



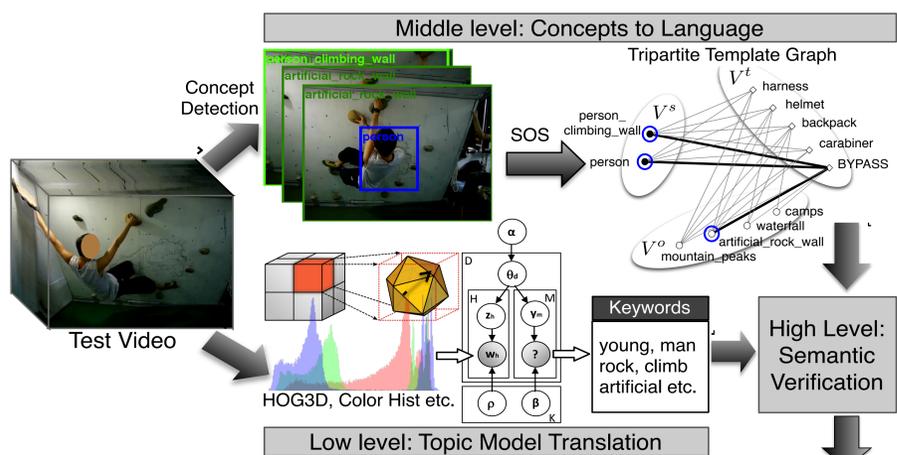
A Thousand Frames in Just a Few Words: Linguistic Description of Videos through Latent Topics and Sparse Object Stitching

Pradipto Das*, Chenliang Xu*, Richard F. Doell and Jason J. Corso
 Department of Computer Science and Engineering - SUNY at Buffalo, Buffalo, NY
 (* denotes equal contribution)

Objective:

- Generate natural language descriptions of a video that incorporate fine-grained information extraction.
- Improve relevance of descriptions in a hybrid framework that leverages bottom-up keyword prediction semantically verified by top-down concept detection and tri-partite template graphs.

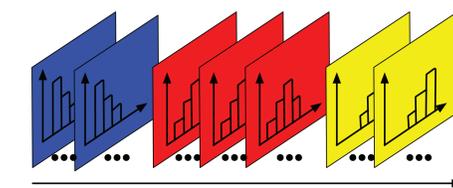
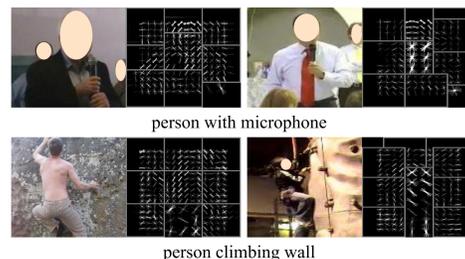
System Overview:



Output from our system: 1. A person in on artificial rock wall 2. A person climbing a wall is on artificial rock wall 3. Person climbs rock wall indoors 4. Young man tries to climb artificial rock wall 5. A man demonstrates how to climb a rock wall

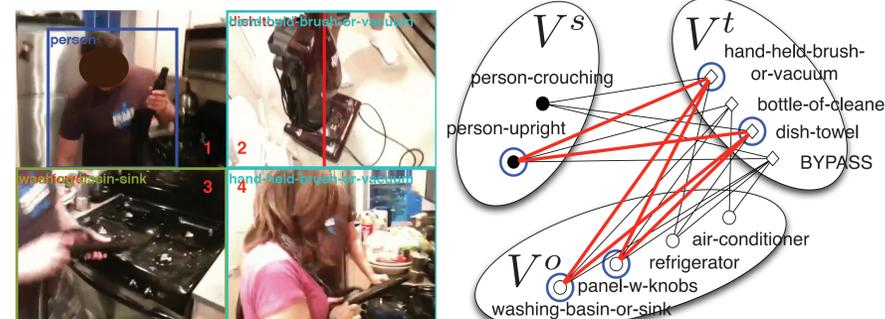
Middle Level: Concepts to Language

- Concept Detectors
 - rich semantics from object, action and scene level
 - reduce visual complexity
 - similar to "visual phrases"
 - deformable parts models are the base detector



- Sparse Object Stitching
 - segment video into a set of concept shots
 - Record distribution of detected concepts per shot
 - avoids need to do expensive dense detection and tracking

- Tripartite Template Graph for sentence generation.



Template Language Output:
 ([a person]) is using ([dish towel] and [hand held brush or vacuum]) to clean ([panel with knobs] and [washing basin or sink])

TRECVID MER12 Dataset:

- Videos are in 5 Events
- Annotation: Human descriptions/synopses & concept bounding boxes
- Train:
 - Concept Detector: 200 Videos/Event
 - Topic Model: 120 Videos/Event (only positive instances)
- Test:
 - 6 Videos/Event (MER12 Test Set for Recounting)

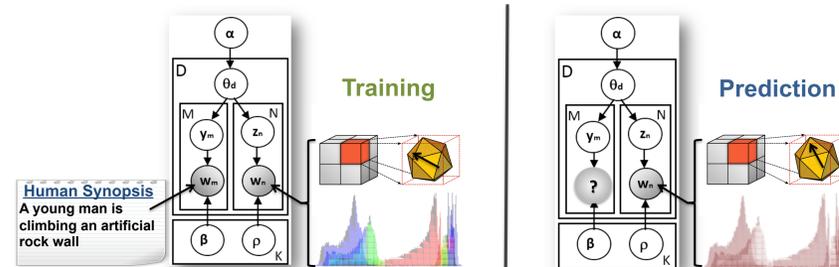
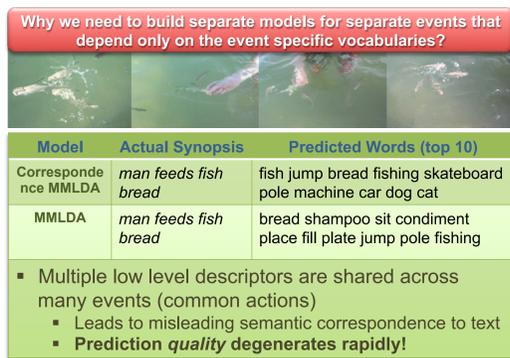
Event	Keywords	Sentences from Our System	Human Synopsis
Cleaning an appliance	refrigerator/OBJ cleans/VERB man/SUBJ-HUMAN clean/VERB blender/OBJ cleaning/VERB woman/SUBJ-HUMAN person/SUBJ-HUMAN stove/OBJ microwave/OBJ sponge/NOUN food/OBJ home/OBJ hose/OBJ oven/OBJ	1. A person is using dish towel and hand held brush or vacuum to clean panel with knobs and washing basin or sink 2. Man cleaning a refrigerator. 3. Man cleans his blender. 4. Woman cleans old food out of refrigerator. 5. Man cleans top of microwave with sponge.	Two standing persons clean a stove top with a vacuum clean with a hose.
Town hall meeting	meeting/VERB town/NOUN hall/OBJ microphone/OBJ talking/VERB people/OBJ podium/OBJ speech/OBJ woman/SUBJ-HUMAN man/SUBJ-HUMAN chairs/NOUN clapping/VERB speaker/VERB questions/VERB giving/VERB	1. A person is speaking to a small group of sitting people and a small group of standing people with board in the back 2. A person is speaking to a small group of standing people with board in the back 3. Man opens town hall meeting. 4. Woman speaks at town meeting. 5. Man gives speech on health care reform at a town hall meeting.	A man talks to a mob of sitting persons who clap at the end of his short speech.
Renovating home	people/SUBJ-HUMAN home/OBJ group/OBJ renovating/VERB working/VERB montage/OBJ stop/VERB motion/OBJ appears/VERB building/VERB floor/OBJ tiles/OBJ floorboards/OTHER man/SUBJ-HUMAN laying/VERB	1. A person is using power drill to renovate a house. 2. A crutching underlay is using power drill to renovate a house. 3. A person is using trowel to renovate a house. 4. man lays out perlay for installing flooring. 5. A man lays a plywood floor in time lapsed video.	Time lapse video of people making a concrete porch with sanders, brooms, vacuums and other tools.
Metal crafts project	metal/OBJ man/SUBJ-HUMAN bending/VERB hammer/VERB piece/OBJ tools/OBJ rods/OBJ hammering/VERB craft/VERB iron/OBJ workshop/OBJ holding/VERB works/VERB steel/OBJ bicycle/OBJ	1. A person is working with pliers. 2. Man hammering metal. 3. Man bending metal in workshop. 4. Man works various pieces of metal. 5. A man works on a metal craft at a workshop.	A man is shaping a star with a hammer.

Automatic Evaluation through ROUGE-1 Metric

Events	Precision				Recall			
	BASELINE	OURS (Low, Middle and High)			BASELINE	OURS (Low, Middle and High)		
Cleaning appliance	20.03	17.52	11.69	10.68	19.16	32.6	35.76	48.15
Renovating home	6.66	15.29	12.55	9.99	7.31	43.41	30.67	49.52
Rock climbing	24.45	16.21	24.52	12.61	44.09	59.22	46.23	65.84
Town hall meeting	17.35	14.41	27.56	13.36	13.8	28.66	45.55	56.44
Metal crafts project	16.73	18.12	31.68	15.63	19.01	41.87	25.87	54.84

Low Level: Topic Model Translation

- Translation model – capture semantic correlations from low level feature codebooks to bag-of-words
- Two families of multimodal (MM) topic models (LDA):
 - Correspondence MMLDA enforces a stronger constraint: Topic sparsity on low level features enforces stronger correspondence to text (Computationally expensive)
 - MMLDA operates more diffusely and focuses on semantic summarization of multiple views based on observation frequencies



High Level: Semantic Verification

- Rank nearest neighbor sentences from the training synopses by a ranking function

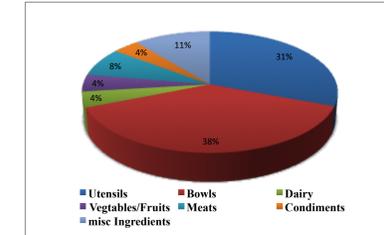
$$r_s = bh(w_1 x_{s_1} + w_2 x_{s_2})$$
 - b (boolean): at least two semantically verified semantically verified
 - h (boolean): at least one human subject
 - x_{s_1} (real): ratio of the total # of matches to the # of words in the sentence
 - x_{s_2} (real): sum of the weights of the predicted words from the topic model
 - Run MMLDA on a vocabulary of training synopses and training concept annotations
 - Semantic Verification: computing # of topic rank inversions for two ranked lists
- $$L_{synopsis} = \left\{ \left\{ k: \sum_{j=1}^V \sum_{m=1}^P p(w_m | \beta_k) \delta(w'_m) \right\} \right\}$$
- $$L_{concept} = \left\{ \left\{ k: \sum_{j=1}^{corrV} \sum_{n=1}^C p(w_n | \rho_k) \delta(w'_n) \right\} \right\}$$
- P : number of word predictions; C : number of positive concept detections

- A person is speaking to a large group of standing people and a small group of standing people with board in the back and a camera man [Sentence from SOS]
- Representative speaks to crowd at town hall meeting
- A man speaks at a town hall meeting about health care
- Man opens town hall meeting care
- People get angry at town hall meeting on health care
- Politician gives speech at town hall meeting



YouCook Dataset:

Cooking videos downloaded from YouTube
 Splits: Train 49 / Test 39
 Annotations:
 -- Human Descriptions from MTurk
 -- Object and Action Bounding Boxes
 ROUGE Benchmark Scoring



Cooking video	Precision Bi/Uni-gram				Recall Bi/Uni-gram			
	BASELINE	OURS			BASELINE	OURS		
Cooking video 1	0.0006	15.47	5.04	24.82	0.0006	19.02	6.81	34.2
Cooking video 2	0.0006	15.47	5.04	24.82	0.0006	19.02	6.81	34.2

Acknowledgements. This work was partially supported by the National Science Foundation CAREER grant (IIS-0845282), the Army Research Office (W911NF-11-1-0090), the DARPA Mind's Eye program (W911NF-10-2-0062), and the Intelligence Advanced Research Projects Activity (IARPA) via Department of Interior National Business Center contract number D11PC20069. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. Disclaimer: The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DOI/NBC, DARPA, ARO, NSF or the U.S. Government. We thank Simon Fraser Univ. & Kitware Inc. for support in feature extraction and the anonymous reviewers for their comments. We also thank the in-house annotators Philip Rosebrough, Yao Li, and Cody Boppert, without whom this project would not have been complete.