

# Unsupervised Discovery of Character Dictionaries in Animation Movies

Krishna Somandepalli, *Member, IEEE*, Naveen Kumar, *Member, IEEE*, Tanaya Guha, *Member, IEEE*, and Shrikanth S. Narayanan, *Fellow, IEEE*

**Abstract**—Automatic content analysis of animation movies can enable an objective understanding of character (actor) representations and their portrayals. It can also help illuminate potential markers of their unconscious biases and their impact. However, multimedia analysis of movie content has predominantly focused on live-action features. A dearth of multimedia research in this field is because of the complexity and heterogeneity in the design of animated characters—an extremely challenging problem to be generalized by a single method or model. In this paper, we address the problem of automatically discovering characters in animation movies as a first step toward automatic character labeling in these media. Movie-specific character dictionaries can act as a powerful first step for subsequent content analysis at scale. We propose an unsupervised approach which requires no prior information about the characters in a movie. We first use a deep neural network-based object detector that is trained on natural images to identify a set of initial character candidates. These candidates are further pruned using saliency constraints and visual object tracking. A character dictionary per movie is then generated from exemplars obtained by clustering these candidates. We are able to identify both anthropomorphic and nonanthropomorphic characters in a dataset of 46 animation movies with varying composition and character design. Our results indicate high precision and recall of the automatically detected characters compared to human-annotated ground truth, demonstrating the generalizability of our approach.

**Index Terms**—Animation movies, deep neural networks, object tracking, saliency, unsupervised clustering, video diarization.

## I. INTRODUCTION

**A**UTOMATIC analysis of movie content is of growing interest in the multimedia research community. One of the driving factors for this research is the large number of movies that are produced, disseminated and consumed annually. Besides being of entertainment value, movies often have an effect

Manuscript received November 28, 2016; revised July 26, 2017; accepted August 4, 2017. This work is based upon work supported by the National Science Foundation. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Xilin Chen. (*Corresponding author: Krishna Somandepalli.*)

K. Somandepalli, N. Kumar, and T. Guha is with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA 90047 USA (e-mail: krishna.somandepalli@gmail.com; knaveen87@gmail.com; tanayaguha@gmail.com).

S. S. Narayanan is with the Signal & Image Processing Institute, University of Southern California, Los Angeles, CA 90089 USA (e-mail: shri@sipi.usc.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TMM.2017.2745712

on certain social and economic aspects, as well as have a global reach and audience.

Researchers have addressed movie content analysis with different objectives and outlooks. Such efforts are often based on efficient indexing and organization of the media content for easy user navigation. They include shot boundary detection for movie segmentation [1], [2], video summarization [3] and abstraction [4]. The study in [2] builds a generative model that incorporates contextual information in order to reorganize interleaved shots into multiple plot threads. Approaches such as in [5] combines the aspects of video summarization, i.e., *who*, *what*, *where* and *when* for a semantic understanding of the movie content and structure. *RoleNet* proposed in [6] examines the movie content from a social network analysis perspective of the movie character roles rather than using audiovisual features. In general, movie content is a rich source of data that includes audio, video and text (dialogs) that enables such multimodal analysis.

Complementary to the aforementioned studies which attempt to achieve a high-level understanding of movies, efforts for a fine-grained (frame level or scene level statistics) analysis of video content have also been emerging. One such application is to quantify the amount of time a character appears on screen in a movie. The study in [7] examined these aspects with respect to gender revealing skewed distributions for the onscreen time of female characters. In order to advance from gender-level statistics to character-level statistics, person identification or character labeling is a crucial step in this direction. We refer to this problem as automatic video diarization – partitioning the video stream into actor-homogeneous segments, i.e., *who appeared*, *when* and for *how long*. Character labeling in live-action TV and movies has been achieved with modest success in [8]–[12]. This is typically performed by clustering the detected faces (e.g., [8]) or by multimodal approaches (e.g., [9], [10]) that model audio and subtitles or scripts alongside the detected faces from video.

It is important to note that all these studies exclusively focus on live-action TV and do not generalize to animated media content. Digital animation movies have contributed to over 10% of the box office market shares in the past decade [13]. Multimedia research in this domain is extremely scarce and technology developed for live-action TV content fails for animated content. Human face detection is the crux of character labeling methods for live action TV. Since human-characters can be uniquely identified by their faces, this method performs adequately well. But, such methods developed for human faces do not work



Fig. 1. Examples illustrating the heterogeneity of animated characters. (a): human-like (Frozen) (b): anthropomorphic (Frozen) (c) and (d): abstract (How to Train your Dragon, and Cars).

84 for the digital animation genre. Animated characters, though  
 85 mostly anthropomorphic (having human characteristics) are not  
 86 always human-like in appearance. They can be fictional animals,  
 87 inanimate objects or abstract in design (see Fig. 1 for a  
 88 few examples).

89 A major obstacle for automating content analysis of animated  
 90 media is the lack of a model that generalizes across different  
 91 characters with varying composition and design. This task be-  
 92 comes extremely complex given that all the characters even  
 93 within a single movie may not share the same structural char-  
 94 acteristics (e.g., human-like and non-human characters from the  
 95 same movie—Fig. 1(a) and (b) from the movie Frozen).

96 In the context of video diarization, when the characters that  
 97 appear on screen are generally not known a priori, a key step  
 98 is to provide a list of characters that form the *who appeared*  
 99 component of the system. We refer to such a list of characters  
 100 specific to each movie as a *character dictionary*. The automatic  
 101 discovery of these character dictionaries is the primary objective  
 102 in this paper. Our overarching goal is to engineer a model for  
 103 animation movie video diarization. With the proposed character  
 104 dictionaries, animation character labeling may be achieved by  
 105 techniques such as [14] that can retrieve frames and shots given  
 106 an object of interest.

107 In content analysis of animated media, researchers have thus  
 108 far focused on problems such as cut detection [15], color-based  
 109 video categorization [16] and movie abstraction [17], [18]. One  
 110 method proposed in [19] performs human-like face detection  
 111 from cartoon images using skin-segmentation techniques. Con-  
 112 sidering the variation in texture, color and shape of animated  
 113 characters in general (as illustrated in Fig. 1), these methods do  
 114 not generalize well. To the best of our knowledge, no work to  
 115 date has specifically addressed the problem of automatic dis-  
 116 covery of characters from animated media in a scalable manner.

117 In contrast to live-action movies, animation movies are com-  
 118 pletely artist generated. Sketches of the characters are designed  
 119 by the artists or the animators, generally referred to as *model*  
 120 *sheets* from which character-specific 3D models are generated.  
 121 Sketch based image retrieval systems such as [20] can be used  
 122 to achieve video diarization when model sheets are available.  
 123 However, model sheets are copyrighted material and mostly  
 124 owned by the animation studio which produced the movie. As

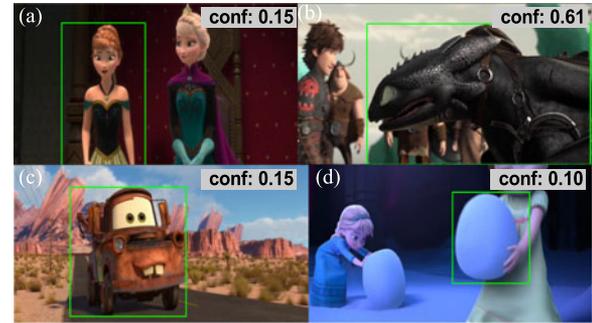


Fig. 2. Character candidates chosen by the Multibox object detector. Conf. indicates the confidence score of the network for the detected object.

such, they are not publicly available and approaches which are  
 based on model sheets will not be scalable for all movies.

125  
 126  
 127 In 1981, Frank Thomas and Ollie Johnston published *The*  
 128 *Illusion of Life* [21]; it outlines a set of twelve basic principles  
 129 of animation. Animators have been using this as a cookbook  
 130 for designing characters in order for the viewers to appreci-  
 131 ate “animation” over mere “movement”. While most of these  
 132 principles aid animators in adding semantic or artistic value  
 133 (e.g., anticipation, exaggeration), a few can be exploited in a  
 134 computer vision context (e.g., *Solid Drawing*: drawing volume  
 135 solidity and illusion of three dimensions; *Staging*: Distinctive  
 136 color, depth of field and positioning in the frame to highlight  
 137 the character). Defining an *animated character* in a complete  
 138 sense would involve delineating abstract concepts such as life  
 139 (or sentience even) from movie content. In this paper, we only  
 140 analyze the video stream from animation movies and leverage  
 141 some of the aforementioned principles of animation as proxies  
 142 to identify the characters.

143 At the outset, we pose our problem as an object detection task  
 144 where any object can be a possible *character candidate*. Ani-  
 145 mation movie frames are comparable with natural photographic  
 146 images, especially in their similarities of depth of field and the  
 147 character presentation in a frame. Additionally, we assume no  
 148 prior models with respect to shape, size, color, or texture for  
 149 these candidates in order for the proposed system to generalize.

150 A few prominent examples of state-of-the-art object detec-  
 151 tion systems include discriminatively trained deformable parts-  
 152 based model (DPM, [22], [23]) and deep neural network (DNN)  
 153 models such as [24]–[26], both of which are supervised and  
 154 trained over a predefined set of object classes. DPMs need  
 155 a carefully designed part-decomposition model of an object  
 156 which makes it unsuitable given the heterogeneity of characters  
 157 within just a single movie. In contrast, DNN-based methods such  
 158 as [24] can detect objects in real-time and outperform DPMs.  
 159 Specifically, DNN models that are saliency-inspired in design  
 160 [25] are of interest for our problem statement. Although super-  
 161 vised with a finite set of object classes, they have been shown to  
 162 detect objects in a *class-agnostic* manner [26] i.e., detect classes  
 163 of objects not used for training the model.

164 Movies in general, portray only a handful of *prominent* char-  
 165 acters. They are more likely to appear frequently in order for  
 166 the viewer to easily comprehend the content and the plot of

the movie. Additionally in movies, characters or the objects-of-interest tend to remain on screen for up to a few seconds depending on the situation. Visual object tracking can be used as an effective method to segment characters locally in time. Several previous works have used tracking as a means to automatically detect a class of objects (e.g., pedestrians, [27]). Object tracking algorithms can be error-prone in a movie video environment because of object deformation, background clutter, changes in illumination, occlusion and lack of a stationary backgrounds. However, visual tracking can minimize the number of detected objects to be considered by accounting for minor deformation or linear motion of the object. Furthermore, tracking also provides time information that can be used for diarization subsequently. For example, in [11], supervisory information available on a profile face is used to learn the appearance of a frontal face from faces tracked in TV series. A reasonable assumption in describing animated character is that the prominent characters are not transient when presented on-screen and appear frequently in the movie. In our method, we use this aspect of character presentation in movies to select character candidates. As a result, the character dictionaries consist of only the frequently occurring characters.

In this paper, we propose a novel approach to automatically discover characters that appear in an animation movie. Our proposed method is unsupervised in the sense that we do not train any aspect of our system with data from animated media content. Furthermore, we use no specific knowledge of the animation style or the physical attributes of the animated characters, thereby ensuring that our system can scale and generalize through the whole spectrum of animation movie content.

The rest of the paper is organized as follows: Section II describes the proposed system for selecting character candidates from an animation movie. In Section III, we present the experiments performed and the creation of an evaluation database. Section IV contains the experimental results and final considerations followed by conclusions and future work in Section V.

## II. METHODS

In this section, we first introduce the different systems that we use to identify and prune the detected objects to obtain a set of possible character candidates. We then use a clustering approach to identify character exemplars that constitute the final character dictionary. The overview of the proposed system is shown in the Fig. 3.

Our animation movie database consisted of forty-six movies, for which we annotated their prominent characters. We then conducted a detailed performance evaluation on eight animation movies which were chosen to represent varying degrees of heterogeneity in character design and composition. The movie-cast data from forty-six movies used for our system evaluation and the output from our system has been released as part of the SAIL Animation Movie character Database (SAIL-AMDb).<sup>1</sup> We have also made the code publicly available.<sup>2</sup>

<sup>1</sup><https://github.com/usc-sail/mica-animation/wiki>

<sup>2</sup><https://github.com/usc-sail/mica-animation>

### A. Coarse Detection of Character Candidates

Animated characters are often designed to have the appearance of a 3D object and characterized by shallow focus where the image plane of the character is in focus while the rest of the frame is out of focus [21]. In other words, they are the salient objects in a given frame. Capitalizing on this, we define a *character candidate* as any object that can be detected by a general-purpose object detector.

We use a pre-trained deep neural network (DNN) called *MultiBox* [25], [26], designed for object detection. Our preliminary experiments with other region proposal networks such as [24] yielded similar results. We chose *MultiBox* since our motivation for using an object detector was only to generate an initial set of potential character candidates.

*MultiBox* is a convolutional neural network (CNN) with an inception-style architecture [28] trained with the full 200-category object detection challenge data set from ImageNet Large Scale Visual Recognition Challenge 2014 (ILSVRC-2014) [29]. This model generates multiple bounding boxes and an associated confidence score that quantifies the network's confidence of each box containing an object. The model has been shown to perform object localization in a *class-agnostic* manner and achieve state-of-the-art performance in object detection tasks [25]. Furthermore, since the network is tailored towards the localization problem, it achieves a scalable representation of multiple salient objects in an image. These features make this model uniquely suitable for our problem. It is important to note that this model is trained with natural images of distinct object classes. Although the authors in [25] have shown that the model generalizes over unseen classes, here we apply the pre-trained DNN for images sampled from animation movies. We refer to this discrepancy as *DNN training bias*. This results in detecting objects that are not characters in a movie (e.g., traffic-light, chair). We refer to such objects as *noisy objects*.

In order to reduce the computational time, we downsample a movie (originally encoded at 23.98 fps) by one frame every 0.42 s (every 10th frame). The resulting frames are input to *MultiBox* [25] to obtain all possible bounding boxes for each image. The confidence score that is returned with each of these boxes was originally optimized in the DNN to match the ground truth object boxes from natural images.

Because of the aforementioned *DNN training bias*, we generally observed lower range of confidence scores for objects detected that were animated characters. We chose to retain objects with a confidence score greater than 0.1. In order to determine this threshold, we randomly sampled 100,000 frames from the movie *Frozen* (2013) in our movie database. We first assumed to have at most five possibly overlapping objects of interest in one frame and obtained the confidence scores for the five most confident objects in each frame. We then examined the distribution of the confidence scores for all the objects detected. We set the confidence threshold to 75th percentile of the distribution of confidence scores which is equal to 0.1002, thus retaining all objects with confidence score greater than 0.1. We apply this confidence threshold for all the movies in our database. A few examples of objects detected and their confidence scores returned by the network are shown in Fig. 2.

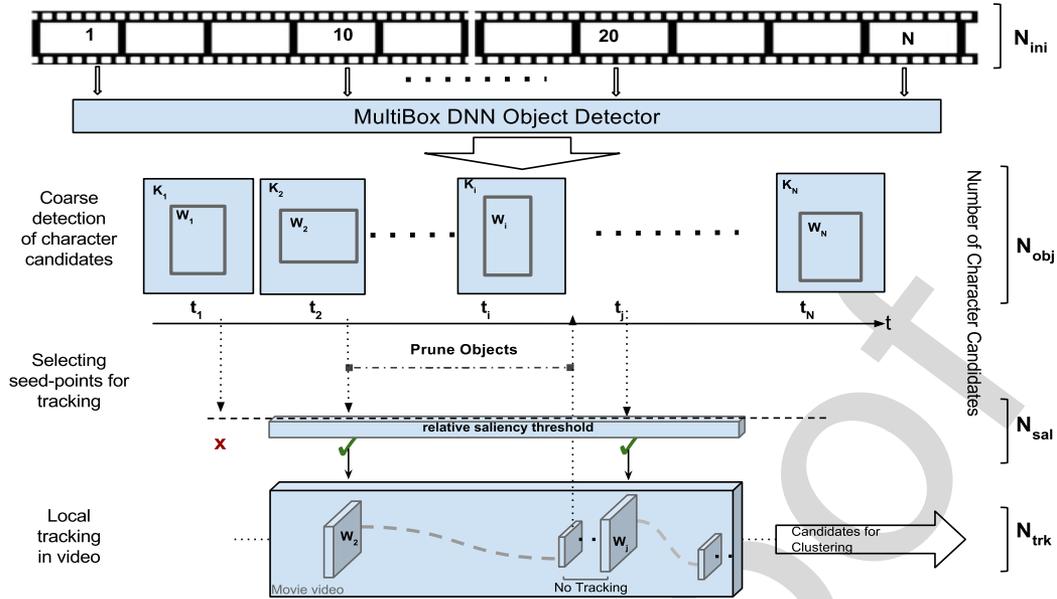


Fig. 3. Schematic diagram of the proposed method.

277 We also computed the area of each bounding box of an object  
 278 relative to the image frame and excluded objects in bounding  
 279 boxes with an area less than 1% or greater than 99% of the entire  
 280 frame. This ensures that very small objects and holistic scenes  
 281 are excluded as character candidates. When multiple objects  
 282 were detected in a single frame, we pruned them to obtain at  
 283 most one object per frame following the approach in [25]. We  
 284 performed non-maximum-suppression with a Jaccard similarity  
 285 [30] threshold of 0.5 and, chose the object with the maximum  
 286 area in that frame. We identified only a single object per frame in  
 287 order to simplify the subsequent step of single-target visual object  
 288 tracking. We refer to a chosen frame containing a character  
 289 candidate as a *candidate frame*.

290 A schematic of the proposed approach is illustrated in Fig. 3.  
 291 Let  $N_{ini}$  be the initial number of images (movie frames) input to *MultiBox* and  $N_{obj}$  be the number of character candidates chosen. We denote the candidate frame  $K$  and the bounding box  $W$  enclosing the object as a set  $M_{obj} = \{(K_{t_i}^{(i)}, W_{t_i}^{(i)}) | i \in [1, N_{obj}]\}$  and  $t_i$  refers to the time (or frame number) in the movie at which the object  $i$  occurs. Qualitative analyses showed that this step captures most of the characters in an animation movie at least once (e.g., images shown in Fig. 2(a)–(c)). However, this set also contains redundant and noisy objects which include non-characters or background objects (e.g., Fig. 2(d)).

### 302 B. Saliency Constraints and Object Tracking

303 In the next phase of our system, we used the saliency of  
 304 the detected object as a constraint to prune the set of character  
 305 candidates obtained in the previous step. We use this pruned  
 306 set of candidate frames as seed-points for tracking. During  
 307 tracking, we do not distinguish camera motion from object motion,  
 308 thereby ensuring that a sufficient condition for a character



Fig. 4. (a) Example for DNN training bias and saliency constraint; (b) Masked regions showing saliency, here *relative saliency score*  $R_s(W_1) = 9.2\%$ .

candidate is its presence on the screen rather than motion (e.g., 309  
 a talking tree). 310

311 1) *Saliency-Constrained Pruning*: As described in  
 312 Section II-A, the DNN training bias may result in choosing  
 313 objects that, although salient, may not be the characters of  
 314 interest (e.g., detected lamp in a scene with two characters—see  
 315 Fig. 4(a)). To quantify this, we use a saliency measure proposed  
 316 in [31] for the character candidate with respect to the entire  
 317 frame. Unsupervised methods that estimate saliency typically  
 318 use pixel-level features such as color, intensity (e.g., [32])  
 319 or background-detection in dynamic scenes (e.g., [33]). In  
 320 contrast, the measure proposed in [31] estimates saliency of  
 321 local areas (instead of pixel level) in static images and requires  
 322 no training. This method uses a kernel-based approach where  
 323 the size of the window relates to the scale of the target objects.  
 324 The saliency of a pixel inside the window is estimated using the  
 325 conditional probability of that pixel drawn from the distribution  
 326 estimated inside that window versus the distribution of the  
 327 surrounding area.

328 We first converted the RGB images to CIELAB color space  
 329 (because of the perceptual uniformity of the CIE color space<sup>3</sup>)  
 330 to estimate a saliency map for the entire candidate frame by

<sup>3</sup><http://www.brucelindbloom.com>

331 choosing window sizes at different scales as described in [34].  
 332 The resulting saliency maps are binarized by setting values  
 333 greater than 0.7 to 1 as recommended in [34]. An example of  
 334 the saliency map is shown in Fig. 4(b). Let  $A_s(W)$  be the area  
 335 of the salient region contained within a bounding box,  $W$  in an  
 336 image frame  $K$ . We define a *relative saliency score*,  $R_s(W)$  of  
 337 an object enclosed by the box  $W$  as the percentage salient area  
 338 it contributes to the frame,  $K$ :

$$R_s(W) = \frac{A_s(W)}{A_s(K)} \times 100. \quad (1)$$

339 We obtained the *relative saliency score*,  $R_s(W_{t_i}^{(i)})$  for every  
 340 character candidate in the set  $\mathbf{M}_{\text{obj}}$  from the *MultiBox* object  
 341 detector. We used a threshold of 10% and retain only those char-  
 342 acter candidates which have a relative saliency score greater than  
 343 this threshold. These candidates are next used as seed-points for  
 344 tracking. This threshold was initially decided based on qualita-  
 345 tive observation. We then conducted additional experiments to  
 346 assess the effect of this threshold parameter as described in the  
 347 Section III-C. The resulting set of *salient* character candidates is  
 348 denoted as  $\mathbf{M}_{\text{sal}} = \{(K_{t_i}^{(i)}, W_{t_i}^{(i)}) | i \in [1, N_{\text{sal}}]\}$ , where  $N_{\text{sal}}$   
 349 is the total number of objects deemed salient after this step  
 350 with  $|\mathbf{M}_{\text{sal}}| \leq |\mathbf{M}_{\text{obj}}|$  where  $|\cdot|$  indicates the cardinality of  
 351 the set.

352 2) *Deformable Object Tracking*: An important property of  
 353 animated characters is their appearance on screen for up to a few  
 354 seconds depending on the context. We utilized this property by  
 355 performing a single-target visual tracking of the salient character  
 356 candidates. Since animated characters are mostly deformable  
 357 bodies, the rigidity assumption that most tracking algorithms  
 358 employ in their motion models (for review, see [35]) does not  
 359 hold. We employ a deformable object tracking algorithm [36]  
 360 which does not impose rigidity assumptions on the object-of-  
 361 interest while tracking.

362 This method first builds a static-appearance model of the ob-  
 363 ject by clustering the key-points into sets of *inliers* (for the  
 364 object body) and *outliers* (for the background) using a dissimi-  
 365 larity measure that quantifies the correspondences between key-  
 366 points. The dissimilarity measure is estimated by computing the  
 367 distance between the initial set of corresponding key-points and  
 368 the transformed version. The model is then adaptively updated in  
 369 time by propagating only the *inlier* correspondences by estimat-  
 370 ing the optical flow of the key-points. The degree of tolerance  
 371 towards the deformation of the object is factored into the model  
 372 by setting a parameter in the tracking algorithm which ensures  
 373 that the cluster of inlier points are spatially localized. We used  
 374 the BRISK [37] features for key-point detection and the param-  
 375 eters were set according to [36] after histogram equalization of  
 376 the images.

377 Tracking every object from the set of salient character candi-  
 378 dates for the full length of the movie is computationally expen-  
 379 sive and may lead to accumulated tracking errors. Hence, we  
 380 performed *local-tracking* in a serial and progressive fashion as  
 381 described in **Algorithm 1**. We refer to the first *candidate frame*  
 382 and the corresponding bounding box for the object of each track  
 383 as a *seed-point*. *Local-tracking* substantially reduced the num-

---

**Algorithm 1: Local Tracking**


---

**Input:** Set of salient character candidates:

$$\mathbf{M}_{\text{sal}} = \{(K_{t_i}^{(i)}, W_{t_i}^{(i)}) | i \in [1, N_{\text{sal}}]\}; \text{ movie, } \mathbf{V}$$

**Output:** Set of track *seed-points*;

$$\mathbf{M}_{\text{trk}} = \{(K_{t_i}^{(i)}, W_{t_i}^{(i)}) | i \in [1, N_{\text{sal}}]\} \text{ and} \\ \text{corresponding track duration } T_i$$

**Parameters:** Track duration threshold:  $\tau$

**while** (  $\mathbf{V}$  open ) **do**

$\mathbf{M}_{\text{trk}} = \{\}$

**while** (  $\mathbf{M}_{\text{sal}} \neq \emptyset$  ) **do**

        Begin tracking at the earliest time frame, i.e.,

$$\{(K_{t_j}^{(j)}, W_{t_j}^{(j)})\} \leftarrow \min_{\forall t_j: j \leq |\mathbf{M}_{\text{sal}}|} \{\mathbf{M}_{\text{sal}}\}$$

        Object tracking lost at  $t_k \geq t_j$

        Track duration,  $T_j \leftarrow t_k - t_j$

**if** (  $T_j > \tau$  ) **then**

            Update tracked seed-points

$$\mathbf{M}_{\text{trk}} \leftarrow \mathbf{M}_{\text{trk}} \cup \{(K_{t_j}^{(j)}, W_{t_j}^{(j)})\}$$

            Prune character candidates

$$\mathbf{M}_{\text{sal}} \leftarrow \mathbf{M}_{\text{sal}} \cap \{(K_{t_m}^{(m)}, W_{t_m}^{(m)}) | \forall m > k\}$$

**else**

$$\mathbf{M}_{\text{sal}} \leftarrow \mathbf{M}_{\text{sal}} \cap \{(K_{t_m}^{(m)}, W_{t_m}^{(m)}) | \forall m > j\}$$

**end**

**end**

$$N_{\text{trk}} = |\mathbf{M}_{\text{trk}}|$$

**end**

---

384 ber of character candidates by eliminating objects that were  
 385 successfully tracked in consecutive frames. As we performed  
 386 single-target visual tracking, this process may also exclude other  
 387 characters that co-occur within a given track. However, since  
 388 *prominent* characters occur quite frequently in a movie, the is-  
 389 sue of losing certain characters was not significantly noted. The  
 390 duration of time for which an object is tracked is used as a  
 391 threshold for retaining objects. We refer to this as the *track du-  
 392 ration threshold*,  $\tau$  and initially set to one frame. This would  
 393 only eliminate the transient and/or spurious object detections.  
 394 Additional experiments varying the  $\tau$  parameter are conducted  
 395 as discussed later. We denote the set of character candidates  
 396 returned after tracking as  $\mathbf{M}_{\text{trk}}$  with  $|\mathbf{M}_{\text{trk}}| = N_{\text{trk}}$  such  
 397 that  $N_{\text{trk}} \leq N_{\text{sal}} \leq N_{\text{obj}}$ . The number of character candidates  
 398 obtained after pruning at each step as a percentage of the initial  
 399 number of input frames is shown in Table II.

### C. Exemplars for Character Representation

400 The character candidates chosen thus far may be redundant  
 401 to some extent, and may contain multiple images with varying  
 402 view-point or segments of the same object. In order to group  
 403 similar objects together, we pose this as an unsupervised clus-  
 404 tering problem with an unknown number of clusters. A suitable  
 405 approach to represent such data is to identify a smaller set of  
 406 samples, referred to as *exemplars*. We use affinity propagation  
 407 (AP) clustering [38] to obtain exemplars which constitute the  
 408 final *character dictionary* for a given movie. AP clustering is  
 409 well suited for this problem because it is deterministic, achieves  
 410

TABLE I  
DETAILS OF THE EVALUATION DATASET

ID	Movie (US Release year)	Duration(mins)	Prominent Characters <sup>†</sup>	Production Studio	Grossing (in \$ millions)
V1	Cars 2 (2011)	107	10 (3)	Pixar	191
V2	Free Birds (2011)	91	11 (4)	Reel FX Creative	55
V3	Frozen (2013)	102	9 (4)	Walt Disney	400
V4	How to Train your Dragon 2 (2014)	102	12 (4)	DreamWorks	177
V5	Shrek Forever After (2010)	93	9 (5)	DreamWorks	238
V6	Tangled (2010)	100	9 (4)	Walt Disney	200
V7	The Lego Movie (2014)	101	12 (3)	Warner Animation	257
V8	Toy Story 3 (2010)	103	18 (9)	Pixar	415

<sup>†</sup> ( ) indicates number of minor characters.

TABLE II  
PERCENTAGE OF INITIAL NUMBER OF OBJECTS AFTER EACH STEP OF PRUNING  
ON THE EVALUATION DATASET

Movie ID	$N_{ini}$	$N_{obj}(\%)$	$N_{sal}(\%)*$	$N_{trk}(\%)+$
V1	15395	19.88	16.99	5.61
V2	13102	14.08	12.50	5.01
V3	14676	9.36	6.83	2.56
V4	14676	9.25	6.32	3.17
V5	13406	10.61	8.06	3.32
V6	14372	9.42	8.22	3.05
V7	14460	9.37	6.96	2.92
V8	14748	11.80	9.79	3.79

\*relative saliency threshold = 10%. + track duration threshold = 1 frame.

a lower clustering error compared to other clustering methods such as k-means [39] and does not require a predetermined number of clusters.

We used the *ImageNet* model proposed in [40] to extract features to cluster the character candidates. Several previous works (e.g., [41]) have shown that feature representations from fully-connected layers in a CNN generalize well for various image recognition tasks. Specifically, we use a 4096-dimensional feature from the second fully connected layer, “FC7” from the ImageNet model which was trained with ILSVRC-2012 [29] competition data.

Because the FC7 features are sparse, we use cosine distance to compute a pairwise similarity matrix,  $\mathbf{S}_{ij}$  between the feature vectors,  $\{\mathbf{v}_i\}$

$$\mathbf{S}_{ij} = \frac{\mathbf{v}_i \mathbf{v}_j^T}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|} \forall i, j \in [1, N_{trk}]. \quad (2)$$

The appearance of most characters is somewhat homogeneous (except for variations in pose and deformation) throughout a movie in terms of shape, color or attire of the character. Leveraging this observation, we also used GIST descriptors [42] for clustering. GIST features provide a low dimensional representation that describes the prominent spatial structure in an image. GIST features have been used for clustering tasks such as scene clustering (e.g., [43]) with some success. We obtained a 960-dimensional GIST descriptor for the character candidates using *pyleargist*<sup>4</sup> package in Python. We then computed negative Euclidean distance between all the candidates from a movie

to form a similarity matrix for clustering. Additionally, we also evaluated the clustering performance of GIST and FC7 features.

We used the AP algorithm proposed in [44] to cluster the similarity matrices obtained from the character candidates. The goal of AP clustering is to choose a character candidate  $j$  to be the exemplar of the  $i$ th candidate. Define *responsibility*  $r(i, j)$ : degree of support that the candidate  $j$  should be the exemplar of  $i$  and *availability*  $a(i, j)$ : degree of support by which the candidate  $i$  should choose  $j$  to be its exemplar. Initialize  $r(i, j), a(i, j) = 0; \forall i, j$  and update responsibility and availability as below:

$$r(i, j) \leftarrow \mathbf{S}_{ij} - \max_{k:k \neq j} (a(k, i) + \mathbf{S}_{ik}) \quad (3)$$

$$a(j, j) \leftarrow \sum_{k:k \neq j} \max[0, r(k, j)] \quad (4)$$

$$a(j, i) \leftarrow \min(0, r(j, j) + \sum_{k:k \notin \{j, i\}} \max[0, r(k, j)]). \quad (5)$$

Introduce a damping factor,  $\lambda \in [0, 1)$  to account for numerical oscillations over iterations in time  $t$

$$r(j, i)_t \leftarrow (1 - \lambda)r(j, i)_t + \lambda r(j, i)_{t-1} \quad (6)$$

$$a(j, i)_t \leftarrow (1 - \lambda)a(j, i)_t + \lambda a(j, i)_{t-1}. \quad (7)$$

Pick  $j$  to be an exemplar of  $i$  if

$$\arg \max_j (r(i, j) + a(j, i)). \quad (8)$$

We set the damping factor,  $\lambda$  which controls the update of  $r(i, j)$  and  $a(i, j)$  in each step to 0.5. Changing this parameter had no effect on the exemplars we obtain. Let  $N_{xmp}$  be the total number of exemplars returned.

AP clustering works well with animation movies since the appearance of most characters (e.g., attire) is consistent within a given movie and the features we used for clustering can capture these attributes. An additional benefit of using AP clustering is that the number of exemplars (i.e., the size of character dictionary) need not be pre-specified. On the other hand, we risk *over-clustering*, i.e., a single character may be represented by multiple exemplars since the features we use are generic and not designed to capture variation in scale, orientation or view-point of a character. This was evident when we performed a second pass of AP clustering on the exemplars obtained here and failed to cluster the *perceptually identical* characters together. In

<sup>4</sup><https://pypi.python.org/pypi/pyleargist>

466 order to penalize for over-clustering, we define an *over-*  
 467 *clustering index* in our performance evaluation measures as  
 468 described in Section III-C.

### 469 III. EXPERIMENTS

470 The problem of identifying character dictionaries for anima-  
 471 tion movies addressed in this paper is unique. Due to the lack  
 472 of existing performance evaluation frameworks for this task,  
 473 we first created a *reference character dictionary* (movie-cast)  
 474 for each movie in our database. We then used these reference  
 475 character dictionaries as ground truth to evaluate the character  
 476 dictionaries output by the proposed method. These reference  
 477 character dictionaries have been made publicly available as a  
 478 part of the SAIL-AMDb<sup>5</sup> along with outputs used for our sys-  
 479 tem evaluation.

#### 480 A. Evaluation Database

481 Our animation movie database consisted of a total of forty-six  
 482 movies produced between 2010–2014. Of the forty-six movies  
 483 available, we chose eight top-grossing movies to evaluate the  
 484 performance of our method in greater detail and to determine  
 485 the best parameter choices for *relative saliency threshold* and  
 486 the *track duration threshold*. The year of release, duration, pro-  
 487 duction company and size of the reference character dictionary  
 488 are shown in Table I. For brevity, we refer to these movies as  
 489 V1–V8.

490 These eight movies were chosen to test the generalizability  
 491 of the proposed system. They represent a diverse set of charac-  
 492 ters in terms of design and composition produced by prominent  
 493 animation studios. These movies include instances of human or  
 494 human-like characters (V3, V5, V6), non-human but anthropo-  
 495 morphic (V3, V5), toy-like (V7, V8) and animals (V2, V4, V5).  
 496 All movies (except V6) include at least one instance of a char-  
 497 acter which is abstract in design. The dataset includes movies  
 498 with varying degrees of illumination, background/environment  
 499 and motion of the characters. For example, V1, V6 and V8 have  
 500 overall higher illumination compared to V3, V4 and V5. The  
 501 movies V1 and V4 have faster moving characters (e.g., drag-  
 502 ons and cars) compared to the others. Quantitative analyses to  
 503 evaluate the diversity of this dataset (e.g., variation in color, il-  
 504 lumination or other characteristics) are beyond the scope of this  
 505 paper (and an objective of our future work).

506 As described in Section II-C, the character dictionary out-  
 507 put by the proposed system for each movie are the exem-  
 508 plars identified by AP clustering. The character candidates  
 509 on which the clustering is performed are obtained by opti-  
 510 mizing two system parameters using a grid search: relative  
 511 saliency threshold and track duration threshold. The settings  
 512 used for the two parameters are  $R_s(X) = \{0, 10, 20, 50, 80, 90\}$   
 513 and  $\tau = \{1, 12, 24, 48, 120\}$ . The values for  $\tau$  (in frames) corre-  
 514 spond to the least possible value (one frame), and approximately  
 515 0.5 s, 1 s, 2 s and 5 s of the movie duration respectively.<sup>6</sup>

#### B. Reference Character Dictionaries

517 We borrow the same definitions for a character as described  
 518 in [45] and [46] to create a movie-specific *reference character*  
 519 *dictionary*. All named characters (speaking and non-speaking)  
 520 displayed on-screen were included. Similar to [46], we first used  
 521 the set of prominent characters as listed by a leading online box-  
 522 office reporting service.<sup>7</sup> The designation of a *minor character*  
 523 available in this resource was retained. This list however, does  
 524 not include non-speaking characters (e.g., dragons). Hence, if a  
 525 character was given a specific name in the movie (as opposed  
 526 to generic names such as a *Spanish ambassador*), we included  
 527 them in the reference. For each of these characters, we obtained  
 528 a representative full-body image from the movie posters or DVD  
 529 covers available online. If the said character was absent in these  
 530 sources, a representative image was manually obtained from  
 531 the internet. The number of prominent characters including the  
 532 number of minor characters are listed in Table I. For annotation  
 533 purposes, all characters in the reference dictionaries are assigned  
 534 a unique ID to preserve character anonymity.

535 We use annotations from Mechanical Turk workers (MTurk; a  
 536 crowdsourcing platform by Amazon Web Services) to compare  
 537 the *proposed* and *reference* character dictionaries. As discussed  
 538 in Section II-C, the exemplars in the proposed dictionaries may  
 539 vary from the representative image used to construct the refer-  
 540 ence. Hence, by using MTurk, we leverage the human perceptual  
 541 ability to match the exemplars with the items in the reference.  
 542 The annotators are instructed to consider an exemplar to be a  
 543 match if 1) it is identifiable regardless to variation in scale, illu-  
 544 mination, orientation or viewpoint or 2) an identifiable segment  
 545 of the reference character is present in the exemplar or 3) if the  
 546 exemplar consists of the said reference character. The annotators  
 547 indicate a match with a unique ID available for every character  
 548 in the *reference*. Furthermore, if an exemplar consists of mul-  
 549 tiple reference characters, the annotators are instructed to list  
 550 all the relevant IDs. Three different annotations were acquired  
 551 for each of the exemplars from unique annotators. In order to  
 552 check for possible confounding factors, additional information  
 553 on whether the annotator had watched the movie prior to anno-  
 554 tating was also collected.

555 We performed an inter-rater reliability analysis to ensure that  
 556 the MTurk annotations were reliable. Since we obtained more  
 557 than two annotations, inter-rater agreement (more specifically,  
 558 inter-annotation agreement) was quantified using Krippendorff’s  
 559 alpha [47] for each movie. The categorical values that were  
 560 used to compute this measure were the unique IDs assigned  
 561 to each character from the reference. Krippendorff’s alpha was  
 562 high for the eight movies used in our system evaluation with  
 563 mean/standard deviation of  $\alpha = 0.81 \pm 0.05$  indicating strong  
 564 agreement. Across all forty-six movies, Krippendorff’s Alpha  
 565 was similarly high (0.82). Furthermore, no difference in agree-  
 566 ment was observed between the set of annotations performed  
 567 by workers who had watched the movie and those who had  
 568 not. Following high agreement, we obtained a single annota-  
 569 tion per exemplar by performing simple majority voting on the

<sup>5</sup><https://goo.gl/WbESbz>

<sup>6</sup>Frame rate for all movies in the dataset was 23.98 fps

<sup>7</sup>[www.boxofficemojo.com](http://www.boxofficemojo.com)

570 three annotations. Three-way ties were resolved with random  
571 assignment.

### 572 C. Performance Evaluation

573 The performance of our method for different experiments was  
574 quantified by comparing the reference character dictionaries  
575 with the output dictionaries from the proposed method. We  
576 refer to the set of exemplars in the proposed dictionary that  
577 were successfully matched to a character in the reference as the  
578 *relevant exemplars* and the remaining as, the *noisy exemplars*.  
579 As described earlier, multiple exemplars can represent a single  
580 character. Therefore, we examine the unique set of character IDs  
581 in the proposed dictionary (*matched characters*) and those never  
582 identified (*missed characters*). Following this, we compute three  
583 measures; precision,  $P$ , recall,  $R$  and F1 score,  $F_1$  as follows:

$$P = \frac{|\{\text{relevant exemplars}\}|}{|\{\text{relevant exemplars}\} \cup \{\text{noisy exemplars}\}|} \quad (9)$$

$$R = \frac{|\{\text{matched characters}\}|}{|\{\text{matched characters}\} \cup \{\text{missed characters}\}|} \quad (10)$$

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}. \quad (11)$$

584 Additionally, we define *over-clustering index* as a measure to  
585 quantify the extent to which multiple exemplars per character  
586 appear in our character dictionaries. In other words, the extent to  
587 which we *over-cluster* the relevant characters. Over-clustering  
588 index for a movie is computed as the median of number of ex-  
589 emplars per character in the set of the relevant exemplars. Since  
590 this metric is defined only over the set of relevant exemplars, it  
591 is independent of precision. It is bounded below by 1 (one exam-  
592 plar per character) and bounded above by  $N_{xmp}$  (all exemplars  
593 represent just one character).

594 In order to compare the clustering performance of GIST and  
595 FC7 features, we measure the *purity* of clustering as described  
596 in [48]. We assign each cluster to the most frequently occurring  
597 character in that cluster. Then, we measure purity by count-  
598 ing the total number of correctly assigned characters, across all  
599 clusters and dividing by the total number of candidates clustered  
600 ( $N_{trk}$ ) as below:

$$\text{purity} = \frac{1}{N_{trk}} \sum_k \max_j |\omega_k \cap c_j| \quad (12)$$

601 where  $\text{purity} \in [0, 1]$ ,  $\Omega = \{\omega_1, \omega_2, \dots, \omega_k\}$  is the set of all clusters  
602 and  $\mathcal{C} = \{c_1, c_2, \dots, c_j\}$  is the set of all relevant exemplars.

603 By our definition of precision (9), a lower value would indicate  
604 that character candidates which are not listed in the refer-  
605 ence were identified as exemplars. These *noisy exemplars* could  
606 either be a result of minor characters not being listed in the refer-  
607 ence or background objects being identified as exemplars. Sim-  
608 ilarly, a high recall (10) would reflect the ability to identify all  
609 the prominent characters at least once. Ideally, recall = 1.0 and  
610 over-clustering index = 1 would indicate that every character  
611 in the reference was detected by exactly one relevant exemplar.  
612 Higher values of the over-clustering index reflect on the failure

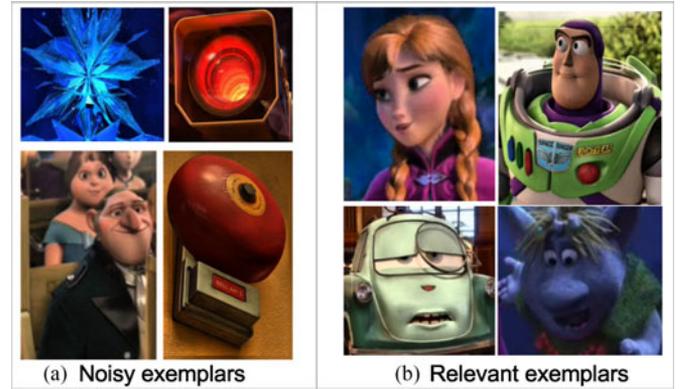


Fig. 5. Examples of noisy and relevant exemplars.

to cluster similar character candidates. This is likely a conse- 613  
614 quence of the features not being invariant to the orientation,  
615 view-point or scale of the character candidates. Complementary  
616 to precision, recall and F1 score which measure the performance  
617 of clustering with respect to a reference, purity (12) measures  
618 the extent to which clusters belonged to a single character, thus  
619 evaluating the features (FC7 versus GIST) used for clustering.

620 The F1 score, precision and recall measures for all eight  
621 movies are averaged for each experiment to determine the best  
622 choice of relative saliency threshold and track duration thresh-  
623 old. These optimal parameters were used to obtain character  
624 dictionaries for the remaining thirty-eight movies in our evalu-  
625 ation dataset.

## 626 IV. RESULTS AND DISCUSSION

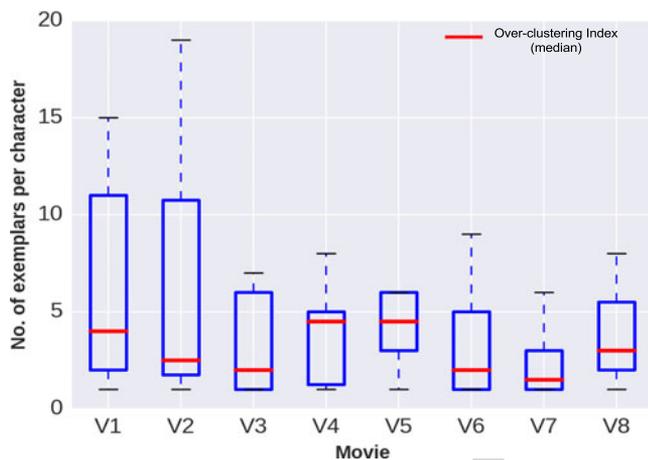
627 A few examples of the *relevant* and *noisy* exemplars from  
628 the proposed character dictionaries are shown in Fig. 5. As de-  
629 scribed earlier, exemplars are categorized as relevant or noisy  
630 based on a reference dictionary constructed for each movie. One  
631 source of noisy exemplars is how we construct these reference  
632 dictionaries. Since the reference consists of only the prominent  
633 characters, it may result in some minor characters being catego-  
634 rized as noisy (See bottom-left image in Fig. 5(a)).

635 The second source of noisy exemplars is the training data used  
636 for the *MultiBox* object detector which comprised only of natural  
637 images. Characters which belong to object classes that the DNN  
638 was trained on tend to get detected more often and consistently  
639 (e.g., traffic lights, bell). The subsequent steps in our method  
640 that use relative saliency score and local-tracking attempt to  
641 eliminate some of these noisy exemplars. However, depending  
642 on the frequency of occurrence or saliency of the character  
643 candidates, they may not always be successfully pruned. Table II  
644 shows the percentage of the input frames pruned at each step.  
645 The proposed character dictionaries for three movies; V1, V2  
646 and V3 are shown in Fig. 10–12 in Appendix B.

647 The precision, recall and F1 score measures that we used to  
648 quantify the performance of our method are shown in Fig. 6.  
649 The relative saliency threshold and track duration threshold  
650 were chosen corresponding to the best F1 score (highlighted  
651 in Fig. 6(a)). These measures were averaged across the eight  
652 movies for each setting of two parameters, relative saliency



Fig. 6. Average (a) F1 score, (b) Precision and (c) Recall for all experiments.

Fig. 7. Distribution of number of exemplars per character in each movie for  $R_s(X) = 10\%$  and  $\tau = 1$ .

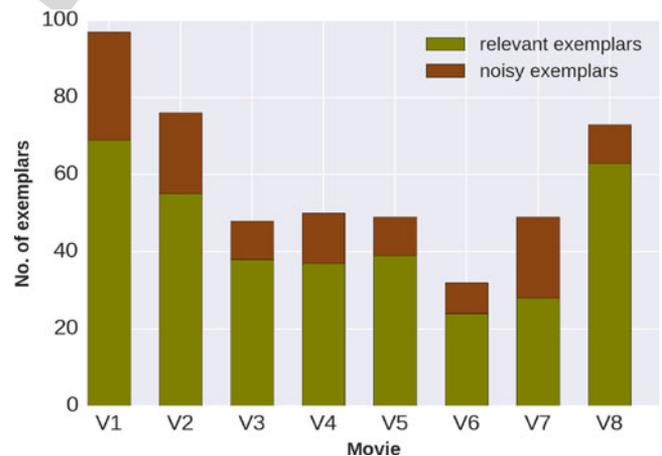
653 threshold,  $R_s(X)$  and track duration threshold,  $\tau$ . Overall, recall is high (over 80% for  $\tau = 1$  and  $R_s(X) = 10\%$ ) which  
 654 indicates that our proposed character dictionaries were able to identify most of the characters in the reference at least once. Precision ranges between 70% and 90% indicating that less than  
 655 one-third of exemplars in our proposed dictionaries are noisy.  
 656

657 We note that the recall measure defined here has to be interpreted alongside over-clustering index; a metric that captures  
 658 the extent to which multiple exemplars represent a single reference character. The distribution of number of relevant exemplars  
 659 per character for the eight movies is shown in Fig. 7. The median number of exemplars per character, i.e., the over-clustering  
 660 index is less than 5 for all the eight movies. As described in Section III-C, this measure lies between 1 and the number of  
 661 exemplars. Here, the number of exemplars range between 35 and 95 (with  $R_s(X) = 10\%$ ;  $\tau = 1$ ) but the over-clustering  
 662 index is less than 5 which reflects on the effective performance of the affinity propagation (AP) algorithm used for clustering.  
 663

664 Additionally, we compared the F1 score and purity of clustering for the eight movies, in order to evaluate the features used  
 665 in clustering, as shown in Table III. Although the F1 scores (computed by comparing the exemplars to the reference) were  
 666  
 667  
 668  
 669  
 670

TABLE III  
F1 SCORE AND PURITY FOR FC7 AND GIST FEATURES USED IN CLUSTERING

Movie ID	FC7 features		GIST descriptors	
	F1 score	Purity	F1 score	Purity
V1	0.691	0.708	0.713	0.414
V2	0.773	0.651	0.769	0.345
V3	0.825	0.842	0.821	0.304
V4	0.764	0.598	0.693	0.322
V5	0.532	0.712	0.653	0.408
V6	0.740	0.677	0.732	0.398
V7	0.732	0.693	0.743	0.438
V8	0.752	0.745	0.799	0.392
Average:	<b>0.726</b>	<b>0.703</b>	<b>0.740</b>	<b>0.378</b>

Fig. 8. Number of relevant and noisy exemplars for each movie with  $R_s(X) = 10\%$  and  $\tau = 1$ .

675 similar between the two descriptors, the clustering purity using  
 676 FC7 features was significantly higher (paired t-test,  $p \ll 0.01$  to  
 677 reject  $H_0 : \mu_0 \leq \mu_1$ ) than that of GIST descriptors. This indicates that FC7 features yield less noisy and more homogeneous  
 678 clusters from AP clustering. Furthermore, FC7 features perform  
 679 better for clustering than GIST features, perhaps because *ImageNet*  
 680 was trained to classify objects robust to variation in the  
 681 background or view-point and occlusions, whereas GIST  
 682 descriptors capture the holistic shape information in an image.  
 683

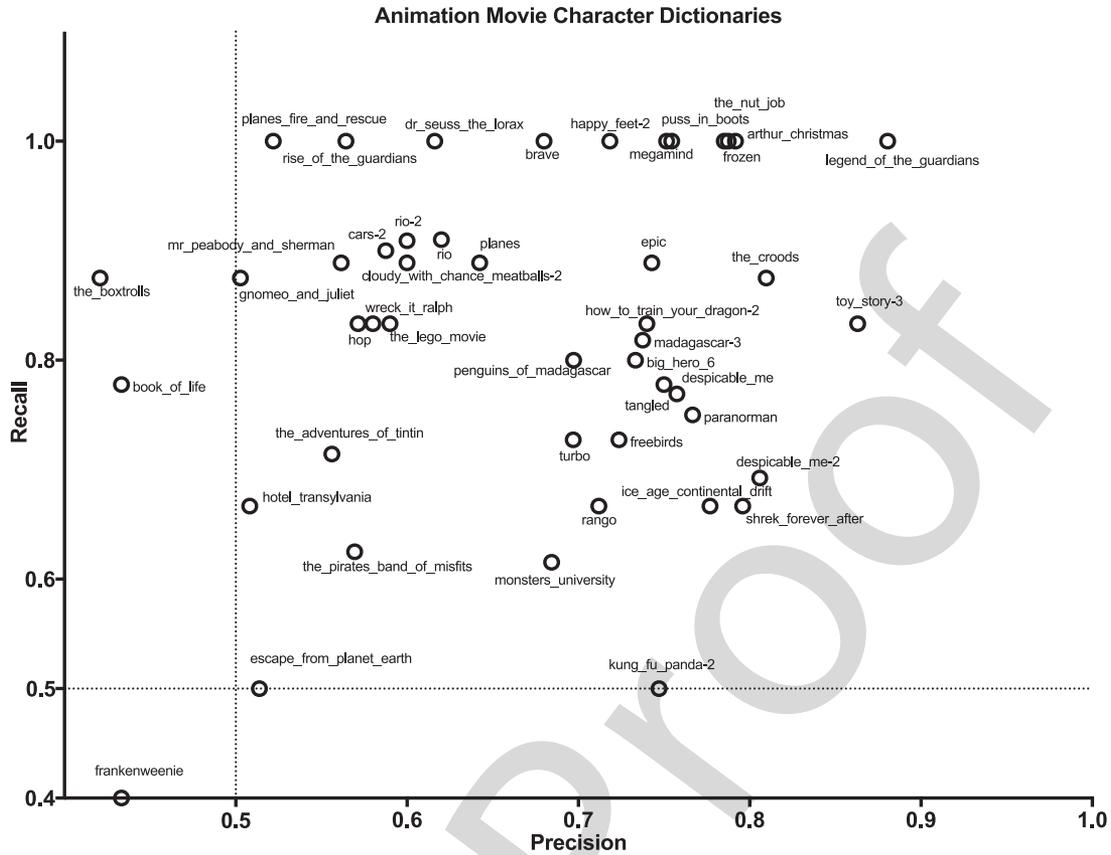


Fig. 9. Precision versus recall for forty-six movies.

684 As shown in Fig. 6(c), recall drops with an increase in  $\tau$  as  
 685 expected. Since, by increasing  $\tau$  we retain only those character  
 686 candidates which remain longer on-screen and do not always  
 687 co-occur with other salient objects. This results in excluding  
 688 some relevant exemplars. In contrast, an increase in precision  
 689 (See Fig. 6(b)) is noticed on increasing  $\tau$  since a few noisy  
 690 character candidates that are infrequent get pruned successfully.  
 691 Relative saliency threshold had the desired effect on the system  
 692 output i.e., increasing  $R_s(X)$  results in an increase in precision.  
 693 However, these gains in precision by increasing  $R_s(X)$  beyond  
 694 10% were not substantial.

695 In order to determine a good choice of the system param-  
 696 eters, we examine F1 score for different combinations of  $R_s(X)$   
 697 and  $\tau$  as shown in Fig. 6(a).  $R_s(X) = 10\%$  and  $\tau = 1$  would  
 698 be the best choice of settings. For these settings, the num-  
 699 ber of relevant and noisy exemplars for the eight movies are  
 700 shown in Fig. 8. It is interesting to note that movies V1 and  
 701 V2 have relatively larger character dictionaries and a higher  
 702 range of number of exemplars per character (See Appendix  
 703 Fig. 10–11). All the characters in the movies, V1 and V2 are  
 704 similar to cars and birds in appearance. The results at a glance  
 705 show that all instances of these characters in different scenes  
 706 were detected in these movies (which include the minor char-  
 707 acters and different appearances of the same character with re-  
 708 spect to view-point). This is likely because both cars and birds  
 709 are among the object classes in ILSVRC-2014 data used to train  
 710 *MultiBox*.

Character dictionaries for the remaining thirty-eight movies  
 were obtained with the choice of  $R_s(X)$  and  $\tau$  determined  
 above. The range of precision was 0.45–0.89 (mean/standard  
 deviation:  $0.66 \pm 0.12$ ) and recall: 0.42–1.0 ( $0.83 \pm 0.16$ ).  
 The range of over-clustering index was 1.5–6.0 ( $3.5 \pm 1.5$ ).  
 See Fig. 9 (Appendix A) for precision and recall measures of  
 all the forty six movies in our dataset. Further error analysis  
 considering different aspects of all the movies (e.g., character  
 design, color, illumination) is warranted and will be a part of  
 our future work.

We note that the movie *Frankenweenie* (2013) which was  
 produced in black and white has the lowest precision and re-  
 call in our dataset. This indicates that color rendering is an  
 important factor since the DNNs we employ were trained with  
 RGB images. The movies, *Boxtrolls* (2014) and *The Book of*  
*Life* (2014) both have a low precision and high recall indicat-  
 ing a larger number of noisy exemplars. On analyzing the errors  
 in these samples, we observed that the local tracking method  
 pruned approximately 42% of the initial character candidates  
 (c.f. the average percentage of candidates pruned by local track-  
 ing for the rest of the movies was 62.23%). This is likely be-  
 cause these movies, unlike the others in the dataset use a *rapid-fire*  
 film editing style which includes fast-action scene cuts and rapidly  
 changing backgrounds which are not ideally suited for visual  
 object tracking.

On the other hand, movies like *Kung Fu Panda 2* (2011) and  
*Escape from Planet Earth* (2013) yield high precision and low



Fig. 10. Character dictionary of the movie Frozen: precision = 0.81, recall = 1.0; over-clustering index=2.

738 recall. This is likely because these movies feature only a small  
 739 number of prominent characters and a larger number of unnamed  
 740 characters which are not included in the reference dictionaries  
 741 that we created. As expected, movies that feature distinct lifelike  
 742 animals or humans, generally performed the best. For example,  
 743 the movie Legend of Guardians (2010) featured only birds and  
 744 The Nut Job (2014) featured animals – both animals and birds  
 745 are included in the set of object categories of the ILSVRC  
 746 datasets.

## 747 V. CONCLUSIONS AND FUTURE WORK

748 In this paper, we proposed an unsupervised method to auto-  
 749 matically create a dictionary of characters from an animation  
 750 movie. We evaluated our method on a set of eight movies with  
 751 diverse character styles and demonstrated high precision and  
 752 recall on a dataset of forty-six movies. We also showed that  
 753 the proposed method generalizes for animation movies at scale.  
 754 These character dictionaries can serve as a powerful tool for  
 755 character labeling to delineate aspects of *who appeared*, *when*  
 756 and for *how long* in a movie (video diarization). We believe  
 757 that our efforts can lay a foundation to provide an impetus for  
 758 multimedia research endeavors specifically involving animated  
 759 media content.

760 One of the drawbacks of the proposed method is that we use  
 761 an object detector that was trained with natural images. We plan  
 762 to address this issue using transfer learning to adapt the existing  
 763 models to specialize the network for detecting characters from  
 764 animation movies. The relevant and noisy exemplars that we  
 765 annotated for the system evaluation can potentially be used for  
 766 these methods. Our future work would also include using the  
 767 relevant exemplars and associated cluster members as a single  
 768 unit to facilitate robust video diarization of animation movies.

### APPENDIX A

769 PRECISION AND RECALL OF OUR PROPOSED METHOD FOR THE  
 770 FORTY-SIX MOVIES

771 The Fig. 9 plots precision vs. recall for all the movies in  
 772 our dataset. The relative saliency threshold and track duration



Fig. 11. Character dictionary of the movie Free Birds: precision = 0.72, recall = 0.72; over-clustering index = 2.5.

threshold was set to 10% and one frame respectively (tuned on 773  
 a subset of 8 movies as described in Section IV). 774

### APPENDIX B EXAMPLES

Fig. 10–12 illustrate the proposed character dictionaries for 775  
 three movies with the settings of relative saliency threshold = 777



Fig. 12. Character dictionary of the movie Cars-2: precision=0.61, recall = 0.9; over-clustering index = 4.

778 10% and track duration threshold = 1 frame. The exemplars  
779 here are arranged in no particular order to maintain their aspect  
780 ratios.

781

## REFERENCES

782 [1] P. P. Mohanta, S. K. Saha, and B. Chanda, "A model-based shot bound-  
783 ary detection technique using frame transition parameters," *IEEE Trans.*  
784 *Multimedia*, vol. 14, no. 1, pp. 223–233, Feb. 2012.  
785 [2] C. Liu, D. Wang, J. Zhu, and B. Zhang, "Learning a contextual multi-  
786 thread model for movie/TV scene segmentation," *IEEE Trans. Multimedia*,  
787 vol. 15, no. 4, pp. 884–897, Jun. 2013.  
788 [3] B. W. Chen, J. C. Wang, and J. F. Wang, "A novel video summarization  
789 based on mining the story-structure and semantic relations among concept  
790 entities," *IEEE Trans. Multimedia*, vol. 11, no. 2, pp. 295–312, Feb. 2009.  
791 [4] Y. Li, S.-H. Lee, C.-H. Yeh, and C. C. J. Kuo, "Techniques for movie  
792 content analysis and skimming: tutorial and overview on video abstraction  
793 techniques," *IEEE Signal Process. Mag.*, vol. 23, no. 2, pp. 79–89,  
794 Mar. 2006.  
795 [5] K. Kurzhals, M. John, F. Heimerl, P. Kuznecov, and D. Weiskopf, "Visual  
796 movie analytics," *IEEE Trans. Multimedia*, vol. 18, no. 11, pp. 2149–2160,  
797 Nov. 2016.  
798 [6] C. Y. Weng, W. T. Chu, and J. L. Wu, "Rolenet: Movie analysis from the  
799 perspective of social networks," *IEEE Trans. Multimedia*, vol. 11, no. 2,  
800 pp. 256–271, Feb. 2009.

[7] T. Guha, C.-W. Huang, N. Kumar, Y. Zhu, and S. S. Narayanan, "Gender  
801 representation in cinematic content: A multimodal approach," in *Proc.*  
802 *ACM Int. Conf. Multimodal Interaction*, 2015, pp. 31–34.  
803 [8] A. Fitzgibbon and A. Zisserman, "On affine invariant clustering and auto-  
804 matic cast listing in movies," *Proc. 7th Eur. Conf. Comput. Vis.*, 2002,  
805 pp. 304–320.  
806 [9] F. Vallet, S. Essid, and J. Carriave, "A multimodal approach to speaker  
807 diarization on TV talk-shows," *IEEE Trans. Multimedia*, vol. 15, no. 3,  
808 pp. 509–520, Apr. 2013.  
809 [10] J. Sivic, M. Everingham, and A. Zisserman, "Who are you?"—Learning  
810 person specific classifiers from video," in *Proc. IEEE Conf. Comput. Vis.*  
811 *Pattern Recogn.*, 2009, pp. 1145–1152.  
812 [11] M. Everingham, J. Sivic, and A. Zisserman, "Hello! My name is...  
813 Buffy"—automatic naming of characters in TV video," in *Proc. Brit.*  
814 *Mach. Vis. Conf.*, 2006.  
815 [12] M. Everingham, J. Sivic, and A. Zisserman, "Taking the bite out of auto-  
816 matic naming of characters in TV video," *Image Vis. Comput.*, vol. 27,  
817 no. 5, pp. 545–559, Apr. 2009.  
818 [13] "Box office history for digital animation" [Online]. Available:  
819 [http://www.the-numbers.com/market/production-method/digital-](http://www.the-numbers.com/market/production-method/digital-animation)  
820 [animation](http://www.the-numbers.com/market/production-method/digital-animation)  
821 Q2  
822 [14] Z. Aghbari, K. Kaneko, and A. Makinouchi, "Content-trajectory approach  
823 for searching video databases," *IEEE Trans. Multimedia*, vol. 5, no. 4,  
824 pp. 516–531, Dec. 2003.  
825 [15] B. Ionescu, V. Buzuloiu, P. Lambert, and D. Coquin, "Improved cut detec-  
826 tion for the segmentation of animation movies," in *Proc. IEEE Int. Conf.*  
827 *Acoust., Speech, Signal Process.*, May 2006, vol. 2, pp. II–II.  
828 [16] B. Ionescu, D. Coquin, P. Lambert, and V. Buzuloiu, "Fuzzy color-based  
829 approach for understanding animated movies content in the indexing task,"  
830 *J. Image Video Process.*, vol. 2008, pp. 8:1–8:17, Jan. 2008.  
831 [17] L. Ott, P. Lambert, B. Ionescu, and D. Coquin, "Animation movie abstrac-  
832 tion: Key frame adaptive selection based on color histogram filtering,"  
833 in *Proc. 14th Int. Conf. Image Anal. Process.*, 2007, pp. 206–211.  
834 [18] B. Ionescu, P. Lambert, D. Coquin, L. Ott, and V. Buzuloiu, "Animation  
835 movies trailer computation," in *Proc. 14th ACM Int. Conf. Multimedia*,  
836 2006, pp. 631–634.  
837 [19] K. Takayama, H. Johan, and T. Nishita, "Face detection and face recog-  
838 nition of cartoon characters using feature extraction," in *Proc. Image,*  
839 *Electron. Visual Comput. Workshop*, 2012, p. 48.  
840 [20] S. Wang, J. Zhang, T. X. Han, and Z. Miao, "Sketch-based image retrieval  
841 through hypothesis-driven object boundary selection with HLR descriptor,"  
842 *IEEE Trans. Multimedia*, vol. 17, no. 7, pp. 1045–1057, Jul. 2015.  
843 [21] *The Illusion of Life: Disney Animation*. New York, NY, USA: Abbeville  
844 Press, 1981.  
845 [22] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan,  
846 "Object detection with discriminatively trained part-based models," *IEEE*  
847 *Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645,  
848 Sep. 2010.  
849 [23] Y. Tang, X. Wang, E. Dellandrea, and L. Chen, "Weakly supervised learn-  
850 ing of deformable part-based models for object detection via region pro-  
851 posals," *IEEE Trans. Multimedia*, vol. 19, no. 2, pp. 393–407, Feb. 2017.  
852 [24] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once:  
853 Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis.*  
854 *Pattern Recogn.*, 2016, pp. 779–788.  
855 [25] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov, "Scalable object  
856 detection using deep neural networks," in *Proc. IEEE Conf. Comput. Vis.*  
857 *Pattern Recogn.*, 2014, pp. 2155–2162.  
858 [26] C. Szegedy, S. E. Reed, D. Erhan, and D. Anguelov, "Scalable, high-  
859 quality object detection," arXiv:1412.1441, 2014.  
860 [27] K. H. Lee and J. N. Hwang, "On-road pedestrian tracking across multiple  
861 driving recorders," *IEEE Trans. Multimedia*, vol. 17, no. 9, pp. 1429–1438,  
862 Sep. 2015.  
863 [28] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. IEEE Conf.*  
864 *Comput. Vis. Pattern Recogn.*, Jun. 2015.  
865 [29] O. Russakovsky *et al.*, "Imagenet large scale visual recognition challenge,"  
866 *Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, Dec. 2015.  
867 [30] P. Jacard, "The distribution of the flora in the alpine zone," *New Phytol-*  
868 *ogist*, vol. 11, pp. 37–50, 1912.  
869 [31] E. Rahtu and J. Heikkilä, "A simple and efficient saliency detector for  
870 background subtraction," in *Proc. 12th Int. Conf. Comput. Vis. Workshops*,  
871 Sep. 2009, pp. 1137–1144.  
872 [32] Y. Hu, X. Xie, W.-Y. Ma, L.-T. Chia, and D. Rajan, "Salient region de-  
873 tection using weighted feature maps based on the human visual attention  
874 model," in *Proc. 5th Pacific Rim Conf. Multimedia Adv. Multimedia Inf.*  
875 *Process.*, 2005, pp. 993–1000.

- 876 [33] V. Mahadevan and N. Vasconcelos, "Background subtraction in highly  
877 dynamic scenes," in *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, Jun.  
878 2008, pp. 1–6.
- 879 [34] E. Rahtu, J. Kannala, M. Salo, and J. Heikkilä, "Segmenting salient objects  
880 from images and videos," *Proc. 11th Eur. Conf. Comput. Vis. Comput. Vis.*,  
881 2010, pp. 366–379.
- 882 [35] M. Kristan *et al.*, "The visual object tracking VOT2015 challenge results,"  
883 in *Proc. IEEE Int. Conf. Comput. Vis. Workshops*, Dec. 2015.
- 884 [36] G. Nebel and R. Pflugfelder, "Clustering of static-adaptive correspondences  
885 for deformable object tracking," *Comput. Vis. Pattern Recogn.*,  
886 Jun. 2015, pp. 2784–2791.
- 887 [37] S. Leutenegger, M. Chli, and R. Y. Siegwart, "Brisk: Binary robust invariant  
888 scalable keypoints," in *Proc. Int. Conf. Comput. Vis.*, Nov. 2011, pp.  
889 2548–2555.
- 890 [38] D. Dueck and B. J. Frey, "Non-metric affinity propagation for unsuper-  
891 vised image categorization," in *Proc. IEEE 11th Int. Conf. Comput. Vis.*,  
892 Oct. 2007, pp. 1–8.
- 893 [39] A. K. Jain, M. N. Murty, and P. J. Flynn, "Data clustering: A review,"  
894 *ACM Comput. Surv.*, vol. 31, no. 3, pp. 264–323, Sep. 1999.
- 895 [40] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification  
896 with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Pro-  
897 cess. Syst.*, 2012, pp. 1106–1114.
- 898 [41] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN features  
899 off-the-shelf: An astounding baseline for recognition," in *Proc. IEEE Conf.  
900 Comput. Vis. Pattern Recogn. Workshops*, 2014, pp. 512–519.
- 901 [42] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic  
902 representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, no. 3,  
903 pp. 145–175, May 2001.
- 904 [43] M. Merler, B. Huang, L. Xie, G. Hua, and A. Natsev, "Semantic model  
905 vectors for complex video event recognition," *IEEE Trans. Multimedia*,  
906 vol. 14, no. 1, pp. 88–101, Feb. 2012.
- 907 [44] B. J. Frey and D. Dueck, "Clustering by passing messages between data  
908 points," *Science*, vol. 315, pp. 972–977, 2007.
- 909 [45] "Comprehensive Annenberg Report on Diversity in Entertainment," 2016.
- 910 [46] S. L. Smith *et al.*, "Media, diversity, & social change initiative," 2016.
- 911 [47] K. Krippendorff, "Reliability in content analysis," *Human Commun. Res.*,  
912 vol. 30, no. 3, pp. 411–433, 2004.
- 913 [48] D. M. Christopher, R. Prabhakar, and S. Hinrich, "Introduction to infor-  
914 mation retrieval," *Introduction Inf. Retrieval*, vol. 151, p. 177, 2008.



**Krishna Somandepalli** received the Master's degree from University of California at Santa Barbara, CA, USA, in electrical and computer engineering. He received the Bachelor's degree in electronics and communication engineering from University Visvesvaraya College of Engineering, Bangalore, India. Following his Master's degree, he worked as an Assistant Research Scientist at NYU Langone Medical Center, New York, NY, USA. His research interests include multimodal analysis with image and signal data. He is currently working toward the Ph.D. degree in the

Signal Analysis and Interpretation Laboratory Group, Department of Electrical Engineering, University of Southern California, Los Angeles, CA, USA.



**Naveen Kumar** received the Ph.D. degree in electrical engineering from the USC Viterbi School of Engineering, where he was a member of the Media Informatics and Content Analysis Group, Signal Analysis and Interpretation Lab. He received the B.Tech. degree in instrumentation engineering from the Indian Institute of Technology, Kharagpur, India, in 2009. He currently works at the Sony PlayStation R&D in San Mateo, CA, USA. His broad research interests include machine learning and signal processing for speech, multimedia and multimodal

applications.



**Tanaya Guha** received the Ph.D. degree in electrical and computer engineering from the University of British Columbia (UBC), Vancouver, BC, Canada. She is an Assistant Professor in the Department of Electrical Engineering, Indian Institute of Technology (IIT) Kanpur. Prior to joining IIT Kanpur, she was a Postdoctoral Fellow at the Signal Analysis and Interpretation Lab (SAIL), University of Southern California.

Her current research interests include social and affective computing, multimedia analysis, and multimodal signal processing. Dr. Guha received Mensa Canada Woodhams Memorial Scholarship, Google Anita Borg Scholarship, and Amazon Grace Hopper celebration scholarship.



**Shrikanth S. Narayanan** (SM'88–M'95–SM'02–F'09) is the Niki & C. L. Max Nikias Chair in Engineering at the University of Southern California (USC), Los Angeles, CA, USA, and holds appointments as a Professor of Electrical Engineering, Computer Science, Linguistics, Psychology, Neuroscience and Pediatrics, the Research Director of the Information Science Institute, and as the Founding Director of the Ming Hsieh Institute. Prior to USC, he was with AT&T Bell Labs and AT&T Research from 1995–2000. At USC, he directs the Signal Analysis

and Interpretation Laboratory (SAIL). His research focuses on human-centered signal and information processing and systems modeling with an interdisciplinary emphasis on speech, audio, language, multimodal and biomedical problems and applications with direct societal relevance. Prof. Narayanan is a Fellow of the National Academy of Inventors, the Acoustical Society of America, the International Speech Communication Association (ISCA) and the American Association for the Advancement of Science (AAAS), and a member of Tau Beta Pi, Phi Kappa Phi, and Eta Kappa Nu. He is the Editor-in-Chief for *IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING*, an Editor for the *Computer Speech and Language Journal* and an Associate Editor for the *APSIPA Transactions On Signal And Information Processing*. He was also previously an Associate Editor of the *IEEE TRANSACTIONS OF SPEECH AND AUDIO PROCESSING* (2000–2004), *IEEE SIGNAL PROCESSING MAGAZINE* (2005–2008), *IEEE TRANSACTIONS ON MULTIMEDIA*, (2008–2011), the *IEEE TRANSACTIONS ON SIGNAL AND INFORMATION PROCESSING OVER NETWORKS* (2014–2015), *IEEE TRANSACTIONS ON AFFECTIVE COMPUTING* (2010–2016), and the *Journal of the Acoustical Society of America* (2009–2017). He received several honors including Best Transactions Paper awards from the IEEE Signal Processing Society in 2005 (with A. Potamianos) and in 2009 (with C. M. Lee) and selection as an IEEE Signal Processing Society Distinguished Lecturer for 2010–2011 and ISCA Distinguished Lecturer for 2015–2016. Papers coauthored with his students have received awards including the 2014 Ten-year Technical Impact Award from ACM ICMI and at several conferences. He has published more than 750 papers and has been granted 17 U.S. patents.

915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941

942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992

- 994 Q1. Author: Please check whether the affiliation of authors is okay as set.  
995 Q2. Author: Please provide year information in Ref. [13].  
996 Q3. Author: Please update Ref. [26], if already published in prints.  
997 Q4. Author: Please provide page range in Refs. [19], [28], and [35].  
998 Q5. Author: Please provide full bibliographic details in Ref. [46].

IEEE Proof

# Unsupervised Discovery of Character Dictionaries in Animation Movies

Krishna Somandepalli, *Member, IEEE*, Naveen Kumar, *Member, IEEE*, Tanaya Guha, *Member, IEEE*, and Shrikanth S. Narayanan, *Fellow, IEEE*

**Abstract**—Automatic content analysis of animation movies can enable an objective understanding of character (actor) representations and their portrayals. It can also help illuminate potential markers of unconscious biases and their impact. However, multimedia analysis of movie content has predominantly focused on live-action features. A dearth of multimedia research in this field is because of the complexity and heterogeneity in the design of animated characters—an extremely challenging problem to be generalized by a single method or model. In this paper, we address the problem of automatically discovering characters in animation movies as a first step toward automatic character labeling in these media. Movie-specific character dictionaries can act as a powerful first step for subsequent content analysis at scale. We propose an unsupervised approach which requires no prior information about the characters in a movie. We first use a deep neural network-based object detector that is trained on natural images to identify a set of initial character candidates. These candidates are further pruned using saliency constraints and visual object tracking. A character dictionary per movie is then generated from exemplars obtained by clustering these candidates. We are able to identify both anthropomorphic and nonanthropomorphic characters in a dataset of 46 animation movies with varying composition and character design. Our results indicate high precision and recall of the automatically detected characters compared to human-annotated ground truth, demonstrating the generalizability of our approach.

**Index Terms**—Animation movies, deep neural networks, object tracking, saliency, unsupervised clustering, video diarization.

## I. INTRODUCTION

**A**UTOMATIC analysis of movie content is of growing interest in the multimedia research community. One of the driving factors for this research is the large number of movies that are produced, disseminated and consumed annually. Besides being of entertainment value, movies often have an effect

Manuscript received November 28, 2016; revised July 26, 2017; accepted August 4, 2017. This work is based upon work supported by the National Science Foundation. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Xilin Chen. (*Corresponding author: Krishna Somandepalli.*)

K. Somandepalli, N. Kumar, and T. Guha is with the Department of Electrical Engineering, University of Southern California, Los Angeles, CA 90047 USA (e-mail: krishna.somandepalli@gmail.com; knaveen87@gmail.com; tanayaguha@gmail.com).

S. S. Narayanan is with the Signal & Image Processing Institute, University of Southern California, Los Angeles, CA 90089 USA (e-mail: shri@sipi.usc.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TMM.2017.2745712

on certain social and economic aspects, as well as have a global reach and audience.

Researchers have addressed movie content analysis with different objectives and outlooks. Such efforts are often based on efficient indexing and organization of the media content for easy user navigation. They include shot boundary detection for movie segmentation [1], [2], video summarization [3] and abstraction [4]. The study in [2] builds a generative model that incorporates contextual information in order to reorganize interleaved shots into multiple plot threads. Approaches such as in [5] combines the aspects of video summarization, i.e., *who*, *what*, *where* and *when* for a semantic understanding of the movie content and structure. *RoleNet* proposed in [6] examines the movie content from a social network analysis perspective of the movie character roles rather than using audiovisual features. In general, movie content is a rich source of data that includes audio, video and text (dialogs) that enables such multimodal analysis.

Complementary to the aforementioned studies which attempt to achieve a high-level understanding of movies, efforts for a fine-grained (frame level or scene level statistics) analysis of video content have also been emerging. One such application is to quantify the amount of time a character appears on screen in a movie. The study in [7] examined these aspects with respect to gender revealing skewed distributions for the onscreen time of female characters. In order to advance from gender-level statistics to character-level statistics, person identification or character labeling is a crucial step in this direction. We refer to this problem as automatic video diarization – partitioning the video stream into actor-homogeneous segments, i.e., *who appeared*, *when* and for *how long*. Character labeling in live-action TV and movies has been achieved with modest success in [8]–[12]. This is typically performed by clustering the detected faces (e.g., [8]) or by multimodal approaches (e.g., [9], [10]) that model audio and subtitles or scripts alongside the detected faces from video.

It is important to note that all these studies exclusively focus on live-action TV and do not generalize to animated media content. Digital animation movies have contributed to over 10% of the box office market shares in the past decade [13]. Multimedia research in this domain is extremely scarce and technology developed for live-action TV content fails for animated content. Human face detection is the crux of character labeling methods for live action TV. Since human-characters can be uniquely identified by their faces, this method performs adequately well. But, such methods developed for human faces do not work

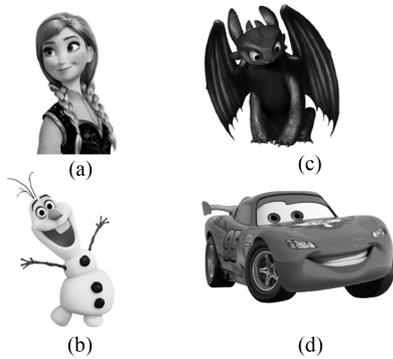


Fig. 1. Examples illustrating the heterogeneity of animated characters. (a): human-like (Frozen) (b): anthropomorphic (Frozen) (c) and (d): abstract (How to Train your Dragon, and Cars).

84 for the digital animation genre. Animated characters, though  
 85 mostly anthropomorphic (having human characteristics) are not  
 86 always human-like in appearance. They can be fictional animals,  
 87 inanimate objects or abstract in design (see Fig. 1 for a  
 88 few examples).

89 A major obstacle for automating content analysis of animated  
 90 media is the lack of a model that generalizes across different  
 91 characters with varying composition and design. This task be-  
 92 comes extremely complex given that all the characters even  
 93 within a single movie may not share the same structural charac-  
 94 teristics (e.g., human-like and non-human characters from the  
 95 same movie—Fig. 1(a) and (b) from the movie Frozen).

96 In the context of video diarization, when the characters that  
 97 appear on screen are generally not known a priori, a key step  
 98 is to provide a list of characters that form the *who appeared*  
 99 component of the system. We refer to such a list of characters  
 100 specific to each movie as a *character dictionary*. The automatic  
 101 discovery of these character dictionaries is the primary objective  
 102 in this paper. Our overarching goal is to engineer a model for  
 103 animation movie video diarization. With the proposed character  
 104 dictionaries, animation character labeling may be achieved by  
 105 techniques such as [14] that can retrieve frames and shots given  
 106 an object of interest.

107 In content analysis of animated media, researchers have thus  
 108 far focused on problems such as cut detection [15], color-based  
 109 video categorization [16] and movie abstraction [17], [18]. One  
 110 method proposed in [19] performs human-like face detection  
 111 from cartoon images using skin-segmentation techniques. Con-  
 112 sidering the variation in texture, color and shape of animated  
 113 characters in general (as illustrated in Fig. 1), these methods do  
 114 not generalize well. To the best of our knowledge, no work to  
 115 date has specifically addressed the problem of automatic dis-  
 116 covery of characters from animated media in a scalable manner.

117 In contrast to live-action movies, animation movies are com-  
 118 pletely artist generated. Sketches of the characters are designed  
 119 by the artists or the animators, generally referred to as *model*  
 120 *sheets* from which character-specific 3D models are generated.  
 121 Sketch based image retrieval systems such as [20] can be used  
 122 to achieve video diarization when model sheets are available.  
 123 However, model sheets are copyrighted material and mostly  
 124 owned by the animation studio which produced the movie. As

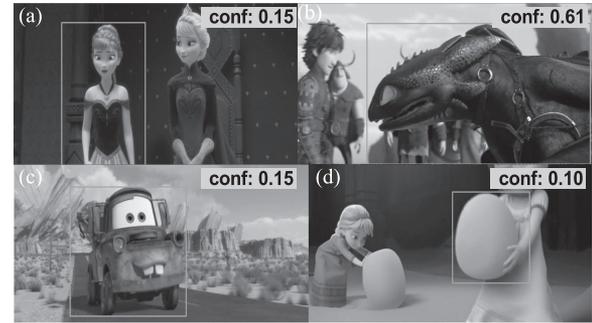


Fig. 2. Character candidates chosen by the Multibox object detector. Conf. indicates the confidence score of the network for the detected object.

such, they are not publicly available and approaches which are  
 based on model sheets will not be scalable for all movies.

125  
 126  
 127 In 1981, Frank Thomas and Ollie Johnston published *The*  
 128 *Illusion of Life* [21]; it outlines a set of twelve basic principles  
 129 of animation. Animators have been using this as a cookbook  
 130 for designing characters in order for the viewers to appreci-  
 131 ate “animation” over mere “movement”. While most of these  
 132 principles aid animators in adding semantic or artistic value  
 133 (e.g., anticipation, exaggeration), a few can be exploited in a  
 134 computer vision context (e.g., *Solid Drawing*: drawing volume  
 135 solidity and illusion of three dimensions; *Staging*: Distinctive  
 136 color, depth of field and positioning in the frame to highlight  
 137 the character). Defining an *animated character* in a complete  
 138 sense would involve delineating abstract concepts such as life  
 139 (or sentience even) from movie content. In this paper, we only  
 140 analyze the video stream from animation movies and leverage  
 141 some of the aforementioned principles of animation as proxies  
 142 to identify the characters.

143 At the outset, we pose our problem as an object detection task  
 144 where any object can be a possible *character candidate*. Ani-  
 145 mation movie frames are comparable with natural photographic  
 146 images, especially in their similarities of depth of field and the  
 147 character presentation in a frame. Additionally, we assume no  
 148 prior models with respect to shape, size, color, or texture for  
 149 these candidates in order for the proposed system to generalize.

150 A few prominent examples of state-of-the-art object detec-  
 151 tion systems include discriminatively trained deformable parts-  
 152 based model (DPM, [22], [23]) and deep neural network (DNN)  
 153 models such as [24]–[26], both of which are supervised and  
 154 trained over a predefined set of object classes. DPMs need  
 155 a carefully designed part-decomposition model of an object  
 156 which makes it unsuitable given the heterogeneity of characters  
 157 within just a single movie. In contrast, DNN-based methods such  
 158 as [24] can detect objects in real-time and outperform DPMs.  
 159 Specifically, DNN models that are saliency-inspired in design  
 160 [25] are of interest for our problem statement. Although super-  
 161 vised with a finite set of object classes, they have been shown to  
 162 detect objects in a *class-agnostic* manner [26] i.e., detect classes  
 163 of objects not used for training the model.

164 Movies in general, portray only a handful of *prominent* char-  
 165 acters. They are more likely to appear frequently in order for  
 166 the viewer to easily comprehend the content and the plot of

the movie. Additionally in movies, characters or the objects-of-interest tend to remain on screen for up to a few seconds depending on the situation. Visual object tracking can be used as an effective method to segment characters locally in time. Several previous works have used tracking as a means to automatically detect a class of objects (e.g., pedestrians, [27]). Object tracking algorithms can be error-prone in a movie video environment because of object deformation, background clutter, changes in illumination, occlusion and lack of a stationary backgrounds. However, visual tracking can minimize the number of detected objects to be considered by accounting for minor deformation or linear motion of the object. Furthermore, tracking also provides time information that can be used for diarization subsequently. For example, in [11], supervisory information available on a profile face is used to learn the appearance of a frontal face from faces tracked in TV series. A reasonable assumption in describing animated character is that the prominent characters are not transient when presented on-screen and appear frequently in the movie. In our method, we use this aspect of character presentation in movies to select character candidates. As a result, the character dictionaries consist of only the frequently occurring characters.

In this paper, we propose a novel approach to automatically discover characters that appear in an animation movie. Our proposed method is unsupervised in the sense that we do not train any aspect of our system with data from animated media content. Furthermore, we use no specific knowledge of the animation style or the physical attributes of the animated characters, thereby ensuring that our system can scale and generalize through the whole spectrum of animation movie content.

The rest of the paper is organized as follows: Section II describes the proposed system for selecting character candidates from an animation movie. In Section III, we present the experiments performed and the creation of an evaluation database. Section IV contains the experimental results and final considerations followed by conclusions and future work in Section V.

## II. METHODS

In this section, we first introduce the different systems that we use to identify and prune the detected objects to obtain a set of possible character candidates. We then use a clustering approach to identify character exemplars that constitute the final character dictionary. The overview of the proposed system is shown in the Fig. 3.

Our animation movie database consisted of forty-six movies, for which we annotated their prominent characters. We then conducted a detailed performance evaluation on eight animation movies which were chosen to represent varying degrees of heterogeneity in character design and composition. The movie-cast data from forty-six movies used for our system evaluation and the output from our system has been released as part of the SAIL Animation Movie character Database (SAIL-AMDb).<sup>1</sup> We have also made the code publicly available.<sup>2</sup>

<sup>1</sup><https://github.com/usc-sail/mica-animation/wiki>

<sup>2</sup><https://github.com/usc-sail/mica-animation>

### A. Coarse Detection of Character Candidates

Animated characters are often designed to have the appearance of a 3D object and characterized by shallow focus where the image plane of the character is in focus while the rest of the frame is out of focus [21]. In other words, they are the salient objects in a given frame. Capitalizing on this, we define a *character candidate* as any object that can be detected by a general-purpose object detector.

We use a pre-trained deep neural network (DNN) called *MultiBox* [25], [26], designed for object detection. Our preliminary experiments with other region proposal networks such as [24] yielded similar results. We chose *MultiBox* since our motivation for using an object detector was only to generate an initial set of potential character candidates.

*MultiBox* is a convolutional neural network (CNN) with an inception-style architecture [28] trained with the full 200-category object detection challenge data set from ImageNet Large Scale Visual Recognition Challenge 2014 (ILSVRC-2014) [29]. This model generates multiple bounding boxes and an associated confidence score that quantifies the network's confidence of each box containing an object. The model has been shown to perform object localization in a *class-agnostic* manner and achieve state-of-the-art performance in object detection tasks [25]. Furthermore, since the network is tailored towards the localization problem, it achieves a scalable representation of multiple salient objects in an image. These features make this model uniquely suitable for our problem. It is important to note that this model is trained with natural images of distinct object classes. Although the authors in [25] have shown that the model generalizes over unseen classes, here we apply the pre-trained DNN for images sampled from animation movies. We refer to this discrepancy as *DNN training bias*. This results in detecting objects that are not characters in a movie (e.g., traffic-light, chair). We refer to such objects as *noisy objects*.

In order to reduce the computational time, we downsample a movie (originally encoded at 23.98 fps) by one frame every 0.42 s (every 10th frame). The resulting frames are input to *MultiBox* [25] to obtain all possible bounding boxes for each image. The confidence score that is returned with each of these boxes was originally optimized in the DNN to match the ground truth object boxes from natural images.

Because of the aforementioned *DNN training bias*, we generally observed lower range of confidence scores for objects detected that were animated characters. We chose to retain objects with a confidence score greater than 0.1. In order to determine this threshold, we randomly sampled 100,000 frames from the movie *Frozen* (2013) in our movie database. We first assumed to have at most five possibly overlapping objects of interest in one frame and obtained the confidence scores for the five most confident objects in each frame. We then examined the distribution of the confidence scores for all the objects detected. We set the confidence threshold to 75th percentile of the distribution of confidence scores which is equal to 0.1002, thus retaining all objects with confidence score greater than 0.1. We apply this confidence threshold for all the movies in our database. A few examples of objects detected and their confidence scores returned by the network are shown in Fig. 2.

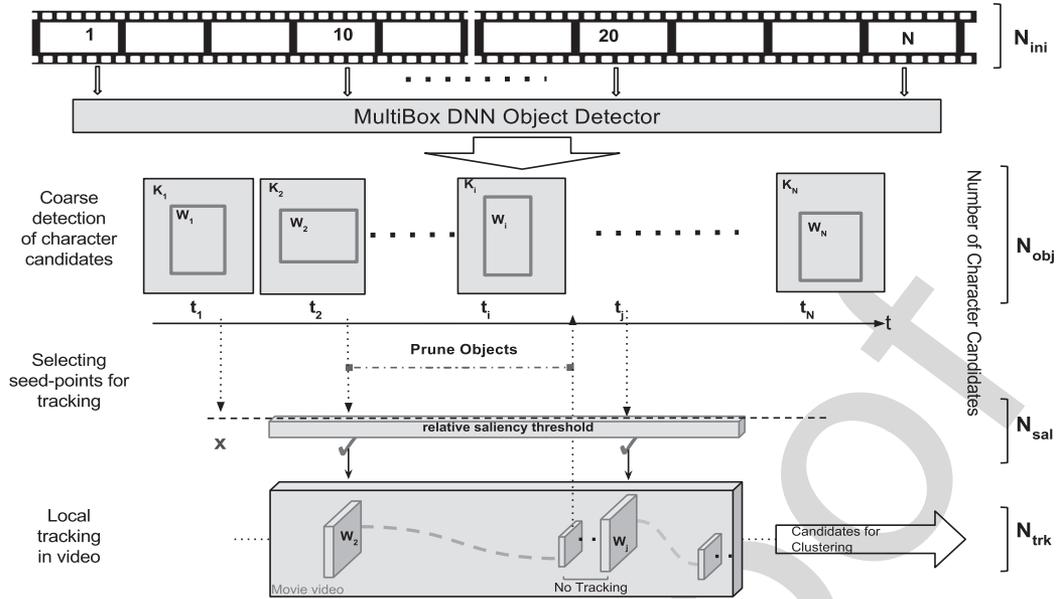


Fig. 3. Schematic diagram of the proposed method.

277 We also computed the area of each bounding box of an object  
 278 relative to the image frame and excluded objects in bounding  
 279 boxes with an area less than 1% or greater than 99% of the entire  
 280 frame. This ensures that very small objects and holistic scenes  
 281 are excluded as character candidates. When multiple objects  
 282 were detected in a single frame, we pruned them to obtain at  
 283 most one object per frame following the approach in [25]. We  
 284 performed non-maximum-suppression with a Jaccard similarity  
 285 [30] threshold of 0.5 and, chose the object with the maximum  
 286 area in that frame. We identified only a single object per frame in  
 287 order to simplify the subsequent step of single-target visual object  
 288 tracking. We refer to a chosen frame containing a character  
 289 candidate as a *candidate frame*.

290 A schematic of the proposed approach is illustrated in Fig. 3.  
 291 Let  $N_{ini}$  be the initial number of images (movie frames) input to *MultiBox*  
 292 and  $N_{obj}$  be the number of character candidates chosen. We denote the candidate frame  $K$   
 293 and the bounding box  $W$  enclosing the object as a set  $M_{obj} = \{(K_{t_i}^{(i)},$   
 294  $W_{t_i}^{(i)}) | i \in [1, N_{obj}]\}$  and  $t_i$  refers to the time (or frame number) in the movie at which the object  $i$  occurs. Qualitative analyses showed that this step captures most of the characters in an animation movie at least once (e.g., images shown in Fig. 2(a)–(c)). However, this set also contains redundant and noisy objects which include non-characters or background objects (e.g., Fig. 2(d)).

### 302 B. Saliency Constraints and Object Tracking

303 In the next phase of our system, we used the saliency of  
 304 the detected object as a constraint to prune the set of character  
 305 candidates obtained in the previous step. We use this pruned  
 306 set of candidate frames as seed-points for tracking. During  
 307 tracking, we do not distinguish camera motion from object motion,  
 308 thereby ensuring that a sufficient condition for a character



Fig. 4. (a) Example for DNN training bias and saliency constraint; (b) Masked regions showing saliency, here *relative saliency score*  $R_s(W_1) = 9.2\%$ .

candidate is its presence on the screen rather than motion (e.g., 309  
 a talking tree). 310

311 1) *Saliency-Constrained Pruning*: As described in  
 312 Section II-A, the DNN training bias may result in choosing  
 313 objects that, although salient, may not be the characters of  
 314 interest (e.g., detected lamp in a scene with two characters—see  
 315 Fig. 4(a)). To quantify this, we use a saliency measure proposed  
 316 in [31] for the character candidate with respect to the entire  
 317 frame. Unsupervised methods that estimate saliency typically  
 318 use pixel-level features such as color, intensity (e.g., [32])  
 319 or background-detection in dynamic scenes (e.g., [33]). In  
 320 contrast, the measure proposed in [31] estimates saliency of  
 321 local areas (instead of pixel level) in static images and requires  
 322 no training. This method uses a kernel-based approach where  
 323 the size of the window relates to the scale of the target objects.  
 324 The saliency of a pixel inside the window is estimated using the  
 325 conditional probability of that pixel drawn from the distribution  
 326 estimated inside that window versus the distribution of the  
 327 surrounding area.

328 We first converted the RGB images to CIELAB color space  
 329 (because of the perceptual uniformity of the CIE color space<sup>3</sup>)  
 330 to estimate a saliency map for the entire candidate frame by

<sup>3</sup><http://www.brucelindbloom.com>

331 choosing window sizes at different scales as described in [34].  
 332 The resulting saliency maps are binarized by setting values  
 333 greater than 0.7 to 1 as recommended in [34]. An example of  
 334 the saliency map is shown in Fig. 4(b). Let  $A_s(W)$  be the area  
 335 of the salient region contained within a bounding box,  $W$  in an  
 336 image frame  $K$ . We define a *relative saliency score*,  $R_s(W)$  of  
 337 an object enclosed by the box  $W$  as the percentage salient area  
 338 it contributes to the frame,  $K$ :

$$R_s(W) = \frac{A_s(W)}{A_s(K)} \times 100. \quad (1)$$

339 We obtained the *relative saliency score*,  $R_s(W_{t_i}^{(i)})$  for every  
 340 character candidate in the set  $\mathbf{M}_{\text{obj}}$  from the *MultiBox* object  
 341 detector. We used a threshold of 10% and retain only those char-  
 342 acter candidates which have a relative saliency score greater than  
 343 this threshold. These candidates are next used as seed-points for  
 344 tracking. This threshold was initially decided based on qualita-  
 345 tive observation. We then conducted additional experiments to  
 346 assess the effect of this threshold parameter as described in the  
 347 Section III-C. The resulting set of *salient* character candidates is  
 348 denoted as  $\mathbf{M}_{\text{sal}} = \{(K_{t_i}^{(i)}, W_{t_i}^{(i)}) | i \in [1, N_{\text{sal}}]\}$ , where  $N_{\text{sal}}$   
 349 is the total number of objects deemed salient after this step  
 350 with  $|\mathbf{M}_{\text{sal}}| \leq |\mathbf{M}_{\text{obj}}|$  where  $|\cdot|$  indicates the cardinality of  
 351 the set.

352 2) *Deformable Object Tracking*: An important property of  
 353 animated characters is their appearance on screen for up to a few  
 354 seconds depending on the context. We utilized this property by  
 355 performing a single-target visual tracking of the salient character  
 356 candidates. Since animated characters are mostly deformable  
 357 bodies, the rigidity assumption that most tracking algorithms  
 358 employ in their motion models (for review, see [35]) does not  
 359 hold. We employ a deformable object tracking algorithm [36]  
 360 which does not impose rigidity assumptions on the object-of-  
 361 interest while tracking.

362 This method first builds a static-appearance model of the ob-  
 363 ject by clustering the key-points into sets of *inliers* (for the  
 364 object body) and *outliers* (for the background) using a dissimi-  
 365 larity measure that quantifies the correspondences between key-  
 366 points. The dissimilarity measure is estimated by computing the  
 367 distance between the initial set of corresponding key-points and  
 368 the transformed version. The model is then adaptively updated in  
 369 time by propagating only the *inlier* correspondences by estimat-  
 370 ing the optical flow of the key-points. The degree of tolerance  
 371 towards the deformation of the object is factored into the model  
 372 by setting a parameter in the tracking algorithm which ensures  
 373 that the cluster of inlier points are spatially localized. We used  
 374 the BRISK [37] features for key-point detection and the param-  
 375 eters were set according to [36] after histogram equalization of  
 376 the images.

377 Tracking every object from the set of salient character candi-  
 378 dates for the full length of the movie is computationally expen-  
 379 sive and may lead to accumulated tracking errors. Hence, we  
 380 performed *local-tracking* in a serial and progressive fashion as  
 381 described in **Algorithm 1**. We refer to the first *candidate frame*  
 382 and the corresponding bounding box for the object of each track  
 383 as a *seed-point*. *Local-tracking* substantially reduced the num-

---

**Algorithm 1: Local Tracking**


---

**Input:** Set of salient character candidates:

$$\mathbf{M}_{\text{sal}} = \{(K_{t_i}^{(i)}, W_{t_i}^{(i)}) | i \in [1, N_{\text{sal}}]\}; \text{ movie, } \mathbf{V}$$

**Output:** Set of track *seed-points*;

$$\mathbf{M}_{\text{trk}} = \{(K_{t_i}^{(i)}, W_{t_i}^{(i)}) | i \in [1, N_{\text{sal}}]\} \text{ and} \\ \text{corresponding track duration } T_i$$

**Parameters:** Track duration threshold:  $\tau$

**while** (  $\mathbf{V}$  open ) **do**

$\mathbf{M}_{\text{trk}} = \{\}$

**while** (  $\mathbf{M}_{\text{sal}} \neq \emptyset$  ) **do**

        Begin tracking at the earliest time frame, i.e.,

$$\{(K_{t_j}^{(j)}, W_{t_j}^{(j)})\} \leftarrow \min_{\forall t_j: j \leq |\mathbf{M}_{\text{sal}}|} \{\mathbf{M}_{\text{sal}}\}$$

        Object tracking lost at  $t_k \geq t_j$

        Track duration,  $T_j \leftarrow t_k - t_j$

**if** (  $T_j > \tau$  ) **then**

            Update tracked seed-points

$$\mathbf{M}_{\text{trk}} \leftarrow \mathbf{M}_{\text{trk}} \cup \{(K_{t_j}^{(j)}, W_{t_j}^{(j)})\}$$

            Prune character candidates

$$\mathbf{M}_{\text{sal}} \leftarrow \mathbf{M}_{\text{sal}} \cap \{(K_{t_m}^{(m)}, W_{t_m}^{(m)}) | \forall m > k\}$$

**else**

$$\mathbf{M}_{\text{sal}} \leftarrow \mathbf{M}_{\text{sal}} \cap \{(K_{t_m}^{(m)}, W_{t_m}^{(m)}) | \forall m > j\}$$

**end**

**end**

$$N_{\text{trk}} = |\mathbf{M}_{\text{trk}}|$$

**end**

---

384 ber of character candidates by eliminating objects that were  
 385 successfully tracked in consecutive frames. As we performed  
 386 single-target visual tracking, this process may also exclude other  
 387 characters that co-occur within a given track. However, since  
 388 *prominent* characters occur quite frequently in a movie, the is-  
 389 sue of losing certain characters was not significantly noted. The  
 390 duration of time for which an object is tracked is used as a  
 391 threshold for retaining objects. We refer to this as the *track du-  
 392 ration threshold*,  $\tau$  and initially set to one frame. This would  
 393 only eliminate the transient and/or spurious object detections.  
 394 Additional experiments varying the  $\tau$  parameter are conducted  
 395 as discussed later. We denote the set of character candidates  
 396 returned after tracking as  $\mathbf{M}_{\text{trk}}$  with  $|\mathbf{M}_{\text{trk}}| = N_{\text{trk}}$  such  
 397 that  $N_{\text{trk}} \leq N_{\text{sal}} \leq N_{\text{obj}}$ . The number of character candidates  
 398 obtained after pruning at each step as a percentage of the initial  
 399 number of input frames is shown in Table II.

### C. Exemplars for Character Representation

400 The character candidates chosen thus far may be redundant  
 401 to some extent, and may contain multiple images with varying  
 402 view-point or segments of the same object. In order to group  
 403 similar objects together, we pose this as an unsupervised clus-  
 404 tering problem with an unknown number of clusters. A suitable  
 405 approach to represent such data is to identify a smaller set of  
 406 samples, referred to as *exemplars*. We use affinity propagation  
 407 (AP) clustering [38] to obtain exemplars which constitute the  
 408 final *character dictionary* for a given movie. AP clustering is  
 409 well suited for this problem because it is deterministic, achieves  
 410

TABLE I  
DETAILS OF THE EVALUATION DATASET

ID	Movie (US Release year)	Duration(mins)	Prominent Characters <sup>†</sup>	Production Studio	Grossing (in \$ millions)
V1	Cars 2 (2011)	107	10 (3)	Pixar	191
V2	Free Birds (2011)	91	11 (4)	Reel FX Creative	55
V3	Frozen (2013)	102	9 (4)	Walt Disney	400
V4	How to Train your Dragon 2 (2014)	102	12 (4)	DreamWorks	177
V5	Shrek Forever After (2010)	93	9 (5)	DreamWorks	238
V6	Tangled (2010)	100	9 (4)	Walt Disney	200
V7	The Lego Movie (2014)	101	12 (3)	Warner Animation	257
V8	Toy Story 3 (2010)	103	18 (9)	Pixar	415

<sup>†</sup> ( ) indicates number of minor characters.

TABLE II  
PERCENTAGE OF INITIAL NUMBER OF OBJECTS AFTER EACH STEP OF PRUNING  
ON THE EVALUATION DATASET

Movie ID	$N_{ini}$	$N_{obj}(\%)$	$N_{sal}(\%)*$	$N_{trk}(\%)+$
V1	15395	19.88	16.99	5.61
V2	13102	14.08	12.50	5.01
V3	14676	9.36	6.83	2.56
V4	14676	9.25	6.32	3.17
V5	13406	10.61	8.06	3.32
V6	14372	9.42	8.22	3.05
V7	14460	9.37	6.96	2.92
V8	14748	11.80	9.79	3.79

\*relative saliency threshold = 10%. + track duration threshold = 1 frame.

a lower clustering error compared to other clustering methods such as k-means [39] and does not require a predetermined number of clusters.

We used the *ImageNet* model proposed in [40] to extract features to cluster the character candidates. Several previous works (e.g., [41]) have shown that feature representations from fully-connected layers in a CNN generalize well for various image recognition tasks. Specifically, we use a 4096-dimensional feature from the second fully connected layer, “FC7” from the ImageNet model which was trained with ILSVRC-2012 [29] competition data.

Because the FC7 features are sparse, we use cosine distance to compute a pairwise similarity matrix,  $\mathbf{S}_{ij}$  between the feature vectors,  $\{\mathbf{v}_i\}$

$$\mathbf{S}_{ij} = \frac{\mathbf{v}_i \mathbf{v}_j^T}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|} \forall i, j \in [1, N_{trk}]. \quad (2)$$

The appearance of most characters is somewhat homogeneous (except for variations in pose and deformation) throughout a movie in terms of shape, color or attire of the character. Leveraging this observation, we also used GIST descriptors [42] for clustering. GIST features provide a low dimensional representation that describes the prominent spatial structure in an image. GIST features have been used for clustering tasks such as scene clustering (e.g., [43]) with some success. We obtained a 960-dimensional GIST descriptor for the character candidates using *pyleargist*<sup>4</sup> package in Python. We then computed negative Euclidean distance between all the candidates from a movie

to form a similarity matrix for clustering. Additionally, we also evaluated the clustering performance of GIST and FC7 features.

We used the AP algorithm proposed in [44] to cluster the similarity matrices obtained from the character candidates. The goal of AP clustering is to choose a character candidate  $j$  to be the exemplar of the  $i$ th candidate. Define *responsibility*  $r(i, j)$ : degree of support that the candidate  $j$  should be the exemplar of  $i$  and *availability*  $a(i, j)$ : degree of support by which the candidate  $i$  should choose  $j$  to be its exemplar. Initialize  $r(i, j), a(i, j) = 0; \forall i, j$  and update responsibility and availability as below:

$$r(i, j) \leftarrow \mathbf{S}_{ij} - \max_{k:k \neq j} (a(k, i) + \mathbf{S}_{ik}) \quad (3)$$

$$a(j, j) \leftarrow \sum_{k:k \neq j} \max[0, r(k, j)] \quad (4)$$

$$a(j, i) \leftarrow \min(0, r(j, j) + \sum_{k:k \notin \{j, i\}} \max[0, r(k, j)]). \quad (5)$$

Introduce a damping factor,  $\lambda \in [0, 1)$  to account for numerical oscillations over iterations in time  $t$

$$r(j, i)_t \leftarrow (1 - \lambda)r(j, i)_t + \lambda r(j, i)_{t-1} \quad (6)$$

$$a(j, i)_t \leftarrow (1 - \lambda)a(j, i)_t + \lambda a(j, i)_{t-1}. \quad (7)$$

Pick  $j$  to be an exemplar of  $i$  if

$$\arg \max_j (r(i, j) + a(j, i)). \quad (8)$$

We set the damping factor,  $\lambda$  which controls the update of  $r(i, j)$  and  $a(i, j)$  in each step to 0.5. Changing this parameter had no effect on the exemplars we obtain. Let  $N_{xmp}$  be the total number of exemplars returned.

AP clustering works well with animation movies since the appearance of most characters (e.g., attire) is consistent within a given movie and the features we used for clustering can capture these attributes. An additional benefit of using AP clustering is that the number of exemplars (i.e., the size of character dictionary) need not be pre-specified. On the other hand, we risk *over-clustering*, i.e., a single character may be represented by multiple exemplars since the features we use are generic and not designed to capture variation in scale, orientation or view-point of a character. This was evident when we performed a second pass of AP clustering on the exemplars obtained here and failed to cluster the *perceptually identical* characters together. In

<sup>4</sup><https://pypi.python.org/pypi/pyleargist>

466 order to penalize for over-clustering, we define an *over-*  
 467 *clustering index* in our performance evaluation measures as  
 468 described in Section III-C.

### 469 III. EXPERIMENTS

470 The problem of identifying character dictionaries for anima-  
 471 tion movies addressed in this paper is unique. Due to the lack  
 472 of existing performance evaluation frameworks for this task,  
 473 we first created a *reference character dictionary* (movie-cast)  
 474 for each movie in our database. We then used these reference  
 475 character dictionaries as ground truth to evaluate the character  
 476 dictionaries output by the proposed method. These reference  
 477 character dictionaries have been made publicly available as a  
 478 part of the SAIL-AMDb<sup>5</sup> along with outputs used for our sys-  
 479 tem evaluation.

#### 480 A. Evaluation Database

481 Our animation movie database consisted of a total of forty-six  
 482 movies produced between 2010–2014. Of the forty-six movies  
 483 available, we chose eight top-grossing movies to evaluate the  
 484 performance of our method in greater detail and to determine  
 485 the best parameter choices for *relative saliency threshold* and  
 486 the *track duration threshold*. The year of release, duration, pro-  
 487 duction company and size of the reference character dictionary  
 488 are shown in Table I. For brevity, we refer to these movies as  
 489 V1–V8.

490 These eight movies were chosen to test the generalizability  
 491 of the proposed system. They represent a diverse set of charac-  
 492 ters in terms of design and composition produced by prominent  
 493 animation studios. These movies include instances of human or  
 494 human-like characters (V3, V5, V6), non-human but anthropo-  
 495 morphic (V3, V5), toy-like (V7, V8) and animals (V2, V4, V5).  
 496 All movies (except V6) include at least one instance of a char-  
 497 acter which is abstract in design. The dataset includes movies  
 498 with varying degrees of illumination, background/environment  
 499 and motion of the characters. For example, V1, V6 and V8 have  
 500 overall higher illumination compared to V3, V4 and V5. The  
 501 movies V1 and V4 have faster moving characters (e.g., drag-  
 502 ons and cars) compared to the others. Quantitative analyses to  
 503 evaluate the diversity of this dataset (e.g., variation in color, il-  
 504 lumination or other characteristics) are beyond the scope of this  
 505 paper (and an objective of our future work).

506 As described in Section II-C, the character dictionary out-  
 507 put by the proposed system for each movie are the exem-  
 508 plars identified by AP clustering. The character candidates  
 509 on which the clustering is performed are obtained by opti-  
 510 mizing two system parameters using a grid search: relative  
 511 saliency threshold and track duration threshold. The settings  
 512 used for the two parameters are  $R_s(X) = \{0, 10, 20, 50, 80, 90\}$   
 513 and  $\tau = \{1, 12, 24, 48, 120\}$ . The values for  $\tau$  (in frames) corre-  
 514 spond to the least possible value (one frame), and approximately  
 515 0.5 s, 1 s, 2 s and 5 s of the movie duration respectively.<sup>6</sup>

#### B. Reference Character Dictionaries

517 We borrow the same definitions for a character as described  
 518 in [45] and [46] to create a movie-specific *reference character*  
 519 *dictionary*. All named characters (speaking and non-speaking)  
 520 displayed on-screen were included. Similar to [46], we first used  
 521 the set of prominent characters as listed by a leading online box-  
 522 office reporting service.<sup>7</sup> The designation of a *minor character*  
 523 available in this resource was retained. This list however, does  
 524 not include non-speaking characters (e.g., dragons). Hence, if a  
 525 character was given a specific name in the movie (as opposed  
 526 to generic names such as a *Spanish ambassador*), we included  
 527 them in the reference. For each of these characters, we obtained  
 528 a representative full-body image from the movie posters or DVD  
 529 covers available online. If the said character was absent in these  
 530 sources, a representative image was manually obtained from  
 531 the internet. The number of prominent characters including the  
 532 number of minor characters are listed in Table I. For annotation  
 533 purposes, all characters in the reference dictionaries are assigned  
 534 a unique ID to preserve character anonymity.

535 We use annotations from Mechanical Turk workers (MTurk; a  
 536 crowdsourcing platform by Amazon Web Services) to compare  
 537 the *proposed* and *reference* character dictionaries. As discussed  
 538 in Section II-C, the exemplars in the proposed dictionaries may  
 539 vary from the representative image used to construct the refer-  
 540 ence. Hence, by using MTurk, we leverage the human perceptual  
 541 ability to match the exemplars with the items in the reference.  
 542 The annotators are instructed to consider an exemplar to be a  
 543 match if 1) it is identifiable regardless to variation in scale, illu-  
 544 mination, orientation or viewpoint or 2) an identifiable segment  
 545 of the reference character is present in the exemplar or 3) if the  
 546 exemplar consists of the said reference character. The annotators  
 547 indicate a match with a unique ID available for every character  
 548 in the *reference*. Furthermore, if an exemplar consists of mul-  
 549 tiple reference characters, the annotators are instructed to list  
 550 all the relevant IDs. Three different annotations were acquired  
 551 for each of the exemplars from unique annotators. In order to  
 552 check for possible confounding factors, additional information  
 553 on whether the annotator had watched the movie prior to anno-  
 554 tating was also collected.

555 We performed an inter-rater reliability analysis to ensure that  
 556 the MTurk annotations were reliable. Since we obtained more  
 557 than two annotations, inter-rater agreement (more specifically,  
 558 inter-annotation agreement) was quantified using Krippendorff’s  
 559 alpha [47] for each movie. The categorical values that were  
 560 used to compute this measure were the unique IDs assigned  
 561 to each character from the reference. Krippendorff’s alpha was  
 562 high for the eight movies used in our system evaluation with  
 563 mean/standard deviation of  $\alpha = 0.81 \pm 0.05$  indicating strong  
 564 agreement. Across all forty-six movies, Krippendorff’s Alpha  
 565 was similarly high (0.82). Furthermore, no difference in agree-  
 566 ment was observed between the set of annotations performed  
 567 by workers who had watched the movie and those who had  
 568 not. Following high agreement, we obtained a single annota-  
 569 tion per exemplar by performing simple majority voting on the

<sup>5</sup><https://goo.gl/WbESbz>

<sup>6</sup>Frame rate for all movies in the dataset was 23.98 fps

<sup>7</sup>[www.boxofficemojo.com](http://www.boxofficemojo.com)

570 three annotations. Three-way ties were resolved with random  
571 assignment.

### 572 C. Performance Evaluation

573 The performance of our method for different experiments was  
574 quantified by comparing the reference character dictionaries  
575 with the output dictionaries from the proposed method. We  
576 refer to the set of exemplars in the proposed dictionary that  
577 were successfully matched to a character in the reference as the  
578 *relevant exemplars* and the remaining as, the *noisy exemplars*.  
579 As described earlier, multiple exemplars can represent a single  
580 character. Therefore, we examine the unique set of character IDs  
581 in the proposed dictionary (*matched characters*) and those never  
582 identified (*missed characters*). Following this, we compute three  
583 measures; precision,  $P$ , recall,  $R$  and F1 score,  $F_1$  as follows:

$$P = \frac{|\{\text{relevant exemplars}\}|}{|\{\text{relevant exemplars}\} \cup \{\text{noisy exemplars}\}|} \quad (9)$$

$$R = \frac{|\{\text{matched characters}\}|}{|\{\text{matched characters}\} \cup \{\text{missed characters}\}|} \quad (10)$$

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}. \quad (11)$$

584 Additionally, we define *over-clustering index* as a measure to  
585 quantify the extent to which multiple exemplars per character  
586 appear in our character dictionaries. In other words, the extent to  
587 which we *over-cluster* the relevant characters. Over-clustering  
588 index for a movie is computed as the median of number of ex-  
589 emplars per character in the set of the relevant exemplars. Since  
590 this metric is defined only over the set of relevant exemplars, it  
591 is independent of precision. It is bounded below by 1 (one exam-  
592 plar per character) and bounded above by  $N_{xmp}$  (all exemplars  
593 represent just one character).

594 In order to compare the clustering performance of GIST and  
595 FC7 features, we measure the *purity* of clustering as described  
596 in [48]. We assign each cluster to the most frequently occurring  
597 character in that cluster. Then, we measure purity by count-  
598 ing the total number of correctly assigned characters, across all  
599 clusters and dividing by the total number of candidates clustered  
600 ( $N_{trk}$ ) as below:

$$\text{purity} = \frac{1}{N_{trk}} \sum_k \max_j |\omega_k \cap c_j| \quad (12)$$

601 where  $\text{purity} \in [0, 1]$ ,  $\Omega = \{\omega_1, \omega_2, \dots, \omega_k\}$  is the set of all clus-  
602 ters and  $\mathcal{C} = \{c_1, c_2, \dots, c_j\}$  is the set of all relevant exemplars.

603 By our definition of precision (9), a lower value would indi-  
604 cate that character candidates which are not listed in the refer-  
605 ence were identified as exemplars. These *noisy exemplars* could  
606 either be a result of minor characters not being listed in the refer-  
607 ence or background objects being identified as exemplars. Sim-  
608 ilarly, a high recall (10) would reflect the ability to identify all  
609 the prominent characters at least once. Ideally, recall = 1.0 and  
610 over-clustering index = 1 would indicate that every character  
611 in the reference was detected by exactly one relevant exemplar.  
612 Higher values of the over-clustering index reflect on the failure

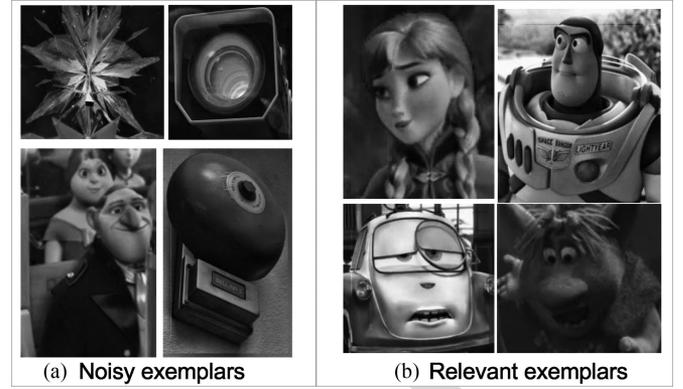


Fig. 5. Examples of noisy and relevant exemplars.

613 to cluster similar character candidates. This is likely a conse-  
614 quence of the features not being invariant to the orientation,  
615 view-point or scale of the character candidates. Complementary  
616 to precision, recall and F1 score which measure the performance  
617 of clustering with respect to a reference, purity (12) measures  
618 the extent to which clusters belonged to a single character, thus  
619 evaluating the features (FC7 versus GIST) used for clustering.

620 The F1 score, precision and recall measures for all eight  
621 movies are averaged for each experiment to determine the best  
622 choice of relative saliency threshold and track duration thresh-  
623 old. These optimal parameters were used to obtain character  
624 dictionaries for the remaining thirty-eight movies in our evalu-  
625 ation dataset.

## 626 IV. RESULTS AND DISCUSSION

627 A few examples of the *relevant* and *noisy* exemplars from  
628 the proposed character dictionaries are shown in Fig. 5. As de-  
629 scribed earlier, exemplars are categorized as relevant or noisy  
630 based on a reference dictionary constructed for each movie. One  
631 source of noisy exemplars is how we construct these reference  
632 dictionaries. Since the reference consists of only the prominent  
633 characters, it may result in some minor characters being catego-  
634 rized as noisy (See bottom-left image in Fig. 5(a)).

635 The second source of noisy exemplars is the training data used  
636 for the *MultiBox* object detector which comprised only of natural  
637 images. Characters which belong to object classes that the DNN  
638 was trained on tend to get detected more often and consistently  
639 (e.g., traffic lights, bell). The subsequent steps in our method  
640 that use relative saliency score and local-tracking attempt to  
641 eliminate some of these noisy exemplars. However, depending  
642 on the frequency of occurrence or saliency of the character  
643 candidates, they may not always be successfully pruned. Table II  
644 shows the percentage of the input frames pruned at each step.  
645 The proposed character dictionaries for three movies; V1, V2  
646 and V3 are shown in Fig. 10–12 in Appendix B.

647 The precision, recall and F1 score measures that we used to  
648 quantify the performance of our method are shown in Fig. 6.  
649 The relative saliency threshold and track duration threshold  
650 were chosen corresponding to the best F1 score (highlighted  
651 in Fig. 6(a)). These measures were averaged across the eight  
652 movies for each setting of two parameters, relative saliency

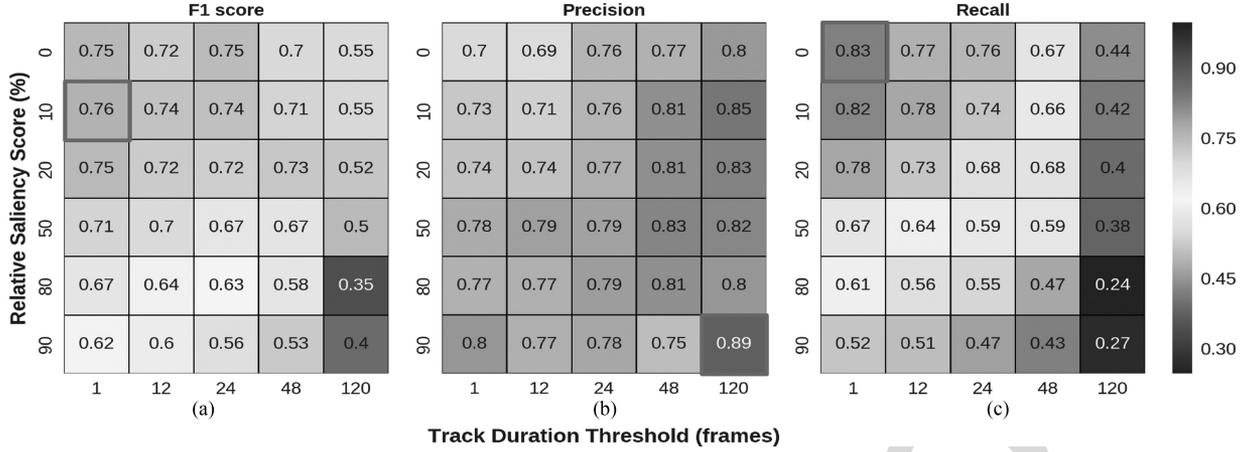
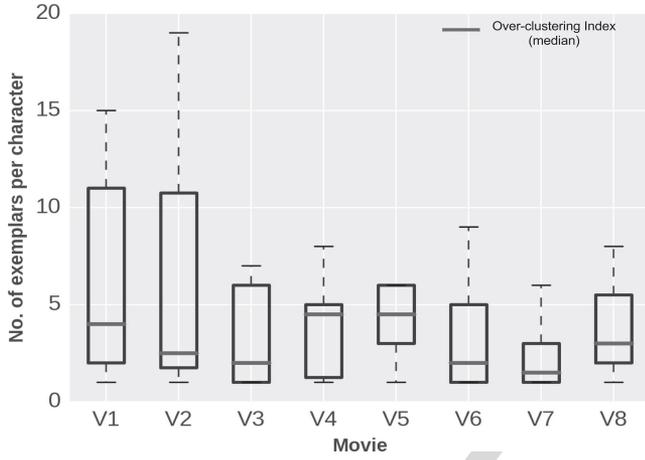


Fig. 6. Average (a) F1 score, (b) Precision and (c) Recall for all experiments.

Fig. 7. Distribution of number of exemplars per character in each movie for  $R_s(X) = 10\%$  and  $\tau = 1$ .

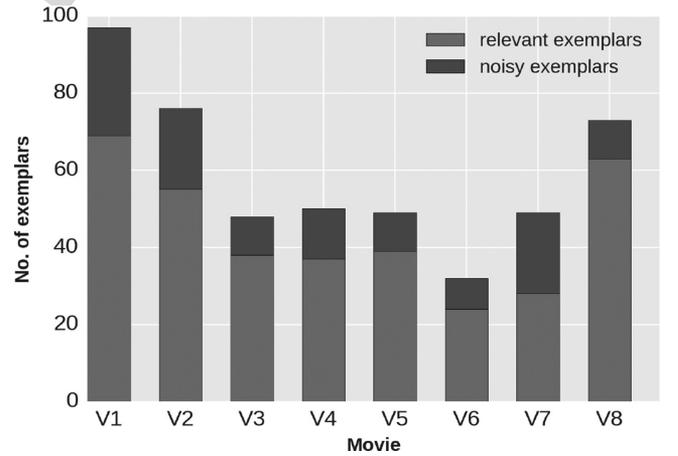
653 threshold,  $R_s(X)$  and track duration threshold,  $\tau$ . Overall, re-  
 654 call is high (over 80% for  $\tau = 1$  and  $R_s(X) = 10\%$ ) which  
 655 indicates that our proposed character dictionaries were able to  
 656 identify most of the characters in the reference at least once. Pre-  
 657 cision ranges between 70% and 90% indicating that less than  
 658 one-third of exemplars in our proposed dictionaries are noisy.

659 We note that the recall measure defined here has to be inter-  
 660 preted alongside over-clustering index; a metric that captures  
 661 the extent to which multiple exemplars represent a single refer-  
 662 ence character. The distribution of number of relevant exemplars  
 663 per character for the eight movies is shown in Fig. 7. The me-  
 664 dian number of exemplars per character, i.e., the over-clustering  
 665 index is less than 5 for all the eight movies. As described in  
 666 Section III-C, this measure lies between 1 and the number of  
 667 exemplars. Here, the number of exemplars range between 35 and  
 668 95 (with  $R_s(X) = 10\%$ ;  $\tau = 1$ ) but the over-clustering in-  
 669 dex is less than 5 which reflects on the effective performance of  
 670 the affinity propagation (AP) algorithm used for clustering.

671 Additionally, we compared the F1 score and purity of cluster-  
 672 ing for the eight movies, in order to evaluate the features used  
 673 in clustering, as shown in Table III. Although the F1 scores  
 674 (computed by comparing the exemplars to the reference) were

TABLE III  
F1 SCORE AND PURITY FOR FC7 AND GIST FEATURES USED IN CLUSTERING

Movie ID	FC7 features		GIST descriptors	
	F1 score	Purity	F1 score	Purity
V1	0.691	0.708	0.713	0.414
V2	0.773	0.651	0.769	0.345
V3	0.825	0.842	0.821	0.304
V4	0.764	0.598	0.693	0.322
V5	0.532	0.712	0.653	0.408
V6	0.740	0.677	0.732	0.398
V7	0.732	0.693	0.743	0.438
V8	0.752	0.745	0.799	0.392
Average:	<b>0.726</b>	<b>0.703</b>	<b>0.740</b>	<b>0.378</b>

Fig. 8. Number of relevant and noisy exemplars for each movie with  $R_s(X) = 10\%$  and  $\tau = 1$ .

675 similar between the two descriptors, the clustering purity using  
 676 FC7 features was significantly higher (paired t-test,  $p \ll 0.01$  to  
 677 reject  $H_0 : \mu_0 \leq \mu_1$ ) than that of GIST descriptors. This indi-  
 678 cates that FC7 features yield less noisy and more homogeneous  
 679 clusters from AP clustering. Furthermore, FC7 features perform  
 680 better for clustering than GIST features, perhaps because *Ima-  
 681 geNet* was trained to classify objects robust to variation in the  
 682 background or view-point and occlusions, whereas GIST  
 683 descriptors capture the holistic shape information in an image.

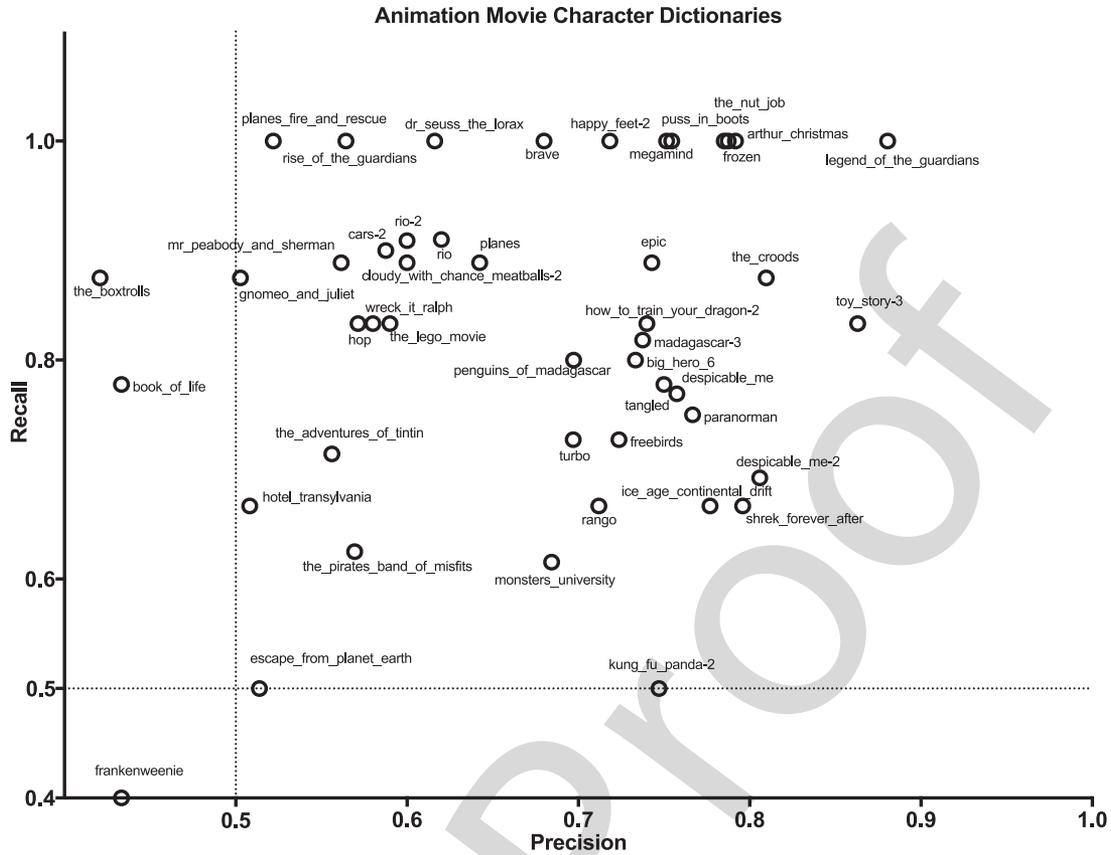


Fig. 9. Precision versus recall for forty-six movies.

684 As shown in Fig. 6(c), recall drops with an increase in  $\tau$  as  
 685 expected. Since, by increasing  $\tau$  we retain only those character  
 686 candidates which remain longer on-screen and do not always  
 687 co-occur with other salient objects. This results in excluding  
 688 some relevant exemplars. In contrast, an increase in precision  
 689 (See Fig. 6(b)) is noticed on increasing  $\tau$  since a few noisy  
 690 character candidates that are infrequent get pruned successfully.  
 691 Relative saliency threshold had the desired effect on the system  
 692 output i.e., increasing  $R_s(X)$  results in an increase in precision.  
 693 However, these gains in precision by increasing  $R_s(X)$  beyond  
 694 10% were not substantial.

695 In order to determine a good choice of the system param-  
 696 eters, we examine F1 score for different combinations of  $R_s(X)$   
 697 and  $\tau$  as shown in Fig. 6(a).  $R_s(X) = 10\%$  and  $\tau = 1$  would  
 698 be the best choice of settings. For these settings, the num-  
 699 ber of relevant and noisy exemplars for the eight movies are  
 700 shown in Fig. 8. It is interesting to note that movies V1 and  
 701 V2 have relatively larger character dictionaries and a higher  
 702 range of number of exemplars per character (See Appendix  
 703 Fig. 10–11). All the characters in the movies, V1 and V2 are  
 704 similar to cars and birds in appearance. The results at a glance  
 705 show that all instances of these characters in different scenes  
 706 were detected in these movies (which include the minor char-  
 707 acters and different appearances of the same character with re-  
 708 spect to view-point). This is likely because both cars and birds  
 709 are among the object classes in ILSVRC-2014 data used to train  
 710 *MultiBox*.

Character dictionaries for the remaining thirty-eight movies  
 were obtained with the choice of  $R_s(X)$  and  $\tau$  determined  
 above. The range of precision was 0.45–0.89 (mean/standard  
 deviation:  $0.66 \pm 0.12$ ) and recall: 0.42–1.0 ( $0.83 \pm 0.16$ ).  
 The range of over-clustering index was 1.5–6.0 ( $3.5 \pm 1.5$ ).  
 See Fig. 9 (Appendix A) for precision and recall measures of  
 all the forty six movies in our dataset. Further error analysis  
 considering different aspects of all the movies (e.g., character  
 design, color, illumination) is warranted and will be a part of  
 our future work.

We note that the movie *Frankenweenie* (2013) which was  
 produced in black and white has the lowest precision and re-  
 call in our dataset. This indicates that color rendering is an  
 important factor since the DNNs we employ were trained with  
 RGB images. The movies, *Boxtrolls* (2014) and *The Book of*  
*Life* (2014) both have a low precision and high recall indicat-  
 ing a larger number of noisy exemplars. On analyzing the errors  
 in these samples, we observed that the local tracking method  
 pruned approximately 42% of the initial character candidates  
 (c.f. the average percentage of candidates pruned by local track-  
 ing for the rest of the movies was 62.23%). This is likely be-  
 cause these movies, unlike the others in the dataset use a *rapid-fire*  
 film editing style which includes fast-action scene cuts and rapidly  
 changing backgrounds which are not ideally suited for visual  
 object tracking.

On the other hand, movies like *Kung Fu Panda 2* (2011) and  
*Escape from Planet Earth* (2013) yield high precision and low



Fig. 10. Character dictionary of the movie Frozen: precision = 0.81, recall = 1.0; over-clustering index=2.

738 recall. This is likely because these movies feature only a small  
 739 number of prominent characters and a larger number of unnamed  
 740 characters which are not included in the reference dictionaries  
 741 that we created. As expected, movies that feature distinct lifelike  
 742 animals or humans, generally performed the best. For example,  
 743 the movie Legend of Guardians (2010) featured only birds and  
 744 The Nut Job (2014) featured animals – both animals and birds  
 745 are included in the set of object categories of the ILSVRC  
 746 datasets.

## 747 V. CONCLUSIONS AND FUTURE WORK

748 In this paper, we proposed an unsupervised method to auto-  
 749 matically create a dictionary of characters from an animation  
 750 movie. We evaluated our method on a set of eight movies with  
 751 diverse character styles and demonstrated high precision and  
 752 recall on a dataset of forty-six movies. We also showed that  
 753 the proposed method generalizes for animation movies at scale.  
 754 These character dictionaries can serve as a powerful tool for  
 755 character labeling to delineate aspects of *who appeared*, *when*  
 756 and for *how long* in a movie (video diarization). We believe  
 757 that our efforts can lay a foundation to provide an impetus for  
 758 multimedia research endeavors specifically involving animated  
 759 media content.

760 One of the drawbacks of the proposed method is that we use  
 761 an object detector that was trained with natural images. We plan  
 762 to address this issue using transfer learning to adapt the existing  
 763 models to specialize the network for detecting characters from  
 764 animation movies. The relevant and noisy exemplars that we  
 765 annotated for the system evaluation can potentially be used for  
 766 these methods. Our future work would also include using the  
 767 relevant exemplars and associated cluster members as a single  
 768 unit to facilitate robust video diarization of animation movies.

### APPENDIX A

769 PRECISION AND RECALL OF OUR PROPOSED METHOD FOR THE  
 770 FORTY-SIX MOVIES

771 The Fig. 9 plots precision vs. recall for all the movies in  
 772 our dataset. The relative saliency threshold and track duration



Fig. 11. Character dictionary of the movie Free Birds: precision = 0.72, recall = 0.72; over-clustering index = 2.5.

threshold was set to 10% and one frame respectively (tuned on 773  
 a subset of 8 movies as described in Section IV). 774

### APPENDIX B EXAMPLES

Fig. 10–12 illustrate the proposed character dictionaries for 775  
 three movies with the settings of relative saliency threshold = 777



Fig. 12. Character dictionary of the movie Cars-2: precision=0.61, recall = 0.9; over-clustering index = 4.

778 10% and track duration threshold = 1 frame. The exemplars  
779 here are arranged in no particular order to maintain their aspect  
780 ratios.

781

## REFERENCES

- 782 [1] P. P. Mohanta, S. K. Saha, and B. Chanda, "A model-based shot bound-  
783 ary detection technique using frame transition parameters," *IEEE Trans.*  
784 *Multimedia*, vol. 14, no. 1, pp. 223–233, Feb. 2012.
- 785 [2] C. Liu, D. Wang, J. Zhu, and B. Zhang, "Learning a contextual multi-  
786 thread model for movie/TV scene segmentation," *IEEE Trans. Multimedia*,  
787 vol. 15, no. 4, pp. 884–897, Jun. 2013.
- 788 [3] B. W. Chen, J. C. Wang, and J. F. Wang, "A novel video summarization  
789 based on mining the story-structure and semantic relations among concept  
790 entities," *IEEE Trans. Multimedia*, vol. 11, no. 2, pp. 295–312, Feb. 2009.
- 791 [4] Y. Li, S.-H. Lee, C.-H. Yeh, and C. C. J. Kuo, "Techniques for movie  
792 content analysis and skimming: tutorial and overview on video abstraction  
793 techniques," *IEEE Signal Process. Mag.*, vol. 23, no. 2, pp. 79–89,  
794 Mar. 2006.
- 795 [5] K. Kurzhals, M. John, F. Heimerl, P. Kuznecov, and D. Weiskopf, "Visual  
796 movie analytics," *IEEE Trans. Multimedia*, vol. 18, no. 11, pp. 2149–2160,  
797 Nov. 2016.
- 798 [6] C. Y. Weng, W. T. Chu, and J. L. Wu, "Rolenet: Movie analysis from the  
799 perspective of social networks," *IEEE Trans. Multimedia*, vol. 11, no. 2,  
800 pp. 256–271, Feb. 2009.
- [7] T. Guha, C.-W. Huang, N. Kumar, Y. Zhu, and S. S. Narayanan, "Gender  
representation in cinematic content: A multimodal approach," in *Proc. ACM Int. Conf. Multimodal Interaction*, 2015, pp. 31–34.
- [8] A. Fitzgibbon and A. Zisserman, "On affine invariant clustering and auto-  
matic cast listing in movies," *Proc. 7th Eur. Conf. Comput. Vis.*, 2002,  
pp. 304–320.
- [9] F. Vallet, S. Essid, and J. Carrive, "A multimodal approach to speaker  
diarization on TV talk-shows," *IEEE Trans. Multimedia*, vol. 15, no. 3,  
pp. 509–520, Apr. 2013.
- [10] J. Sivic, M. Everingham, and A. Zisserman, "Who are you?"—Learning  
person specific classifiers from video," in *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 2009, pp. 1145–1152.
- [11] M. Everingham, J. Sivic, and A. Zisserman, "Hello! My name is...  
Buffy"—automatic naming of characters in TV video," in *Proc. Brit. Mach. Vis. Conf.*, 2006.
- [12] M. Everingham, J. Sivic, and A. Zisserman, "Taking the bite out of auto-  
mated naming of characters in TV video," *Image Vis. Comput.*, vol. 27,  
no. 5, pp. 545–559, Apr. 2009.
- [13] "Box office history for digital animation" [Online]. Available:  
[http://www.the-numbers.com/market/production-method/digital-  
animation](http://www.the-numbers.com/market/production-method/digital-animation)
- [14] Z. Aghbari, K. Kaneko, and A. Makinouchi, "Content-trajectory approach  
for searching video databases," *IEEE Trans. Multimedia*, vol. 5, no. 4,  
pp. 516–531, Dec. 2003.
- [15] B. Ionescu, V. Buzuloiu, P. Lambert, and D. Coquin, "Improved cut detec-  
tion for the segmentation of animation movies," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, May 2006, vol. 2, pp. II–II.
- [16] B. Ionescu, D. Coquin, P. Lambert, and V. Buzuloiu, "Fuzzy color-based  
approach for understanding animated movies content in the indexing task,"  
*J. Image Video Process.*, vol. 2008, pp. 8:1–8:17, Jan. 2008.
- [17] L. Ott, P. Lambert, B. Ionescu, and D. Coquin, "Animation movie abstrac-  
tion: Key frame adaptive selection based on color histogram filtering,"  
in *Proc. 14th Int. Conf. Image Anal. Process.*, 2007, pp. 206–211.
- [18] B. Ionescu, P. Lambert, D. Coquin, L. Ott, and V. Buzuloiu, "Animation  
movies trailer computation," in *Proc. 14th ACM Int. Conf. Multimedia*,  
2006, pp. 631–634.
- [19] K. Takayama, H. Johan, and T. Nishita, "Face detection and face recog-  
nition of cartoon characters using feature extraction," in *Proc. Image, Electron. Visual Comput. Workshop*, 2012, p. 48.
- [20] S. Wang, J. Zhang, T. X. Han, and Z. Miao, "Sketch-based image retrieval  
through hypothesis-driven object boundary selection with HLR descriptor,"  
*IEEE Trans. Multimedia*, vol. 17, no. 7, pp. 1045–1057, Jul. 2015.
- [21] *The Illusion of Life: Disney Animation*. New York, NY, USA: Abbeville  
Press, 1981.
- [22] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan,  
"Object detection with discriminatively trained part-based models," *IEEE  
Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645,  
Sep. 2010.
- [23] Y. Tang, X. Wang, E. Dellandrea, and L. Chen, "Weakly supervised learn-  
ing of deformable part-based models for object detection via region pro-  
posals," *IEEE Trans. Multimedia*, vol. 19, no. 2, pp. 393–407, Feb. 2017.
- [24] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once:  
Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 2016, pp. 779–788.
- [25] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov, "Scalable object  
detection using deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, 2014, pp. 2155–2162.
- [26] C. Szegedy, S. E. Reed, D. Erhan, and D. Anguelov, "Scalable, high-  
quality object detection," arXiv:1412.1441, 2014.
- [27] K. H. Lee and J. N. Hwang, "On-road pedestrian tracking across multiple  
driving recorders," *IEEE Trans. Multimedia*, vol. 17, no. 9, pp. 1429–1438,  
Sep. 2015.
- [28] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, Jun. 2015.
- [29] O. Russakovsky *et al.*, "Imagenet large scale visual recognition challenge,"  
*Int. J. Comput. Vis.*, vol. 115, no. 3, pp. 211–252, Dec. 2015.
- [30] P. Jacard, "The distribution of the flora in the alpine zone," *New Phytologist*, vol. 11, pp. 37–50, 1912.
- [31] E. Rahtu and J. Heikkilä, "A simple and efficient saliency detector for  
background subtraction," in *Proc. 12th Int. Conf. Comput. Vis. Workshops*,  
Sep. 2009, pp. 1137–1144.
- [32] Y. Hu, X. Xie, W.-Y. Ma, L.-T. Chia, and D. Rajan, "Salient region de-  
tection using weighted feature maps based on the human visual attention  
model," in *Proc. 5th Pacific Rim Conf. Multimedia Adv. Multimedia Inf. Process.*, 2005, pp. 993–1000.

- 876 [33] V. Mahadevan and N. Vasconcelos, "Background subtraction in highly  
877 dynamic scenes," in *Proc. IEEE Conf. Comput. Vis. Pattern Recogn.*, Jun.  
878 2008, pp. 1–6.
- 879 [34] E. Rahtu, J. Kannala, M. Salo, and J. Heikkilä, "Segmenting salient objects  
880 from images and videos," *Proc. 11th Eur. Conf. Comput. Vis. Comput. Vis.*,  
881 2010, pp. 366–379.
- 882 [35] M. Kristan *et al.*, "The visual object tracking VOT2015 challenge results,"  
883 in *Proc. IEEE Int. Conf. Comput. Vis. Workshops*, Dec. 2015.
- 884 [36] G. Nebehay and R. Pflugfelder, "Clustering of static-adaptive correspondences  
885 for deformable object tracking," *Comput. Vis. Pattern Recogn.*,  
886 Jun. 2015, pp. 2784–2791.
- 887 [37] S. Leutenegger, M. Chli, and R. Y. Siegwart, "Brisk: Binary robust invariant  
888 scalable keypoints," in *Proc. Int. Conf. Comput. Vis.*, Nov. 2011, pp.  
889 2548–2555.
- 890 [38] D. Dueck and B. J. Frey, "Non-metric affinity propagation for unsuper-  
891 vised image categorization," in *Proc. IEEE 11th Int. Conf. Comput. Vis.*,  
892 Oct. 2007, pp. 1–8.
- 893 [39] A. K. Jain, M. N. Murty, and P. J. Flynn, "Data clustering: A review,"  
894 *ACM Comput. Surv.*, vol. 31, no. 3, pp. 264–323, Sep. 1999.
- 895 [40] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification  
896 with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Pro-  
897 cess. Syst.*, 2012, pp. 1106–1114.
- 898 [41] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN features  
899 off-the-shelf: An astounding baseline for recognition," in *Proc. IEEE Conf.  
900 Comput. Vis. Pattern Recogn. Workshops*, 2014, pp. 512–519.
- 901 [42] A. Oliva and A. Torralba, "Modeling the shape of the scene: A holistic  
902 representation of the spatial envelope," *Int. J. Comput. Vis.*, vol. 42, no. 3,  
903 pp. 145–175, May 2001.
- 904 [43] M. Merler, B. Huang, L. Xie, G. Hua, and A. Natsev, "Semantic model  
905 vectors for complex video event recognition," *IEEE Trans. Multimedia*,  
906 vol. 14, no. 1, pp. 88–101, Feb. 2012.
- 907 [44] B. J. Frey and D. Dueck, "Clustering by passing messages between data  
908 points," *Science*, vol. 315, pp. 972–977, 2007.
- 909 [45] "Comprehensive Annenberg Report on Diversity in Entertainment," 2016.
- 910 [46] S. L. Smith *et al.*, "Media, diversity, & social change initiative," 2016.
- 911 [47] K. Krippendorff, "Reliability in content analysis," *Human Commun. Res.*,  
912 vol. 30, no. 3, pp. 411–433, 2004.
- 913 [48] D. M. Christopher, R. Prabhakar, and S. Hinrich, "Introduction to infor-  
914 mation retrieval," *Introduction Inf. Retrieval*, vol. 151, p. 177, 2008.



**Krishna Somandepalli** received the Master's degree from University of California at Santa Barbara, CA, USA, in electrical and computer engineering. He received the Bachelor's degree in electronics and communication engineering from University Visvesvaraya College of Engineering, Bangalore, India. Following his Master's degree, he worked as an Assistant Research Scientist at NYU Langone Medical Center, New York, NY, USA. His research interests include multimodal analysis with image and signal data. He is currently working toward the Ph.D. degree in the

Signal Analysis and Interpretation Laboratory Group, Department of Electrical Engineering, University of Southern California, Los Angeles, CA, USA.



**Naveen Kumar** received the Ph.D. degree in electrical engineering from the USC Viterbi School of Engineering, where he was a member of the Media Informatics and Content Analysis Group, Signal Analysis and Interpretation Lab. He received the B.Tech. degree in instrumentation engineering from the Indian Institute of Technology, Kharagpur, India, in 2009. He currently works at the Sony PlayStation R&D in San Mateo, CA, USA. His broad research interests include machine learning and signal processing for speech, multimedia and multimodal

applications.



**Tanaya Guha** received the Ph.D. degree in electrical and computer engineering from the University of British Columbia (UBC), Vancouver, BC, Canada. She is an Assistant Professor in the Department of Electrical Engineering, Indian Institute of Technology (IIT) Kanpur. Prior to joining IIT Kanpur, she was a Postdoctoral Fellow at the Signal Analysis and Interpretation Lab (SAIL), University of Southern California.

Her current research interests include social and affective computing, multimedia analysis, and multimodal signal processing. Dr. Guha received Mensa Canada Woodhams Memorial Scholarship, Google Anita Borg Scholarship, and Amazon Grace Hopper celebration scholarship.



**Shrikanth S. Narayanan** (SM'88–M'95–SM'02–F'09) is the Niki & C. L. Max Nikias Chair in Engineering at the University of Southern California (USC), Los Angeles, CA, USA, and holds appointments as a Professor of Electrical Engineering, Computer Science, Linguistics, Psychology, Neuroscience and Pediatrics, the Research Director of the Information Science Institute, and as the Founding Director of the Ming Hsieh Institute. Prior to USC, he was with AT&T Bell Labs and AT&T Research from 1995–2000. At USC, he directs the Signal Analysis

and Interpretation Laboratory (SAIL). His research focuses on human-centered signal and information processing and systems modeling with an interdisciplinary emphasis on speech, audio, language, multimodal and biomedical problems and applications with direct societal relevance. Prof. Narayanan is a Fellow of the National Academy of Inventors, the Acoustical Society of America, the International Speech Communication Association (ISCA) and the American Association for the Advancement of Science (AAAS), and a member of Tau Beta Pi, Phi Kappa Phi, and Eta Kappa Nu. He is the Editor-in-Chief for *IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING*, an Editor for the *Computer Speech and Language Journal* and an Associate Editor for the *APSIPA Transactions On Signal And Information Processing*. He was also previously an Associate Editor of the *IEEE TRANSACTIONS OF SPEECH AND AUDIO PROCESSING* (2000–2004), *IEEE SIGNAL PROCESSING MAGAZINE* (2005–2008), *IEEE TRANSACTIONS ON MULTIMEDIA*, (2008–2011), the *IEEE TRANSACTIONS ON SIGNAL AND INFORMATION PROCESSING OVER NETWORKS* (2014–2015), *IEEE TRANSACTIONS ON AFFECTIVE COMPUTING* (2010–2016), and the *Journal of the Acoustical Society of America* (2009–2017). He received several honors including Best Transactions Paper awards from the IEEE Signal Processing Society in 2005 (with A. Potamianos) and in 2009 (with C. M. Lee) and selection as an IEEE Signal Processing Society Distinguished Lecturer for 2010–2011 and ISCA Distinguished Lecturer for 2015–2016. Papers coauthored with his students have received awards including the 2014 Ten-year Technical Impact Award from ACM ICMI and at several conferences. He has published more than 750 papers and has been granted 17 U.S. patents.

Q5  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941

942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992

- 994 Q1. Author: Please check whether the affiliation of authors is okay as set.  
995 Q2. Author: Please provide year information in Ref. [13].  
996 Q3. Author: Please update Ref. [26], if already published in prints.  
997 Q4. Author: Please provide page range in Refs. [19], [28], and [35].  
998 Q5. Author: Please provide full bibliographic details in Ref. [46].

IEEE Proof