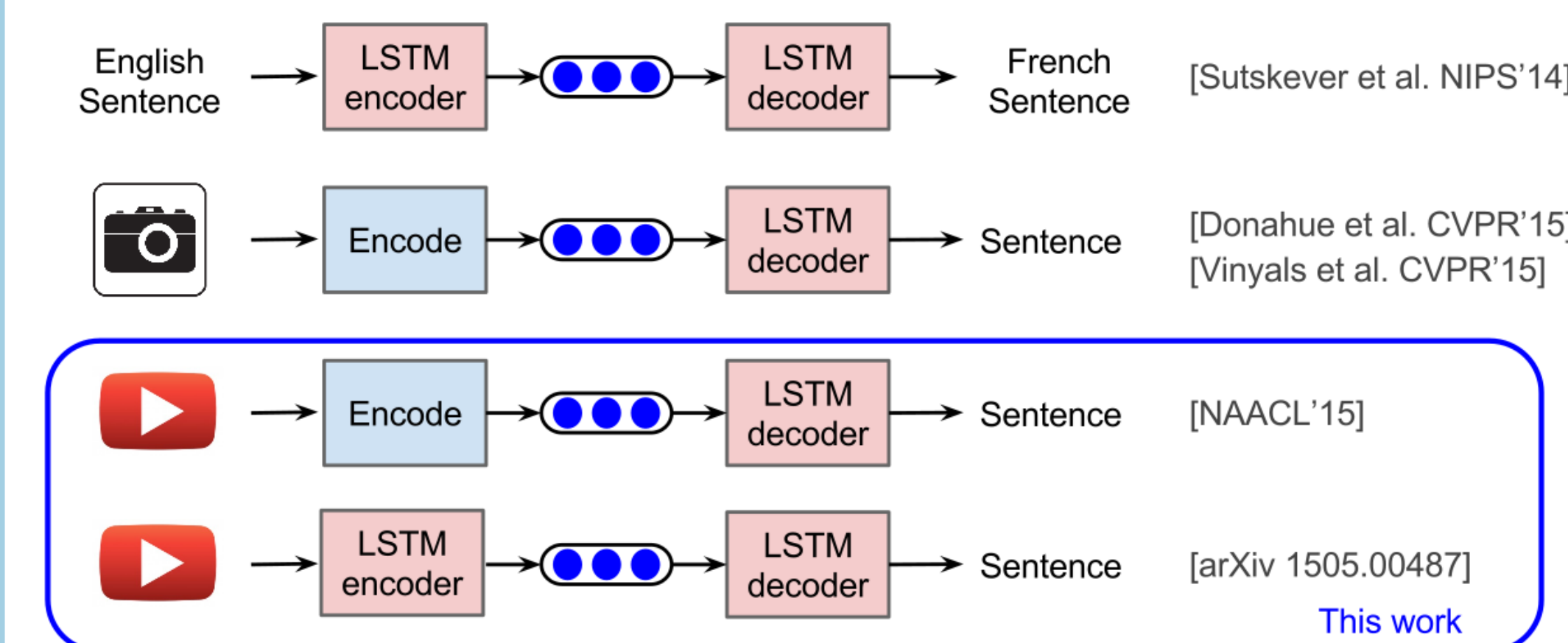**We demonstrate our approach on a large, realistic collection of YouTube videos and movies.**



A woman is cooking onions.
Someone is cooking in a pan.
someone preparing something
a person coking.
recipe for katsu curry

A man is sitting and playing a guitar
A man is playing guitar
Street artists play guitar.
A man is playing a guitar.
a lady is playing the guitar.

A girl is ballet dancing.
A girl is dancing on a stage.
A girl is performing as a ballerina.
A woman dances.

A train is rolling by.
A train passes by Mount Fuji.
A bullet train zooms through the countryside.
A train is coming down the tracks.

**(a) YouTube Video corpus**



**DVS:** Abby gets in the basket.

**Script:** After a moment a frazzled Abby pops up in his place.

Mike leans over and sees how high they are.

Mike looks down to see – they are now fifteen feet above the ground.

Abby clasps her hands around his face and kisses him passionately. For the first time in her life, she stops thinking and grabs Mike and kisses the hell out of him.
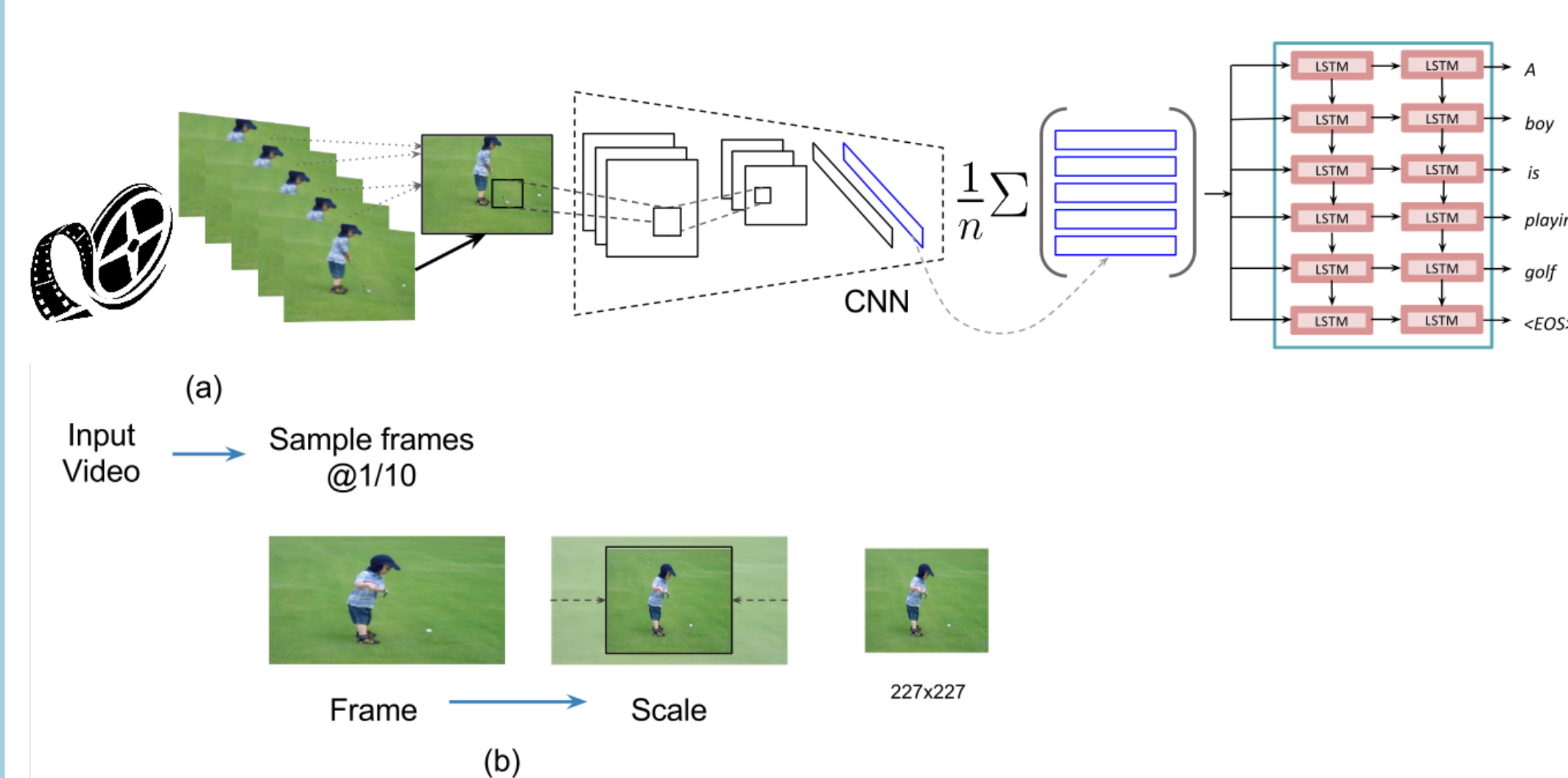
**(b) MPII Movie Description Dataset**

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 large movie description corpora like the Montreal and MPII movie description datasets, see (b) above.

## INSIGHT



English Sentence → LSTM encoder → LSTM decoder → French Sentence [Sutskever et al. NIPS'14]

→ Encode → LSTM decoder → Sentence [Donahue et al. CVPR'15] [Vinyals et al. CVPR'15]

→ Encode → LSTM decoder → Sentence [NAACL'15]

→ LSTM encoder → LSTM decoder → Sentence [arXiv 1505.00487]

This work

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



CNN

Input Video   (a)   Sample frames @1/10

Frame → Scale → 227x227
(b)

The image is forward propagated through a convolutional neural network. The activations of the fully connected layer just before the classification is considered as the image feature which is then mean-pooled (as in 1) or is directly provided as input to the LSTM network in the sequence to sequence (S2VT) models (in 2).



CNN

fc7: 4096 dimension "feature vector"

Forward propagate
Output: "fc7" features
(activations before classification layer)



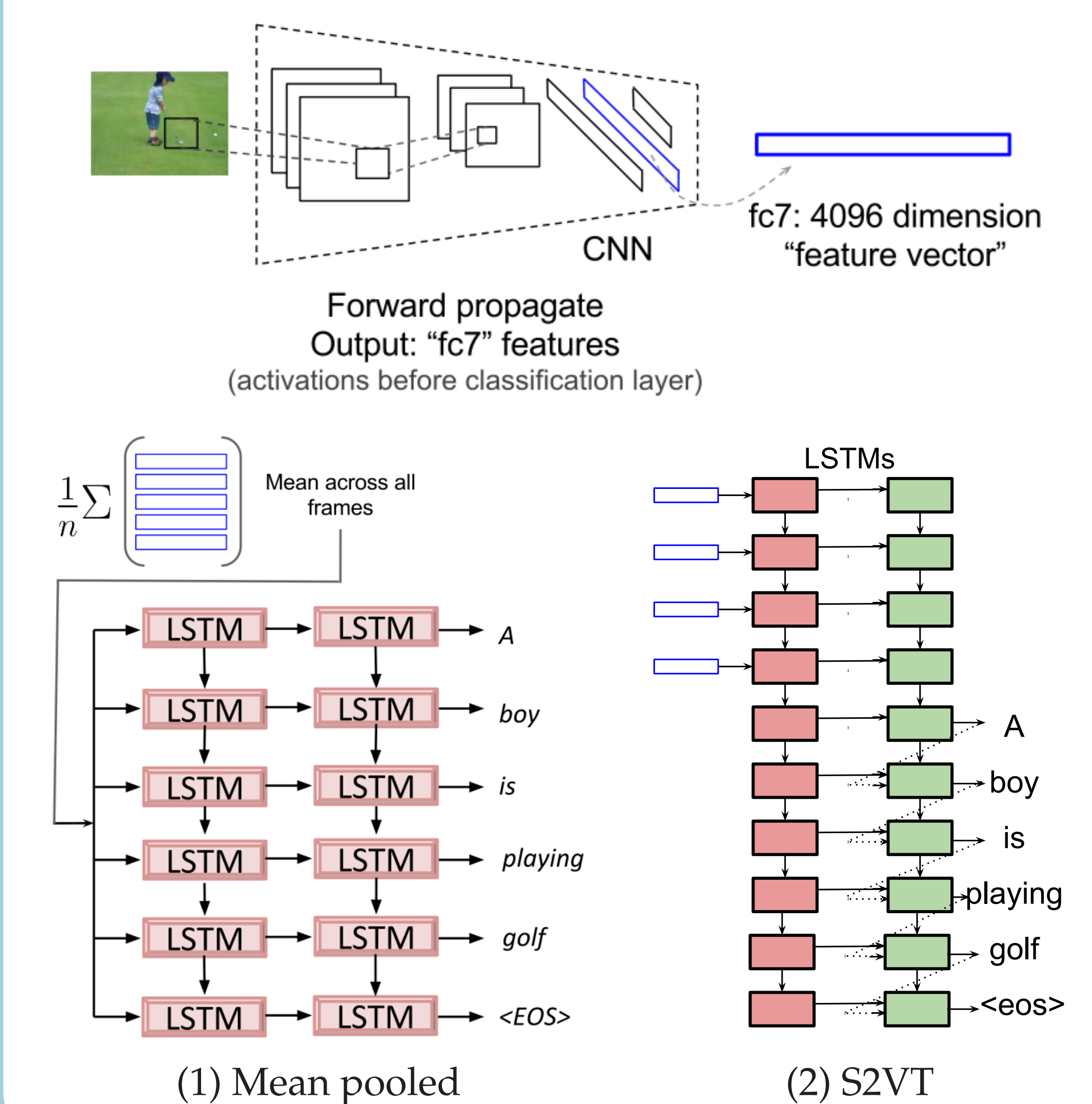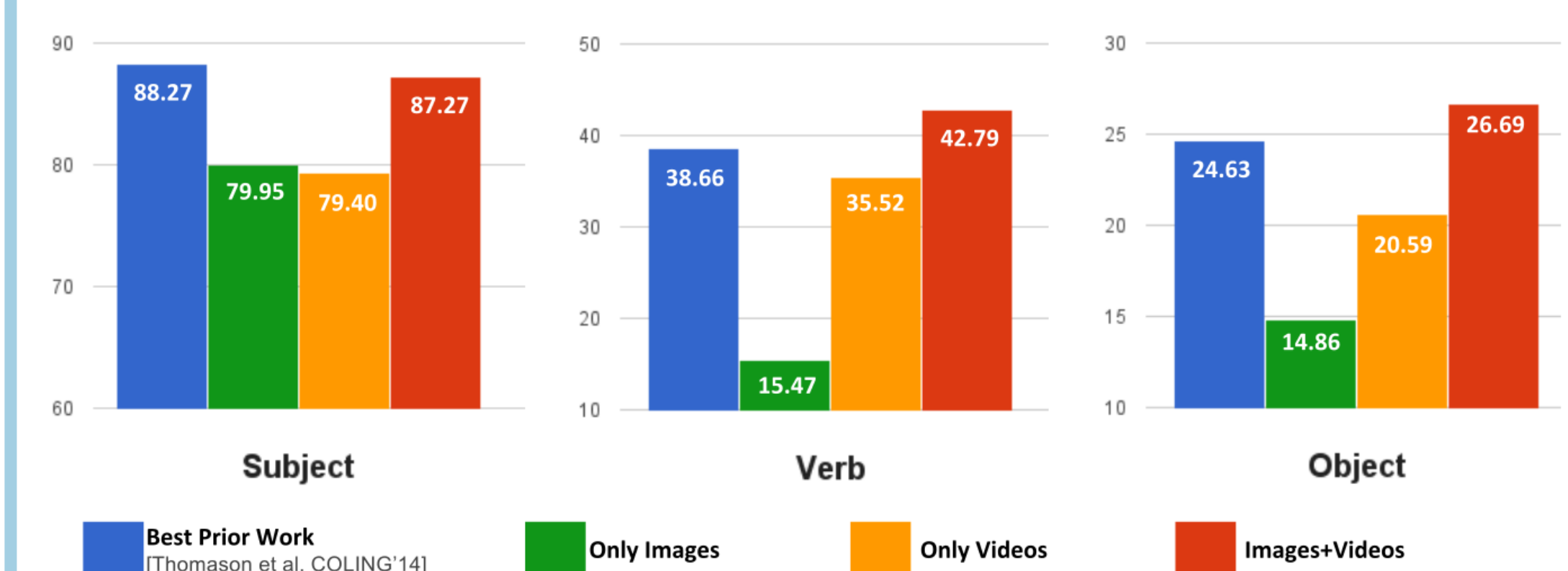Mean across all frames

LSTMs

(1) Mean pooled   (2) S2VT

## IMAGE PRE-TRAINING

Annotated video data is scarce. We can optionally initialize the weights of our network by pre-training on image-caption datasets. Flickr30k and MSCOCO datasets together have over 150,000 images and 750,000 caption (5 descriptions per image).



# Training videos - 1300

Flickr30k - 30,000 images, 150,000 descriptions

MSCOCO - 120,000 images, 600,000 descriptions

## SVO ACCURACY

**Accuracy of Subject, Verb and Object of the Mean-Pooled models.**



| | Subject | Verb | Object |
|---|---|---|---|
| Best Prior Work [Thomason et al. COLING'14] | 88.27 | 38.66 | 24.63 |
| Only Images | 79.95 | 15.47 | 14.86 |
| Only Videos | 79.40 | 35.52 | 20.59 |
| Images+Videos | 87.27 | 42.79 | 26.69 |

The subject, verb and object are extracted from the generated sentence and compared against all valid subjects, verbs, and objects amongst the ground truth descriptions.

## SENTENCE EVALUATION

We use the machine translation metric METEOR to compare the quality of the generated description against the multiple ground truth reference sentences.
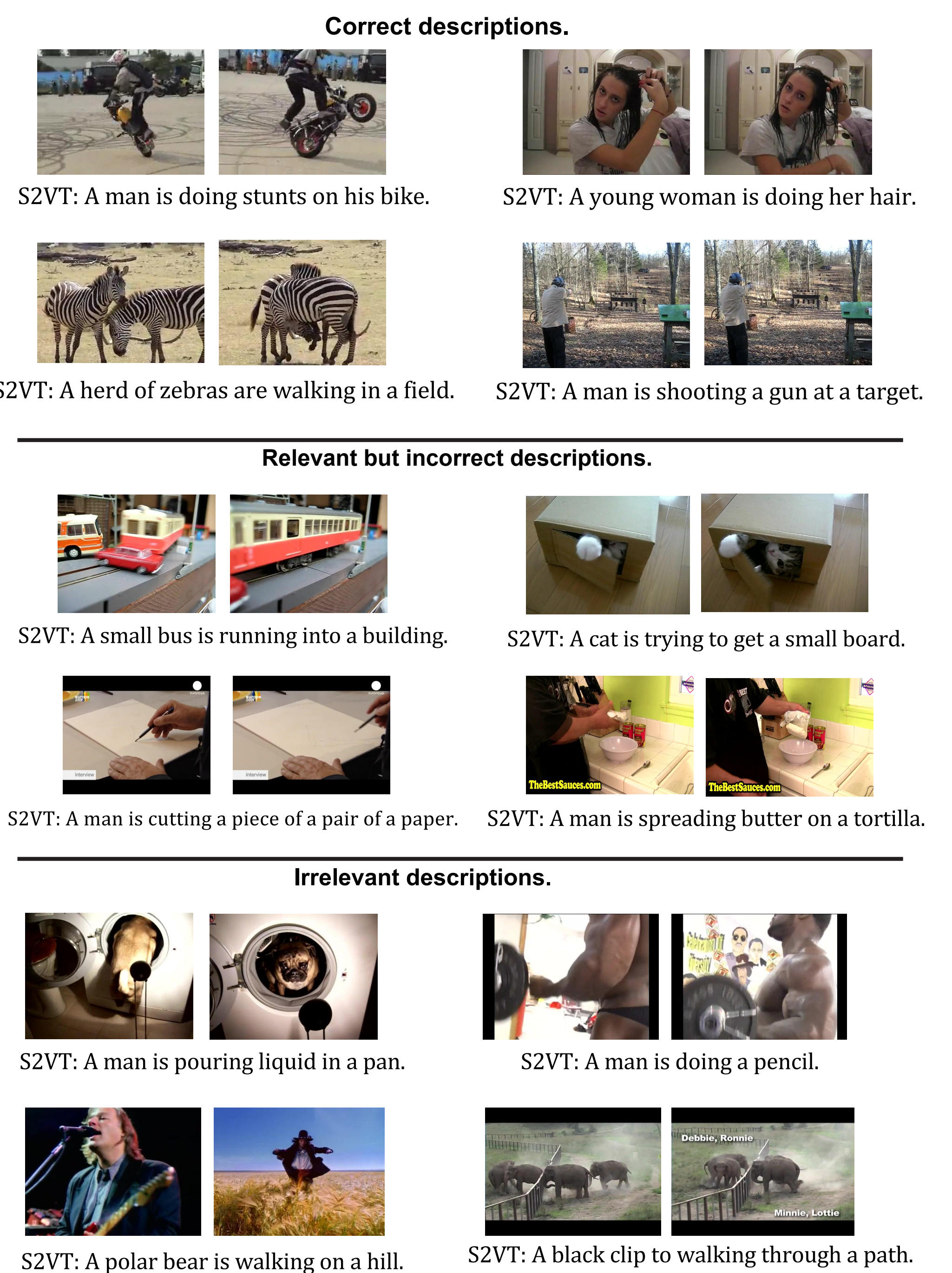
| Model | METEOR |
|---|---|
| FGM (best prior work) [1] | 23.9 |
| Mean pool | |
| - AlexNet [3] | 26.9 |
| - VGG | 27.7 |
| - AlexNet COCO pre-trained [3] | 29.1 |
| - GNet [4] | 28.7 |
| Soft-attention | |
| - GoogleNet [4] | 29.0 |
| - GoogleNet + 3D-CNN [4] | 29.6 |
| S2VT [2] | |
| - RGB (AlexNet) | 27.9 |
| - Flow (AlexNet) | 24.3 |
| - RGB (VGG) | 29.2 |
| - RGB (VGG) + Flow (AlexNet) | 29.8 |

## MOVIE DATASET RESULTS

| On MPII Movie Corpus | METEOR |
|---|---|
| SMT (best variant) | 5.6 |
| S2VT: RGB (VGG) [ours] | 6.3 |

| On Montreal M-VAD Corpus | METEOR |
|---|---|
| Soft-attention (GNet + 3D-CNN) [4]* | 4.1 |
| S2VT: RGB (VGG) [ours] | |
| - trained on M-VAD | 5.6 |
| - trained on MPII-MD & M-VAD | 6.7 |

## QUALITATIVE RESULTS

**Correct descriptions.**



S2VT: A man is doing stunts on his bike.

S2VT: A young woman is doing her hair.

S2VT: A herd of zebras are walking in a field.

S2VT: A man is shooting a gun at a target.

**Relevant but incorrect descriptions.**



S2VT: A small bus is running into a building.

S2VT: A cat is trying to get a small board.

S2VT: A man is cutting a piece of a pair of a paper.

S2VT: A man is spreading butter on a tortilla.

**Irrelevant descriptions.**



S2VT: A man is pouring liquid in a pan.

S2VT: A man is doing a pencil.

S2VT: A polar bear is walking on a hill.

S2VT: A black clip to walking through a path.

## REFERENCES

[1] J. Thomason, S. Venugopalan, S. Guadarrama, K. Saenko, and R. J. Mooney. Integrating language and vision to generate natural language descriptions of videos in the wild. In *COLING*, 2014.

[2] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, K. Saenko, and T. Darrell. Sequence to sequence – video to text. *arXiv:1505.00487*, 2015.

[3] S. Venugopalan, H. Xu, J. Donahue, M. Rohrbach, R. Mooney, and K. Saenko. Translating videos to natural language using deep recurrent neural networks. *NAACL-HLT*, 2015.

[4] L. Yao, A. Torabi, K. Cho, N. Ballas, C. Pal, H. Larochelle, and A. Courville. Describing videos by exploiting temporal structure. *arXiv:1502.08029v4*, 2015.