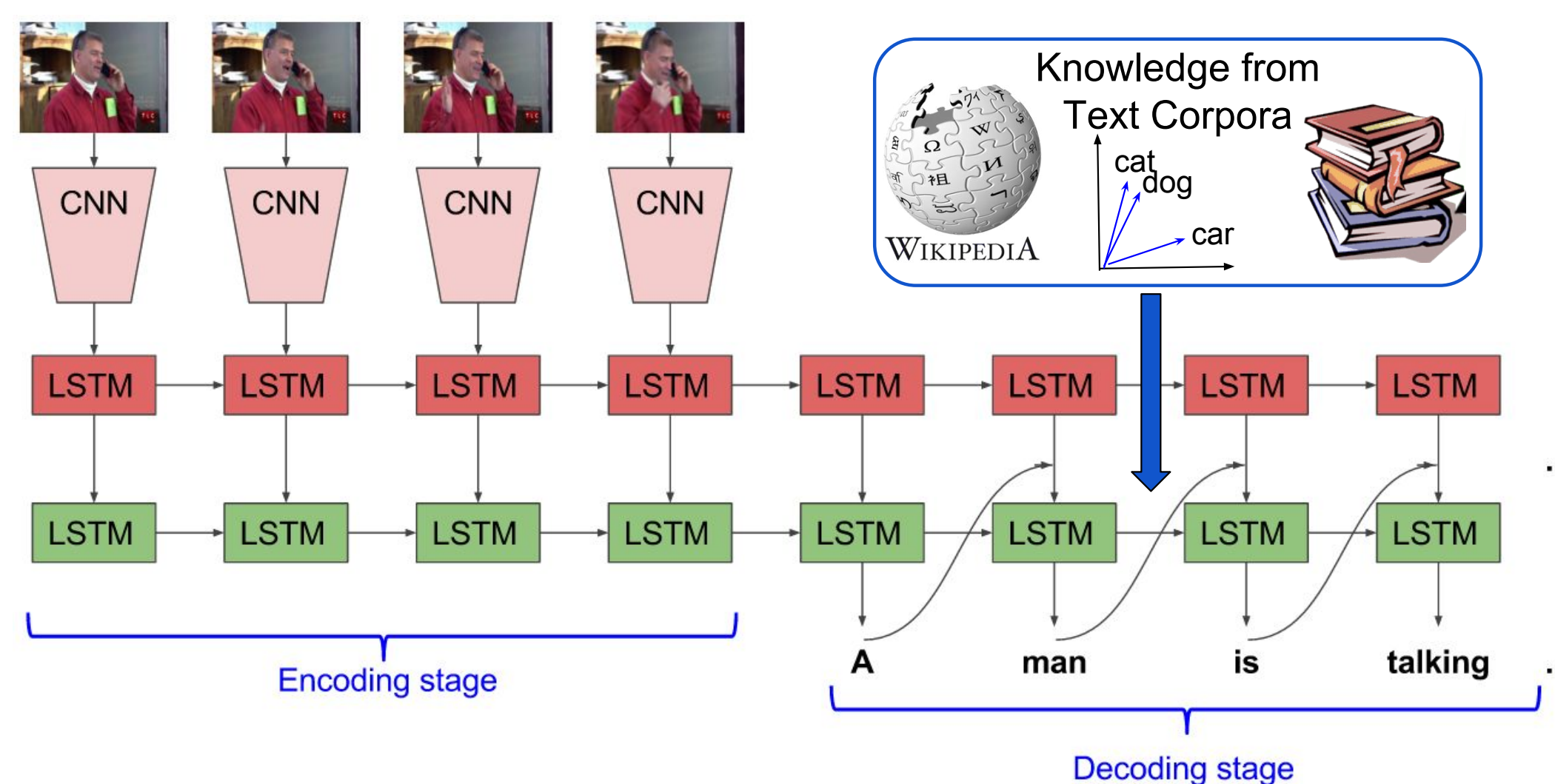


# Improving LSTM-based Video Description with Linguistic Knowledge Mined from Text

Subhashini Venugopalan<sup>1</sup>, Lisa Anne Hendricks<sup>2</sup>, Raymond Mooney<sup>1</sup>, Kate Saenko<sup>3</sup>  
<sup>1</sup> UT-Austin <sup>2</sup> UC-Berkeley <sup>3</sup> UMass-Lowell

## GOALS

Incorporate linguistic knowledge mined from external text to improve Video Description. [1]



We propose to add linguistic knowledge from external text corpora to enhance the quality of LSTM-based video description networks.

## DATASETS

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



(a) YouTube Video corpus



**DVS:** Abby gets in the basket.  
**Script:** After a moment a frazzled Abby pops up in his place.  
 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.

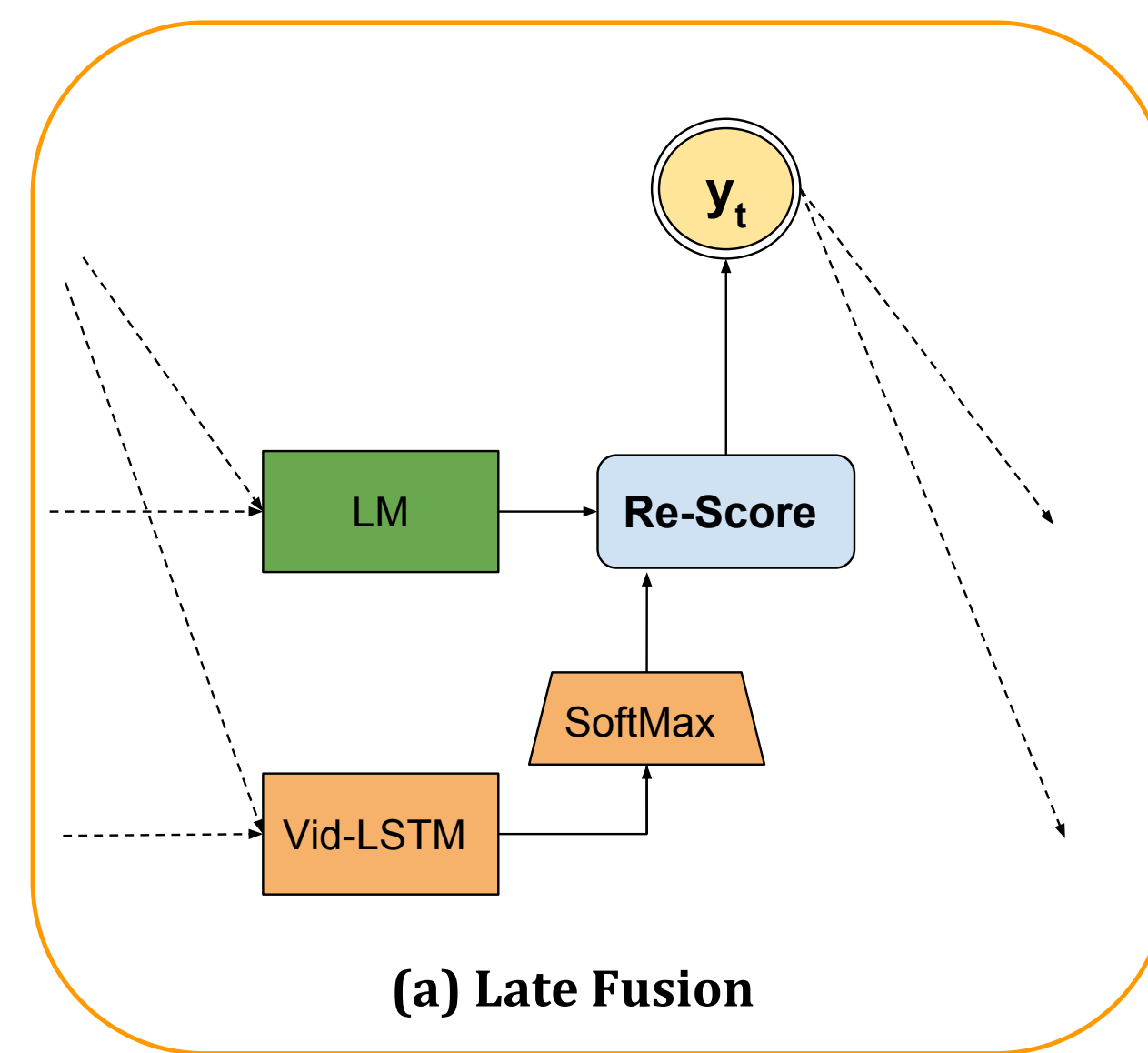
## REFERENCES

- [1] S. Venugopalan, L. A. Hendricks, R. Mooney, and K. Saenko. Improving lstm-based video description with linguistic knowledge mined from text. *arXiv:1604.01729, EMNLP*, 2016.
- [2] S. Venugopalan, L. A. Hendricks, M. Rohrbach, R. Mooney, T. Darrell, and K. Saenko. Captioning images with diverse objects. *arXiv:1606.07770*, 2016.

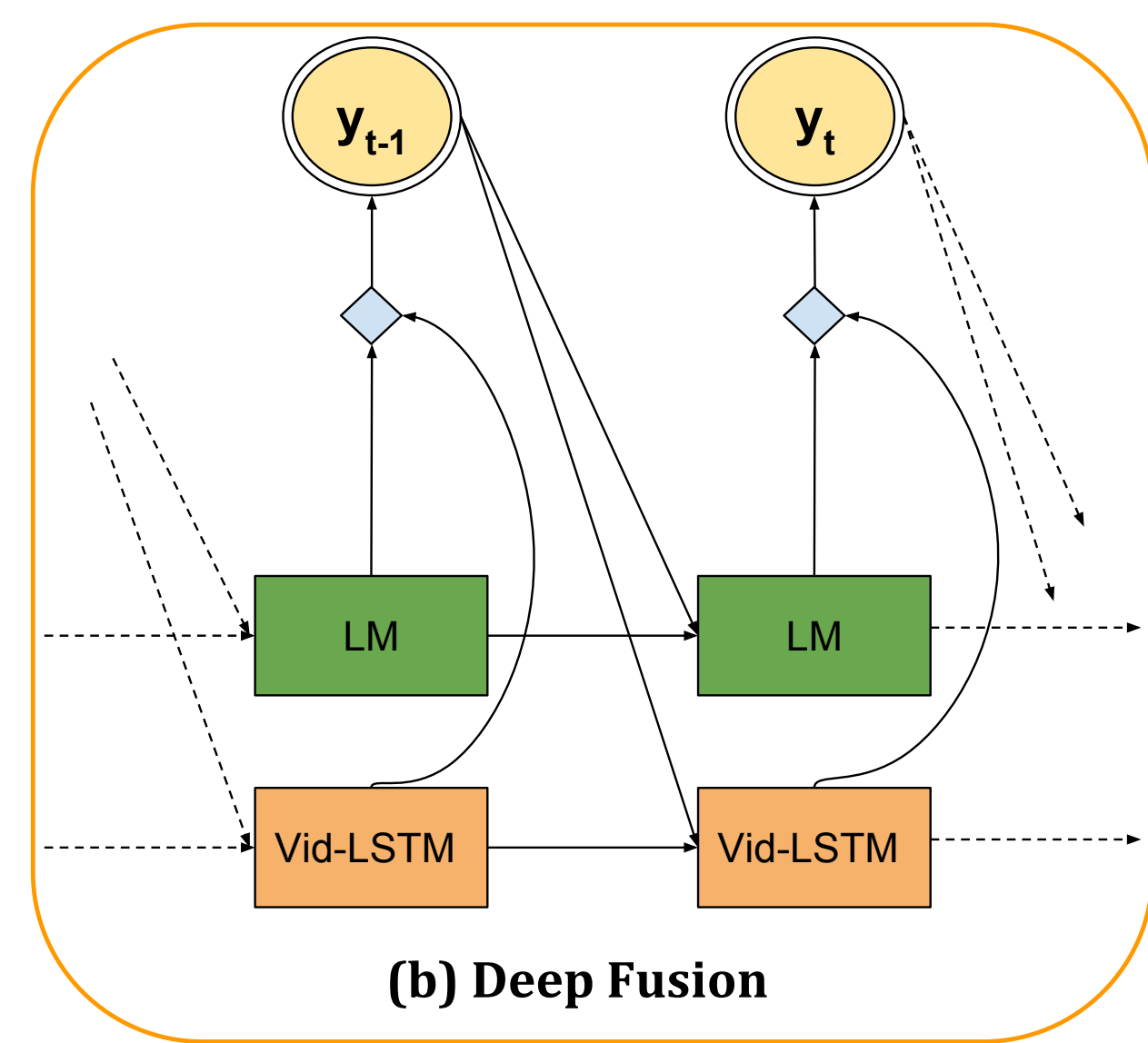
## TECHNIQUES

We propose multiple approaches to integrate language models trained on unannotated text corpora with existing video-captioning systems.

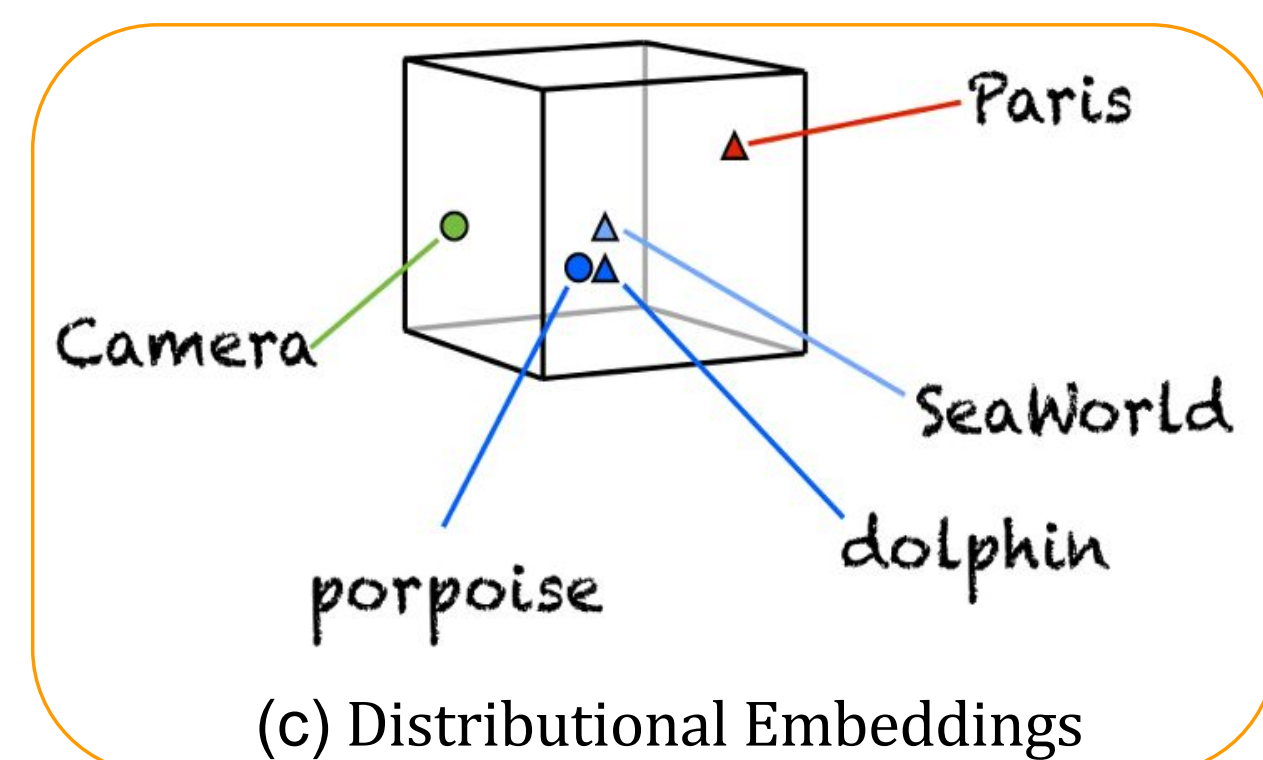
Our **early fusion** approach, simply initializes the weights of the video-captioning LSTM with an LSTM language model trained on large external text corpora.



Our **late fusion** approach combines the language model and the video captioning model by re-scoring the softmax output of the video-to-text model.



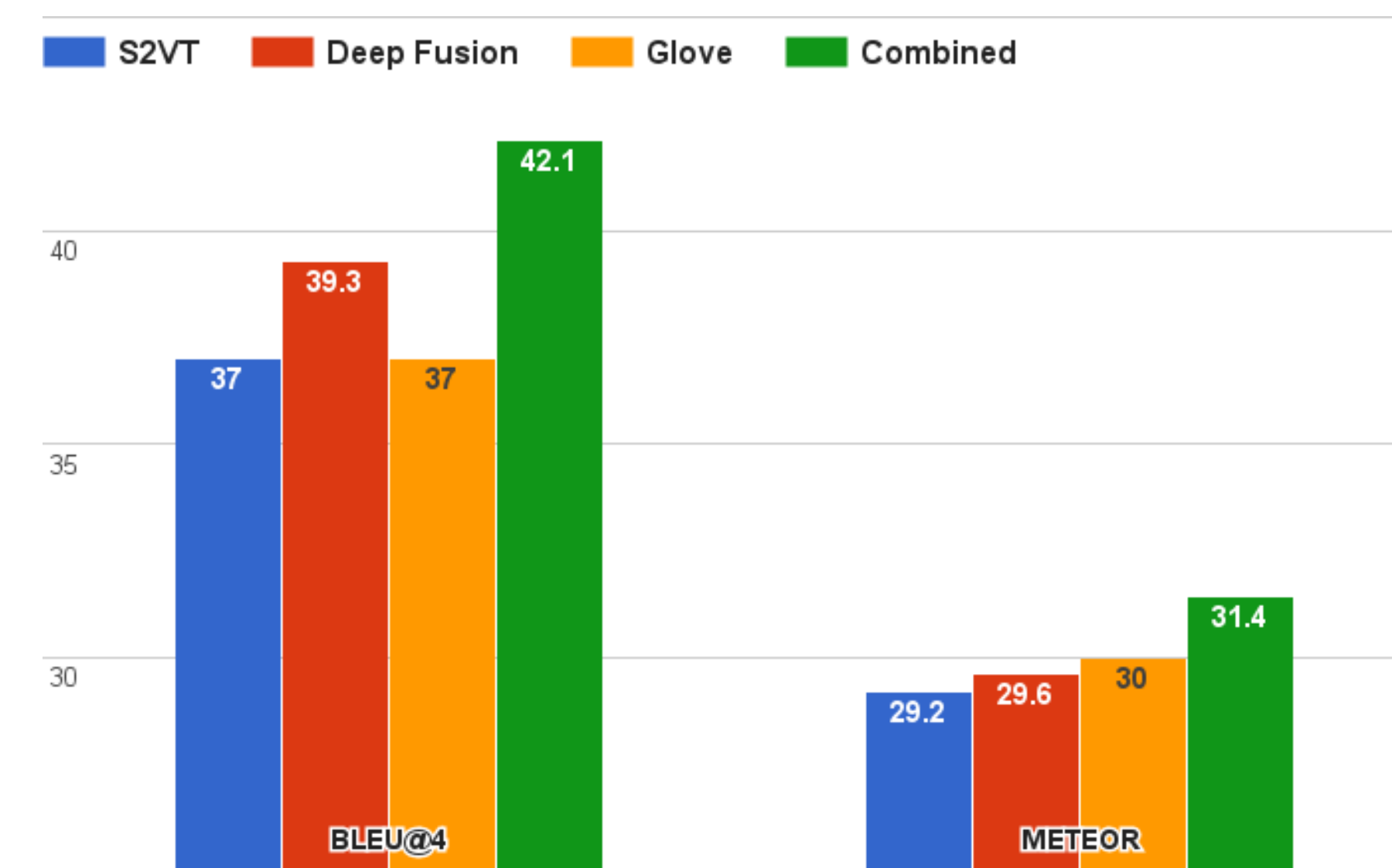
Our **deep fusion** model learns jointly from the hidden representations of the Language Model and S2VT video-to-text model (Vid-LSTM). So, the video description network is trained with the knowledge of the LSTM LM.



Additionally, instead of the typical “one-hot” representations for words, we use distributional word embeddings. Embeddings such as word2vec and Glove are dense vector-space representations of words that capture fine-grained semantic and syntactic regularities

## EVALUATIONS

We evaluate our models on automatic machine translation metrics - METEOR and BLEU. We also obtain human judgements on ‘Relevance’ and ‘Grammar’.



Evaluations using automated metrics on the Youtube video dataset.

## RESULTS - YOUTUBE

Model	METEOR	B-4	Relevance	Grammar
S2VT	29.2	37.0	2.06	3.76
Early Fusion	29.6	37.6	-	-
Late Fusion	29.4	37.2	-	-
Deep Fusion	29.6	39.3	-	-
Glove	30.0	37.0	-	-
Glove+Deep	30.3	38.1	2.12	4.05*
- Web Corpus	30.3	38.8	2.21*	4.17*
Ensemble	<b>31.4</b>	<b>42.1</b>	<b>2.24*</b>	<b>4.20*</b>
Groundtruth			4.52	4.47

Youtube dataset: METEOR and BLEU@4 in %, and human ratings (1-5) on relevance and grammar. Best results in bold, \* indicates significant over S2VT.

## RESULTS - MOVIES

Model	MPII-MD		M-VAD	
	METEOR	Grammar	METEOR	Grammar
S2VT <sup>†</sup>	6.5	2.6	6.6	2.2
Early Fusion	6.7	-	6.8	-
Late Fusion	6.5	-	6.7	-
Deep Fusion	6.8	-	6.8	-
Glove	6.7	3.9*	6.7	3.1*
Glove+Deep	6.8	<b>4.1*</b>	6.7	<b>3.3*</b>

Movie dataset: METEOR in %, and human ratings (1-5) on grammar (higher is better). Best results in bold, \* indicates significant over S2VT.

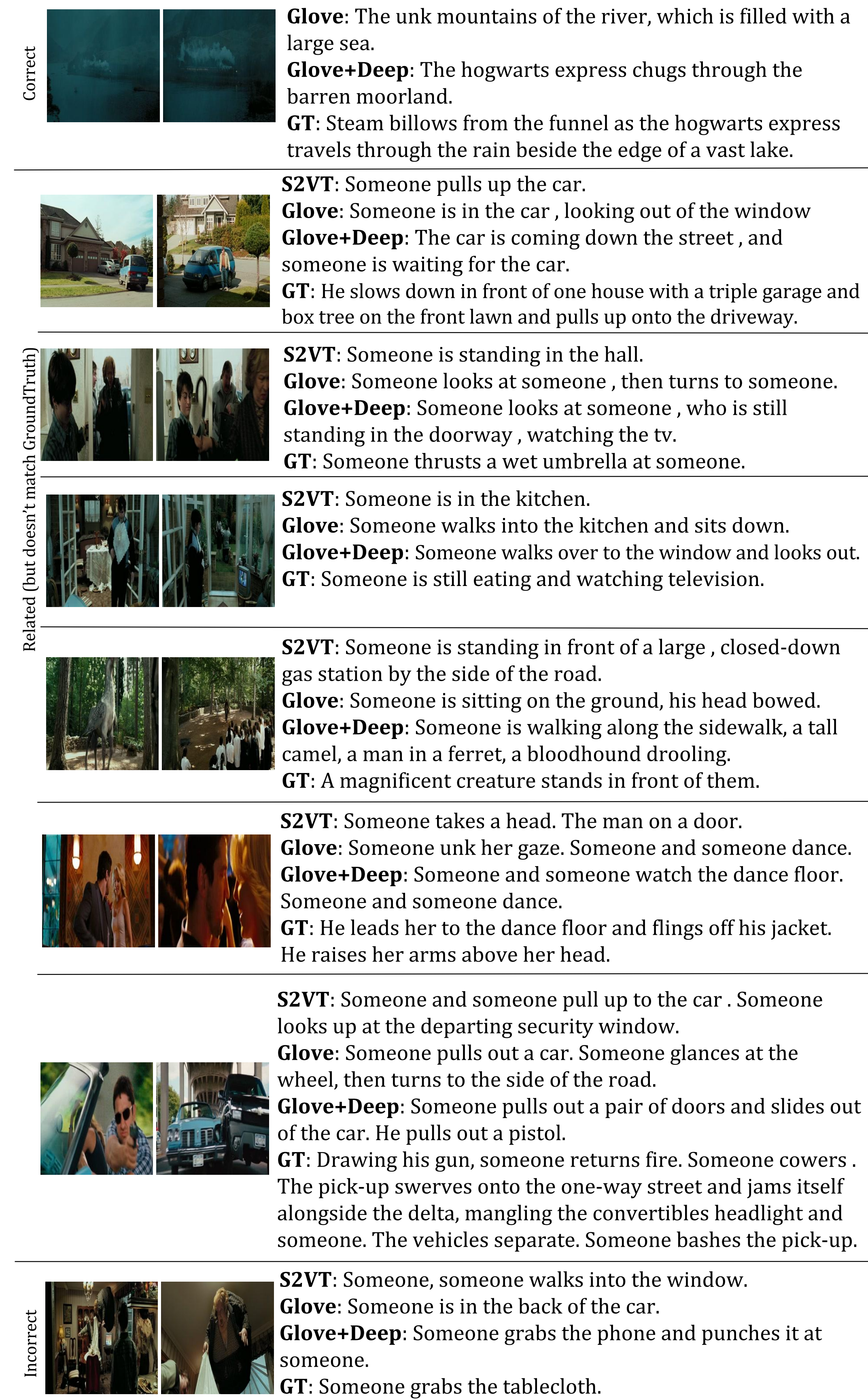
## CODE

Project Page with Code and Examples

[http://vsubhashini.github.io/language\\_fusion.html](http://vsubhashini.github.io/language_fusion.html)

## EXAMPLES

Results on the Movie Description Corpora.



## FUTURE DIRECTIONS

We also show that knowledge from external sources can be used to generate captions for obejts without paired image-caption data. [2]

