

000  
001  
002003 

# Seeing What You're Told: Sentence-Guided Activity Recognition In Video

004  
005  
006  
007  
008  
009  
010  
011012 Anonymous CVPR submission  
013014 Paper ID 1701  
015016 

## Abstract

017 We present a system that demonstrates how the compositional  
018 structure of events, in concert with the compositional  
019 structure of language, can interplay with the underlying  
020 focusing mechanisms in video action recognition, thereby  
021 providing a medium, not only for top-down and bottom-up  
022 integration, but also for multi-modal integration between  
023 vision and language. We show how the roles played by par-  
024 ticipants (nouns), their characteristics (adjectives), the ac-  
025 tions performed (verbs), the manner of such actions (ad-  
026 verbs), and changing spatial relations between participants  
027 (prepositions) in the form of whole sentential descriptions  
028 mediated by a grammar, guides the activity-recognition pro-  
029 cess. Further, the utility and expressiveness of our frame-  
030 work is demonstrated by performing three separate tasks  
031 in the domain of multi-activity videos: sentence-guided fo-  
032 cus of attention, generation of sentential descriptions of  
033 video, and query-based video search, simply by leveraging  
034 the framework in different manners.

035 

## 1. Introduction

036 The ability to describe the observed world in natural lan-  
037 guage is a quintessential component of human intelligence.  
038 A particular feature of this ability is the use of rich sen-  
039 tences, involving the composition of multiple nouns, ad-  
040 jjectives, verbs, adverbs, and prepositions, to describe not just  
041 static objects and scenes, but also events that unfold over  
042 time. Furthermore, this ability appears to be learned by vir-  
043 tually all children. The deep semantic information learned  
044 is multi-purpose: it supports comprehension, generation,  
045 and inference. In this work, we investigate the intuition,  
046 and the precise means and mechanisms that will enable us  
047 to support such ability in the domain of activity recognition  
048 in multi-activity videos.

049 Suppose we wanted to recognize an occurrence of an  
050 event described by the sentence *The ball bounced*, in a  
051 video. Nominally, we would need to detect the *ball* and  
052 its position in the field of view in each frame and determine  
053 that the sequence of such detections satisfied the require-  
054 ments of *bounce*. The sequence of such object detections  
055 and their corresponding positions over time constitutes a  
056

057 track for that object. In this view, the semantics of an in-  
058 transitive verb like *bounce* would be formulated as a unary  
059 predicate over object tracks. Recognizing occurrences of  
060 events described by sentences containing transitive verbs,  
061 like *The person approached the ball*, would require detect-  
062 ing and tracking two objects, the *person* and the *ball* con-  
063 strained by a binary predicate.

064 In an ideal world, event recognition would proceed in a  
065 purely feed-forward fashion: robust and unambiguous ob-  
066 ject detection and tracking followed by application of the  
067 semantic predicates on the recovered tracks. However, the  
068 current state-of-the-art in computer vision is far from this  
069 ideal. Object detection alone is highly unreliable. The best  
070 current average-precision scores on PASCAL VOC hover  
071 around 40%-50% [3]. As a result, object detectors suf-  
072 fer from both false positives and false negatives. One way  
073 around this is to use detection-based tracking [17], where  
074 one biases the detector to overgenerate, alleviating the prob-  
075 lem of false negatives, and uses a different mechanism to  
076 select among the overgenerated detections to alleviate the  
077 problem of false positives. One such mechanism selects de-  
078 tections that are temporally coherent, *i.e.* the track motion  
079 being consistent with optical flow. Barbu *et al.* [2] proposed  
080 an alternate mechanism that selected detections for a track  
081 that satisfied a unary predicate such as one would construct  
082 for an intransitive verb like *bounce*. We significantly ex-  
083 tend that approach, selecting detections for multiple tracks  
084 that collectively satisfy a complex multi-argument predicate  
085 representing the semantics of an entire sentence. That pred-  
086 icate is constructed as a conjunction of predicates represent-  
087 ing the semantics of individual words in that sentence. For  
088 example, given the sentence *The person to the left of the*  
089 *chair approached the trash can*, we construct a logical form.

$$\begin{aligned} & \text{PERSON}(P) \wedge \text{TOTHELEFTOF}(P, Q) \wedge \text{CHAIR}(Q) \\ & \wedge \text{APPROACH}(P, R) \wedge \text{TRASHCAN}(R) \end{aligned}$$

090 Our tracker is able to simultaneously construct three  
091 tracks *P*, *Q*, and *R*, selecting out detections for each, in an  
092 optimal fashion that simultaneously optimizes a joint mea-  
093 sure of detection score and temporal coherence while also  
094 satisfying the above conjunction of predicates. We obtain  
095 the aforementioned detections by employing a state-of-the-  
096

108 art object detector [5], where we train a model for each object (*e.g. person, chair, etc.*), which when applied to an image, produces axis-aligned bounding boxes with associated scores indicating strength of detection.  
 109  
 110  
 111  
 112

113 We represent the semantics of lexical items like *person*,  
 114 *to the left of*, *chair*, *approach*, and *trash can* with predicates  
 115 over tracks like  $\text{PERSON}(P)$ ,  $\text{TOTHELEFTOF}(P, Q)$ ,  
 116  $\text{CHAIR}(Q)$ ,  $\text{APPROACH}(P, R)$ , and  $\text{TRASHCAN}(R)$ . These  
 117 predicates are in turn represented as regular expressions (*i.e.*  
 118 finite state recognizers or FSMs) over features extracted  
 119 from the sequence of detection positions, shapes, and sizes  
 120 as well as their temporal derivatives. For example, the predi-  
 121 cate  $\text{TOTHELEFTOF}(P, Q)$  might be a single state FSM  
 122 where, on a frame-by-frame basis, the centers of the de-  
 123 tections for  $P$  are constrained to have a lower  $x$ -coordinate  
 124 than the centers of the detections for  $Q$ . The actual formula-  
 125 tion of the predicates (Table 2) is far more complex to deal  
 126 with noise and variance in real-world video. What is central  
 127 is that the semantics of *all* parts of speech, namely nouns,  
 128 adjectives, verbs, adverbs, and prepositions (both those that  
 129 describe spatial-relations and those that describe motion), is  
 130 uniformly represented by the same mechanism: predicates  
 131 over tracks formulated as finite state recognizers over fea-  
 132 tures extracted from the detections in those tracks.

133 We refer to this capacity as the *Sentence Tracker*, which  
 134 is a function  $\mathcal{S} : (D, \Phi) \mapsto (\tau, Z)$ , that takes as input an  
 135 overgenerated set  $D$  of detections along with a complex  
 136 sentential predicate  $\Phi$  and produces a score  $\tau$  together with  
 137 a set  $Z$  of tracks that satisfy  $\Phi$  while optimizing a linear  
 138 combination of detection scores and temporal coherence.  
 139 This can be used for three distinct purposes:

140 **focus of attention** One can apply the sentence tracker to  
 141 the same video  $D$ , that depicts multiple simultaneous  
 142 events taking place in the field of view with different  
 143 participants, with two different sentences  $\Phi_1$  and  $\Phi_2$ .  
 144 In other words, one can compute  $(\tau_1, Z_1) = \mathcal{S}(D, \Phi_1)$   
 145 and  $(\tau_2, Z_2) = \mathcal{S}(D, \Phi_2)$  to yield two different sets  
 146 of tracks  $Z_1$  and  $Z_2$  corresponding to the different sets  
 147 of participants in the different events described by  $\Phi_1$   
 148 and  $\Phi_2$ . We demonstrate this in section 4.1.

149 **generation** One can take a video  $D$  as input and systematically  
 150 search the space of all possible  $\Phi$  that correspond  
 151 to sentences that can be generated by a context-free gram-  
 152 mar and find that sentence that corresponds to the  $\Phi^*$  for  
 153 which  $(\tau^*, Z^*) = \mathcal{S}(D, \Phi^*)$  yields the maximal  $\tau^*$ . This  
 154 can be used to generate a sentence that describes an input  
 155 video  $D$ . We demonstrate this in section 4.2.

156 **retrieval** One can take a collection  $\mathcal{D} = \{D_1, \dots, D_n\}$  of  
 157 videos (or a single long video chopped into short clips)  
 158 along with a sentential query  $\Phi$ , compute  $(\tau_i, Z_i) =$   
 159  $\mathcal{S}(D_i, \Phi)$  for each  $D_i$ , and find the clip  $D_i$  with maxi-  
 160 mal score  $\tau_i$ . This can be used to perform sentence-based  
 161 video search. We demonstrate this in section 4.3.

162 However, we first present the two central algorithmic con-  
 163 tributions of this work. In section 2 we present the de-  
 164 tails of the sentence tracker, the mechanism for efficiently  
 165 constraining several parallel detection-based trackers, one  
 166 for each participant, with a conjunction of finite state rec-  
 167 ognizers. In section 3 we present lexical semantics for  
 168 a small vocabulary of 17 lexical items (5 nouns, 2 adjec-  
 169 tives, 4 verbs, 2 adverbs, 2 spatial-relation prepositions, and  
 170 2 motion prepositions) all formulated as finite state rec-  
 171 ognizers over features extracted from detections produced by  
 172 an object detector, together with compositional semantics  
 173 that maps a sentence to a semantic formula  $\Phi$  constructed  
 174 from these finite state recognizers where the object tracks  
 175 are assigned to arguments of these recognizers.

## 2. The Sentence Tracker

176 Barbu *et al.* [2] address the issue of selecting detec-  
 177 tions for a track that simultaneously satisfies a temporal-  
 178 coherence measure and a single predicate corresponding to  
 179 an intransitive verb such as *bounce*. Doing so constitutes the  
 180 integration of top-down high-level information, in the form  
 181 of an event model, with bottom-up low-level information in  
 182 the form of object detectors. We provide a short review of  
 183 the relevant material in that work to introduce notation and  
 184 provide the basis for our exposition of the sentence tracker.

$$\max_{j^1, \dots, j^T} \sum_{t=1}^T f(b_{j^t}^t) + \sum_{t=2}^T g(b_{j^{t-1}}^{t-1}, b_{j^t}^t) \quad (1)$$

185 The first component is a detection-based tracker. For a given  
 186 video with  $T$  frames, let  $j$  be the index of a detection and  $b_j^t$   
 187 be a particular detection in frame  $t$  with score  $f(b_j^t)$ . A  
 188 sequence  $\langle j^1, \dots, j^T \rangle$  of detection indices, one for each  
 189 frame  $t$ , denotes a track comprising detections  $b_{j^t}^t$ . We seek  
 190 a track that maximizes a linear combination of aggregate  
 191 detection score, summing  $f(b_j^t)$  over all frames, and a mea-  
 192 sure of temporal coherence, as formulated in Eq. 1. The  
 193 temporal coherence measure aggregates a local measure  $g$   
 194 computed between pairs of adjacent frames, taken to be the  
 195 negative Euclidean distance between the center of  $b_{j^t}^t$  and  
 196 the forward-projected center of  $b_{j^{t-1}}^{t-1}$  computed with opti-  
 197 cal flow. Eq. 1 can be computed in polynomial time using  
 198 dynamic-programming with the Viterbi [15] algorithm. It  
 199 does so by formulating a lattice, whose rows are indexed  
 200 by  $j$  and whose columns are indexed by  $t$ , where the node  
 201 at row  $j$  and column  $t$  is the detection  $b_j^t$ . Finding a track  
 202 thus reduces to finding a path through this lattice.

$$\max_{k^1, \dots, k^T} \sum_{t=1}^T h(k^t, b_{j^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t) \quad (2)$$

203 The second component recognizes events with hidden  
 204 Markov models (HMMs), by finding a maximum *a posteriori*  
 205 probability (MAP) estimate of an event model given  
 206 a track. This is computed as shown in Eq. 2, where  $k^t$  de-  
 207 notes the state for frame  $t$ ,  $h(k, b)$  denotes the log proba-  
 208 bility of generating a detection  $b$  conditioned on being in

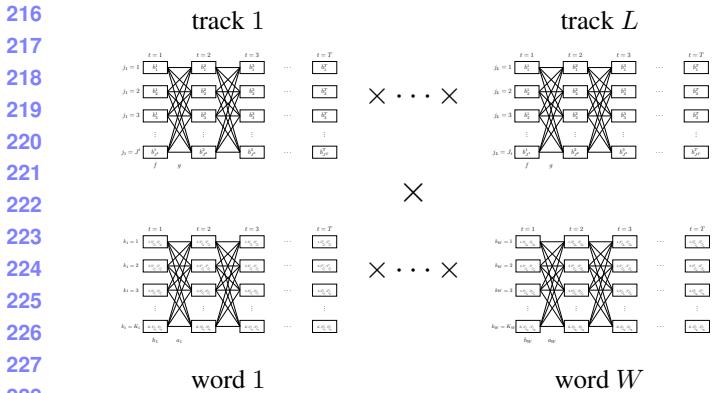


Figure 1. The cross-product lattice used by the sentence tracker, consisting of  $L$  tracking lattices and  $W$  event-model lattices.

state  $k$ ,  $a(k', k)$  denotes the log probability of transitioning from state  $k'$  to  $k$ , and  $j^t$  denotes the index of the detection produced by the tracker in frame  $t$ . This can also be computed in polynomial time using the Viterbi algorithm. Doing so induces a lattice, whose rows are indexed by  $k$  and whose columns are indexed by  $t$ .

The two components, detection-based tracking and event recognition, can be combined by combining the cost functions from Eq. 1 and Eq. 2 to yield a unified cost function

$$\max_{\substack{j_1^1, \dots, j_T^1 \\ k^1, \dots, k^T}} \sum_{t=1}^T f(b_{j^t}^t) + \sum_{t=2}^T g(b_{j^{t-1}}^t, b_{j^t}^t) + \sum_{t=1}^T h(k^t, b_{j^t}^t) + \sum_{t=2}^T a(k^{t-1}, k^t)$$

that computes the joint MAP estimate of the best possible track and the best possible state sequence. This is done by replacing the  $j^t$  in Eq. 2 with  $j^t$ , allowing the joint maximization over detection and state sequences. This too can be computed in polynomial time with the Viterbi algorithm, finding the optimal path through a cross-product lattice where each node represents a detection paired with an event-model state. This formulation combines a single tracker lattice with a single event model, constraining the detection-based tracker to find a track that is not only temporally coherent but also satisfies the event model. This can be used to select that *ball* track from a video that contains multiple balls that exhibits the motion characteristics of an intransitive verb such as *bounce*.

One would expect that encoding the semantics of a complex sentence such as *The person to the right of the chair quickly carried the red object towards the trash can*, which involves nouns, adjectives, verbs, adverbs, and spatial-relation and motion prepositions, would provide substantially more mutual constraint on the *collection* of tracks for the participants than a single intransitive verb would constrain a single track. We thus extend the approach described above by incorporating a complex multi-argument predicate that represents the semantics of an entire sentence in-

stead of one that only represents the semantics of a single intransitive verb. This involves formulating the semantics of other parts of speech, in addition to intransitive verbs, also as HMMs. We then construct a large cross-product lattice, illustrated in Fig. 1, to support  $L$  tracks and  $W$  words. Each node in this cross-product lattice represents  $L$  detections and the states for  $W$  words. To support  $L$  tracks, we subindex each detection index  $j$  as  $j_l$  for track  $l$ . Similarly, to support  $W$  words, we subindex each state index  $k$  as  $k_w$  for word  $w$  and the HMM parameters  $h$  and  $a$  for word  $w$  as  $h_w$  and  $a_w$ . The argument-to-track mappings  $\theta_w^1$  and  $\theta_w^2$  specify the tracks that fill arguments 1 and 2 (where necessary) of word  $w$  respectively. We then seek a path through this cross-product lattice that optimizes

$$\max_{\substack{j_1^1, \dots, j_L^T \\ j_L^1, \dots, j_L^T \\ k_1^1, \dots, k_W^T \\ k_W^1, \dots, k_W^T}} \sum_{l=1}^L \sum_{t=1}^T f(b_{j_l^t}^t) + \sum_{t=2}^T g(b_{j_l^{t-1}}^t, b_{j_l^t}^t) + \sum_{w=1}^W \sum_{t=1}^T h_w(k_w^t, b_{j_{\theta_w^1}^t}^t, b_{j_{\theta_w^2}^t}^t) + \sum_{t=2}^T a_w(k_w^{t-1}, k_w^t)$$

This can also be computed in polynomial time using the Viterbi algorithm. This describes a method by which the function  $S(D, \Phi) \mapsto (\tau, Z)$ , discussed earlier, can be computed, where  $D$  is the collection of detections  $b_j^t$  and  $Z$  is the collection of tracks  $j_l^t$ .

### 3. Natural-Language Semantics

The sentence tracker uniformly represents the semantics of words in all parts of speech, namely nouns, adjectives, verbs, adverbs, and prepositions (both those that describe spatial relations and those that describe motion), as HMMs. Finite state recognizers (FSMs) are a special case of HMMs where the transition matrices  $a$  and the output models  $h$  are 0/1. Here, we formulate the semantics of a small fragment of English consisting of 17 lexical items (5 nouns, 2 adjectives, 4 verbs, 2 adverbs, 2 spatial-relation prepositions, and 2 motion prepositions), by hand, as FSMs. We do so to focus on what once can do with this approach, namely take sentences as input and focus the attention of a tracker, take video as input and produce sentential descriptions as output, and perform content-based video retrieval given a sentential input query, as discussed in Section 4. It is particularly enlightening that the FSMs we use are perspicuous and clearly encode pretheoretic human intuitions about the semantics of these words. But nothing turns on the use of hand-coded FSMs. Our framework, as described above, supports HMMs. A companion submission describes a method by which one can automatically learn such HMMs for the lexicon, grammar, and corpus discussed in this paper.

Nouns (*e.g. person*) may be represented by constructing static FSMs over discrete features, such as detector class. Adjectives (*e.g. red, tall, and big*) may be represented as

|     |  |     |
|-----|--|-----|
| 324 | $S \rightarrow NP VP$  | 378 |
| 325 | $NP \rightarrow D [A] N [PP]$  | 379 |
| 326 | $D \rightarrow an   the$   | 380 |
| 327 | $A \rightarrow blue   red$   | 381 |
| 328 | $N \rightarrow person   backpack   trash can   chair   object$   | 382 |
| 329 | $PP \rightarrow P NP$  | 383 |
| 330 | $P \rightarrow to the left of   to the right of$   | 384 |
| 331 | $VP \rightarrow V NP [ADV] [PPM]$  | 385 |
| 332 | $V \rightarrow picked up   put down   carried   approached$  | 386 |
| 333 | $ADV \rightarrow quickly   slowly$   | 387 |
| 334 | $PPM \rightarrow PM NP$  | 388 |
| 335 | $PM \rightarrow towards   away from$   | 389 |
| 336 | (a)  | 390 |
| 337 | $to the left of = (agent patient) (referent)$  | 391 |
| 338 | $to the right of = (agent patient) (referent)$   | 392 |
| 339 | $picked up = (agent) (patient)$  | 393 |
| 340 | $put down = (agent) (patient)$   | 394 |
| 341 | $carried = (agent) (patient)$  | 395 |
| 342 | $approached = (agent) (goal)$  | 396 |
| 343 | $towards = (agent patient) (goal)$   | 397 |
| 344 | $away from = (agent patient) (source)$   | 398 |
| 345 | $other = (agent patient referent goal source)$   | 399 |
| 346 | (b)  | 400 |
| 347 | Table 1. (a) The grammar for our lexicon of 19 lexical entries (2 determiners, 2 adjectives, 5 nouns, 2 spatial relations, 4 verbs, 2 adverbs, and 2 motion prepositions). Note that the grammar allows for infinite recursion in the noun phrase. (b) The theta grid, specifying the number of arguments and roles such arguments refer to. (c) A selection of sentences drawn from the grammar based on which we collected multiple videos for our corpus. | 401 |
| 348 |  | 402 |
| 349 |  | 403 |
| 350 |  | 404 |
| 351 |  | 405 |
| 352 |  | 406 |
| 353 |  | 407 |
| 354 |  | 408 |
| 355 |  | 409 |
| 356 |  | 410 |
| 357 |  | 411 |
| 358 |  | 412 |
| 359 |  | 413 |
| 360 |  | 414 |
| 361 |  | 415 |
| 362 |  | 416 |
| 363 |  | 417 |
| 364 |  | 418 |
| 365 |  | 419 |
| 366 |  | 420 |
| 367 |  | 421 |
| 368 |  | 422 |
| 369 |  | 423 |
| 370 |  | 424 |
| 371 |  | 425 |
| 372 |  | 426 |
| 373 |  | 427 |
| 374 |  | 428 |
| 375 |  | 429 |
| 376 |  | 430 |
| 377 |  | 431 |
|     | (c)  |     |
|     | 1a. <i>The backpack approached the trash can.</i>  |     |
|     | b. <i>The chair approached the trash can.</i>  |     |
|     | 2a. <i>The red object approached the chair.</i>  |     |
|     | b. <i>The blue object approached the chair.</i>  |     |
|     | 3a. <i>The person to the left of the trash can put down an object.</i>   |     |
|     | b. <i>The person to the right of the trash can put down an object.</i>   |     |
|     | 4a. <i>The person put down the trash can.</i>  |     |
|     | b. <i>The person put down the backpack.</i>  |     |
|     | 5a. <i>The person carried the red object.</i>  |     |
|     | b. <i>The person carried the blue object.</i>  |     |
|     | 6a. <i>The person picked up an object to the left of the trash can.</i>  |     |
|     | b. <i>The person picked up an object to the right of the trash can.</i>  |     |
|     | 7a. <i>The person picked up an object.</i>   |     |
|     | b. <i>The person put down an object.</i>   |     |
|     | 8a. <i>The person picked up an object quickly.</i>   |     |
|     | b. <i>The person picked up an object slowly.</i>   |     |
|     | 9a. <i>The person carried an object towards the trash can.</i>   |     |
|     | b. <i>The person carried an object away from the trash can.</i>  |     |
|     | 10. <i>The backpack approached the chair.</i>  |     |
|     | 11. <i>The red object approached the trash can.</i>  |     |
|     | 12. <i>The person put down the chair.</i>  |     |

Table 1. (a) The grammar for our lexicon of 19 lexical entries (2 determiners, 2 adjectives, 5 nouns, 2 spatial relations, 4 verbs, 2 adverbs, and 2 motion prepositions). Note that the grammar allows for infinite recursion in the noun phrase. (b) The theta grid, specifying the number of arguments and roles such arguments refer to. (c) A selection of sentences drawn from the grammar based on which we collected multiple videos for our corpus.

static FSMs that describe select properties of the detections for a single participant, such as color, shape, or size, independent of other features of the overall event. Intransitive verbs (*e.g. bounce*) may be represented as FSMs that describe the changing motion characteristics of a single participant, such as *moving downward* followed by *moving upward*. Transitive verbs (*e.g. approach*) may be represented as FSMs that describe the changing relative motion characteristics of two participants, such as *moving closer*. Adverbs (*e.g. slowly* and *quickly*) may be represented by FSMs that describe the velocity of a single participant, independent of the direction of motion. Spatial-relation prepositions (*e.g. to the left of*) may be represented as static FSMs that describe the relative position of two participants. Motion prepositions (*e.g. towards* and *away from*) may be represented as FSMs that describe the changing relative position of two participants. As is often the case, even simple static properties, such as detector class, object color, shape, and size, spatial relations, and direction of motion, might hold only for a portion of an event. We handle such temporal uncertainty by incorporating garbage states into the FSMs that always accept and do not affect the scores computed. This also allows for alignment between multiple words in a temporal interval during a longer aggregate event. We formulate the FSMs for specifying the word meanings as regular expressions over predicates computed from detections. The particular set of regular expressions and associated predicates that are used in the experiments are given in Table 2. The predicates are formulated around a number of primitive functions. The function *avgFlow(b)* computes a vector that represents the average optical flow

inside the detection *b*. The functions *x(b)*, *model(b)*, and *hue(b)* return the *x*-coordinate of the center of *b*, its object class, and the average hue of the pixels inside *b* respectively. The function *fwdProj(b)* displaces *b* by the average optical flow inside *b*. The functions *∠* and *angleSep* determine the angular component of a given vector and angular distance between two angular arguments respectively. The function *normal* computes a normal unit vector for a given vector. The argument *v* to NOJITTER denotes a specified direction represented as a 2D unit vector in that direction. Regular expressions are formulated around predicates as atoms. A given regular expression must be formed solely from output models of the same arity and denotes an FSM with a  $-\infty/0$  transition matrix. We use a new regular-expression operator,  $R^{[n]} \triangleq (R [\text{TRUE}])^{\{n\}}$  to indicate that *R* must be repeated at least *n* times but can optionally have a single frame of noise between each repetition. This allows for some flexibility in the models.

A sentence may describe an activity involving multiple tracks, where different (collections of) tracks fill the arguments of different words. This gives rise to the requirement of compositional semantics: dealing with the mappings from arguments to tracks. Given a sentence, say *The person to the right of the chair picked up the backpack*, argument-to-track assignment is a function  $\mathcal{T}(\Lambda, \Gamma, \Psi) \mapsto (\Phi)$ , that takes, as input, a sentence  $\Lambda$  and a grammar  $\Gamma$ , along with a specification of the argument arity and role types  $\Psi$  for the words in the lexicon and produces a formula  $\Phi$  that specifies which tracks fill which arguments of which predicate instances for the words in the sentence. Such a function, applied to our example sentence with the grammar  $\Gamma$  as

| 432 | Constants                                     | Simple Predicates   | Complex Predicates   | 486 |
|-----|---|---|--|-----|
| 433 | $x\text{BOUNDARY} \triangleq 300\text{PX}$    | $\text{NOJITTER}(b, v) \triangleq \ \text{avgFlow}(b) \cdot v\  \leq \Delta\text{JUMP}$                       | $\text{STATIONARYCLOSE}(b_1, b_2) \triangleq \text{STATIONARY}(b_1) \wedge \text{STATIONARY}(b_2) \wedge \neg\text{ALIKE}(b_1, b_2) \wedge \text{CLOSE}(b_1, b_2)$ | 487 |
| 434 | $\text{NEXTTo} \triangleq 50\text{PX}$        | $\text{ALIKE}(b_1, b_2) \triangleq \text{model}(b_1) = \text{model}(b_2)$                                     | $\text{STATIONARYFAR}(b_1, b_2) \triangleq \text{STATIONARY}(b_1) \wedge \text{STATIONARY}(b_2) \wedge \neg\text{ALIKE}(b_1, b_2) \wedge \text{FAR}(b_1, b_2)$     | 488 |
| 435 | $\Delta\text{STATIC} \triangleq 6\text{PX}$   | $\text{FAR}(b_1, b_2) \triangleq  x(b_1) - x(b_2)  \geq x\text{BOUNDARY}$                                     | $\text{CLOSER}(b_1, b_2) \triangleq  x(b_1) - x(b_2)  >  x(\text{fwdProj}(b_1)) - x(b_2)  + \Delta\text{CLOSING}$  | 489 |
| 436 | $\Delta\text{JUMP} \triangleq 30\text{PX}$    | $\text{CLOSE}(b_1, b_2) \triangleq  x(b_1) - x(b_2)  < x\text{BOUNDARY}$                                      | $\text{FARTHER}(b_1, b_2) \triangleq  x(b_1) - x(b_2)  <  x(\text{fwdProj}(b_1)) - x(b_2)  + \Delta\text{CLOSING}$   | 490 |
| 437 | $\Delta\text{QUICK} \triangleq 80\text{PX}$   | $\text{LEFT}(b_1, b_2) \triangleq 0 < x(b_2) - x(b_1) \leq \text{NEXTTo}$                                     | $\text{MOVECLOSER}(b_1, b_2) \triangleq \text{NOJITTER}(b_1, (0, 1)) \wedge \text{NOJITTER}(b_2, (0, 1)) \wedge \text{CLOSER}(b_1, b_2)$                           | 491 |
| 438 | $\Delta\text{SLOW} \triangleq 30\text{PX}$    | $\text{RIGHT}(b_1, b_2) \triangleq 0 < x(b_1) - x(b_2) \leq \text{NEXTTo}$                                    | $\text{MOVEFARTHER}(b_1, b_2) \triangleq \text{NOJITTER}(b_1, (0, 1)) \wedge \text{NOJITTER}(b_2, (0, 1)) \wedge \text{FARTHER}(b_1, b_2)$                         | 492 |
| 439 | $\Delta\text{CLOSING} \triangleq 10\text{PX}$ | $\text{HASCOLOUR}(b, \text{hue}) \triangleq \text{angleSep}(\text{hue}(b), \text{hue}) \leq \Delta\text{HUE}$ | $\text{ALONGDIR}(b, v) \triangleq \text{angleSep}(\angle\text{avgFlow}(b), \angle v) < \Delta\text{DIRECTION} \wedge \neg\text{STATIONARY}(b)$                     | 493 |
| 440 | $\Delta\text{DIRECTION} \triangleq 30^\circ$  | $\text{STATIONARY}(b) \triangleq \ \text{avgFlow}(b)\  \leq \Delta\text{STATIC}$                              | $\text{MOVINGDIR}(b, v) \triangleq \text{ALONGDIR}(b, v) \wedge \text{NOJITTER}(b, \text{normal}(v))$  | 494 |
| 441 | $\Delta\text{HUE} \triangleq 30^\circ$        | $\text{QUICK}(b) \triangleq \ \text{avgFlow}(b)\  \geq \Delta\text{QUICK}$                                    | $\text{APPROACHING}(b_1, b_2) \triangleq \neg\text{ALIKE}(b_1, b_2) \wedge \text{STATIONARY}(b_2) \wedge \text{MOVECLOSER}(b_1, b_2)$                              | 495 |
| 442 |   | $\text{SLOW}(b) \triangleq \ \text{avgFlow}(b)\  \geq \Delta\text{SLOW}$                                      | $\text{DEPARTING}(b_1, b_2) \triangleq \neg\text{ALIKE}(b_1, b_2) \wedge \text{STATIONARY}(b_2) \wedge \text{MOVEFARTHER}(b_1, b_2)$                               | 496 |
| 443 |   | $\text{ISPERSON}(b) \triangleq \text{model}(b) = \text{person}$   | $\text{PICKINGUP}(b_1, b_2) \triangleq \neg\text{ALIKE}(b_1, b_2) \wedge \text{STATIONARY}(b_1) \wedge \text{MOVINGDIR}(b_2, (0, 1))$                              | 497 |
| 444 |   | $\text{ISBACKPACK}(b) \triangleq \text{model}(b) = \text{backpack}$   | $\text{PUTTINGDOWN}(b_1, b_2) \triangleq \neg\text{ALIKE}(b_1, b_2) \wedge \text{STATIONARY}(b_1) \wedge \text{MOVINGDIR}(b_2, (0, -1))$                           | 498 |
| 445 |   | $\text{ISTRASHCAN}(b) \triangleq \text{model}(b) = \text{trashcan}$   | $\text{CARRY}(b_1, b_2) \triangleq \text{MOVINGDIR}(b_1, v) \wedge \text{MOVINGDIR}(b_2, v)$   | 499 |
| 446 |   | $\text{ISCHAIR}(b) \triangleq \text{model}(b) = \text{chair}$   | $\text{CARRYING}(b_1, b_2) \triangleq \text{CARRY}(b_1, b_2, (0, 1)) \vee \text{CARRY}(b_1, b_2, (0, -1))$   | 500 |
| 447 |   | $\text{ISBLUE}(b) \triangleq \text{HASCOLOUR}(b, 225^\circ)$  |  | 501 |
| 448 |   | $\text{ISRED}(b) \triangleq \text{HASCOLOUR}(b, 0^\circ)$   |  | 502 |
| 449 |   |   |  | 503 |
| 450 |   |   |  | 504 |
| 451 |   |   |  | 505 |
| 452 |   |   |  | 506 |
| 453 |   |   |  | 507 |
| 454 |   |   |  | 508 |
| 455 |   |   |  | 509 |
| 456 |   |   |  | 510 |
| 457 |   |   |  | 511 |
| 458 |   |   |  | 512 |
| 459 |   |   |  | 513 |
| 460 |   |   |  | 514 |
| 461 |   |   |  | 515 |
| 462 |   |   |  | 516 |
| 463 |   |   |  | 517 |
| 464 |   |   |  | 518 |
| 465 |   |   |  | 519 |
| 466 |   |   |  | 520 |
| 467 |   |   |  | 521 |
| 468 |   |   |  | 522 |
| 469 |   |   |  | 523 |
| 470 |   |   |  | 524 |
| 471 |   |   |  | 525 |
| 472 |   |   |  | 526 |
| 473 |   |   |  | 527 |
| 474 |   |   |  | 528 |
| 475 |   |   |  | 529 |
| 476 |   |   |  | 530 |
| 477 |   |   |  | 531 |
| 478 |   |   |  | 532 |
| 479 |   |   |  | 533 |
| 480 |   |   |  | 534 |
| 481 |   |   |  | 535 |
| 482 |   |   |  | 536 |
| 483 |   |   |  | 537 |
| 484 |   |   |  | 538 |
| 485 |   |   |  | 539 |

Table 2. The finite-state recognizers corresponding to the lexicon in Table 1(a).

specified in Table 1(a) and *theta grid*  $\Psi$ , as specified in Table 1(b), would produce the following formula.

$$\begin{aligned} \text{PERSON}(P) \wedge \text{TOTHERIGHTOF}(P, Q) \wedge \text{CHAIR}(Q) \\ \wedge \text{PICKEDUP}(P, R) \wedge \text{BACKPACK}(R) \end{aligned}$$

To do so, we first construct a parse tree of the sentence  $\Lambda$  given the grammar  $\Gamma$ , using a recursive-descent parser, producing a parse tree. Such a parse tree encodes in its structure, the dependency relationships between different parts of speech as specified by the grammar. For each word, we then determine from the parse tree, which words in the sentence are determined to be its *dependents* in the sense of *government*, and how many such *dependents* exist, from the theta grid specified in Table 1(b). For example, the dependents of *to the right of* are determined to be *person* and *chair*, filling its first and second arguments respectively. Moreover, we determine a consistent assignment of roles, one of agent, patient, source, goal, and referent, for each participant track that fills the word arguments, from the allowed roles specified for that word and argument in the theta grid. Here,  $P$ ,  $Q$ , and  $R$  are participants that play the agent, referent, and patient roles respectively.

## 4. Experimental Evaluation

The sentence tracker supports three distinct capabilities. It can take sentences as input and focus the attention of a tracker, it can take video as input and produce sentential descriptions as output, and it can perform content-based video retrieval given a sentential input query. To evaluate these, we filmed a corpus of 94 short videos, of varying length, in 3 different outdoor environments. The camera was moved for each video so that the varying background precluded unanticipated confounds. These videos, filmed

with a variety of actors, each depicted one or more of the 21 sentences from Table 1(c). The depiction, from video to video, varied in scene layout and the actor(s) performing the event. The corpus was carefully constructed in a number of ways. First, many videos depict more than one sentence. In particular, many videos depict simultaneous distinct events. Second, each sentence is depicted by multiple videos. Third the corpus was constructed with minimal pairs: pairs of videos whose depicted sentences differ in exactly one word. These minimal pairs are indicated as the ‘a’ and ‘b’ variants of sentences 1–9 in Table 1(c). That varying word was carefully chosen to span all parts of speech and all sentential positions: sentence 1 varies subject noun, sentence 2 varies subject adjective, sentence 3 varies subject preposition, sentence 4 varies object noun, sentence 5 varies object adjective, sentence 6 varies object preposition, sentence 7 varies verb, sentence 8 varies adverb, and sentence 9 varies motion preposition. We filmed our own corpus as we are unaware of any existing corpora that exhibit the above properties. We annotated each of the 21 sentences with ground truth judgments for each of the 21 sentences, indicating whether the given clip depicted the given sentence. This set of 1974 judgments was used for the following analyses.

### 4.1. Focus of Attention

Tracking is traditionally performed using cues from motion, object detection, or manual initialization on an object of interest. However, in the case of a cluttered scene involving multiple activities occurring simultaneously, there can be many moving objects, many instances of the same object class, and perhaps even multiple simultaneously occurring instances of the same event class. This presents a significant

540 obstacle to the efficacy of existing methods in such scenarios.  
 541 To alleviate this problem, one can decide which objects  
 542 to track based on which ones participate in a target event.  
 543

544 The sentence tracker can focus its attention on just those  
 545 objects that participate in an event specified by a sentential  
 546 description. Such a description can differentiate between  
 547 different simultaneous events taking place between many  
 548 moving objects in the scene using descriptions constructed  
 549 out of a variety of parts of speech: nouns to specify object  
 550 class, adjectives to specify object properties, verbs to  
 551 specify events, adverbs to specify motion properties, and  
 552 prepositions to specify (changing) spatial relations between  
 553 objects. Furthermore, such a sentential description can even  
 554 differentiate which objects to track based on the role that  
 555 they play in an event: agent, patient, source, goal, or referent.  
 556 Fig. 2 demonstrates this ability: different tracks are  
 557 produced for the same video that depicts multiple simultaneous  
 558 events when focused with different sentences.  
 559

560 We further evaluated this ability on all 9 minimal pairs,  
 561 collectively applied to all 24 suitable videos in our corpus.  
 562 For 21 of these, both sentences in the minimal pair yielded  
 563 tracks deemed to be correct depictions. We include example  
 564 videos for all 9 minimal pairs in the supplementary material.

## 565 4.2. Generation

566 Much of the prior work on generating sentences to de-  
 567 scribe images [4, 7, 8, 12, 13, 18] and video [1, 6, 9, 10, 16]  
 568 uses special-purpose natural-language-generation methods.  
 569 We can instead use the ability of the sentence tracker to  
 570 score a sentence paired with a video as a general-purpose  
 571 natural-language generator by searching for the highest-  
 572 scoring sentence for a given video. However, this has a  
 573 problem. Since  $h$  and  $a$  are log probabilities,  $g$  is a nega-  
 574 tive Euclidean distance, and we constrain  $f$  to be nega-  
 575 tive, scores decrease with longer word strings and greater  
 576 numbers of tracks that result from longer word strings. So  
 577 we don't actually search for the highest-scoring sentence,  
 578 which would bias the process towards short sentences. In-  
 579 stead we seek complex sentences that are true of the video  
 580 as they are more informative.

581 Nominally, this search process would be intractable since  
 582 the space of possible sentences can be huge and even infinite.  
 583 However, we can use beam search to get an approximate  
 584 answer. This is possible because the sentence tracker  
 585 can score any collection of words, not just complete phrases  
 586 or sentences. We can select the  $k$  top-scoring single-word  
 587 strings and then repeatedly extend the  $k$  top-scoring  $n$ -word  
 588 strings, by one word, to select the  $k$  top-scoring  $n + 1$ -word  
 589 strings, subject to the constraint that these  $n + 1$ -word  
 590 strings can be extended to grammatical sentences by inser-  
 591 tion of additional words. Thus we terminate the search pro-  
 592 cess when the *contraction threshold*, the ratio between the  
 593 score of an expanded string and the score of the string it ex-  
 panded from, exceeds a specified value and the string being

594 expanded is a complete sentence. This contraction thresh-  
 595 old controls complexity of the generated sentence.  
 596

597 When restricted to FSMs,  $h$  and  $a$  will be 0/1 which be-  
 598 come  $-\infty/0$  in log space. Thus increase in the number of  
 599 words can only decrease a score to  $-\infty$ , meaning that a  
 600 string of words is no-longer true of a video. Since we seek  
 601 true sentences, we terminate the above beam search process  
 602 before the score goes to  $-\infty$ . In this case, there is no ap-  
 603 proximation: a beam search maintaining all  $n$ -word strings  
 604 with finite score yields the highest-scoring sentence before  
 605 the contraction threshold is met.  
 606

607 To evaluate this approach, we searched the space of sen-  
 608 tences in the grammar in Table 1(a) to find the best true sen-  
 609 tence for each of the 94 videos in our corpus. Note that the  
 610 grammar generates an infinite number of sentences due to  
 611 recursion in NP. Even restricting the grammar to eliminate  
 612 NP recursion yields a space of 147,123,874,800 sentences.  
 613 Despite not restricting the grammar in this fashion, we are  
 614 able to effectively find good descriptions of the videos. We  
 615 computed the accuracy of the sentence tracker in generat-  
 616 ing descriptions for all 94 videos in our corpus for multiple  
 617 contraction thresholds. Accuracy was computed as the per-  
 618 centage of the 94 videos for which the sentence tracker pro-  
 619 duced descriptions that were deemed to be true. Contraction  
 620 thresholds of 0.95, 0.90, and 0.85 yielded accuracies of  
 621 63.82%, 69.14%, and 64.89% respectively. We demon-  
 622 strate examples of this approach in Fig. 3. The supplemen-  
 623 tary material contains additional examples.  
 624

## 625 4.3. Retrieval

626 The availability of vast video corpora, such as on  
 627 YouTube, has created a rapidly growing demand for  
 628 content-based video search and retrieval. The existing sys-  
 629 tems, however, only provide a means to search via human-  
 630 provided captions. The inefficacy of such an approach is  
 631 evident. Attempting to search for even simple queries such  
 632 as *pick up* or *put down* yields surprisingly poor results, let  
 633 alone searching for more complex queries such as *person*  
 634 *approached horse*. Furthermore, prior work on content-  
 635 based video-retrieval systems like Sivic and Zisserman [14]  
 636 search only for objects and like Laptev *et al.* [11] search  
 637 only for events. Even combining such to support conjunc-  
 638 tive queries for videos with specified collections of objects  
 639 jointly with a specified event, would not effectively rule out  
 640 videos where the specified objects did not play a role in  
 641 the event or played different roles in the event. For exam-  
 642 ple, it could not rule out a video depicting a person jump-  
 643 ing next to a stationary ball for a query *ball bounce* or dis-  
 644 tinguish between the queries *person approached horse* and  
 645 *horse approached person*. The sentence tracker exhibits the  
 646 ability to serve as the basis of a much better video search  
 647 and retrieval tool, one that performs content-based search  
 648 with complex sentential queries to find precise semantically  
 649 relevant clips, as demonstrated in Fig. 4.  
 650

A photograph showing two men in a park-like setting. One man in a red t-shirt and white shorts stands by a blue trash bin. The other man, wearing a blue long-sleeved shirt and dark pants, is bent over, reaching for a red backpack lying on the grass. A blue rectangular bounding box is drawn around the red backpack, and a red rectangular bounding box is drawn around the head of the man bending over.

*The person picked up an object.*

The person put down the object.

*The person put down an object.*

Figure 2. Sentence-guided focus of attention: different sets of tracks for the same video produced under guidance of different sentences.

A photograph showing a man standing next to a blue trash can in a park. A red box highlights a backpack on the ground to the left, and a green box highlights the trash can itself. The background shows trees and a paved area.

*The backpack approached the trash can.*

A photograph showing a man with a beard and a blue shirt standing next to a blue trash can on a paved area. A red bounding box encloses the man, and a green bounding box encloses the trash can. A blue line extends from the bottom left corner of the image to the trash can.

*The person to the left of the trash can put down the chair.*

Figure 3. Generation of sentential descriptions: constructing the highest-scoring sentence for each video that is generated by the grammar in Table 1(a), by means of a beam search.

To evaluate this approach, we scored every video in our corpus against every sentence in Table 1(c), rank ordering the videos for each sentence, yielding the following statistics over the 1974 scores.

|   |         |
|---|---------|
| chance that a random video depicts a given sentence               | 13.12%  |
| top-scoring video depicts the given sentence                      | 85.71%  |
| at least 1 of the top 3 scoring videos depicts the given sentence | 100.00% |

The judgment of whether a video depicted a given sentence was made using our annotation. We conducted an additional evaluation with this annotation. One can threshold the sentence-tracker score to yield a binary predicate on video-sentence pairs. We performed 4-fold cross validation on our corpus, selecting the threshold for each fold that maximized accuracy of this predicate, relative to the annotation, on 75% of the videos and evaluating the accuracy with this selected threshold on the remaining 25%. This yielded an average accuracy of 91.74%.

## 5. Conclusion

We have presented a novel framework that utilizes the compositional structure of events and the compositional

The figure consists of three side-by-side photographs of a scene in a park. In each photo, a man in a red t-shirt and white shorts is bending over a blue trash can. A second man in a dark blue t-shirt and dark pants stands to his right, holding a brown backpack. A red dashed box highlights the area around the trash can. In the first panel, the man in red has just placed the backpack on the ground next to the trash can. In the second panel, he is in the middle of putting it down. In the third panel, he has finished putting it down and is standing upright again.

The same video produced under guidance of different sentences.

structure of language to drive a semantically meaningful and targeted approach towards activity recognition. This multimodal framework integrates low-level visual components, such as object detectors, with high-level semantic information in the form of sentential descriptions in natural language. This is facilitated by the shared structure of detection-based tracking, which incorporates the low-level object-detector components, and of finite-state recognizers, which incorporate the semantics of the words in a lexicon.

We demonstrated the utility and expressiveness of our framework by performing three separate tasks on our corpus, requiring no training or annotation, simply by leveraging our framework in different manners. The first, sentence-guided focus of attention, showcases the ability to focus the attention of a tracker on the activity described in a sentence, indicating the capability to identify such subtle distinctions as between *The person picked up the chair to the left of the trash can* and *The person picked up the chair to the right of the trash can*. The second, generation of sentential description of video, showcases the ability to produce a complex description of a video, involving multiple parts of speech, by performing an efficient search for the best description

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809



The person carried an object away from the trash can.



The person picked up an object to the left of the trash can.

Figure 4. Sentential-query-based video search: returning the best-scoring video, in a corpus of 94 videos, for a given sentence.

though the space of all possible descriptions. The final task, query-based video search, showcases the ability to perform content-based video search and retrieval, allowing for such distinctions as between *The person approached the trash can* and *The trash can approached the person*.

## References

- [1] A. Barbu, A. Bridge, Z. Burchill, D. Coroian, S. Dickinson, S. Fidler, A. Michaux, S. Mussman, N. Siddharth, D. Salvi, L. Schmidt, J. Shangguan, J. M. Siskind, J. Waggoner, S. Wang, J. Wei, Y. Yin, and Z. Zhang. Video in sentences out. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence*, pages 102–112, Aug. 2012. **6**
- [2] A. Barbu, N. Siddharth, A. Michaux, and J. M. Siskind. Simultaneous object detection, tracking, and event recognition. *Advances in Cognitive Systems*, 2:203–220, Dec. 2012. **1, 2**
- [3] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes (VOC) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010. **1**
- [4] A. Farhadi, M. Hejrati, M. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth. Every picture tells a story: Generating sentences from images. In *Proceedings of the European Conference on Computer Vision*, pages 15–29, 2010. **6**
- [5] P. F. Felzenszwalb, R. B. Girshick, and D. McAllester. Cascade object detection with deformable part models. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2241–2248, 2010. **2**
- [6] C. Fernández Tena, P. Baiget, X. Roca, and J. González. Natural language descriptions of human behavior from video sequences. In *Advances in Artificial Intelligence*, pages 279–292, 2007. **6**
- [7] A. Gupta, Y. Verma, and C. Jawahar. Choosing linguistics over vision to describe images. In *Proceedings of the Twenty-Sixth National Conference on Artificial Intelligence*, 2012. **6**
- [8] L. Jie, B. Caputo, and V. Ferrari. Who's doing what: Joint modeling of names and verbs for simultaneous face and pose annotation. In *Proceedings of the Neural Information Processing Systems Conference*, 2009. **6**
- [9] M. U. G. Khan and Y. Gotoh. Describing video contents in natural language. In *Proceedings of the Workshop on Innovative Hybrid Approaches to the Processing of Textual Data*, pages 27–35, 2012. **6**
- [10] A. Kojima, T. Tamura, and K. Fukunaga. Natural language description of human activities from video images based on concept hierarchy of actions. *International Journal of Computer Vision*, 50(2):171–184, Nov. 2002. **6**
- [11] I. Laptev, M. Marszałek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2008. **6**
- [12] P. Li and J. Ma. What is happening in a still picture? In *First Asian Conference on Pattern Recognition*, pages 32–36, Nov. 2011. **6**
- [13] M. Mitchell, J. Dodge, A. Goyal, K. Yamaguchi, K. Stratos, X. Han, A. Mensch, A. C. Berg, T. L. Berg, and H. Daumé III. Midge: Generating image descriptions from computer vision detections. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, 2012. **6**
- [14] J. Sivic and A. Zisserman. Video Google: a text retrieval approach to object matching in videos. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, pages 1470–1477, 2003. **6**
- [15] A. J. Viterbi. Convolutional codes and their performance in communication systems. *IEEE Transactions on Communication*, 19:751–772, Oct. 1971. **2**
- [16] Z. Wang, G. Guan, Y. Qiu, L. Zhuo, and D. Feng. Semantic context based refinement for news video annotation. *Multimedia Tools and Applications*, pages 1–21, 2012. **6**
- [17] J. K. Wolf, A. M. Viterbi, and G. S. Dixon. Finding the best set of K paths through a trellis with application to multitarget tracking. *IEEE Transactions on Aerospace and Electronic Systems*, 25(2):287–296, Mar. 1989. **1**
- [18] Y. Yang, C. L. Teo, H. Daumé III, and Y. Aloimonos. Corpus-guided sentence generation of natural images. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 444–454, 2011. **6**

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863